

# Variational Data Assimilation and Model Learning Jerome Monnier

#### ▶ To cite this version:

Jerome Monnier. Variational Data Assimilation and Model Learning. Master. France. 2021. hal- $03040047 \mathrm{v2}$ 

# HAL Id: hal-03040047 https://hal.science/hal-03040047v2

Submitted on 23 Nov 2021 (v2), last revised 12 Apr 2024 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

#### **INSA** Toulouse

**Department of Applied Mathematics** 

# Variational Data Assimilation and Model Learning

To assimilate datasets into physically-based models, to efficiently compute a model-output gradient, to calibrate models / identify model parameters, to learn model terms

by Jérôme Monnier

Fall 2021

Keywords: inverse problems, model-based feedback control, PDE models, optimisation, gradient-based algorithms, adjoint method, best estimate, error covariances, Earth sciences.





INSTITUT DE MATHÉMATIQUES de TOULOUSE



# Contents

| 1        | Fundamentals & Practical description                       |  |    |  |  |  |
|----------|--|--|----|--|--|--|
|          | 1.1  | Data Assimilation at a glance  | 6  |  |  |  |
|          | 1.2  | Direct vs Inverse modeling   | 8  |  |  |  |
|          |  | 1.2.1 Direct, inverse  | 8  |  |  |  |
|          |  | 1.2.2 Examples   | 9  |  |  |  |
|          |  | 1.2.3 Ill-posedness, well-posedness  | 11 |  |  |  |
|          | 1.3  | Least-square solutions of regression problems                              | 12 |  |  |  |
|          |  | 1.3.1 Linear least-square problems   | 12 |  |  |  |
|          |  | 1.3.2 Ill-posed inverse problem & Tikhonov regularization                  | 14 |  |  |  |
|          |  | 1.3.3 Non linear least-square problems                                     | 16 |  |  |  |
|          |  | 1.3.4 Rank deficience and SVD *  | 19 |  |  |  |
|          | 1.4  | From the control of dynamical systems to data assimilation                 | 21 |  |  |  |
|          |  | 1.4.1 Connections between Optimal Control and Variational DataAssimilation | 21 |  |  |  |
|          |  | 1.4.2 Example: optimal dynamic trajectory of a vehicle                     | 22 |  |  |  |
|          | 1.5  | Connections with Artificial Neural Networks in Machine Learning            | 24 |  |  |  |
|          | 1.6  | The Programming Practical in Earth sciences                                | 27 |  |  |  |
|          |  | 1.6.1 The direct model: river flow dynamics                                | 27 |  |  |  |
|          |  | 1.6.2 The spatial hydrology inverse problem                                | 28 |  |  |  |
| <b>2</b> | Optimal Control of ODE: the Linear-Quadratic (LQ) case $3$ |  |    |  |  |  |
| -        | 2.1  | Introduction   | 32 |  |  |  |
|          | $\frac{2.2}{2.2}$  | Illustrative example   | 32 |  |  |  |
|          |  | 2.2.1 The controlled system  | 33 |  |  |  |
|          |  | 2.2.2 Different states of the system according to the control values       | 33 |  |  |  |
|          | 2.3  | Controllability *  | 35 |  |  |  |
|          | -  | 2.3.1 Reachable sets   | 35 |  |  |  |
|          |  | 2.3.2 Controllability of autonomous linear systems: Kalman's condition     | 37 |  |  |  |
|          | 2.4  | The Linear-Quadratic (LQ) problem  | 37 |  |  |  |
|          |  | 2.4.1 The Linear model   | 38 |  |  |  |
|          |  | 2.4.2 The Quadratic cost function  | 39 |  |  |  |
|          | 2.5  | The Pontryagin principle & the Hamiltonian                                 | 40 |  |  |  |
|          |  | 2.5.1 Existence and uniqueness of the solution                             | 40 |  |  |  |
|          |  | 2.5.2 The Pontryagin principle in the LQ case                              | 42 |  |  |  |
|          |  | 2.5.3 The Hamiltonian  | 43 |  |  |  |

|   |              | 2.5.4 Examples   | 45 |
|---|--------------|--|----|
|   | 2.6          | Feedback law and the Riccati equation *  | 45 |
|   | 2.7          | A few remarks  | 47 |
|   | 2.8          | Numerical methods for optimal control problems   | 47 |
|   |              | 2.8.1 Two classes of numerical methods   | 47 |
|   |              | 2.8.2 Direct methods   | 47 |
|   |              | 2.8.3 Indirect methods: boundary value problem and shooting methods                        | 48 |
|   |              | 2.8.4 Direct vs indirect methods   | 49 |
|   | 2.9          | The fundamentals equations at a glance   | 50 |
|   |              |  | 00 |
| 3 | VDA          | A: optimal control of PDE models, adjoint equations  | 51 |
|   | 3.1          | The direct model   | 53 |
|   |              | 3.1.1 General non-linear stationary direct model   | 53 |
|   |              | 3.1.2 Examples   | 53 |
|   | 3.2          | The objective and cost functions: misfit to data   | 55 |
|   | 3.3          | VDA formulation  | 56 |
|   | 3.4          | On the difficulty to compute the gradient for large size control variables                 | 57 |
|   |              | 3.4.1 Global minimization vs local one   | 57 |
|   |              | 3.4.2 Relationship between the differential $j'(\cdot)$ and the gradient $\nabla j(\cdot)$ | 58 |
|   |              | 3.4.3 How to compute the cost function gradient ?  | 58 |
|   | 3.5          | Mathematical purposes *  | 59 |
|   |              | 3.5.1 Differentiability of the cost function   | 59 |
|   |              | 3.5.2 Existence and uniqueness of the optimal control in the LQ case                       | 60 |
|   |              | 3.5.3 LQ problems vs real-world problems   | 61 |
|   | 3.6          | Deriving the optimality system from the Lagrangian   | 62 |
|   |              | 3.6.1 Weak forms of the equations  | 62 |
|   |              | 3.6.2 The Lagrangian & the optimality system   | 63 |
|   | 3.7          | Computing the cost function gradient without adjoint model                                 | 64 |
|   | 0            | 371 The Finite Difference gradient   | 64 |
|   |              | 372 The straightforward "gradient" expression  | 65 |
|   |              | 373 LTM-based gradient   | 65 |
|   | 38           | Bigorous derivation of the adjoint equations and the corresponding gradient ex-            | 00 |
|   | 0.0          | pression   | 66 |
|   |              | 3.8.1 The adjoint-based gradient theorem   | 67 |
|   |              | 3.8.2 The optimality system  | 69 |
|   |              | 3.8.3 Cradient components: in weak or classical forms ? *                                  | 70 |
|   | 3.0          | The VDA algorithm  | 70 |
|   | 3.5<br>3.10  | The fundamental equations at a glance  | 74 |
|   | 3.10<br>3.11 | Another example: control of an electic membrane deformation *                              | 74 |
|   | 3.11<br>3.10 | How to access your adjoint computational code & gradient *                                 | 70 |
|   | J.1Z         | 2 12 1 The gapler product test   | 11 |
|   |              | 2.12.1 The scalar product test   | 11 |
|   |              | 5.12.2 The gradient test   | 18 |
| 4 | Eau          | ivalences between VDA, the BLUE and Bayesian inferences                                    | 81 |
| I | 4 1          | The Best Linear Unbiased Estimator (BLUE) and VDA  | 83 |
|   | 1.1          |  | 00 |

|          |     | 4.1.1 The least-square solutions depend on the metrics   | 33 |
|----------|-----|--|----|
|          |     | 4.1.2 The BLUE, scalar case  | 34 |
|          |     | 4.1.3 Equivalence with a VDA solution  | 35 |
|          |     | 4.1.4 The BLUE (vectorial case) & resulting optimal metric   | 37 |
|          |     | 4.1.5 Sequential filters formalism   | 38 |
|          | 4.2 | The Bayesian view  | 39 |
|          | 4.3 | Conclusion: in the LQ case, VDA, the Maximum A-Posteriori (MAP) estimation   |    |
|          |     | and BLUE can be equivalent   | 93 |
|          | 4.4 | Another example of BLUE and VDA solution   | )5 |
| <b>5</b> | VD. | A for time-dependent PDE systems (parabolic, hyperbolic) 9   | )9 |
|          | 5.1 | The direct model and cost function   | )1 |
|          |     | 5.1.1 Classical PDE models   | )1 |
|          |     | 5.1.2 The general (and formal) direct model $\ldots \ldots \ldots$ | )3 |
|          |     | 5.1.3 The cost function $\ldots \ldots \ldots$                     | )4 |
|          | 5.2 | The optimization problem & regularization terms  | )5 |
|          |     | 5.2.1 The optimization problem $\ldots \ldots \ldots$              | )5 |
|          |     | 5.2.2 Regularization terms $\ldots \ldots \ldots$                  | )6 |
|          | 5.3 | Derivation of the optimality equations (in continuous form)  | )6 |
|          |     | 5.3.1 The Linear Tangent Model (LTM) and the resulting gradient 10   | )6 |
|          |     | 5.3.2 The adjoint model and the gradient expression  | )8 |
|          | 5.4 | Derivation of the optimality system in a discrete formalism  | 10 |
|          | 5.5 | $Algorithm(s) \dots \dots$   | 11 |
|          |     | 5.5.1 The 4D-var algorithm $\ldots \ldots \ldots$                  | 11 |
|          |     | 5.5.2 A few remarks $\ldots \ldots \ldots$                         | 12 |
|          | 5.6 | The fundamental equations at a glance  | 15 |
|          | 5.7 | Examples in Earth Sciences *   | 16 |
|          |     | 5.7.1 Twin experiments concept   | 16 |
|          |     | 5.7.2 Identification of the topography in a 2d shallow-water model 11  | 16 |
|          |     | 5.7.3 Identification of inflow boundary conditions in a 2d shallow-water model.  |    |
|          |     | Flood flows context  | 17 |
|          | 5.8 | Exercises: optimality systems of classical PDE models *  | 20 |
|          |     | 5.8.1 Viscous Burgers' equation  | 20 |
|          |     | 5.8.2 Diffusion equation with non constant coefficients  | 21 |
| 6        | Dic | tionary-based model learning 12  | 23 |
|          | 6.1 | Learning ODE terms from a dictionary   | 24 |
|          |     | 6.1.1 Basic principle  | 24 |
|          |     | 6.1.2 The inverse problem $\ldots \ldots \ldots$                   | 25 |
|          |     | 6.1.3 The optimisation problem $\ldots \ldots \ldots$              | 29 |
|          |     | 6.1.4 The minimisation algorithms: LASSO, STRidge  | 32 |
|          |     | 6.1.5 When the goal is to identify model parameters $*$  | 33 |
|          | 6.2 | Numerical experiments  | 34 |
|          |     | 6.2.1 The direct model: non-linear dynamics of a spring  | 34 |
|          |     | 6.2.2 The inverse problem & strategy to determine $\lambda$  | 38 |

|              |      | 6.2.3 Numerical results  |
|--------------|------|--|
|              |      | 6.2.4 Programming practical  |
|              | 6.3  | Learning PDE model terms   |
|              |      | 6.3.1 Basic formalism for a generic scalar PDE equation  |
|              |      | 6.3.2 Case of the viscous Burgers equation   |
|              | 6.4  | Appendix   |
|              |      | 6.4.1 Notations matricielles & normes  |
|              |      | 6.4.2 Algorithmes d'optimisation   |
| Ap           | pen  | dices 153  |
| $\mathbf{A}$ | A fe | w mathematical and optimization recalls 15   |
|              | A.1  | Differential calculus, functional analysis   |
|              |      | A.1.1 Differentials  |
|              |      | A.1.2 Green formulas   |
|              |      | A.1.3 The Riesz-Fréchet theorem  |
|              |      | A.1.4 The adjoint operator   |
|              | A.2  | Differentiable optimization  |
|              |      | A.2.1 Convexity $\ldots \ldots \ldots$ |
|              |      | A.2.2 Optimality conditions  |
|              |      | A.2.3 Lagrangian and saddle-point  |
|              | A.3  | Descent algorithms and Newton-like methods   |
|              |      | A.3.1 Gradient type methods  |
|              |      | A.3.2 The linear search (step descent)   |
|              |      | A.3.3 Newton method (2nd order)  |
|              |      | A.3.4 Gauss-Newton method  |
|              |      | A.3.5 Quasi-Newton method  |
|              | A.4  | Basic notations of statistics  |
| в            | Algo | orithmic - Automatic Differentiation * 16  |
|              | B.1  | What adjoint code?   |
|              | B.2  | From mathematical differentiation to source code differentiation   |
|              | B.3  | Automatic differentiation in short   |
|              | B.4  | Exercices  |
|              | B.5  | Recycle your linear solver   |
|              |      | B.5.1 How to adjoint-ize your linear solver ?  |
|              |      | B.5.2 A simple example   |
|              | B.6  | Complement: MPI instructions   |
|              | B.7  | Complement: optimize your memory. Checkpointing  |
|              | B.8  | Practical work Tapenade  |
|              | B.9  | Example of a software architecture designed for Automatic Differentiation 194  |
| $\mathbf{C}$ | Wes  | ther forecasting & VDA algorithm variants *  |
| -            | C.1  | Data for weather forecasting   |
|              | C.2  | The incremental VDA (4D-Var) algorithm   |
|              | -    |  |

#### Bibliography

1

<sup>&</sup>lt;sup>1</sup>The sections indicated with a \* are "to go further sections". These sections can be skipped as a first reading, or if you are not interested in deeper mathematical basis, mathematical proofs.

#### WhereTo & How

#### Language

At INSA Toulouse, this course is tough in GlobeEngliche excepted if *all* attendants understand French sufficiently well.

#### Goal of this course

- To formulate the Variation Data Assimilation (VDA) equations aiming at fusing a physicalbased model, databases and prior information (e.g. covariance operators).

- To design optimal control algorithms for systems governed by a PDE or an ODE (or even SDE).

- To compute gradients by deriving the adjoint equations.

- To formulate real-like numerical modelling problems by combining "optimally": mathematical - physical equations (PDE models mainly) and databases containing measurements of the phenomena.

- To perform model calibration, parameter identification, local sensitivity analysis by assimilating the available data into the model.

- To identify (learn) model terms from datasets.

- To understand the link between the VDA approach, the Bayesian approach and the sequential filters approaches (Kalman's filter).

- To known to assess computational codes including the adjoint code and the resulting gradient.

- To generate an adjoint code by employing "Automatic Differentiation" (source-to-source code transformation such as Tapenade tool) (also called "backward propagation gradient" in machine learning).

At the end, the students are supposed to be able :

- To set up a data assimilation inverse formulation e.g. in (geophysical) fluid mechanics, structural mechanics, biology etc in presence of databases, measurements.

- To compute efficiently the gradient with a large number of parameters by deriving the adjoint model.

- To design the complete optimisation process,

- To perform local sensitivities analysis, to calibrate the model, to estimate uncertain parameters by assimilating the data into the physical-based model (PDE but also potentially SDE).

- To learn model terms (ODE or PDE) from datasets.

#### Content

Fundamentals: inverse problems, least-squares solutions (linear and not), SVD, optimisation, regularisation of optimisation problems.

Examples of inverse problems, parameter identification, model calibration, data assimilation for prediction.

Basics of optimal control:

- ODE systems: the Linear-Quadratic (LQ) case, Pontryagin's principle, the Hamiltonian.

- Non-linear PDE systems (steady and unsteady): gradient computation, adjoint equations, optimality system, the Lagrangian.

- Variational Data Assimilation (VDA) formulation: cost function gradient, gradient-based optimisation, regularisation terms, connections with prior probabilistic model - covariances operators, connection with the Kalman's filter and the Bayesian view, 4D-var algorithm(s).

Examples and exercises.

A programming practical: assimilation of altimetry type data into a river flow model (non-linear advection-diffusion model).

Implementation in Python (or FreeFem++).

Algorithmic Differentiation, Automatic differentiation (source-to-source e.g. with Tapenade software); link with the "backward propagation gradient".

The sections indicated with a \* are "to go further sections". These sections can be skipped as a first reading, or if you are not interested in deeper mathematical basis, mathematical proofs.

**Keywords** PDE models, uncertainties, data basis, variational calculus, optimal control, adjoint equations, variational data assimilation, sensitivities, prior probabilistic information, covariance operators, algorithmic differentiation.

**Background required** Differential calculus, PDE models, numerical analysis, optimisation, programming, basics of probability.

(Basics knowledge of neural networks - deep learning enables to make various links between these two complimentary approaches).

#### J. Monnier

**Duration** At INSA Toulouse, this course is scheduled within 9 weeks, 25H of plenary classes, exercise sessions and a couple of sessions to complete the programming practical. The expected home work duration is 25H too.

**Examinations** A few  $\approx$  45mn long exams : on-line exercises with multiple-choice questionnaires .

1 programming practical work with a (short) written report + the source code.

# Chapter 1

# Fundamentals & Practical description

The outline of this chapter is as follows.

#### Contents

| 1.1 | Data  | A Assimilation at a glance  | 6         |
|-----|-------|---|-----------|
| 1.2 | Dire  | ect vs Inverse modeling   | 8         |
|     | 1.2.1 | Direct, inverse   | 8         |
|     | 1.2.2 | Examples  | 9         |
|     | 1.2.3 | Ill-posedness, well-posedness   | 11        |
| 1.3 | Leas  | st-square solutions of regression problems  | 12        |
|     | 1.3.1 | Linear least-square problems  | 12        |
|     | 1.3.2 | Ill-posed inverse problem & Tikhonov regularization $\ldots \ldots \ldots \ldots$ | 14        |
|     | 1.3.3 | Non linear least-square problems  | 16        |
|     | 1.3.4 | Rank deficience and SVD *   | 19        |
| 1.4 | From  | n the control of dynamical systems to data assimilation $\ldots$                  | <b>21</b> |
|     | 1.4.1 | Connections between Optimal Control and Variational DataAssimilation              | 21        |
|     | 1.4.2 | Example: optimal dynamic trajectory of a vehicle                                  | 22        |
| 1.5 | Con   | nections with Artificial Neural Networks in Machine Learning                      | <b>24</b> |
| 1.6 | The   | Programming Practical in Earth sciences   | 27        |
|     | 1.6.1 | The direct model: river flow dynamics   | 27        |
|     | 1.6.2 | The spatial hydrology inverse problem   | 28        |

# 1.1 Data Assimilation at a glance

Data Assimilation (DA) is the science combining "at best" the knowledge of a system that is: 1) the mathematical model(s), 2) the observations (measurements / datasets), 3) potential priors such as statistical errors or a-priori probabilistic laws.

The mathematical models are generally Partial Derivatives Equation (PDE) systems. It may be a Stochastic Differential Equation (SDE) too; SDE are not considered in the present course.

The observations may be in-situ or remote sensed e.g. satellite or drones images. They can constitute large size databases ("big data").

DA may be seen as physically-based processes solving inverse problems. Variational Data Assimilation (VDA) constitutes one of the two types of DA methods. VDA relies on an optimization process. The present course mainly focuses on VDA.

Common features with Artificial Neural Networks (ANN) employed for regression problems can be established, not in the way the methods work but on the goals and the results.



Figure 1.1: Case of a time-dependent model, one of the two goals of DA: setting up a good predictive model. (Image extracted from a ECMWF course).

Assimilating data into models, calibrating models: what for ? When data are assimilated into models, some model "parameters" e.g. an effective fluid viscosity, or even the IC, a BC value of the BVP, are estimated (calibrated).

One calibrates models from past measurements (observations) to:

• 1) identify uncertain-unknown parameters values of the model (*identification / calibration process*),

#### J. Monnier

• 2) run the (calibrated) model for *prediction*.

Indeed, after calibration, the model may be supposed to be more accurate than if not, for more or less long horizon...

Building up a model able to describe "at best" available observations is cheering before using it to predict, estimate unobserved values.

However, a robust calibration process does not fully ensure that the resulting model is reliable for prediction. Typically a calibrated model should not be used out of range of its calibration value set. This remark fully applied to supervised learning in machine learning too.



Figure 1.2: Principle of VDA: control some uncertain input model parameters to make fit the model outputs with data.

#### Two classes of Data Assimilation (DA) methods

DA refers to two types of methods, see e.g. [14]:

1) A purely stochastic approach: *sequential filters*. For example, the classical Kalman filter which computes the Best Linear Unbiased Estimate (BLUE). One may mention Ensemble Kalman Filterq (EnKF) too.

2) The VDA approach based on the *optimal control theory*. It is a deterministic approach since optimizing a cost function while respecting a physical-based model constraint.

However a-priori stochastic information can be introduced in the process through the employed metrics.

A cost function measuring the misfit between the model output and the measurements is minimized.

Each approach presents advantages and drawbacks. The present course focuses on the VDA approach which is based on the adjoint equations.

The sequential filter approaches are not studied in details here. However, it will be shown that in the Linear Quadratic (LQ) case, both approaches are equivalent.

In real-world problems, data are generally heterogeneous, multiscale in space and time. Moreover, data are generally very sparse and/or unfrequent. As a consequence, they measure very partially only the modeled phenomena.

Recall that *numerical simulation*, based on a mathematical model, is a fundamental step for a large range of industrial processes - conceptions - designs or physical - geophysical analysis - descriptions - predictions. Numerical simulation replaces more and more (real) experiments. Experiments are often built up for validation of the (virtual) numerical results only.

**Some history** *Data Assimilation* is a relatively recent science.

The *Kalman filter* (sequential approach) has been developed in the 60's by R. Kalman<sup>1</sup>. These sequential filters have been first (?) used during the NASA Apollo program to better estimate the trajectories.

Least-square method is a standard method to approximate a solution of over-determined systems. Historically, it has been elaborated by J.C. Gauss (1777-1855) and A.-M. Legendre (1752-1833).

J.C. Gauss calculated, at 24 years old (...), a correct prediction of an asteroid trajectory based on past observations.

Somehow, VDA solves *non-linear least-square problems* with an underlying physically-based model: the physically-based model constitutes a constraint of the optimization problem.

The historical application of DA and more particularly VDA, is the weather forecast problem (atmosphere dynamics). This extremely complex problem is briefly presented in appendix.

# 1.2 Direct vs Inverse modeling

#### 1.2.1 Direct, inverse

In the common sense, the term "model" denotes a *direct* model (also called "forward" model).

A direct problem based on the model operator  $\mathcal{M}(.)$  consists to find (compute) a solution y given the input parameter k: find  $y, y = \mathcal{M}(k)$ .

 $<sup>^{1}</sup>$ R. Kálmán, Hungarian-born American electrical engineer and mathematician, inventor of the Kalman filters



Figure 1.3: Representation of a direct model with its input "parameter" k (a-priori vectorial) and its output variable y (its "response"); and its inverse counterpart.

The inverse problem consists to find (compute) k, given y i.e. solve:  $k = \mathcal{M}^{-1}(y)$ .

#### 1.2.2 Examples

**Example 1)** (Parameter identification).

Find a polynomial p(x) of degree n, of coefficients  $(k_0, \dots, k_n)$  that fit given values  $(y_1, \dots, y_n)$  at some points  $(x_1, \dots, x_n)$ .

This Lagrange interpolation problem is the inverse problem of the following direct problem: calculate the given polynomial p(x) at the points  $(x_1, ..., x_n)$ .

This example is a basic data assimilation problem, aiming at *identifying parameters* of the "model" p(k; x).

**Example 2)** (Model term identification) (from [26]). Given a real symmetric  $n \times n$  matrix M and n real numbers  $(\lambda_1, \dots, \lambda_n)$ , find a diagonal matrix D such that (M + D) has the eigenvalues  $(\lambda_1, \dots, \lambda_n)$ .

This problem is inverse to the following direct problem: compute the eigenvalues of the given matrix (M + D).

This inverse problem aims at *identifying "model" terms*.

**Example 3)** (Model Learning).

Given the first few terms  $(x_1, \dots, x_m)$  of a sequence, find the law of formation of the sequence; that is find  $x_n$  for all n.

The corresponding direct problem is to evaluate the sequence  $(x_n)_n$  given the law of formation (the model).

From [26]: this inverse problem is used on intelligence tests. In a mathematical point of view, such inverse problems have many solutions (that is why theirs use on intelligence tests has seemingly been criticized).

This example may be seen as a *model learning* example.

#### **Example 4)** (PDE based DA problem).

Let us consider a diffusion phenomena in a material. The non homogeneous diffusivity (it may be a conductivity) of the material (eg. biological tissue) is denoted by  $\nu(x)$ . The model is the following.

Given  $\nu(x)$  in the domain  $\Omega$ , find the scalar quantity u(x,t) (e.g. an electrical field or wave intensity) satisfying:

$$\begin{cases} \partial_t u(x,t) - div(\nu(x)\nabla u(x,t)) = 0 & \text{in } \Omega \times ]0, T[\\ u(x,0) = u_0(x) & \text{in } \Omega \\ u(x,t) = u_d(x,t) & \text{in } \partial\Omega \times ]0, T[ \end{cases}$$
(1.1)

The initial condition  $u_0$  and the value  $u_d$  at boundary are given. The direct problem consists to solve this classical Boundary Value Problem (BVP).

The inverse problem is as follows.

One acquires measurements of the field u(x,t) and the flux  $(\nu(x)\partial_n u(x,t))$  on the boundary  $\partial\Omega$ .

Given these boundary measurements, one seek to determine the unknown-uncertain diffusivity coefficient  $\nu(x)$  (eg. a conductivity) in the domain  $\Omega$ .

This inverse problem is the "impedance tomography" problem.

In the context of Electrical Impedance Tomography (EIT), the inverse problem aims at recovering conductivity (and permittivity) inside a body from surface measurements of electrical currents and potentials.

It is a potentially ill-posed problem, depending on the assumptions, see e.g. [L. Borcea, Inv. Problems (2002)].

EIT problem is still active research problem; it still posed challenging questions for mathematicians, numerician and experimentalists. "EIT is a noninvasive type of medical imaging in which the electrical conductivity, permittivity, and impedance of a part of the body is inferred from surface electrode measurements and used to form a tomographic image of that part" (from Wikipedia page). See eg. Fig. 1.4.

The reader may consult e.g. [26] or [3] to read other standard instructive inverse problems in various technological contexts; the mathematical problems involves integral equations, PDE systems etc.



Figure 1.4: Electrical Impedance Tomography (EIT) for Cardio-Pulmonary Monitoring: voltage measurements around the thorax using an EIT system with 16 electrodes. (R) Image extracted from C. Putensen et al., J. Clinical Medecine (2019).

#### 1.2.3 Ill-posedness, well-posedness

In the Hadamard sense<sup>2</sup>, a model is well-posed if the following conditions hold:

i) it admits an unique solution,

ii) the solution depends continuously on the data or input parameters.

The first condition i) is entirely related to the space the solution is sough in.

The second condition ii) is a stability condition. This property is crucial too.

Indeed, if this condition is not satisfied then any measurement or numerical error generates large errors in the model response i.e. a highly perturbed solution.

In all the sequel, it will be assumed that the direct model is well-posed; this is a necessary condition to go further !...

This assumption implies in particular that the "model operator" M defined as follows, is continuous.

$$\mathcal{M}: k \in \mathcal{K} \mapsto y \in \mathcal{Y}$$
(1.2)

Note that thanks to the open mapping theorem, see e.g. [11]: If  $\mathcal{M}$  is linear and continuous with  $\mathcal{X}$  and  $\mathcal{Y}$  Banach spaces, then the inverse model operator  $\mathcal{M}^{-1}$  is continuous

*Exercise.* Propose a general PDE-based example which is trivially well-posed. (*Hint : mind to the linear elliptic BVP, coercitive in a Hilbert space V, with the Lax-Milgram theory applying*).  $\Box$ 

In real-life modeling, even if the direct problem is well-posed, the inverse problem is very often severely ill-posed !...

*Ill-posed inverse problems* are somehow the opposite of well-posed direct problems. Direct problems are usually the way that can be solved easily (compared to the inverse problem). Actually, direct and inverse models are back-and-forth transmission of information between k and y. Roughly, "direct problem involves going from cause to effect, whereas the inverse problem attempts to go from the effects to the cause" (from [2]).

Many inverse problems in science and engineering consist to determine properties of some unmeasurable quantities (in space and/or in time).

The information, data are far to be complete or even accurate. This lack of observations, data (both in terms of quantity and quality) is one of the source of the encountered difficulties to solve the inverse problems.

The measurements of Y are almost always incomplete (and inaccurate). Also, the direct models are often non-linear.

In real-life modeling, inverse problems are as common as direct problems. The inverse problem classes and techniques are extremely wide. Inverse modeling is a fast growing topic, requiring both stochastic and deterministic skills.

The historical and common applications fields are in geosciences (e.g. weather forecast, oil reservoirs, hydrology), neutronic (nuclear power plants), inverse scattering (e.g. seismology) and imaging (e.g. medical imaging, tomography, image restoration).

Data Assimilation (DA) problems are particular class of inverse problems. This course focuses on *Variational Data Assimilation* (VDA) which is based on calculus variations and optimal control.

VDA solutions can be seen as *least-square solutions with physical constraints*.

Then, let us recall basic properties of least-square solutions and Singular Value Decomposition (SVD) in the case of simple inverse problems.

## **1.3** Least-square solutions of regression problems

#### **1.3.1** Linear least-square problems

Let assume that we have m measurements  $(z_i)_{1 \le i \le m}$  we seek to describe by a "model". To do so, we choose to consider the following linear model with n unknown parameters  $(k_i)_{1 \le i \le n}$ :

$$\begin{array}{rcl}
 a_{11}k_1 + \dots + a_{1n}k_n &= z_1 \\
 \dots &= \dots \\
 a_{m1}k_1 + \dots + a_{mn}k_n &= z_m
\end{array}$$

We denote:  $A = (a_{ij})_{1 \le i \le m, 1 \le j \le n}$  the chosen linear transformation (the linear model is given). A is a  $m \times n$ -matrix,  $z \in \mathbb{R}^m$  the observation vector and  $k \in \mathbb{R}^n$  the unknown input parameter vector. This is a *identification parameter problem*. This problem reads as: Given A, find  $k \in \mathbb{R}^n$  such that: Ak = z.

In the case there is as much parameters  $k_i$  as observations  $z_k$  (yes, that sounds weird...), i.e. n = m, the model is well-posed if and only if A is a full-rank matrix. In this case, it exists a unique set of parameters k describing exactly the data.

Still in the case n = m but if the model is ill-posed in the sense rank(A) < n, then it exists solutions but they are not unique.

In this case, the kernel of A, Ker(A), contains non zero vectors y such that: Ay = 0. If  $k^*$  is a solution then  $(k^* + y)$  with  $y \in Ker(A)$  is also solution.

In practice there is no reason to have  $n = m \dots$ 

In the case n > m i.e. in the unlikely case there is more input parameters than observations, the system is under-determined. Generally, it exists solutions but they are non-unique.

In the case n < m i.e. in the usual case there is less input parameters than observations, the system is over-determined. A-priori it does not exist any solution fitting all data.

Indeed, with m input parameters, it can exist a unique solution making fit the observations; but what about the extra (n-m) "constraint equations" ?



Figure 1.5: Linear least-square problem (with dense observations-measurements !). (L) The most simple case: linear regression; two parameters to identify. (R) An other simple case (polynomial, degree 3).

Least-square solutions. Instead of seeking a solution satisfying each equation above (i.e. a solution making fit exactly all the observations), it is interesting to find a solution satisfying "at best" the system; in other words, a solution k minimizing the norm ||Ak - z||.

Of course the choice of different norms will lead to different solutions...

An easy choice of norm is the Euclidian norm  $\|.\|_2$  since it is associated to a *scalar product* on contrary to the norms  $\|.\|_1$  and  $\|.\|_\infty$  for example.

Then the problem becomes: find  $k^* \in \mathbb{R}^n$  such that

$$j(k^*) = \min_{\mathbb{R}^n} j(k)$$
 with  $j(k) = \frac{1}{2} ||Ak - z||_{2,m}^2$ 

It is an unconstrained optimization problem in a convex set. The functional j reads:  $j(k) = \frac{1}{2}(A^TAk, k) - (Ak, z) + \frac{1}{2}(z, z)$ . j is quadratic, convex since  $A^TA$  is positive, hence j admits a minimum in  $\mathbb{R}^n$ 

Furthermore, if  $A^T A$  is definite (it means  $n \leq m$  and A is full rank, dim(Im(A) = m) then the solution is unique.

The gradient of j reads:  $\nabla j(k) = A^T A k - A^T z$ . Then, the solution k has to satisfy the necessary optimal condition:

$$A^T A k = A^T z$$

This is the normal equations.

Full rank case. In the case A full rank (i.e.  $\dim(\operatorname{Im}(A)) = m$ ), then the solution is unique since  $A^T A$  is symmetric positive definite. Even if A sparse, then  $A^T A$  is a-priori non sparse; also  $K_2(A^T A) = (K_2(A))^2$ , hence the normal equations can be an ill-conditioned system. Then, a good algorithm to solve the normal equations is a-priori not the Cholesky algorithm, especially if m large. Remembering that the 2-norm is preserved under orthogonal transformations, a better option is to solve the following equivalent system:

$$\min_{k \in \mathbb{R}^n} \frac{1}{2} \| QAk - Qz \|_{2,m}^2$$

with Q an orthogonal matrix.

By performing QR-factorizations (Householder's method), the resulting linear system to be inverted is a triangular  $n \times n$  system, with its original conditioning number  $K_2(A)$  preserved.

#### 1.3.2 Ill-posed inverse problem & Tikhonov regularization

It has been mentioned above that in the case A is full rank  $(\dim(Im(A)) = m)$  then the least-square solution (the solution of the normal equations) is unique.

However, even if A is full rank, the system may be ill-conditioned e.g. the smallest eigenvalue is extremely small. This is a frequent case in real-world problems !

In order to better define the problem or to select a solution with desirable properties, a "regularization term" is introduced.

The most classical way to regularize the inverse problem consists to compute  $k^* \in \mathbb{R}^n$  such that  $j(k^*) = \min_{\mathbb{R}^n} j(k)$  with

$$j(k) = \frac{1}{2} \|Ak - z\|_{2,m}^2 + \alpha_{reg} \frac{1}{2} \|Ck\|_{2,n}^2$$
(1.3)

with C the regularization matrix and  $\alpha_{reg}$  a positive scalar (the weight) making the balance between the two terms.

The added regularization term is a convex, differentiable term. It is the so-called *Tikhonov regularization method*.

(A. Tikhonov, Russian mathematician, 1906-1993). This method seems to be due to Phillips and Miller too.

The solution of the corresponding normal equations reads:

$$k^* = (A^T A + \alpha_{reg}^2 C^T C)^{-1} \cdot A^T z$$

The choice C = Id provides the least-square solution with a minimal norm.

In real-world modeling problem, the regularization operator C is often defined as a differential operator



Figure 1.6: A fitting problem with/without Tikhonov regularization term. (L) The dataset to fit "at best" (= in the LS sense). (M) A LS optimal fitting curve (without regularization). (R) The same class of LS solution but with a regularization term. (Image extracted from a MIT course).

(e.g. the gradient operator) or a Fourier operator (e.g. filtering high frequencies).

Indeed, it turns out that numerous real-world direct models (phenomena) have the effect of low-pass filters. As a consequence, the inverse operator acts as a high-pass filter. Indeed, the eigenvalues (or singular values) are the largest in the reverse mapping when they are the smallest in the forward mapping.

Amplifying high frequency is not a desirable feature since it amplifies noise and uncertainties (e.g. of data). As a consequence, regularization operators play a crucial role in inverse formulations.

The weight parameter value  $\alpha_{reg}$  is empirically set since its 'optimal' value is a-priori not known. Various methods are proposed to select a value; let us cite the discrepancy principle (Morozov's principle), cross-validation techniques and the L-curve method. The latter is briefly presented in next paragraph.

In statistics, the *Tikhonov regularization* is called the *ridge regression* method. In Machine Learning (ML), it corresponds to the so-called *weight decay* method.

Other choice of norm than  $\|.\|_2$ , compress sensing Other regularization terms can be considered; typically the 1-norm. In this case, the regularization term reads:  $\alpha_{reg} \frac{1}{2} \|Cx\|_{1,n}$ .

This term is still convex but it is not differentiable therefore leading to a non-differential convex optimization problem.

Obviously this is no longer a least square problem; this is the so-called *LASSO* problem.

The 1-norm is highly interesting for its *compress sensing* features.

**On multi-objective optimization** The introduction of regularization term naturally leads to a bi-objective optimization, see (1.3). As already mentioned, the weight parameter value  $\alpha_{reg}$  has to be a-priori set. Obviously the computed optimal solution depends on  $\alpha_{reg}$ .

No magic criteria nor method exist to determine an optimal value of  $\alpha_{reg}$ .

The weight coefficient  $\alpha_{reg}$  has to respect a "good" trade-off between the two terms  $j_{misfit}(\cdot)$  and  $j_{reg}(\cdot)$ .

An expertise of the modeled system is likely the best way to determine an optimal value of  $\alpha_{reg}$ . The so-called L-curve helps to visualize this trade-off, see Fig. 1.7.



Figure 1.7: Bi-objective optimization: an example. The optimal computed optimal solution depends on the weight parameter  $\alpha_{reg}$  (here denoted by  $\lambda$ ). (L) The obtained minimal value of each cost function terms  $j_{misfit}(k^*)$  and  $j_{reg}(k^*)$  (here equal to ||Dk||) vs  $\alpha_{reg}$ . (R) The L-curve: values of  $j_{misfit}(k^*)$  vs values of  $j_{reg}(k^*)$  for different  $\alpha_{reg}$  values.

#### **1.3.3** Non linear least-square problems

Let us consider the same problem as previously: set up a model describing "at best" m available data  $z_i$ ,  $1 \le i \le m$  but using a *non-linear* model. In this case, the corresponding least-square formulation reads:

$$j(k^*) = \min_{\mathbb{R}^n} j(k) \text{ with } j(k) = \frac{1}{2} \|M(k)\|_{2,m}^2$$
(1.4)

with M defined from  $\mathbb{R}^n$  into  $\mathbb{R}^m$ , M non linear in k. Again, it is an unconstrained optimization problem in a convex set.

In the previous linear case, M was equal to: M(k) = Ak - z.

Note that neither the existence of a solution, nor its uniqueness is guaranteed without (strong) assumptions. Non-linear problems are much more difficult than linear ones. Bellow, and more generally in this course, one consider heuristic computational algorithms only.

#### Examples

**Example 1) The location problem (eg by GPS or Galileo).** The location problem by Global Navigation System Satellites (GNSS) consists to estimate  $x, x \in \mathbb{R}^3$  (the location value). The measurements are distances  $d_i^{obs}$  from given locations  $a_i$  with  $i = 1, \dots, m$  (the *m* satellites):

$$d_i^{obs} = ||a_i - x^{exact}||_{2,\mathbb{R}^3} + \varepsilon_i \quad i = 1, ..., m$$

where  $\varepsilon$  denotes a noise (the inacurracy of the measurements).

#### J. Monnier

The least-square problem to be solved is:

$$\min_{x \in \mathbb{R}^3} j(x) \tag{1.5}$$

with:  $j : \mathbb{R}^3 \to \mathbb{R}$ ,

$$j(k) = \|\|a_i - k\|_{2,\mathbb{R}^3} - d_i^{obs}\|_{2,\mathbb{R}^m}^2 = \sum_{i=1}^m \left(\|a_i - k\|_{2,\mathbb{R}^3} - d_i^{obs}\right)^2$$

The functional  $j(\cdot)$  is non convex, see Fig. 1.8.



Figure 1.8: Non-linear least square problems: the location problem. (Image extracted from Vandenberghe's UCLA course). (L) The graph of j(k), non convex functional in k. (R) Isovalues of j(k), data, exact solution and the least-square one.

**Example 2) A standard fitting problem** The goal is to fit at best a (dense) dataset; the regression problem plotted in Fig. 1.9(L).

The chosen functional ("model") is:  $j(k_1, \dots, k_4) = k_1 \exp(k_2 x) \cos(k_3 x + k_4)$  with the four parameters  $k_i$ ,  $i = 1, \dots, 4$ , to identify.

**Example 3) A classification problem** A basic way to solve a binary classification problem is to solve the non-linear least square problem defined by:

$$j(k_1, \cdots, k_p) = \sum_{i=1}^m \left( \Phi(k_1 f_1(x_i) + \cdots + k_p f_p(x_i)) - d_i \right)^2$$

where the functions  $f_j(\cdot)$  are given.

 $\Phi(u)$  is he sigmoidal function;  $\Phi(u) = \frac{\exp(u) - \exp(-u)}{\exp(u) + \exp(-u)}$ .  $\Phi$  is a differentiable function approximating the sign function, Fig. 1.9(R).



Figure 1.9: Non-linear least square problems: examples. (L) A fitting problem (with dense data); 4 parameters to identify. (R) The (differentiable) sigmoid function employed to solve binary classification problems.

#### Derivatives & optimality conditions

Let us write the necessary first order optimality condition for the general non-linear square problem (1.4). This optimality condition reads:

$$\nabla j(k) = 0$$

We have:  $j(k) = \frac{1}{2}M^T(k)M(k)$ . Let us denote the Jacobian of M as follows:

$$DM(k) = \left(\frac{\partial M_i}{\partial k_j}(k)\right)_{1 \le i \le m, 1 \le j \le m}$$

and the Hessian for each model component  $M_i(k), 1 \leq i \leq m$ , as follows:

$$D^{2}M_{i}(k) = \left(\frac{\partial^{2}M_{i}}{\partial k_{l}\partial k_{j}}(k)\right)_{1 \leq l \leq m, 1 \leq j \leq n}$$

Then the gradient and the Hessian of j read:

$$\nabla j(k) = DM^T(k)M(k)$$
$$H_j(k) \equiv D^2 j(k) = DM^T(k)DM(k) + \sum_{i=1}^m M_i(k)D^2M_i(k)$$

Let us remark that in the *previous linear case*, the gradient and the Hessian read:

$$\nabla j(k) = A^T M(k) = A^T (Ak - z)$$
 and  $H_j(k) = A^T A$ 

since the term  $D^2 M_i$  in the Hessian expression vanishs.

#### Gauss-Newton method, Levenberg-Marquardt method

The Newton method applied to the optimality condition  $\nabla j(k) = 0$  reads at each iteration:

$$H_j(k^n) \cdot \delta k = -\nabla j(k^n); \quad k^{n+1} = k^n + \delta k$$

Newton's method requires the computation and the inversion of the Hessian of j. For complex real-world models, this feature is often prohibitive because too complex or too CPU-time consuming.

The principle of the *Gauss-Newton method* is to consider in the Newton method the linearized Hessian (i.e. its approximation omitting the second order term):

$$H_i(k) \approx DM^T(k)DM(k)$$

It gives at each iteration:

$$DM^{T}(k^{n})DM(k^{n})\cdot\delta x = -DM^{T}(k^{n})M(k^{n})$$
(1.6)

Then the linear system to be inverted is symmetric positive, and definite if  $DM(k^n)$  is full rank. The Gauss-Newton method is observed to be very efficient if  $DM^T(k^n)$  is full rank and if M(k) small when close to the solution. On the contrary it becomes inefficient if these two conditions are not

satisfied and/or if the model is locally "far to be linear" i.e. if the term  $\sum_{i=1}^{m} M_i(k) D^2 M_i(k)$  is dominant.

Finally, a good alternative method to solve non-linear least square problems is the *Levenberg-Marquardt* algorithm, see e.g. [3].

This algorithm is somehow a damped version of the Gauss-Newton method. It can be seen as a combination of a descent algorithm in a trust region and the Gauss-Newton algorithm.

In summary, the most classical algorithms to solve "explicit" least square problems, linear or not, have been recalled.

VDA algorithms are based on optimal control methods; the latter are based on the minimization of a criteria j too; somehow, in the least square sense...

However, on the contrary to the "explicit" least square problems above, the unknown k is based on an underlying model; the physically-based model is an additional constraint to the optimization problem. As a consequence, the optimal solution  $k^*$  cannot be explicitly determine as in the cases above.

#### 1.3.4 Rank deficience and SVD \*

This is a "To go further section".

In the case A rank deficient  $(\operatorname{rank}(A) = r < m)$ , solutions exist but are not unique. Recall that if k is solution then  $(k + y), \forall y \in Ker(A)$ , is solution too. In this case, the functional j is convex, but not strictly convex since  $A^T A$  is not definite.

Nevertheless, it can be easily proved that it exists an unique solution  $k^*$  to the normal equations in  $Ker^T(A)^3$ .

The right mathematical tool to analyze such ill-posed linear inverse problem is the Singular Value Decomposition (SVD).

 $<sup>^{3}</sup>$ If necessary, the reader should consult his favourite course on the topic, or any excellent book / course available online.

Recalls on the SVD. Given a rectangular  $m \times n$ -matrix A, rank(A) = r < m, the SVD of A reads:

$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$$

where:

.  $\Sigma$  is the  $r \times r$ -diagonal matrix containing the singular values  $\sigma_i$  of A:

$$(\sigma_i)^2 = \lambda_i (A^T A) = \lambda_i (A A^T)$$
 for all  $\lambda_i \neq 0, \ 1 \le i \le r$   
 $0 < \sigma_r \le \dots \le \sigma_1$ 

- .  $U = (U_1, ..., U_r), m \times r$ -matrix, contains the unit orthogonal eigenvector basis of  $AA^T$ :  $(AA^T)U = U\Sigma^2$  with  $U^TU = \mathcal{I}_r$ ,
- .  $V = (V_1, ..., V_r), n \times r$ -matrix, contains the unit orthogonal eigenvector basis of  $A^T A$ :  $(A^T A)V = V\Sigma^2$  with  $V^T V = \mathcal{I}_r$ .

The vectors of U constitute an unit orthogonal basis of  $Im(A) \subset \mathbb{R}^m$ , while the vectors of V constitute an unit orthogonal basis of  $Ker^T(A) \subset \mathbb{R}^n$ .

From the SVD, a *pseudo-inverse* (or generalized inverse) of the rectangular matrix A can be defined as follows:

$$A^{-1} = V\Sigma^{-1}U^T$$

The least-square solution in terms of SVD. Let us recall that since A is rank deficient, rank(A) = r < m, the calibration problem admits solutions but non-unique. Nevertheless, it exists an unique solution  $k^*$  presenting a minimal norm  $||x||_2$  (and satisfying of course the calibration problem "Ak = z at best"). In other words, the following problem has an unique solution  $k^*$ :

$$\min_{k \in \mathbb{R}^n \text{ s.t. } A^T A k = A^T z} \|k\|_2$$

Furthermore, it can be proved that  $^4$ ,  $^5$ :

$$k^* = A^{-1}z = V\Sigma^{-1}U^T z = \sum_{i=1}^r (\sigma_i)^{-1}u_i^T v_i z$$

with the singular values  $\sigma_i$  ordered as mentioned previously i.e. decreasing gradually to 0.

The expression of  $k^*$  shows that the more the number of singular values are taken into account (from 1 to r at maximum), the more the data errors are amplified. In practice, while performing a SVD analysis of an ill-posed linear inverse problem, a key-point is to consider the "right number" of singular values; not too much, not too few...

<sup>&</sup>lt;sup>4</sup>Keypoint for the proof: show that  $(A^T A)A^{-1}z = A^T z$ .

<sup>&</sup>lt;sup>5</sup>The reader can find a very nice illustrated summary of the SVD and their application to inverse problems in the course online written by H.E. Krogstad, IMF, Norway.

# 1.4 From the control of dynamical systems to data assimilation

Let us present now the connection between the optimal control of differential equations and the socalled Variational Data Assimilation (VDA) approach.



Figure 1.10: Step 1: Direct modeling and comparison of the numerical model outputs with measurements-observations-data. Step 2: feedback to the model and so on

## 1.4.1 Connections between Optimal Control and Variational DataAssimilation

The control theory aims at analyzing dynamical systems modified by a command (the control).

The two objectives of the control theory are the followings:

A) Acting on the control in order to bring the system state to a final state given (if possible). It is a *controllability* problem.

The goal is generally to stabilize the system in order to make it insensitive to perturbation. This approach is not considered in the present course since not directly related to Data Assimilation.

B) Defining a control u such that a criteria (the cost function) j is minimal.

The cost function j depends on the control and the state of the system: it is an *optimal control* problem.

Generally the control must satisfy constraints. Optimal control is a science between the automatic science and applied mathematics.

The Variational Data Assimilation (VDA) approach is based the optimal control of differential equations (the physically-based model). In the next chapters, the optimal control of systems will be studied. First, the mathematical models will be Ordinary Differential Equation (ODE, dynamical systems), see Chapter 2. One focuses on the important *Linear-Quadratic (LQ) case*: the model is linear, the cost function is quadratic.

Next, the model is an elliptic Partial Derivatives Equations (steady-state PDE models), see Chapter 3. Finally, the considered model is an unsteady PDE models (parabolic or hyperbolic PDE systems), see Chapter 5.

VDA approach can be applied to EDS too; this mathematical equation type is not addressed in the present course.

Recall that all these classes of models represent an extremely wide range of systems encountered in engineering (and in academic researches) e.g. in fluid mechanics (geophysical or not), structural mechanics, micro-electronics, nano-technologies, biological systems, coupled multi-physic systems etc.

Calculus of variations deals with the minimization of functionals (i.e. mappings from a set of functions to  $\mathbb{R}$ ; functionals are often definite integrals involving functions and their derivatives). Optimal control theory is somehow an extension of the calculus of variations: it is mathematical optimization problems with an underlying model (the differential equation). The goal is to derive a control policy.

To correctly attempt this course, students will need to be confortable with the basics of differential calculus. Recalls and exercises are proposed in Appendix.

Historically, optimal control appeared after the second world war with applications in aeronautics (missile guidance). A key point of the optimal control is the *Pontryagin minimum principle* (L. Pontryagin (blind) Russian mathematician (1908-1988)). The latter provides a necessary condition of optimality for an ODE system.

#### 1.4.2 Example: optimal dynamic trajectory of a vehicle

Let us consider a vehicle (robot, space vehicle) represented by x(t). x represents its position and its velocity. The goal is to control the trajectory by acting on a command u(t), its engine power and/or its wheel direction for example.

We have  $x: [0,T] \to \mathbb{R}^n$ , n = 2d, and  $u: [0,T] \to \mathbb{R}^m$ .

The mechanical laws provide the vehicle trajectory as follows:

$$\begin{cases} \text{Given } u(t) \text{ and } s(t), \text{find } x(t) \text{ which satisfies:} \\ x'(t) = A(t)x(t) + B(t)u(t) + s(t) \text{ for } t \in I = [0, T] \\ \text{with the I. C.: } x(0) = x_0. \end{cases}$$
(1.7)

where A, B are two linear mappings defined from I = [0, T] onto  $M_{n,n}(\mathbb{R})$ ,  $M_{n,m}(\mathbb{R})$ , s(t) is an external force (source term given in  $\mathbb{R}^n$ ). For a sake of simplicity, the model considered is here a 1st order linear ODE. Given u(t), this problem has an unique solution x(t).

Next, the *inverse problem* (optimal control problem here) is as follows.

Given a target trajectory z(t), can we find (and compute) a control such that the vehicle is as close as possible to z(t), moreover by consuming a minimum of power ?

To formulate this question in a mathematical point of view, the following cost function is defined:

$$j(u) = \frac{1}{2} \int_0^T \|x^u(t) - z(t)\|_W^2 dt + \frac{1}{2} \int_0^T \|u(t)\|_U^2 dt$$

#### J. Monnier

with x the unique solution of the (direct) model, given u.

W and U are two symmetrical positive matrices, U definite, thus defining semi-norm and norm respectively.

The optimal control problem reads:

$$\min_{u(t)\in K} j(u)$$

where K is a convex closed subset of  $\mathbb{R}^m$ , including constraints on the control values (e.g. the power is bounded etc).

The cost function j depends on u explicitly through its second term, but also through x(t) in its first term: this is why it is a *optimal control problem*. The model operator  $\mathcal{M}$  is defined as:  $\mathcal{M}(u) = x^u$ .

Note that the optimization formulation is bi-objective.

Since the model is linear, the cost function is strongly convex; this result is proved in next chapter. This problem is a *Linear-Quadratic (LQ) optimal control problem*. LQ problems are studied in Chapter 2.

A Python code solving a 1d version of this problem, by a basic direct method, is proposed on the course Moodle page.

#### Example of VDA based on a PDE model

Please read the proposed programming practical subject in next section.

#### The VDA course outline

The present DA course follows the outline below.

Chapter 3 aims at deriving the optimality systems for elliptic PDEs (a-priori non-linear) and focus on the computational aspects. A result of existence and uniqueness of the control is presented in the LQ PDE case. We present in detail how to solve numerically the optimality system. The final goal is the *identification* of the uncertain / unknown input parameters; in other words a better *calibration* of the model.

In Chapter 5, the method is extended to unsteady PDEs systems (either parabolic or hyperbolic). Calculations derived in these cases are more formal. The classical computational algorithms is presented; it is sometimes the "4D-var" algorithm.

In Chapter 4, the equivalence between the present VDA approach and the Best Linear Unbiased estimator (BLUE), also the sequential Kalman filter, also the most probable solution (Bayesian view) are demonstrated in simple examples. These equivalence are valid for Linear-Quadratic (LQ) problem that is the model is linear, the observation operator is linear therefore the cost function is strictly convex and it admits an unique minimum.

A few references The reader can find many references related to the present multi-topic course. Concerning VDA, let us cite a few pioneer works: Y. Sasaki in the  $\sim 60-70$ 's [38, 39], F.X. Le Dimet - O. Talagrand [27]; see also [21, 15, 32] plus the numerous historical cited references therein. The course attendants may consult the books [7, 8], the following historical online course [10] (and [19] if more interested in sequential filters). Up-to-date surveys are available in numerous excellent PhD thesis manuscripts too.

#### Acknowledgments

The author is very grateful to all his (current and former) great collaborators: PhD students, research engineers and post-doctoral researchers. We've worked together on VDA, see the references indicated on the computational software webpage [34]. As a consequence, all these colleagues have greatly contributed to the course author's understanding and to the material presented in this manuscript. Moreover, the programming practicals have been elaborated in collaboration with some of these colleagues; the computational code has been developed by one of them, see the programming practical description for details.

# 1.5 Connections with Artificial Neural Networks in Machine Learning

Supervised learning applied to physically-based data can be an alternative approach to build up an effective model and/or to solve identification problems.

In this sense, ML may be seen as "Data Assimilation" not relying on a physical model (e.g. equations based on law conservations), but relying on the empirically built Artificial Neural Network (ANN). In

#### J. Monnier

other words, the ANN is the model.

The principle feature of an ANN-based model is its architecture. The latter defines the multi-scale and non-linear resolution features.

ANNs (with their numerous versions like CNN, RNN etc) enable to remarkably find nonlinear trends between data, if trained on a sufficiently marge amount of data, and if the prediction are performed sufficiently close to the learning data subset....



Figure 1.11: An Artificial Neural Network (ANN) with a single hidden layer. (Image extracted from Wikipedia page).

Let us summarize very briefly how an ANN can be built up, trained from (large) datasets, and finally how it may be employed for prediction.

A large training dataset  $\mathcal{D}$  describing the targeted phenomena has to be given.  $\mathcal{D}$  contain "examples" (samples)  $(I_i, O_i)^{obs}$ ,  $i = 1, \dots, m$  with  $I_i$  the *i*-th input variable and  $O_i$  the corresponding output. In the present course context and if adopting the same notations, the operator  $I \mapsto O$  may be either the underlying direct model or the underlying inverse model, depending on the point of view. That is, one may consider (depending on the point of view):

$$\mathcal{M}(\mathbf{k};\cdot): I \mapsto O_{(\mathbf{k};I)} \text{ or } \mathcal{M}^{-1}(\mathbf{k};\cdot): I \mapsto O_{(\mathbf{k};I)}$$

The learning stage aims at producing the desired output for each input. In other words, it consists to built up the *unknown underlying model*  $\mathcal{M}$  (or  $\mathcal{M}^{-1}$ , see above).

In ML terminology, the mapping operator above is the called the *estimator*.

The network consists of connections between perceptrons. Each connection is assigned to a weight value  $k_{i,j}$ , measuring the connection importance.

The estimator  $\mathcal{M}$  (or  $\mathcal{M}^{-1}$ , see above) is determined as the minimizer of the following differentiable

cost function (called loss function in ML):

$$j_{\mathcal{D}}(\mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \left( O_i^{obs} - \mathcal{M}(\mathbf{k}; I_i^{obs}) \right)^2 = \| O^{obs} - \mathcal{M}(\mathbf{k}; I^{obs}) \|_{2,m,h}^2$$
(1.8)

The optimization variable **k** is constituted by the ANN weight parameters. Each weight  $k_j$  is a matrix of dimension  $n_{out} \times n_{in}$ ,  $n_{\Box}$  the connections number.

The resulting optimization problem is of very large size:  $\mathbf{k}$  is a very high dimension vector. The optimization problem reads:

$$\min_{\mathbf{k}} j_{\mathcal{D}}(\mathbf{k})$$

The training step consists to solve this optimization problem i.e. to identify the optimal values of the ANN parameters  $k_j$ , j = 1, ..., Nlay. Nlay is eg. equal to the number of layers. ANN are efficient if built up from numerous hidden layers: this is *deep learning*.

In the end, the "ANN model" is determined by its architecture and its weight parameters  $(k_1, \dots, k_{Nlay})$ , see Fig. 1.11.

For each perceptron, the most efficient activation function is the rectified linear unit (ReLU) function. The "hyper-parameters" of the algorithm (learning rate, decay rate, dropout probability) can be seen as "priors" of the model. They are experimentally chosen. The selected values are those providing the minimal value of  $j(\cdot)$ .

To minimize  $j_{\mathcal{D}}(\mathbf{k})$  a first-order gradient-based stochastic optimization is employed e.g. the classical Adam method.

The cost gradient is computed by the so-called "back-propagation" procedure. It is the application of the derivative of the (numerous) function compositions. This important technical step can be performed using "Automatic Differentiation".

After the training phase, one obtains an empirical effective model, also generally not understood.

The latter may be employed as a predictor, however the ANN model thus built is very likely valid within the learning range values only !

Nowadays, an ANN can been coded in Python very easily e.g. using PyTorch - Mpi4Py libraries, TensorFlow etc.

### **1.6** The Programming Practical in Earth sciences

Below is presented a moderately complex example. This example will be developed along the course. A Python code solving the problem and assessing the employed VDA method is provided on the INSA Moodle page. This programming practical has been initially designed in collaboration with Thibault Malou during his PhD at INSA-IMT& CLS corp. in 2018-21.



Figure 1.12: The real-world inverse problem arising in global hydrology: identifying river features (unobserved bathymetry b(x), friction parametrization K(x; h(x, t)) and discharges values Q(x, t)) from sparse large scale altimetry measurements. The latter are the water surface elevation H(x, t) and the river width w(x, t).

#### **1.6.1** The direct model: river flow dynamics

River flow dynamics is a subject widely studied especially for a very long time, in particular with the Saint-Venant equations (also known as the 1D shallow water equations). (Barré de Saint-Venant, 1797-1886, French "ingénieur", physicist and mathematician).

The (single) dimensional number of the Saint-Venant equations is the Froude number denoted by Fr. At large scale, river flows presents in great majority low Froude numbers :  $Fr \ll 0.3$  and lower. This leads to subcritical flows, no shock can occur.

The river flow is here modeled in 1D. The river geometry is defined as indicated on Fig. 1.13. The wetted cross-sections A(x,t)  $(m^2)$  are represented by rectangles of width w(x,t) (m).

We denote by: H(x,t) the free surface elevation, b(x) the bathymetry elevation, h(x,t) the water depth.

The subcritical flow hypothesis  $(Fr < 1, Fr = \frac{u}{\sqrt{gh}})$  implies that the mass and momentum equations can be combined into a single equation; this leads to the so-called *diffusive wave equation* ("lubrication theory").

In its steady-state version, the present simple flow model reads as follows:

$$-\Lambda(H,b;x)\partial_{xx}^2 H(x) + \partial_x H(x) = \partial_x b(x) \text{ in } (0,L)$$
(1.9)

with :

$$\Lambda(H,b;x) = \frac{(H(x) - b(x))}{|\partial_x H(x)|}$$
(1.10)

The direct model consists to determine the water flow height H(x), solution of (1.9) accompanied Dirichlet boundary conditions, given the river bathymetry b(x).

Any flow dynamics model requires the bathymetry to be performed (therefore simulating the flow dynamics). Thus, this spatially-distributed parameter of the model is a crucial data. However, bathymetry are for a large majority of rivers in the world uncertain or even unknown.



Figure 1.13: River geometry and notations (in the unsteady case). Image from T. Malou.

#### 1.6.2 The spatial hydrology inverse problem

Spatial hydrology combines the historical hydrology science with large scale measurements from space. Spatial altimetry measures free surface elevations  $H^{obs}$  at some points along the river, with  $\pm 10$ cm uncertainty at 1 km scale (with a few days frequency).
## J. Monnier

Then the inverse problem we propose to tackle is the following.

Given a few measurements of the free surface elevation  $H^{obs}$ , estimate the bathymetry b(x) such that the flow model outputs H(x) fit with the observations  $H^{obs}$ .

# Chapter 2

# Optimal Control of ODE: the Linear-Quadratic (LQ) case

In this chapter, the optimal control of systems governed by an Ordinary Differential Equation is briefly studied; we refer for example to [18, 42] for a full course on the subject. The present chapter follows closely the presentation done in [42].

The basic and important case is studied: the case with a linear ODE model and a quadratic cost function; this the so-called "LQ case". A proof of existence and uniqueness of the optimal solution is given. Next its characterization using the Pontryagin's minimum principle and the Hamiltonian is presented.

Also the notion of controllability and the Kalman's condition are briefly presented (these concepts may be skipped).

The outline of this chapter is as follows<sup>1</sup>

#### Contents

| 2.1 | Intro | oduction   | 32        |
|-----|-------|--|-----------|
| 2.2 | Illus | trative example  | <b>32</b> |
|     | 2.2.1 | The controlled system  | 33        |
|     | 2.2.2 | Different states of the system according to the control values $\ldots$ .  | 33        |
| 2.3 | Con   | trollability *   | <b>35</b> |
|     | 2.3.1 | Reachable sets   | 35        |
|     | 2.3.2 | Controllability of autonomous linear systems: Kalman's condition $\ . \ .$ | 37        |
| 2.4 | The   | Linear-Quadratic (LQ) problem  | 37        |
|     | 2.4.1 | The Linear model $\ldots$  | 38        |
|     | 2.4.2 | The Quadratic cost function  | 39        |

<sup>1</sup>Recall that the sections indicated with a \* are "to go further sections". They can be skipped in a first reading or if the reader is not particularly interested in deeper mathematical basis, mathematical proofs.

| <b>2.5</b> | The  | Pontryagin principle & the Hamiltonian  | 40        |  |  |
|------------|--|---|-----------|--|--|
|            | 2.5.1  | Existence and uniqueness of the solution  | 40        |  |  |
|            | 2.5.2  | The Pontryagin principle in the LQ case   | 42        |  |  |
|            | 2.5.3  | The Hamiltonian   | 43        |  |  |
|            | 2.5.4  | Examples  | 45        |  |  |
| 2.6        | Feed   | back law and the Riccati equation $*$   | <b>45</b> |  |  |
| 2.7        | A fe   | w remarks   | 47        |  |  |
| 2.8        | 2.8 Numerical methods for optimal control problems |   |           |  |  |
|            | 2.8.1  | Two classes of numerical methods $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$                | 47        |  |  |
|            | 2.8.2  | Direct methods  | 47        |  |  |
|            | 2.8.3  | Indirect methods: boundary value problem and shooting methods $~$   | 48        |  |  |
|            | 2.8.4  | Direct vs indirect methods $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$ | 49        |  |  |
| 2.9        | 2.9 The fundamentals equations at a glance         |   |           |  |  |

# 2.1 Introduction

In order to illustrate what is an optimal control problem let us give a simple example (extracted from a Wikipedia webpage).

Consider a car traveling on a straight line through a hilly road. The question is, how should the driver press the accelerator pedal in order to minimize the total traveling time? The control law refers specifically to the way in which the driver presses the accelerator and shifts the gears. The "system" or "direct model" is the car dynamics, and the optimality criterion is the minimization of the total traveling time. Control problems usually include additional constraints. For example the amount of available fuel might be limited, the speed limits must be respected, etc.

The "cost function" is a mathematical functional giving the traveling time as a function of the speed, initial conditions and parameters of the system.

Another optimal control problem can be to find the way to drive the car so as to minimize its fuel consumption, given that it must complete a given course in a time not exceeding some amount. Yet another control problem can be to minimize the total monetary cost of completing the trip, given assumed monetary prices for time and fuel.

Let us notice that the optimal control problem may have multiple solutions (i.e. the solution may not be unique), if it exists...

# 2.2 Illustrative example

The aim of this section is simply to illustrate (or recall) how the response of an apparently gentle ODE equation can widely differ depending on the source term (or the applied command).

#### 2.2.1 The controlled system

This classical example is treated in detail in [42], we refer to this book for more details



Figure 2.1: The spring example.

Let us consider a spring - mass system (see Figure 2.1). The mass m is submitted to a force f(x), which supposed to be equal to: $-[k_1(x-L) + k_2(x-L)^3]$ . L is the spring length at rest,  $k_*$  are the spring stiffness parameters.

We will apply an external force:  $u(t)\vec{i}$  (depending on time t).

Given the external force, the spring position x(t) is described by the equation:

$$mx''(t) + k_1(x(t) - L) + k_2(x(t) - L)^3 = u(t)$$
 for  $t \ge 0$ 

It is the *direct model*. It is an ordinary differential equation, linear if  $k_2 = 0$ , non-linear if not. The initial condition,  $x(0) = x_0, x'(0) = y_0$ , is given.

The problem we consider is as follows. Given an initial state  $(x_0, y_0)$ , find the "best" external force u(t) such that x(t) = L in a minimal time, under the constraint  $u(t) \leq 1$  (the external force is bounded).

u(t) is the *control* of the problem.

For sake of simplicity (and without any change of the nature of problem), we set:  $m = 1, k_1 = 1, L = 0$ . Then in the phase space (x, x'), the direct model reads:

$$X'(t) = AX(t) + f(X(t)) + Bu(t) , \quad X(0) = X0.$$

where X(t) = (x, x')(t) is the *state* of the system,

$$A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} , \quad B = (0, 1)^T , \quad f = (0, -k_2 x^3)^T$$

The model is a 1st order differential system, linear if  $k_2 = 0$ , non-linear if not.

## 2.2.2 Different states of the system according to the control values

1) In the case u(t) = 0 (no external action), the solution of the system satisfies:

$$x''(t) + x(t) + 2x(t)^3 = 0$$

It is a particular case of the Duffing equation. Its solution x(t) satisfies:  $x(t)^2 + x(t)^4 + x'(t)^2 = cste$ . All solutions are periodic and can be represented by an algebraic curve. The phase diagram and the



Figure 2.2: The spring example. State of the system without control, from [42].

trajectory are plotted on Fig. 2.2.

2) In the case u(t) = -x'(t), i.e. we try to damp the spring, the model reads:

$$x''(t) + x(t) + 2x(t)^3 + x'(t) = 0$$

Its numerical solution is computed, the corresponding phase diagram and the trajectory are plotted on Fig. 2.3.



Figure 2.3: The spring example. State of the system with damp control, from [42].

Using the Lyapunov theory, it can be shown that the origin is asymptotically stable. The spring position and the velocity reach the equilibrium position in infinite time, not in finite time. Therefore this control do not satisfies the original problem.

3) Let us set  $u(t) = -(x(t)^2 - 1)x'(t)$ . With this particular control expression, the model reads:

$$x''(t) + x(t) + 2x(t)^{3} + (x(t)^{2} - 1)x'(t) = 0$$

It is a particular case of the Van der Pol equation.

Two different solutions are computed and plotted on Fig. 2.4 (phase diagram and trajectories).

Using the Lyapunov theory, it can be proved that there exists a periodic solution which is attractive (Figure 2.4), [42].



Figure 2.4: The spring example. State of the system with a given control (Van der Pol like equation), from [42].

These three examples simply illustrate the wide range of behaviors which can be obtained from a very simple differential system by changing the control term.

# 2.3 Controllability \*

This is a "To go further" section...

This section presents the basic concept of controllability of a dynamical system and the Kalman's condition. These concepts are important in a context of an automatic course, and more secondary in the present Data Assimilation concept. Therefore this section may be skipped if the reader wants to focus on VDA only.

## 2.3.1 Reachable sets

Let  $x_u(t)$  be the solution of the ODE system corresponding to a given control u(t) ( $x_u$  is supposed to exist and to be unique).

**Definition 2.1.** The set of reachable states from the initial state  $x_0$  in time T > 0, is defined by:

 $Acc(x_0, T) = \{x_u(T), with \ u \in L^{\infty}([0, T], \Omega)\}$ 

with  $\Omega$  a compact subset of  $\mathbb{R}^m$ . We set:  $Acc(x_0, 0) = x_0$ .

**Theorem 2.2.** Let us consider the following 1st order linear ODE system in  $\mathbb{R}^n$ :

$$x'(t) = A(t)x(t) + B(t)u(t) + r(t)$$

Then  $Acc(x_0, t)$  is compact, convex, and varies continuously with  $t, t \in [0, T]$ .

Partial proof. We prove only that  $Acc(x_0, t)$  is convex in the case  $u \in L^{\infty}([0, T], \Omega)$  with  $\Omega$  convex. We refer to [42] for the proof of the remaining results.

Let  $\Omega$  be convex. Let  $x_1^t, x_2^t \in Acc(x_0, t)$ ; we denote by  $u_i$  the corresponding controls. We have:

$$x_i(0) = x_0$$
; ;  $x'_i(s) = A(s)x_i(s) + B(s)u_i(s) + r(s)$ ;  $i = 1, 2$ 

We denote:  $x_i^t \equiv x_i(t)$ . Let  $\lambda \in [0, 1]$ . We seek to prove that:  $(\lambda x_1^t + (1 - \lambda)x_2^t) \in Acc(x_0, t)$ . We solve the 1st order linear ODE, we obtain:

$$x_i^t \equiv x_i(t) = M(t)x_0 + M(t)\int_0^t M(s)^{-1}B(s)u_i(s)ds$$

with  $M(t) \in M_{n,n}(\mathbb{R})$  such that: M'(t) = A(t)M(t), M(0) = Id. Note that if A(t) = A constant then  $M(t) = \exp tA$ .

We set:  $u(t) = \lambda u_1(t) + (1 - \lambda)u_2(t)$ . Since  $\Omega$  is convex,  $u(t) \in \Omega$ ; also  $u \in L^{\infty}([0, T], \Omega)$ . Let x(t) be the solution associated to u(t). We have  $x(t) \in Acc(x_0, t)$  and:

$$x(t) = M(t)x_0 + M(t)\int_0^t M(s)^{-1}B(s)u(s)ds$$

In other respect,

$$\lambda x_1^t + (1-\lambda)x_2^t = M(t)x_0 + \lambda[M(t)\int_0^t M(s)^{-1}B(s)u_1(s)ds] + (1-\lambda)[M(t)\int_0^t M(s)^{-1}B(s)u_2(s)ds]$$
  
= x(t)

Thus  $(\lambda x_1^t + (1 - \lambda) x_2^t) \in Acc(x_0, t)$  and  $Acc(x_0, t)$  is convex.



Figure 2.5: Reachable set in the linear case, from [42].

#### A more simple case

Let us consider: r = 0 and  $x_0 = 0$ . The system becomes:

$$x'(t) = A(t)x(t) + B(t)u(t); \quad x(0) = x_0$$

and its (unique) solution reads:

$$x(t) = M(t) \int_0^t M(s)^{-1} B(s) u(s) ds$$

Thus x is linear in u. Then we have

**Proposition 2.3.** Let us consider the system: x'(t) = A(t)x(t) + B(t)u(t);  $x(0) = x_0$ , and  $\Omega = \mathbb{R}^m$ . Then for all t > 0, the reachable set Acc(0,t) is a vectorial subspace of  $\mathbb{R}^n$ . If we assume in addition B(t) = B constant, then for any  $0 < t_1 < t_2$ , we have:  $Acc(0,t_1) \subset Acc(0,t_2)$ .

We refer to [42] for the proof.

# 2.3.2 Controllability of autonomous linear systems: Kalman's condition

In the following we consider a first order linear autonomous ODE system in  $\mathbb{R}^n$ :

$$x'(t) = Ax(t) + Bu(t) + r(t)$$

with A and B independent of t. Furthermore, we set  $\Omega = \mathbb{R}^m$  (i.e. we do not constraint the control variable).

We say that the system is controllable at any time T if  $Acc(x_0, T) = \mathbb{R}^n$ ; it means that for any  $x_0$  and  $x_1$  in  $\mathbb{R}^n$ , there exists a control u such that  $x(0) = x_0$  and  $x(T) = x_1$ .

**Theorem 2.4.** Under the assumptions above, the first order linear autonomous ODE system is controllable at any time T if and only if the rank of matrix  $C = (B, AB, ..., A^{n-1}B)$  is equal to n.

The matrix C is called the Kalman matrix and the condition rank(C) = n is called the Kalman controllability condition.

Since the Kalman's condition does not depend neither on T nor on  $x_0$ , then the first order linear autonomous system is controllable at any time from any initial condition. (R. Kalman is a hongarian-american mathematician-electrical engineer, born in 1930).

Sketch of the proof. First, we proves that rank(C) = n if and only if the following linear application  $\Phi$  is surjective.

$$\Phi: L^{\infty}([0,T],\mathbb{R}^m) \mapsto \mathbb{R}^n ; \quad u \mapsto \int_0^T \exp\left(((T-t)A)Bu(t)dt\right)$$

Thus  $\Phi: L^{\infty}([0,T], \mathbb{R}^m) = \mathbb{R}^n$ . In other respect, given u, we have:

$$x(T) = \exp{(TA)x_0} + \int_0^T \exp{((T-t)A)(Bu(t) + r(t))dt}$$

Thus the reachable states is:

$$Acc(x_0,T) = \exp{(TA)x_0} + \int_0^T \exp{((T-t)A)r(t)dt} + \Phi(L^\infty) = \mathbb{R}^n$$

Therefore, the system is controllable.

We refer to [42] for the full proof of the theorem

These results give an overview of notions of controllability for *linear systems*. We refer e.g. to [42] for more details and the study of non-linear systems.

# 2.4 The Linear-Quadratic (LQ) problem

This section is important in all the following since it introduces the reference inverse problem: the Linear-Quadratic (LQ) problem.

## 2.4.1 The Linear model

Let A, B and s be three mappings defined from I = [0, T] into  $M_{n,n}(\mathbb{R}), M_{n,m}(\mathbb{R})$  and  $\mathbb{R}^n$  respectively. The three mappings are assumed to be bounded i.e.  $L^{\infty}(I)$  (this assumption could be relaxed since locally integrable would be sufficient). We consider the following linear 1st order ODE.

$$\begin{cases} \text{Given } u(t), \text{ find } x(t) \text{ such that:} \\ x'(t) = A(t)x(t) + B(t)u(t) + s(t) \text{ for } t \in I = [0, T] \\ \text{with the initial condition: } x(0) = x_0. \end{cases}$$

$$(2.1)$$

The function u(t) is assumed to be in  $L^{\infty}(I)$ .

In other words, we consider a phenomena which can be modelled by this linear ODE. (2.1) is the *direct* model, x(t) is the *the state* of the system, and u(t) will be the *control* of the system.

Existence and uniqueness of the solution. We known (classical theorem of existence) that (2.1) has one and only one solution x(t), x(t) continuous from I into  $\mathbb{R}^n$ .

Let us remark that we known an explicit expression of x in an integral form; the unique solution x(t) depends on u(t) of course.

The general optimal control problem reads as follows: find a control u(t) minimizing a given criteria (cost function) j(u).

Often, additional constraint on the state are imposed; e.g. x(t) must go from  $x_0$  (the initial state) to  $x_1$  in finite time.

A controllability problem would read as follows: given  $x_1 \in \mathbb{R}^n$ , find a control u(t) such that x(t) goes from  $x_0$  (at initial time) to  $x_1$  in finite time.

Courses of "Automatic" or "Optimal command" may focus on the controllability problems for linear systems. In the present course, we need to consider optimal control problems only, not controllability problems.

**Remark 2.5.** - Since the previous spring academic example is modeled by a linear ODE in the case  $k_2 = 0$ , the results which follow will apply to such a case.

- Among the historical optimal control problems, let us cite the so-called brachistochrone problem (end of the 17th century, from the greek words "brakhisto" (shorter) and "chronos" (time)) where the system is a simple 1d gravitational dynamic without friction, and j is time, from  $x_0$  to  $x_1$ .

Let us define the "Model Operator"  $\mathcal{M}(u)$  as:

$$\mathcal{M}: u(t) \mapsto x^u(t) M = L^2(I, \mathbb{R}^m) \to C^0(I, \mathbb{R}^n)$$

We have the important property:

**Proposition 2.6.** The model operator  $\mathcal{M}(u(t))$  is affine for all t in [0, T].

Proof. The proof is similar to Theorem 2.2 one; however let us recall its principles.

Let x(t) be the (unique) solution associated to u(t):  $x(t) = \mathcal{M}(u(t))$ .

In the present case, the model is a linear EDO; therefore one can write an explicit expression of its solution.

Indeed, we have:

$$x(t) = M(t)x_0 + M(t)\int_0^t M(s)^{-1}B(s)u(s)ds$$

with  $M(t) \in M_{n,n}(\mathbb{R})$  such that: M'(t) = A(t)M(t), M(0) = Id. Note that if A(t) = A constant then  $M(t) = \exp tA$ . The result follows straightforwardly from the explicit expression for

The result follows straightforwardly from the explicit expression fo x(t).

## 2.4.2 The Quadratic cost function

The choice of the cost function to be minimized is part of the problem definition. Since it will be minimized, convexity properties are expected.

Moreover if using computational gradient-based methods, differentiability properties will be required too. As a consequence, quadratic cost functionals are almost always chosen. Somehow, this set the inverse control problem as a least-square problem.

Recall that the quadratic property implies differentiability and strict convexity.

A typical generic "quadratic" objective function reads as follows:

$$J(u;x) = \frac{1}{2} \int_0^T \|x(t)\|_W^2 dt + \frac{1}{2} \int_0^T \|u(t)\|_U^2 dt + \frac{1}{2} \|x(T)\|_Q^2$$
(2.2)

The cost function j(u) is next defined from the objective function J(u; x) by:

$$j(u) = J(u; x^u)$$
(2.3)

with  $x^u$  the unique solution of the (linear) model, given u.

The three terms are the time averaged cost of the state, the control and the final state (in the metrics W, U and Q respectively).

This functional j is a multi-objective cost functional.

Note that the objective function J(u; x) is quadratic in its primal variables (u; x), but j(u) is not quadratic in u. However, it will be demonstrated later that j(u) is strictly convex (in u) since the model operator  $\mathcal{M}$  is affine.

The operators (matrices) Q, W and U are given in  $M_{n,n}(\mathbb{R})$ ,  $M_{n,n}(\mathbb{R})$  and  $M_{m,m}(\mathbb{R})$  respectively; they define some (semi-)norms. They are symmetric positive; moreover U is definite.

Note that Q, W must be symmetric positive, but definite is not necessary. Indeed, for example, one can have the cost function minimal for a vanishing control.

Recall that:  $\|\cdot\|_{\square}^2 = < \Box \cdot, \cdot >.$ 

Since the observation functional is quadratic and the model is linear, the natural functional space for control variable is  $M = L^2([0, T], \mathbb{R}^m)$ .

Let us point out that the a-priori natural space  $C^0([0,T],\mathbb{R}^m)$  is not a Hilbert space...

Let us recall that for a matrix M symmetric positive definite, in vertu of Courant-Fischer's theorem (Rayleigh's quotient), there exists a constant c > 0 such that:

$$\forall v \in \mathbb{R}^m, \quad \|v\|_M^2 \ge c\|v\|^2$$

In other words, the (linear) operator M is coercive, uniformly in time.

**The Linear-Quadratic (LQ) optimal control problem** The considered optimal control problem is a Linear-Quadratic (LQ) optimal control problem. It reads as follows: Given  $x_0$  and T, find u(t) minimizing the "quadratic" cost function j(u) defined by (2.3), with x(t) solution of the linear EDO (2.1).

**Remark 2.7.** The linear-quadratic problem seems a-priori idealistic; it is not so much. Indeed, observation functions are often defined quadratic in view to minimize differentiable functionals. Moreover, to better understand more complex non-linear problems, a good understanding of simplified linear model is often a necessary step.

In the LQ case, many instructive properties of the system can be written, both in a mathematical and numerical point of view.

For a sake of simplicity, in the direct model we consider the source term s(t) = 0.

# 2.5 The Pontryagin principle & the Hamiltonian

In this section, the fundamental concepts and results are introduced.

## 2.5.1 Existence and uniqueness of the solution

First, let us prove the existence and uniqueness of the optimal control solution in the LQ case.

**Theorem 2.8.** There exists an unique solution  $u \in M$  minimizing j(u), j(u) defined by (2.3), with the "constraint" x(t) solution of (2.1). In other words, it exists a unique optimal control u(t) and a corresponding trajectory x(t) to the LQ problem.

#### J. Monnier

*Proof.* A) Existence ("To go further").

It is based on the convergence of minimizing sequence (calcul of variations, D. Hilbert, 1900 approx.). Step 1). The cost function is bounded by below:  $\inf\{j(u), u \in M\} > -\infty$  since  $j(u) \ge 0$ . There exists a minimizing sequence  $(u_n)$  defined for all  $t \in [0, T]$ ; i.e. a sequence such that:

$$\lim_{n \to +\infty} j(u_n) = \inf\{j(u), u \in M\}$$

(As a matter of fact,  $\forall n \in \mathbb{N}$ ,  $\exists v_n$  such that:  $m \leq j(v_n) < m + \frac{1}{n}$ ).

Step 2) There exists  $\alpha > 0$  such that:  $j(u) \ge \alpha ||u||_M^2$ . Thus, the minimizing sequence  $(u_n)$  is bounded in M. Hence there exists a sub-sequence  $(u_{n_k})$  which converges weakly to a control u in M:

$$u_{n_k} \to u \text{ in } L^2(I)$$

Step 3) Let us denote by  $x_n$  (resp. x) the state associated to  $u_n$  (resp. u). The system (2.1) is a first order linear O.D. E. (and with s(t) = 0); we known an explicit expression of the solution:

$$\forall t, \ x_n(t) = M(t)x_0 + M(t) \int_0^t M(s)^{-1} B(s) u_n(s) ds$$
(2.4)

with  $M(t) \in M_{n,n}(\mathbb{R})$  such that: M'(t) = A(t)M(t), M(0) = Id. (If A(t) = A constant then  $M(t) = \exp(tA)$ ).

Similarly, we have:  $\forall t, x(t) = M(t)x_0 + M(t) \int_0^t M(s)^{-1} B(s)u(s) ds$ . Thus, we obtain that the sequence  $(x_n)_n$  converge to x.

Passing to the limit in (2.4), we obtain the existence of  $x_u$ , solution corresponding to u.

Step 4) It remains to prove that u minimizes j. Since  $u_n \rightarrow u$  in  $L^2$ , since j is continuous hence lower semi-continuous, we have (by definition):

$$j(u) \le \lim_n \inf j(u_n)$$

and necessarily  $j(u) = \inf_{v \in M} j(v)$ .



Figure 2.6: A function j lower semi-continuous at  $u_0$ : for u close to  $u_0$ , j(u) is either close to  $u_0$  or lower than  $u_0$ .

In other words, u(t) minimizes j(u)  $(j(u) = \min_{v \in M} j(v))$  and the corresponding state (trajectory)  $x_u$  is an optimal trajectory.

B) Uniqueness. Let us prove that j(u) is strictly convex i.e.  $\forall (u_1, u_2) \in M^2, \forall t \in ]0, 1[$ ,

$$j(tu_1 + (1-t)u_2) < tj(u_1) + (1-t)j(u_2)$$

unless  $u_1 = u_2$ .

For all t,  $||u(t)||_U$  is a norm hence convex but not strictly convex... Proof of this assertion: the triangle inequality.

However, in a Hilbert space, the square of a norm (eg  $||u(t)||_U^2$ ) is strictly convex, see e.g. [6] chap. 10, p118.

In vertu of Lemma 2.6, the model operator  $\mathcal{M}$  is affine therefore convex. In other respects  $\|.\|_W$  and  $\|.\|_Q$  are semi-norms hence convex.

Finally the cost function j(u) is strictly convex and the uniqueness follows straightforwardly.

Note that the strict convexity of  $j(\cdot)$  is due to the quadratic term  $||u(t)||_U^2$ . It may be viewed as a Tykhonov regularization term.

Exercice 2.9. Prove the uniqueness

Solution: Let  $u_1$  and  $u_2$  be such that:  $j(u_k) = \inf_{v \in M} j(v), \ k = 1, 2.$ We have:  $j(tu_1 + (1-t)u_2) < tj(u_1) + (1-t)j(u_2).$ Hence:  $j(tu_1 + (1-t)u_2) < \inf_{v \in M} j(v)$  unless  $u_1 = u_2$ , which must be the case.

**Remark 2.10.** In the autonomous case (A and B constant), we have:

$$||x'(t)|| \le ||A|| ||x(t)|| + ||B|| ||u(t)|| \le cst(||x(t)||^2 + ||u(t)||^2)$$

Then, if all hypothesis are satisfied in  $I = [0, +\infty$  then x'(t) is in  $L^1(I)$  and necessarily the minimizing trajectory x(t) tends to 0 when t tends to  $+\infty$ .

### 2.5.2 The Pontryagin principle in the LQ case

(L. Pontryagin, russian mathematician, 1908-1988).

In the case of a non-linear state equation, the cost function is a-priori non-convex and the Pontryagin minimum principle (also called maximum principle) is a necessary condition of optimality. In the LQ case, the Pontryagin minimum principle is a necessary and sufficient condition of optimality. It introduces the Hamiltonian<sup>2</sup>, the adjoint equation, and the Hamiltonian minimisation/maximisation.

In the present LQ case, the Pontryagin maximum principle reads:

<sup>&</sup>lt;sup>2</sup>The so-called "Hamiltonian" in control theory is a particular case of the Hamiltonian in mechanics; it is inspired by the Lagrangian you have studied in (differentiable) optimization.

**Theorem 2.11.** The trajectory x(t) associated with the control u(t), is optimal for the LQ problem (system (2.1)(2.3)) if there exists an adjoint vector p(t) which satisfies:

$$p'(t) = -p(t)A(t) + x(t)^T W(t) \text{ for almost } t \in [0, T]$$
(2.5)

with the final condition:  $p^T(T) = -Qx(T)$ . (p is here by convention a line vector). Furthermore, an optimal control u expression reads:

$$u(t) = U(t)^{-1}B(t)^{T}p(t)^{T} \text{ for almost } t \in [0, T]$$
(2.6)

The proof is based on 'calculus of variations'. The calculations are similar to those for the PDE case, see Theorem 3.21 in next chapter. The proof will be detailed in the case of general non-linear elliptic PDE case in next chapter. As a consequence we here refer to [42] for the present proof.  $\Box$ 

Note the theorem provides an expression of the optimal control, denoted by  $u^*$ , depending on the adjoint p, but not explicitly on the state x. However, p depends on x through the adjoint equation. An explicit expression of the optimal control  $u^*$  in function of x will be presented later: this is the so-called *feedback law*.

#### Remark 2.12.

If the term in x(T) in the cost function is more general and reads g(x(T)), with g any function defined from  $\mathbb{R}^n$  into  $\mathbb{R}$ ,  $C^1$  and convex. That is:

$$j(u) = \frac{1}{2} \int_0^T \|x(t)\|_W^2 dt + \frac{1}{2} \int_0^T \|u(t)\|_U^2 dt + g(x(T))$$
(2.7)

Then the final condition of the adjoint reads:

$$p(T) = -\nabla g(x(T))$$
(2.8)

#### Remark 2.13.

a) In case of an infinite time interval  $(T = +\infty)$ , the final condition becomes:  $\lim_{t \to \infty} p(t) = 0$ . b) If the system contains a non vanishing source term s(t), the same result holds.

## 2.5.3 The Hamiltonian

(W. Hamilton, 1805-1865, Irish physicist, astronomer and mathematician). In this section, the crucial concept of Hamiltonian is introduced.

Let us recall the direct model with no source term (s(t) = 0):

$$\begin{cases} \text{Given } u(t), \text{ find } x(t) \text{ such that:} \\ x'(t) = A(t)x(t) + B(t)u(t) \text{ for } t \in (0,T) \\ \text{with the initial condition: } x(0) = x_0. \end{cases}$$
(2.9)

and the cost function expression:

$$j(u) = \frac{1}{2} \int_0^T \|x(t)\|_W^2 dt + \frac{1}{2} \int_0^T \|u(t)\|_U^2 dt + \frac{1}{2} \|x(T)\|_Q^2$$
(2.10)

In this case the (control) Hamiltonian H reads:

$$H(x, p, u) = p(Ax + Bu) - \frac{1}{2}(||x||_W^2 + ||u||_U^2)$$
(2.11)

with  $H: \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ .

In short, the necessary and sufficient condition of the LQ problem solution (Theorem 2.11) reads as: the Hamiltonian must be maximized.

Then Theorem 2.11 re-reads: the following *relations* have to be satisfied.

$$\begin{cases} x'(t) = \partial_p H(x, p, u) = Ax + Bu \\ -p'(t) = \partial_x H(x, p, u) = pA - x^T W \\ 0 = \partial_u H(x, p, u) & \text{since } pB - u^T U = 0 \end{cases}$$
(2.12)

with the final condition:  $p(T) = -\nabla g(x(T)) = -x(T)^T Q$ . (Recall, p is a line vector in  $\mathbb{R}^n$ ).

This is the so-called *optimality system*.

The first two equations are the state equation and the adjoint state equations (with the final condition depending on x(T)). It is the *Hamiltonian equations*. They are accompanied by the last equation which is the optimality condition on u (necessary and sufficient condition in the present LQ case). In general these three equations are fully coupled.

#### The Hamiltonian: the conserved quantity in time.

Indeed, for all t:

$$d_t H(x, p, u) = \partial_x H(x, p, u) x'(t) + \partial_p H(x, p, u) p'(t) + \partial_u H(x, p, u) u'(t)$$
  
=  $-p'(t) x'(t) + x'(t) p'(t) + 0$   
=  $0$ 

Therefore the mapping  $t \mapsto H(x, p, u)(t)$  is constant.

In some contexts, the Hamiltonian denotes the energy of the system; in such a case the state and adjoint equations constitute the so-called Hamiltonian equations.

Some links can be done between the Hamiltonian and the Lagrangian (functional introduced in next section). The reader may consult e.g. the on-line course [18] presenting the concepts of "Hamiltonian mechanics" and "Lagrangian mechanics".

### 2.5.4 Examples

#### The most simple Linear-Quadratic problem

Let us consider the trivial linear state equation :

$$x'(t) = u(t)$$
;  $x(0) = x_0$  given

Given the final time T, we seek to minimize the "quadratic" cost function :

$$j(u) = \int_0^T [\alpha_1 x^2(t) + \alpha_2 u^2(t)] dt$$

#### Exercice 2.14.

- 1) Write the optimality system.
- 2) Give an explicit expression of the optimal control.
- 3) Give an expression of the corresponding optimal trajectory.

#### A simple vehicle control problem

Let us consider the control of the 1d trajectory of a vehicle (without friction forces). We consider the following idealistic model:

$$x''(t) = -x'(t) + u(t) ; \quad x(0) = x'(0) = 0$$

Given the final time T, we seek to maximize the distance travelled with a minimum of energy consumed. This goal can be translated by defining one of the two following cost function to be minimized:

$$j_1(u) = -\alpha_1 x(T) + \alpha_2 \int_0^T u^2(t) dt$$
$$j_2(u) = -\alpha_1 x^2(T) + \alpha_2 \int_0^T u^2(t) dt$$

with  $\alpha_*$  given coefficients modeling the weight we want to give to each term. For a sake of simplicity, we set:  $\alpha_1 = \alpha_2 = 1/2$ .

#### Exercice 2.15.

1) Write the optimality system in both cases:  $j_1$  and  $j_2$ .

2) In the case  $j_1$ , give the explicit expression of the optimal control.

Next, indicate how to deduce the corresponding optimal trajectory.

Note that this problem is numerically solved in a programming practical. See the course Moodle page for details.

# 2.6 Feedback law and the Riccati equation \*

"To go further" section.

Riccati family (father and son), italian mathematicians, 18th century.

The optimality condition derived above gives an expression of the optimal control u(t) in function of the adjoint solution p(t). In some (simple) cases (e.g. the exercises above), it is possible to integrate explicitly the adjoint equation, resulting to an expression of u(t) depending on the state x(t), that is the so-called *feedback law or closed-loop control*. The latter is the law sough in automatic.

In general, the expression of u in function of x is far to be trivial... But in the LQ case, the feedback law is known: it is obtained by solving the Riccati's equation presented below.

Such feedback laws are not required for VDA; therefore the result below may be skipped for readers interested in VDA only.

### Feedback law in the LQ case

In the LQ case we have:

**Theorem 2.16.** Under the assumptions of the existence - uniqueness Theorem 2.8, the (unique) optimal control u writes as a feedback function (closed-loop control):

$$u(t) = K(t)x(t)$$
 with  $K(t) = U(t)^{-1}B^{T}(t)E(t)$ 

where the matrix E(t),  $E(t) \in M_n(\mathbb{R})$ , is solution of the Riccati equation:

$$E'(t) = W(t) - A^{T}(t)E(t) - E(t)A(t) - E(t)B(t)U^{-1}(t)B^{T}(t)E(t) \quad \forall t \in (0,T)$$

with the final condition: E(T) = -Q.

For all t, the matrix E(t) is symmetric. Furthermore, since Q and W are positive definite, E(t) is positive definite too.

We refer to [42] for the proof.

This result gives an expression of the optimal control u in function of the state x.



Figure 2.7: Optimal control of a system: closed loop

# 2.7 A few remarks

What happens in non-linear cases ? In practice, optimal control problems are often non-linear, thus the results presented for the LQ problem do not apply directly, in particular the analytic expression of the optimal control. Nevertheless, a good understanding of the LQ solution structure, the strict convexity of the resulting cost function in particular, is useful to tackle non-linear problems. Then, to solve a *non-linear optimal control problem* (in the sense computing an optimal control u and the optimal trajectory  $x^u$ ), the option consists to employ *numerical methods*.

**Connection between the control of ODEs and PDEs** For Partial Derivatives Equations (PDE) systems (see next chapter), the Pontryagin principle does not apply anymore and the feedback laws are generally not known. These questions are still an active mathematical research field and some recent results exist (eg. the Riccati-like equations for the Navier-Stokes fluid flows equations at low Reynolds number etc).

Nevertheless, for PDE systems, it remains possible to write equations characterizing the optimal control solution: it is the *optimality system* based on the *adjoint equations*. This approach is developed for non-linear elliptic PDE systems in next chapter.

# 2.8 Numerical methods for optimal control problems

We present very briefly in this section two types of approaches to numerically solve optimal control problems.

## 2.8.1 Two classes of numerical methods

Basically, it exists two classes of numerical methods in optimal control: direct methods and indirect methods.

- . Direct methods consist in discretizing the state x and the control u, and to reduce the problem to a standard optimization problem. The approach is next based on (nonlinear) programming algorithms such as e.g. the Sequential Quadratic Programming (SQP) algorithm.
- . The indirect methods consist in numerically solving a problem resulting from the maximum principle, that is the necessary first-order optimality conditions based on the Hamiltonian. The resulting problem is a boundary value problem which is solved by a "shooting" method.

Methods mixing the two classes of approaches exist too.

# 2.8.2 Direct methods

Direct methods consist to write the optimal control problem equations in a discrete form, next to solve the optimization problem by standard differentiable optimization algorithms.

Let us write in a discrete form the general LQ problem (2.1)(2.3). We consider a time grid  $(t_0, \dots, t_N)$  of constant time step h. Then in its discrete form, the problem to be solved is:

$$\operatorname{argmin}_{(u_1,\cdots,u_M)} j(u_1,\cdots,u_M) \tag{2.13}$$

with:  $j(u_h) = J(u_h; x_h^{u_h}), u_h = (u_1, \dots, u_M).$  $u_i$  denotes an approximation of  $u(t_i)$  and  $x_h^{u_h}$  the numerical solution of the state equation given  $u_h$ . The numerical state  $x_h$  is obtained by performing a numerical scheme such as eg Runge-Kutta 4. For a sake of simplicity, let us consider the basic explicit Euler method. In this case, the discrete state system reads:

$$x_{i+1} = h_i(1+A_i)x_i + h_i B_i u_i \text{ for } i = 0 \cdots (N-1); \quad x_0 = x_0$$
(2.14)

Finally, the finite dimensional optimization problem reads:

 $\min_{u_h = (u_1, \cdots, u_M)} j(u_h)$  under the N constraints (2.14) on  $x_h$ (2.15)and under potential (equality-inequality) constraints on  $u_h$ .

This a finite dimensional optimization problem with constraints. This can be solved by e.g. classical Sequential Quadratic Programming (SQP) methods. Numerous variants of formulations can be considered.

**Recall on the SQP methods** Sequential Quadratic Programming (SQP) methods denote iterative methods for constrained nonlinear optimization problems. The objective function and the constraints are supposed to be  $C^2$ .

The SQP methods principle is as follows. One solves a sequence of optimization subproblems; each of them optimizes a quadratic representation of the objective function, under the constraints which are linearized.

If the problem is unconstrained, then the method reduces to the Newton method : the optimum is such that it makes vanish the gradient.

If the problem has equality constraints and no inequality constraints, then the method is equivalent to apply Newton's method to the first-order optimality conditions (KKT condition).

#### Indirect methods: boundary value problem and shooting meth-2.8.3ods

Let us consider the LQ problem (2.1)(2.3); the final time T is given. The Pontryagin principle states that the unique optimal control  $u^*$  solving this problem makes satisfy the optimality system (3.30).

The shooting method principle By setting z(t) = (x(t), p(t)), the two first equation of the optimality system (the state and the adjoint equations) can be written as a first order dynamical system :

$$z'(t) = F(u(t); z(t))$$

with the I.C.  $z(0) = (x_0, \cdot)$  and the Final Condition (F.C.)  $z(T) = (\cdot, p_T)$ . One can write the IC and FC as: R(z(0), z(T)) = 0, see e.g. [42] for details.

Then, the problem to be solved reads:

$$\begin{cases} z'(t) = F(u(t); z(t)) \\ R(z(0), z(T)) = 0 \end{cases}$$
(2.16)

This is a *boundary value problem*, and not a "simple" Cauchy problem. Let us denote by  $z(z_0;t)$  the solution of the Cauchy problem:  $z'(t) = F(u(t); z(t)), z(0) = z_0$ . We set:  $G(z_0) = R(z_0, z(z_0;T))$ . Then, Problem (2.16) consists to solve:

$$G(z_0) = 0 (2.17)$$

This problem can be numerically solved by the Newton-Raphson method.

The multiple step version An efficient version of the shooting method above consists to split the time interval [0, T] into sub-intervalls  $[t_p, t_{p+1}]$  and compute the values  $z(t_p)$  at the start of each sub-intervall. Conditions of continuity at eeach sub-interval boundary have to be imposed. The multiple shooting method is more stable than the basic method above.

The reader may refer to [42] and references therein for more details on shooting type methods.

## 2.8.4 Direct vs indirect methods

Let us briefly present the pros and cons of each approach. For details, the author may refer to e.g. [42] and references therein.

| Méthodes directes                    | Méthodes indirectes                        |
|--------------------------------------|--|
| mise en oeuvre simple,               | connaissance a priori de la structure      |
| sans connaissance a priori           | de la trajectoire optimale                 |
| peu sensibles au choix de la         | très sensibles au choix de la              |
| condition initiale                   | condition initiale                         |
| facilité de la prise en compte       | difficulté théorique de la prise en compte |
| de contraintes sur l'état            | de contraintes sur l'état                  |
| contrôles (globalement) optimaux     | contrôles (localement) optimaux            |
| en boucle fermée                     | en boucle ouverte                          |
| précision numérique basse ou moyenne | très grande précision numérique            |
| efficaces en basse dimension         | efficaces en toute dimension               |
| gourmandise en mémoire               | calculs parallélisables                    |
| problème des minima locaux           | petit domaine de convergence               |

Figure 2.8: Pros and cons of direct methods vs indirects methods (shooting methods) to numerically solve an optimal control problem. Extracted from [42]

# 2.9 The fundamentals equations at a glance

#### The Linear model

The linear model (without source term s(t)) reads:

$$\begin{cases} \text{Given } u(t), \text{ find } x(t) \text{ such that:} \\ x'(t) = A(t)x(t) + B(t)u(t) \text{ for } t \in I = [0, T] \\ \text{with the initial condition: } x(0) = x_0. \end{cases}$$

$$(2.18)$$

The model operator  $\mathcal{M}(u)$  is defined as:  $\mathcal{M}(u) = x^u$ . This operator  $\mathcal{M}(u(t))$  is here affine in u(t), for all t in [0, T].

**The cost functional** j(u) is chosen quadratic. It is defined from the "observation" functional J(u;x) as follows:  $j(u) = J(u;x^u)$ ,  $x^u$  the unique solution of the (linear) model, given u. We have here:

$$j(u(t)) = J(u(t), x^{u}(t)) = \frac{1}{2} \int_{0}^{T} \|x^{u}(t)\|_{W}^{2} dt + \frac{1}{2} \int_{0}^{T} \|u(t)\|_{U}^{2} dt + \frac{1}{2} \|x^{u}(T)\|_{Q}^{2}$$
(2.19)

#### The LQ optimal control problem is:

Given the I.C.  $x_0$  and the final time T, find u(t) such that  $x^u(t)$ , the unique solution of the state equation (2.18), minimizes the cost function j(u) defined by (2.19).

The Hamiltonian is the conserved quantity in time. Its expressions is:

$$H(x, p, u) = p(Ax + Bu) - \frac{1}{2}(||x||_W^2 + ||u||_U^2)$$
(2.20)

with  $H : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ , u(t) the control, x(t) the state of the system and p(t) the adjoint state, solution of the adjoint model.

The adjoint model reads:

$$p'(t) = -p(t)A(t) + x(t)^T W(t) \text{ for } t \in [T, 0]$$
(2.21)

with the final condition:  $p^T(T) = -Qx(T)$ . (p is a line vector).

#### The Pontryagin maximum principle states that if the control u(t) is defined as:

$$u(t) = U(t)^{-1}B(t)^T p(t)^T$$
 for almost  $t \in [0, T]$  (2.22)

then the state  $x^{u}(t)$  associated to this control u(t), is optimal for the LQ problem.

The Pontryagin maximum principle re-reads: the equations below have to be satisfied.

$$\begin{cases} x'(t) = \partial_p H(x, p, u) = Ax + Bu \\ -p'(t) = \partial_x H(x, p, u) = pA - x^T W \\ 0 = \partial_u H(x, p, u) \iff pB - u^T U = 0 \end{cases}$$
(2.23)

with the final condition:  $p(T) = -\nabla g(x(T)) = -x(T)^T Q$ . This is the **optimality system**.

# Chapter 3

# VDA: optimal control of PDE models, adjoint equations

This chapter presents optimal control basis of PDE systems leading to the Variational Data Assimilation (VDA) formulation. The presentation is focused on the computational aspects: it mainly addresses the derivation of the first order optimality system based on the adjoint equations. Two presentations of how to derive the adjoint equations are proposed: a formal one based on the introduction of the Lagrangian and a more rigorously one. The calculations are derived for a general non-linear elliptic PDE model therefore valid for very wide range of models representing stationary physical phenomena.

The global computational control algorithm (leading to VDA) is highlighted. The algorithm enables the fitting of the model with data. Examples are presented.

Some results and developments require relatively high mathematical skills; this is the case when addressing the differentiability of the PDE solution (with respect to the control variable) and when addressing the existence and uniqueness of the optimal control in the LQ case (PDE version). These mathematical developments are gathered in a dedicated section untitled "mathematical purposes".

Before reading the present chapter, the reader is invited to revise basic but crucial concepts in functional analysis, optimization and gradient-based minimization algorithms. Few of theses basic concepts are shortly recalled in Appendix A. Also, few exercises extracted from the literature are proposed. The reader has to be comfortable with all these basic concepts before going further.

In practice, the numerical resolution of the optimality system enables the *identification of model input parameters* (which are uncertain or unknown) and/or a better *model calibration*.

The outline of this chapter is as follows<sup>1</sup>.

## Contents

| 3.1        | The    | direct model $\ldots$                  | 53        |
|------------|--------|---|-----------|
|            | 3.1.1  | General non-linear stationary direct model $\ldots \ldots \ldots \ldots \ldots$   | 53        |
|            | 3.1.2  | Examples  | 53        |
| <b>3.2</b> | The    | objective and cost functions: misfit to data  | 55        |
| 3.3        | VDA    | $ \  \   {\rm formulation}  . \ . \ . \ . \ . \ . \ . \ . \ . \ .$  | 56        |
| 3.4        | On t   | the difficulty to compute the gradient for large size control   |           |
|            | varia  | ables   | 57        |
|            | 3.4.1  | Global minimization vs local one  | 57        |
|            | 3.4.2  | Relationship between the differential $j'(\cdot)$ and the gradient $\nabla j(\cdot)$                                    | 58        |
|            | 3.4.3  | How to compute the cost function gradient ?   | 58        |
| 3.5        | Mat    | hematical purposes $*$  | 59        |
|            | 3.5.1  | Differentiability of the cost function  | 59        |
|            | 3.5.2  | Existence and uniqueness of the optimal control in the LQ case  | 60        |
|            | 3.5.3  | LQ problems vs real-world problems $\hdots$   | 61        |
| 3.6        | Deri   | ving the optimality system from the Lagrangian  | 62        |
|            | 3.6.1  | Weak forms of the equations $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$                     | 62        |
|            | 3.6.2  | The Lagrangian & the optimality system  | 63        |
| 3.7        | Com    | puting the cost function gradient without adjoint model   | 64        |
|            | 3.7.1  | The Finite Difference gradient  | 64        |
|            | 3.7.2  | The straightforward "gradient" expression   | 65        |
|            | 3.7.3  | LTM-based gradient  | 65        |
| 3.8        | Rigo   | rous derivation of the adjoint equations and the correspond-  |           |
|            | ing g  | $gradient expression \ldots \ldots$ | 66        |
|            | 3.8.1  | The adjoint-based gradient theorem $\hfill \ldots \ldots \ldots \ldots \ldots \ldots \ldots$                            | 67        |
|            | 3.8.2  | The optimality system   | 69        |
|            | 3.8.3  | Gradient components: in weak or classical forms ? *   | 70        |
| 3.9        | The    | VDA algorithm   | 71        |
| 3.10       | ) The  | fundamental equations at a glance   | <b>74</b> |
| 3.1        | l Anot | ther example: control of an elastic membrane deformation $*$ .  | 76        |
| 3.12       | 2 How  | to assess your adjoint computational code & gradient *  | 77        |
|            | 3.12.1 | The scalar product test $\ldots \ldots \ldots$           | 77        |
|            | 3.12.2 | The gradient test   | 78        |
|            |        |   |           |

<sup>&</sup>lt;sup>1</sup>Recall that the sections indicated with a \* are "to go further sections". They can be skipped in a first reading or if the reader is not particularly interested in deeper mathematical basis, mathematical proofs.

# 3.1 The direct model

# 3.1.1 General non-linear stationary direct model

Let  $\Omega$  be a bounded domain ( $\Omega$  Lipschitz). Let U, a Hilbert space, be the controls space (U Banach space only is potentially enough, but we consider it here Hilbert). Let V, a Hilbert space, be the states space.

We consider the following *state equation* (called *direct model* too):

$$\begin{cases} \text{Given } u \in U, \text{ find } y \in V \text{ such that:} \\ A(u; y) = B(u) \text{ in } \Omega \\ \text{with boundary conditions on } \partial\Omega \end{cases}$$
(3.1)

where A is an elliptic operator (with respect to the state y), defined from  $U \times V$  into V' (dual of V).

 $A(\cdot; \cdot)$  is a priori non-linear, both with respect to the parameter u and with respect to the state y.

B is defined from U into  $V^{\prime}.$ 

**Assumption 3.1.** The state equation (3.1) is well posed in the sense that it has an unique solution  $y \in V$ , and this solution is continuous with respect to parameters (in particular with respect to u).

In the linear case (A elliptic operator linear with respect to y) the Lax-Milgram theorem might be the right framework to prove the existence-uniqueness, while the Mint-Browner theorem is useful for a large class of non-linear cases.

## Optimal control terminology: distributed control, boundary control

If the control u appears in the "bulk" (i.e. in the equation posed in  $\Omega$ ) then one say it is a *distributed control*.

If u appears on the boundary conditions only, then one say it is a *boundary control*.

# 3.1.2 Examples

Two simple linear examples based on second order elliptic operators are as follows.

## Example 1)

$$A(u; y) = -div(\lambda \nabla y)$$
 and  $B(u) = u$ 

with mixed boundary conditions: y = 0 on  $\Gamma_0$ ;  $-\lambda \partial_n y = \varphi$  on  $\partial \Omega / \Gamma_0$ , with  $\varphi$  given.  $\lambda \in L^{\infty}(\Omega), \lambda > 0$  a.e. The corresponding functional spaces are :  $U = L^2(\Omega), V = H^1_{\Gamma_0}(\Omega).$ 

The control u is spatially distributed, it is the source term in the Right Hand Side (RHS).

The state equation models a diffusion phenomena e.g. heat diffusion in a structure, concentration in a fluid, the elastic deformation of a membrane under external force f = u, the electrostatic field in a conducting media, etc.

Here, u is a distributed control since defined in  $\Omega$ ; it represents an external force (since in the RHS).

#### Exercice 3.2.

a) Write the state equation (and recall the adequate functional spaces). Prove that it has one and only one (weak) solution in V. b) Prove that the unique solution y is continuous with respect to u i.e. the operator  $\pi : u \in U \mapsto y^u \in V$  is continuous.

#### Solution.

a) In vertu of Lax-Milgram theorem.

b) The inequalities of continuity and coercivity give the result.

#### Example 2)

A non linear version of the previous example is as follows.

$$A(u; y) = -div(\lambda(y)\nabla y)$$
 and  $B(u) = u$ 

with mixed boundary conditions (independent of u). In this (classical) case, the PDE is non linear due to the term  $\lambda(u)$ .

#### Example 3)

$$A(u; y) = -div(\lambda \nabla y)$$
 and  $B(u) = f$ 

A and B are independent on the control u, however the system is closed with the following mixed boundary conditions:

y = 0 on  $\Gamma_0$ ;  $-\lambda \partial_n y = u$  on  $\partial \Omega / \Gamma_0$ , with  $(\lambda, f)$  given in the adequate functional spaces.

The correct functional spaces are :  $U = L^2(\Omega), V = H^1_{\Gamma_0}(\Omega).$ 

In this case u is a boundary control, it represents the flux at boundary.

As previously, the corresponding state equation has one and only solution in V (in vertu of Lax-Milgram theorem), the inequalities of continuity and coercivity give the continuity of y with respect to u.

# 3.2 The objective and cost functions: misfit to data

If seeking to make fit "at best" (in the least-square sense) the model outputs with the available data  $z_d$ , it is natural to measure the discrepancy between both.

The objective function J(u; y) is supposed to be of class  $C^1$ . A natural choice is to define quadratic J(u; y) in the following form:

$$\|y - z_d\|_O^2 + \|u\|_B^2 \tag{3.2}$$

The operators O and B are symmetrical positive definite therefore defining norms (potentially semi-definite).

For example,  $\|.\|_B^2 = (B_{\cdot}, .)_U : \forall v \in U, \ (Bv, v)_U^2 \ge c \|v\|_U^2, \ c > 0.$ 

The data (measurements, observations)  $z^{obs}$ , are not necessarily the state of the system y. For example, one measures a signal intensity and frequency to actually measure a temperature. In this case, it is necessary to define an *observation operator* C which maps the state of the system y onto the observation space  $\mathcal{O}$ :

$$C: y \in V \mapsto z \in \mathcal{O} \tag{3.3}$$

( $\mathcal{O}$  is supposed to be a Hilbert space).

The observation operator may be a complex multi-scale model e.g. if the goal is to measure a 3D fluid dynamic from a sequence of 2D satellite images.

Then, we consider the objective function  $J, J: U \times V \to \mathbb{R}$ , decomposed as follows:

$$J(u;y) = J_{obs}(y) + \alpha_{reg} J_{reg}(u)$$
(3.4)

with  $J_{obs}$  the term of misfit with the observations and  $J_{reg}$  a regularization term. The scalar coefficient  $\alpha_{reg}$  is a given (positive) weight, balancing the two objective functions terms.

A typical expressions of  $J_{obs}$  is as follows:

$$J_{obs}(y) = \left\| C(y) - z_d \right\|_{O}^{2}$$
(3.5)

where:

- C is the (a-priori non-linear) observation operator,

-  $z_d$  a given "data-observation-measurement".

This term measures the discrepancy between the data / target  $z_d$  and the computed solution y, in the observation space  $\mathcal{O}$  (and using the metric  $\|\cdot\|_O$ ).

Classical quadratic regularization terms read as follows:

$$J_{reg}(u) = \left\| u - u_b \right\|_B^2 \tag{3.6}$$

with  $u_b$  is a "targeted" value of the control  $uu_b$  is called the "background" value or the "first guess" value too.

This Tykhonov regularization term  $||u - u_b||_B^2$  penalizes the discrepancy between the control u and the "background value"  $u_b$ , in norm B.

As a consequence,  $u_b$  is a (strong) *prior* of the inverse problem...

Another option is to define  $J_{req}$  as follows:

$$J_{reg}(u) = \left\| D^p u \right\|^2 \tag{3.7}$$

with  $D^p$  (with p = 1 or 2 in practice) denoting a differential operator eg the gradient or the Laplace operator.

This term  $||D^p u||^2$  enforces to find the optimal solution u in a regular (therefore smaller) space. Eg for p = 1,  $||D_p v||^2 = ||\nabla v||^2$ . Moreover recall that:  $\int ||\nabla u||^2 dx$  corresponds to the energy of the Laplace operator  $\Delta u$ ;  $\int ||\Delta u||^2 dx$  corresponds to the energy of the bi-Laplacian operator  $\Delta^2 u$ . Each term integral requiring/imposing a certain amount regularity.

Note that we have:  $J_{obs}: V \to \mathbb{R}$  and  $J_{reg}: U \to \mathbb{R}$ .

Next, the cost function to be minimized is defined from the observation function J as follows:

j(u) = J(u; y(u))where y(u) (also denoted  $y^u$ ) is the (unique) solution of the direct model (3.1). (3.8)

We have:  $j: U \to \mathbb{R}$ .

Note that the norm O has to be clarified; it will be defined from a covariance matrix (operator) modeling a-priori knowledge we may have on the solution. If it is simply defined as a diagonal matrix, the diagonal coefficients represent the a-priori confidence we have on each observation. This point is developed in next chapter.

In the following, one need to define the "model operator":

$$\mathcal{M} : u \in U \mapsto y^u \in V$$
  
with  $y^u$  the (unique) solution of the direct model (3.1). (3.9)

# 3.3 VDA formulation

Variational Data Assimilation (VDA) simply denotes the optimal control of a PDE system with the cost function measuring the discrepancy between data and the model outputs. Let  $U_{ad}$ , subset of U, be the admissible control set. Recall the observation functional:

$$J(u;y) = \left\| C(y) - z_d \right\|_O^2 + \alpha_{reg} \left\| u - u_b \right\|_B^2$$

The optimal control problem reads:

$$\begin{cases} \text{Find } u^* \in U_{ad} \text{ such that:} \\ j(u^*) = \min_{U_{ad}} j(u) \\ \text{with the cost function } j \text{ defined by } (3.8) . \end{cases}$$
(3.10)

Let us point out again that j depends on the unknown parameter (the control) u through the "model operator"  $\mathcal{M}(u) = y^u$ ;  $y^u$  the (unique) solution of the state equation (3.1).

Problem (3.10) can be re-read as:

$$\begin{cases} \text{Minimize } j(u) = J(u; y^u) \text{ in } U_{ad} \\ \text{under the "model constraint" (3.1)} \end{cases}$$
(3.11)

In other words, the problem is an optimization problem under the constraint "model must be satisfied". This point of view naturally leads to introduce the Lagrangian, see next paragraph.

**Exercice 3.3.** Let j be the cost function defined by  $j(u) = J(u; y^u)$  with J defined by (3.4). a) Write a sufficient condition to have j of class  $C^1$ .

b) Write an expression of the differential  $j'(u).\delta u$ , for all  $\delta u$ 

c) Is the cost function j(u) strictly convex ?

Discuss the answer depending on the observation operator C and the regularization term form. Hint: linear + quadratic implies strictly convex.

Note a few exercises on differential calculus is proposed in the supplementary material of this course.

**Exercice 3.4.** Propose the VDA formulation of a real-life like problem (e.g. based on the simple Laplace equation a simple advection-diffusion model).

# 3.4 On the difficulty to compute the gradient for large size control variables

## 3.4.1 Global minimization vs local one

To solve the optimization problem (3.11) few approaches are a-priori possible, depending either if the CPU time  $T_{cpu}^{j}$  required to evaluate the cost function j is reasonable or not...

If  $T_{cpu}^{j}$  is small (let say in fractions of seconds using your laptop or a super-computer, whatever), then one can adopt a global optimization approach based on stochastic algorithms (e.g. Monte-Carlo), or heuristic methods (e.g. genetic algorithms) or surface response approaches.

If the state equation is a PDE system then generally  $T_{cpu}^{j}$  is not reasonable enough. It may require minutes and even much more...

In this case, global optimization is not worth considering. Then, one adopt *local minimization* approaches based on *algorithms of descent* (see the recall sections in Appendix).

Therefore one need to calculate the expression of the differential  $j'(u).\delta u$  for any direction  $\delta u$ , in view to compute its discrete expression i.e. the cost function gradient.

To do so, of course, the functional  $j(\cdot)$  must be differentiable, see the previous section.

# 3.4.2 Relationship between the differential $j'(\cdot)$ and the gradient $\nabla j(\cdot)$

In the computational context, the state equation (direct model) has been discretized using an adequate numerical method (e.g. finite differences, finite elements, finite volumes).

Let us denote  $U_h$  the discrete control space with  $dim(U_h) = m$ . It means there is m discrete control variables.

The gradient  $\nabla j(u)$  to be computed,  $\nabla j(u) \in \mathbb{R}^m$ , is related to the differential j'(u) by the relation:

$$\langle \nabla j(u), \delta u \rangle_{\mathbb{R}^m} = j'(u) \cdot \delta u \text{ for all } \delta u \in U_h \subset \mathbb{R}^m$$

$$(3.12)$$

### 3.4.3 How to compute the cost function gradient ?

The size of the gradient  $\nabla j(u)$  equals m. Descent algorithms require scalar products of the form  $\langle \nabla j(u), \delta u \rangle$  for a potentially large number of directions  $\delta u$ ; typically m i.e. for each gradient component.

Then, the following question arises:

How to compute the scalar values  $\langle \nabla j(u), \delta u \rangle$  for a large number of directions  $\delta u$ ? (Case m large).

In the continuous formalism,  $j'(u) = \frac{dj}{du}(u) \in \mathcal{L}(U; \mathbb{R})$ . Then, the question re-reads as: Given  $u_0 \in U_{ad}$  and  $\delta u \in U$ , how to evaluate efficiently the differential value  $\frac{dj}{du}(u_0).\delta u$  for any  $\delta u \in U ~?$ 

This key question is addressed in a next section with the introduction of the *adjoint model*.

**Remark 3.5.** In the case the control variables includes few components e.g.  $u = (u_1, u_2)$ , then the "gradient"  $j'(u) = \frac{dj}{du}(u)$  reads as (recall it is an element of  $\mathcal{L}(U;\mathbb{R})$ ):

$$"j'(u) = \frac{dj}{du}(u) = \nabla j(u) = \left(\frac{\partial j}{\partial u_1}(u), \frac{\partial j}{\partial u_2}(u)\right)^T,$$
  
and :  $j'(u) \cdot \delta u = \frac{\partial j}{\partial u_1}(u) \cdot \delta u_1 + \frac{\partial j}{\partial u_2}(u) \cdot \delta u_2.$ 

**Exercice 3.6.** 1) Write the VDA formulation of your practical. 2) Write an expression of  $j'(u).\delta u$ 

# 3.5 Mathematical purposes \*

This is a "To go further section".

## 3.5.1 Differentiability of the cost function

In the following, we will need to differentiate the cost function j (with respect to its unique variable u). Thus, the following question is of main interest in order to address the optimal control problem:

Is the cost function (continuously) differentiable?

This question of differentiability is potentially difficult to answer for non-linear systems. For non-linear hyperbolic system in particular, the solution may be even not continuous with respect to the control variable u...

A useful result to address this question of differentiability is the *implicit function theorem*.

**Theorem 3.7.** (Implicit function theorem) Let us assume that: i) the operator A and L in the state equation (3.1) are  $C^1$  (i.e.  $A \in C^1(U \times V)$  and  $B \in C^1(U)$ ), ii) the linearized problem is well-posed (i.e. given  $(u_0, y_0)$  the linearized operator  $\partial_u A(u_0; y_0)$  is an isomorphism from V into V').

Then, the implicit function  $\mathcal{M} : u \mapsto y^u$ , with  $y^u$  is the (unique) solution of the state equation, is locally  $C^1$  (locally means it exists a neighborhood of  $u_0$  such that).

In short, in view to apply the implicit function theorem we need to verify that the state operators are  $C^1$  and the linearized problem is well-posed.

The implicit function is in fact the "model operator"  $\mathcal{M}(u)$ , see (3.9).

Example 3) We set:  $A(u; y) = -u \Delta y$ , B(u) = f, with mixed boundary conditions: y = 0 on  $\Gamma_0$ ;  $-u\partial_n y = \varphi$  on  $\partial\Omega/\Gamma_0$ , with  $\varphi$  given.

Here, u is a distributed control, it is the diffusivity coefficient of the material.

#### Exercice 3.8.

a) Write the corresponding state equation, and prove that it has we and only (weak) solution in V for u given in  $L^{\infty}(\overline{\Omega}), u > 0$  a.e..

b) Prove that the unique solution y is continuous and differentiable with respect to u (in the right functional space).

# 3.5.2 Existence and uniqueness of the optimal control in the LQ case

We refer to [28, 29] for basics on this subject.

#### LQ problem case

In the Linear-Quadratic (LQ) case, with the existence and uniqueness of the optimal control hold.

**Theorem 3.9.** Let us assume that the state equation linear and coercive. Let us consider the cost function defined from (3.4)-(3.6), and  $U_{ad}$  a closed convex subset of U. It exists a unique solution at the optimal control problem (3.10).

The state equation is said to be coercive in V if for all  $y \in V$ , for all  $u \in U$ , there exists  $\alpha > 0$  such that:

$$a(u; y, y) = \langle A(u; y), y \rangle_{V' \times V} \ge \alpha ||y||^2$$

Proof. The proof is similar to those of Theorem 2.8. Let us consider the expression J as in (3.4)-(3.6) with  $u_b = 0$ .

Step 0) First, let us reformulate the cost function expression as follows:

$$j(u) = \|C(y^u - y^0) + Cy^0 - z_d\|^2 + \|u\|_B^2$$

We set:

$$\pi(u,v) = (C(y^u - y^0), C(y^v - y^0))_Z + (Bu, v) \text{ and } D(v) = (z_d - Cy^0, C(y^v - y^0))$$

Since the mapping model  $\mathcal{M}$  is affine and continuous, the form  $\pi$  is bilinear symmetric in U and the form L is linear continuous in U.

Furthermore, the form  $\pi$  is coercive in U:

$$\forall u \in U, \ \pi(u, u) \ge c_0 ||u||^2, \ c_0 > 0$$

#### J. Monnier

The cost function reads:

$$j(u) = \pi(u, u) - 2L(u) + ||z_d - y^0||^2$$

Thus j is continuous and satisfies:

$$j(u) \ge c_0 \|u\|^2 - c_1 \|u\|$$
(3.13)

A) Proof of existence. It is based on the convergence of minimizing sequence (calculus of variations, D. Hilbert, 1900 a.c. approx.).

Step 1). Let  $(u_n)$  be a minimizing sequence:

$$j(u_n) \to_n \inf\{j(u), u \in U_{ad}\}$$

$$(3.14)$$

From (3.13)(3.14), we obtain:

 $||u_n|| \leq constant$ 

Hence there exists a sub-sequence  $(u_{n_k})$  which converges weakly to a control u in  $U_{ad}$ :

 $u_{n_k} \rightarrow u$  in  $U_{ad}$ 

Step 2).  $U_{ad}$  is a closed convex subset hence weakly closed. Hence  $u \in U_{ad}$ .

Step 3). Since the cost function j(u) is continuous (lower semi-continuous would be enough), we have:  $j(u) = \min_{v \in U_{ad}} j(v)$ . In other words, u is solution of the optimal control problem.

B) Uniqueness. The bilinear form  $v \mapsto \pi(v, v)$  is coercitive hence the cost function is strictly convex. Then, the uniqueness is a straightforward consequence of the strict convexity of j(u).

**Exercice 3.10.** Detail the proof of uniqueness. (It is similar to the previous proof for an ODE).

## 3.5.3 LQ problems vs real-world problems

Most of the control problems are not LQ problem since the model is generally not linear.

Moreover even if the direct model is linear then one almost never measure the model output everywhere in  $\Omega$ , that is at the numerical model grid scale ...

Indeed, in real-world problems, data  $z_d$  are often available at some locations only, or at larger scale than the model solution scale (eg satellite measurements of local dynamics flows or camera measurements inside a complex multi-scale material).

As a consequence, without imposing a regularization term (that is  $\alpha_{reg} = 0$ ), the objective function simplifies as:

$$J(y) \equiv J(u; y) = \left\| C(y) - z_d \right\|_O^2$$

In this case, the uniqueness of a global minimum is not guaranteed anymore.

One generally have the existence (without always being able to prove it) but the cost function may be not strictly convex in the whole set  $U_{ad}$ .

In practice, the cost function is often "weakly" convex, with "nearly flat valleys", see eg. Fig. 3.1. This situation leads to ill-conditioned minimization problem.

The introduction of the quadratric  $J_{reg}$  term nacts as a "convexification" of J. Note that if defining  $J_{reg}$  from a "background value"  $u_b$ , this introduces a very strong prior on the sough solution  $u^*$ ...

Other regularization treatment then adding a back ground value term  $||u - u_b||^2$  may be necessary. A potential solution to this difficulty is to define metrics, physically-consistent, if one knowns some... This point is briefly addressed in a next chapter.



Figure 3.1: Tikhonov regularization. (L) A typical "slightly" convex functional  $j_{misfit}(\cdot)$ : illconditioned minimisation problem. (R) Regularized functional  $j(.) = (j_{misfit}(.) + \alpha_{reg}j_{reg}(.))$ , with  $j_{reg}$  defined from a *prior* value  $k_0$ .

# 3.6 Deriving the optimality system from the Lagrangian

As already mentioned, the optimization problem (3.11) may be viewed as a standard differentiable optimization problem with the model (3.1) being an equality constraint.

#### **3.6.1** Weak forms of the equations

First, let us write the weak (variational) form of the direct model. We set:

$$a(u; y, z) : U \times V \times V \to \mathbb{R} ; \quad a(u; y, z) = \langle A(u; y), z \rangle_{V' \times V}$$
$$b(u; z) : U \times V \to \mathbb{R} ; \quad b(u; z) = \langle B(u), z \rangle_{V' \times V}$$

Let us point out that both  $z \mapsto a(.;.,z)$  and  $z \mapsto b(.;z)$  are linear (i.e. linear with respect to the test function z).

Furthermore, if the direct model is linear (i.e. with respect to its unknown y) then the form a(.; y, z) is bilinear; it is not if not.

The weak formulation of the *direct model* (the state equation) is as follows:

$$\begin{cases} \text{Given } u \in U, \text{ find } y \in V \text{ such that:} \\ a(u; y, z) = b(u; z) \text{ for all } z \in V \end{cases}$$
(3.15)

An advantage to consider the weak form of the equations is to naturally taking into account the boundary conditions in the single equation (3.15).

Next, it will be more convenient to derive the correct adjoint equations.

By using the Riez-Frechet representation theorem, (3.15) is equivalent to:

$$A(u;y) = B(u) \text{ in } V' \tag{3.16}$$

**Exercice 3.11.** Apply this general expression to a simple BVP of your choice.

## 3.6.2 The Lagrangian & the optimality system

The optimization problem (3.11) may be viewed as a (differentiable) optimization problem with the model (3.1) being an equality constraint.

Then, it is natural to write the corresponding Lagrangian  $\mathcal{L}$ :

$$\mathcal{L}(u; y, p) = J(u; y) - \langle A(u; y) - B(u), p \rangle_{V' \times V}$$

with p the Lagrangian multiplier.

p is a dual variable belonging to the state space V.

Let us calculate the stationary point of the Lagrangian. The necessary optimality condition  $\nabla \mathcal{L}(u; y, p) = 0$  reads:

$$\begin{cases} \partial_{u}\mathcal{L}(u;y,p).\delta u = 0 \quad \forall \delta u \in U \\ \partial_{y}\mathcal{L}(u;y,p).\delta y = 0 \quad \forall \delta y \in V \\ \partial_{p}\mathcal{L}(u;y,p).\delta p = 0 \quad \forall \delta p \in V \end{cases}$$
(3.17)

We have:  $\mathcal{L}(u; y, p) = J(u; y) - a(u; y, p) + b(u; p)$ . The 1st equation of (3.17) reads:

$$\partial_u J(u; y) \cdot \delta u - [\partial_u a(u; y, p) - \partial_u b(u; p)] \cdot \delta u = 0 \quad \forall \delta u \in U$$

It will be shown in next section that this equation is the necessary condition "the gradient must vanish".

Indeed it will be rigorously shown later that:

$$j'(u) = \partial_u J(u; y) - [\partial_u a(u; y, p) - \partial_u b(u; p)]$$
(3.18)

The 2nd equation of (3.17) provides the equation:

$$\partial_y a(u; y, p) . \delta y = \partial_y J(u, y) . \delta y \quad \forall \delta y \in V$$

The 3rd equation of (3.17) provides the state equation (the direct model):

$$a(u; y, \delta p) = b(u; \delta p) \quad \forall \delta p \in V$$

It will be shown in next section that the 2nd equation is the so-called adjoint equation.

The equations system above constitutes the so-called "optimality system".

**Exercice 3.12.** Re-derive these equations by yourself.

**Remark 3.13.** Note that from the particular decomposition of J(u; y), see (3.4), then one has:

$$\partial_y J(u, y) \cdot \delta y = J'_{obs}(y) \cdot \delta y \text{ and } \partial_u J(u; y) \cdot \delta u = J'_{rea}(u) \cdot \delta u$$

$$(3.19)$$

**Exercice 3.14.** Apply the general expressions above to the case of your practical (or any other BVP of your choice).

# 3.7 Computing the cost function gradient without adjoint model

The goal is the numerical resolution of the optimization problem (3.11). As already mentioned, if the dimension of the discrete control variables  $u_h$  is large, this very likely requires to employ descent algorithms therefore to compute the cost function gradient  $\nabla j(u_h)$ .

#### 3.7.1 The Finite Difference gradient

The historical method (and the most simple one) consists to approximate the gradient values using finite differences.

Let  $U_h$  be the discrete control space,  $\dim(U_h) = m$ . Then, given  $\delta u \in \mathbb{R}^m$ , the gradient in the direction  $\delta u$  is evaluated as follows:

$$j'(u).\delta u \approx \frac{j(u+\varepsilon\delta u)-j(u)}{\varepsilon}$$
 at order 1 in  $\varepsilon$  (3.20)

Then, the computation of the vector  $\nabla j(u)$  requires (m+1) evaluation of j, therefore (m+1) resolutions of the direct model...

If  $T_{cpu}^{j}$  is not small and if *m* not small (typically in geophysical flows,  $m \approx 10^{6}$  or more, see next chapter), then the present finite difference method is not possible.

The finite difference approach is simple to implement, it is a non-intrusive method. Indeed, it is enough to run the direct model as much as required. It is possible if m tiny. On the contrary if m large it is not possible.

Another drawback of this approach is its weak accuracy, depending on the empirical choice of  $\varepsilon$ .

An other approach would be to solve the so-called Linear Tangent Model (LTM) but excepted for a very few control parameters, the computational time becomes prohibitive (see next section). Then, the good approach is to solve the adjoint model (see later).
Exercice 3.15. Propose an algorithm to compute the gradient of your practical by FD.

### 3.7.2 The straightforward "gradient" expression

Firstly let us recall the following differential expression :

**Lemma 3.16.** Let  $u_0$  in U, for all  $\delta u \in U$ , we have:

$$j'(u_0).\delta u = \frac{\partial J}{\partial u}(u_0; y).\delta u + \frac{\partial J}{\partial y}(u_0; y).w^{\delta u}$$
(3.21)

where  $w^{\delta u}$  denotes the derivative of the state y with respect to u in the direction  $\delta u$ :

$$w^{\delta u} = \frac{dy}{du}(u_0).\delta u \tag{3.22}$$

Proof. Straightforward calculation (done in exercise).

### 3.7.3 LTM-based gradient

It is worth noticing that  $w^{\delta u} = \frac{dy}{du}(u_0).\delta u$  can be obtained by simply deriving the direct model. This provides the so-called *Linear Tangent Model* (LTM).

#### Derivation of the LTM

The latter reads as follows:

| ĺ | Given $u \in U$ and $y^u$ the corresponding solution of the state equation (3.1),<br>given a direction $\delta u \in U$ , find $w \in V$ such that:           |        |
|---|---|--------|
| Ì | $\frac{\partial A}{\partial y}(u; y^u).w = \left[\frac{\partial B}{\partial u}(u) - \frac{\partial A}{\partial u}(u; y^u)\right].\delta u \text{ in } \Omega$ | (3.23) |
| l | with corresponding linearized boundary conditions on $\partial \Omega$  |        |

Proof. Straightforward calculation (done in exercise).

The weak formulation of the LTM is as follows:

$$\begin{cases} \text{Given } u \in U \text{ and } y^u \text{ solution of } (3.15), \\ \text{given } \delta u \in U, \text{ find } w \in V \text{ such that:} \\ \frac{\partial a}{\partial y}(u; y^u, z).w = [\frac{\partial b}{\partial u}(u; z) - \frac{\partial a}{\partial u}(u; y^u, z)].\delta u \text{ for all } z \in V \end{cases}$$
(3.24)

Solving the LTM provides  $w^{\delta u} = \frac{dy}{du}(u).\delta u$ , that is the derivative of the state y with respect to the control u in the direction  $\delta u$ . Recall that:  $\mathcal{M}: u \in U \mapsto y^u \in V$ . Therefore:  $\frac{dy}{du}(u) \in \mathcal{L}(U; V)$ , and  $w^{\delta u} \in V$ .

Exercice 3.17. Write the LTM of the BVP of your choice; both the weak and the classical forms.

Exercice 3.18. Write the LTM of your practical; both the weak and the classical forms.

**Remark 3.19.** It is assumed that the linearized tangent model is well-posed. Let us remark that if we have proved existence of solutions to the non-linear model, one likely had to prove that the linearized model is well-posed.

Moreover, if the implicit function theorem holds (and has been applied to prove the differentiability of the state with respect to the control), then the linearized problem must be well-posed.

Nevertheless for real-like non-linear problems, even the linearized model analysis can be non straightforward at all...

### Resulting sensitivity maps

Let us point out that the linearized tangent model provides  $w^{\delta u}$  that is the *local sensitivity of the state*  $y^u$  with respect to the control u, at "point" u in the direction  $\delta u$ .

In a modeling context, the resulting sensitivity map (it is distributed values) constitue rich information to better understand the model and/or the modelled phenomena

**Exercice 3.20.** In your programming practical, write a formal procedure which plot the "sensitivity map" in this context.  $\Box$ 

### FD gradient vs LTM-based gradient

The LTM has to be solved m times to obtain  $w^{\delta u}$  in each direction  $\delta u$ .

Therefore if m is large and the CPU time for each resolution is large than the LTM approach to compute  $w^{\delta u}$  is prohibitive.

However if possible, it is preferable to compute the gradient using the LTM compared to the finite difference method since its accuracy does not depend on the arbitrary setting of  $\varepsilon$ .

# 3.8 Rigorous derivation of the adjoint equations and the corresponding gradient expression

Previously, the optimality system has been formally derived by introducing the Lagrangian and by calculating necessary conditions of stationary points. One of the resulting equation was new; it was the so-called adjoint equation.

Below the adjoint model is rigorously derived.

These equations are a mathematical trick enabling the gradient computation by solving one (1) extra system only; to be compared to (m + 1) resolutions if using the finite difference approach or the LTM !

### 3.8.1 The adjoint-based gradient theorem

In this section the expression of the adjoint equations and the corresponding gradient expression are derived for the general direct model aforementioned.

Let us recall that given F an element of  $U' = \mathcal{L}(U, \mathbb{R})$  (U Banach space), we denote indifferently:  $F(u) \cdot v \equiv \langle F(u), v \rangle_{U' \times U} \quad \forall v \in U; \langle ., . \rangle_{U' \times U}$  denotes the duality product. The required basic concepts of functional analysis and differential calculus are recalled in Appendix A.

#### The fundamental result

We have:

Theorem 3.21. It is assumed that:

i) the state equation (3.15) is well-posed,

ii) its unique solution  $y^u$  is  $C^1$  with respect to u,

iii) the objective function J(u; y) is  $C^1$ ,

iv) the linearized tangent model (3.24) is well-posed.

Then, the cost function j(u) is  $C^1$ , and:

$$j'(u) = \frac{\partial J}{\partial u}(u; y^u) - \left[\frac{\partial a}{\partial u}(u; y^u, p^u) - \frac{\partial b}{\partial u}(u; p^u)\right] \text{ in } \mathcal{L}(U; \mathbb{R})$$
(3.25)

with  $\partial_u J(u; y^u) = \alpha_{reg} J'_{reg}(u)$  if considering the particular decomposition (3.4) of J(u; y).  $y^u$  is the unique solution of the state equation (3.15), and  $p^u$  is solution of the adjoint equation:

$$\begin{cases} Given \ u \ and \ y^u \ the \ unique \ solution \ of \ (3.15),\\ find \ p \in V \ satisfying:\\ \frac{\partial a}{\partial y}(u; y^u, p).z = \frac{\partial J}{\partial y}(u, y^u).z \quad \forall z \in V \end{cases}$$
(3.26)

Its solution  $p^u$  (the adjoint state) exists and is unique. One has  $\partial_y J(u; y^u) = J'_{obs}(y)$  if considering the particular decomposition (3.4) of J(u; y).

Proof.

First, under the present assumptions, the implicit function theorem applies and it gives the differentiability of the state with respect to u; i.e. the "model operator"  $\mathcal{M} : u \in U \mapsto y^u \in V$  is  $C^1$ . Then  $j: U \to \mathbb{R}$  can be differentiated. We have:  $j'(u) \in \mathcal{L}(U; \mathbb{R})$ , and :

$$\langle j'(u), \delta u \rangle_{U' \times U} = \langle \frac{\partial J}{\partial u}(u; y^u), \delta u \rangle_{U' \times U} + \langle \frac{\partial J}{\partial y}(u; y^u), w^{\delta u} \rangle_{V' \times V} \quad \forall \delta u \in U$$
 (3.27)

(see Lemma (3.21)) with  $\langle ., . \rangle_{U' \times U}$ ,  $\langle ., . \rangle_{V' \times V}$  the corresponding duality products. For sake of clarity we may denote too:  $j'(u) \cdot \delta u = \langle j'(u), \delta u \rangle_{U' \times U}$ .

The LTM reads:

$$<\frac{\partial A}{\partial y}(u;y^{u}).w^{\delta u}, z>_{V'\times V} = <\frac{\partial L}{\partial u}(u).\delta u, z>_{V'\times V} - <\frac{\partial A}{\partial u}(u;y^{u}).\delta u, z>_{V'\times V} \quad \forall z \in V$$
(3.28)

with  $w^{\delta u}$  defined by (3.22). By adding the two equations above, we obtain:

$$< j'(u), \delta u >_{U' \times U} = < \frac{\partial J}{\partial u}(u; y^u), \delta u >_{U' \times U} - < [\frac{\partial A}{\partial u}(u; y^u) - \frac{\partial L}{\partial u}(u)]. \delta u, z >_{V' \times V} + < [\frac{\partial J}{\partial y}(u; y^u) - (\frac{\partial A}{\partial y})^*(u; y^u).z], w^{\delta u} >_{V' \times V} \quad \forall \delta u \in U \ \forall z \in V$$

where  $(\partial_y A)^*$  is the adjoint operator of the linearized direct model operator  $\partial_y A$ .

In the expression of  $j'(u).\delta u$  above, we want to make vanish the term in  $w^{\delta u}$ . To do so, let us try to find  $p^u \in V$  satisfying:

$$< \left(\frac{\partial A}{\partial y}(u; y^{u})\right)^{*} \cdot p^{u}, w >_{V' \times V} = < \frac{\partial J}{\partial y}(u, y^{u}), w >_{V' \times V} \quad \forall w \in V$$

$$(3.29)$$

The linearized problem is well-posed therefore the operator  $\partial_y A(u; y^u)$  is an isomorphism from V into V', and its adjoint operator  $(\partial_y A)^*(u; y^u)$  is an isomorphism from V into V' too, see e.g. [11]. As a consequence, the equation above is well-posed. This is the so-called adjoint equation.

By noticing that:  $\langle (\frac{\partial A}{\partial y}(u; y^u))^* . p^u, w \rangle_{V' \times V} = \langle p^u, \frac{\partial A}{\partial y}(u; y^u) \cdot w \rangle_{V' \times V}$ , we obtain the following expression:

$$\langle j'(u), \delta u \rangle_{U' \times U} = \langle \frac{\partial J}{\partial u}(u; y^u), \delta u \rangle_{U' \times U} - \langle [\frac{\partial A}{\partial u}(u; y^u) - \frac{\partial L}{\partial u}(u)]. \delta u, p^u \rangle_{V' \times V} \quad \forall \delta u \in U$$

Therefore the result.

#### A few important remarks

- The expression of j'(u) in the direction  $\delta u$  obtained in Theorem 3.21 does not depend anymore on  $w^{\delta u}$ . Indeed the expression of j'(u) is now explicit with respect to the direction  $\delta u$ .

Thus, if solving the direct model plus the adjoint model then *all components of the gradient* follow i.e. the gradient in all directions !

Let us recall that after discretization (in the finite dimension space  $U_h$ ), we have:

$$\nabla j(u) \in \mathbb{R}^m, \quad \langle \nabla j(u), \delta u \rangle_{\mathbb{R}^m} = j'(u) \cdot \delta u \text{ for all } \delta u \in U_h \subset \mathbb{R}^m$$

- By construction, the adjoint model is linear, whatever if the direct model is linear or not. Recall the adjoint is the adjoint operator of the linearized direct operator.

- Let us point out that excepted for few particular cases (eg if the operator A is self-adjoint,  $A^* = A$ , of course), the adjoint model has in the general case no physical meaning.

- If the direct operator is *self-adjoint*, in other words if a(u, v) is bilinear symmetric then the adjoint operator is the same as the direct operator (but not the right-hand side). Indeed, in such a case, we have:

$$\partial_y a(u; y^u, p) \cdot z = a(u; z, p) = a(u; p, z)$$

Only the source term (RHS) and the boundary conditions differ from the state equation. Then the differential operator, hence the numerical method and numerical solver, are the same.

### Case of non-homogeneous Dirichlet boundary conditions

Let us consider the condition:  $y = y_d$  on  $\Gamma_d \subset \partial \Omega$ . Then, the direct model reads:

$$\begin{cases} Find \ y \in V_t = V_0 \oplus y_d \text{ such that } :\\ a(u; y, z) = b(u; z) \text{ for all } z \in V_0 \end{cases}$$

where  $V_t$ , affine subspace, is the Dirichlet translation of  $V_0$  ( $V_0$  subspace of V Hilbert). A typical example for a second order linear elliptic equation is :  $V = H^1(\Omega)$ ,  $V_0 = \{z \in V, z = 0 \text{ on } \Gamma_d\}$ and  $V_t = V_0 \oplus y_d = \{z \in V, z = y_d \text{ on } \Gamma_d\}$ .

Then the question is : What the non-homogeneous Dirichlet boundary conditions becomes when defining the LTM hence the adjoint model ?

The answer is : the non-homogeneous Dirichlet condition in the direct model becomes the corresponding homogeneous condition in the linear tangent and adjoint models.

Let us show this statement in the linear case.

The direct model, if linear in y, re-reads as follows:

$$\begin{cases} \text{ Find } y_0 \in V_0 \text{ such that }:\\ a(u; y_0, z) = b(u; z) - a(u; \tilde{y}_d, z) = \tilde{b}(u; z) \text{ for all } z \in V_0 \end{cases}$$

with  $y_0 = y - \tilde{y}_d$ ,  $\tilde{y}_d$  being a raising from  $y_d$  on  $\Gamma_d$  onto the whole domain  $\Omega$ . Then, following the proof of Theorem 3.21, it is easy to notice that the corresponding boundary condition in the LTM is homogeneous, hence the same for the adjoint model.

### 3.8.2 The optimality system

The optimality system denotes the set of equations characterizing the optimal control solution, that is: the state equation, the adjoint state equation and the first order necessary optimality condition.

In the case,  $K = U_{ad} = V$  and if considering the particular decomposition (3.4) of J(u; y), the optimality system reads as follows.

The optimal control solution u of Problem (3.10) has to satisfy:

$$\begin{cases} a(u; y^{u}, z) = b(u; z) \quad \forall z \in V \\ a^{*}((u, y^{u}); p, z) \equiv \partial_{y} a(u; y^{u}, p) \cdot z = J'_{obs}(y^{u}) \cdot z \quad \forall z \in V \\ j'(u) \cdot \delta u = 0 = \alpha_{reg} J'_{reg}(u) \cdot \delta u - [\partial_{u} a(u; y^{u}, p^{u}) - \partial_{u} b(u; p^{u})] \cdot \delta u \quad \forall \delta u \in U \end{cases}$$

$$(3.30)$$

The optimality system (3.30) is nothing else than the stationary point conditions (3.17) of the Lagrangian. The adjoint state p is the lagrangian multiplier associated to the "-model constraint".

**Exercice 3.22.** Write the optimality system which characterizes the solution u of the optimal control problem of a basic BVP and inverse problem of your choice.

**Exercice 3.23.** Write the optimality system which characterizes the solution u of the optimal control problem of your practical. Detail both the weak and the classical forms.

**Remark 3.24.** \* ("To go further"). The adjoint equations for a coupled system. If the direct model is composed by two PDEs equations weakly coupled then the adjoint system is composed by the corresponding adjoint equations weakly coupled too by in the reverse way.

If the direct model is composed by two PDEs equations coupled (fully) then the adjoint system is composed by the corresponding adjoint equations (fully) coupled too.

We do not develop more in detail this section since no time enough.

### 3.8.3 Gradient components: in weak or classical forms ? \*

This a "to go further paragraph".

We have written above the state, linear tangent and adjoint equations in their weak forms for few reasons. First, it is the right framework to study the solutions in the weak sense. Second, deriving the correct adjoint equations in the weak form may be more clear since the weak forms include naturally the boundary conditions. Third, it is the natural framework in the case of the finite element method. Nevertheless, in practice if not considering the finite element method, deriving the equations in their weak forms is not compulsory. Once the reader is confortable with these equation derivations process, it is possible to write all the equations sough directly in the classical forms (or going back from the weak to the classical forms !).

Let us point out that in next chapter a re-derivation of the equations is done in the semi-discrete form (time-dependent problems).

**Gradient discretization in the case of FEM** In the case of Finite Element discretization, an extra question remains concerning the discretization of the gradient expression (3.25). Recall  $j'(u) \in \mathcal{L}(U; \mathbb{R})$ . Let us denote  $\{\psi_i^u\}_{1 \le i \le m}$  the finite element basis of the control variable u. Then the ith component of the (discretized) gradient is naturally defined as follows :

$$\partial_i j(u) = \langle j'(u), \psi_i^u \rangle$$

**Example** Let us consider the direct model :  $-\Delta y = u$  (the control is the RHS) with homogeneous Dirichlet boundary conditions on  $\partial \Omega$ . Let us consider the cost function :

$$j(u) = \frac{1}{2} \int_{\Omega} (y^u)^2 dx + \frac{1}{2} \int_{\Omega} u^2 dx$$

If using finite element discretization, then the ith component of the gradient reads :

$$\partial_i j(u) = \int_{\Omega} p^u \psi_i^u dx + \int_{\Omega} u \psi_i^u dx \quad , \ 1 \le i \le m$$

with  $p^u$  the (discrete) adjoint state function (to be decomposed in its own finite element basis) and u to be decomposed in its finite element basis too.

Finite difference case Now let us consider the same problem but using finite difference schemes (and the same discretization for u and p, with m point values). The ith component of the (discrete) gradient reads :

$$\partial_i j(u) = p_i^u + u_i , \ 1 \le i \le m$$

with  $p_i^u$  the ith value of the (discrete) adjoint state.

In the case of *finite element discretization*, there is a choice to make for the gradient definition. In the example above, the choice would read:

$$\int_{\Omega} p^u \psi_i^u dx \text{ vs } p_i^u$$

These two possible expressions do not present the same properties, in particular concerning their dependence on the local mesh element size.

### 3.9 The VDA algorithm

In practice, it is classical to solve the optimization problem (3.10) and the optimality system (3.30) by an iterative descent algorithm therefore performing a local minimization (on the contrary to global optimization methods).

Given a first guess  $u_0$ , we seek  $(u^m)_m$  that decreases the cost function using a first order descent algorithm, preferably a Quasi-Newton method e.g. BFGS. The algorithm is as follows, see Fig. (3.2):

- 1) compute the cost function using the direct model,
- 2) compute the gradient using the adjoint model,
- 3) given the current value  $u^n$ , the current values  $j(u^n)$  and  $\nabla j(u^n)$ , compute a new iteration  $u^{n+1}$  such that:

$$j(u^{n+1}) < j(u^n)$$

\*) Iterate until convergence

**Remark 3.25.** In order to solve the optimality system (3.30), the 'natural' method would be use the Newton algorithm. This algorithm is very efficient (second order, if converging...) but it requires the computation of the Hessian  $H_j$  of j. Generally, the computation of  $H_j$  is either very complex and/or CPU time consuming. That is why in many cases, we solve the 1st order optimality system (3.30) only.



Figure 3.2: VDA algorithm: optimal control of the PDE system (identification process)

The background term in the cost function: an important prior, for the best or worse... If the cost function term  $J_{obs}$  is "weakly" convex only, that is presenting "nearly flat valleys", the first guess value of the iterative algorithm likely determine the local minimum the descent algorithm will converge to.

Moreover, if adding  $J_{reg}$  (Tykhonov regularization term) defined from  $u_b$ , this will enforce the minimization algorithm to converge to a value partially defined by  $u_b$ ...

Note that other important priors are the definition of the norms O and N.

These (semi-)norms may be defined from covariance matrices (symmetric linear positive definite operators) modeling a-priori knowledge we may have on the solution, see later.



Figure 3.3: The first guess: a potentially influencing prior

### J. Monnier

### 3.10 The fundamental equations at a glance

The non-linear stationary PDE model. The direct model reads:

$$\begin{cases} \text{Given } u \in U, \text{ find } y \in V \text{ such that:} \\ A(u; y) = B(u) \text{ in } \Omega \\ \text{with boundary conditions on } \partial\Omega \end{cases}$$
(3.31)

It is supposed to be well-posed.

The model operator  $\mathcal{M}(u)$  is defined as:  $\mathcal{M}(u) = y^u$ . This operator  $\mathcal{M}(\cdot)$  is potentially non linear.

The cost functional j(u) is defined from the "observation" functional J(u; y) as follows:

$$j(u) = J(u; y^u) \tag{3.32}$$

where  $y^u$  denotes the unique solution of the direct model, given u.

The observation functional J(u; y) is chosen quadratic, decomposed as follows:

$$J(u;y) = J_{obs}(y) + \alpha_{reg} J_{reg}(u)$$
(3.33)

with:

$$J_{obs}(y) = \frac{1}{2} \left\| C(y) - z_d \right\|_R^2$$
(3.34)

$$J_{reg}(u) = \frac{1}{2} \left\| u - u_b \right\|_B^2 \text{ or e.g. } \frac{1}{2} \left\| D^p u \right\|_0^2 \text{ with } p = 1 \text{ or } 2.$$
(3.35)

The optimization problem reads:

$$\begin{cases} \text{Minimize } j(u) \text{ in } U_{ad} \text{ under the "model constraint"} \\ \text{since } j(u) = J(u; y^u) \text{ with } y^u = \mathcal{M}(u). \end{cases}$$
(3.36)

The adjoint model reads:

$$(\partial_y A)^T(u; y^u) \cdot p = J'_{obs}(y^u) \text{ in } \Omega$$
(3.37)

In weak form:

$$a^*((u, y^u); p, z) \equiv \partial_y a(u; y^u, p) \cdot z = J'_{obs}(y^u) \cdot z \quad \forall z \in V$$
(3.38)

### The resulting gradient expression.

The differential of j reads: for all  $\delta u \in U$ ,

$$j'(u) \cdot \delta u = \alpha_{reg} J'_{reg}(u) \cdot \delta u - [\partial_u A(u; y^u) \cdot \delta u - B'(u) \cdot \delta u] \cdot p^u$$
(3.39)

In weak form:

$$j'(u) \cdot \delta u = \alpha_{reg} J'_{reg}(u) \cdot \delta u - \partial_u a(u; y^u, p^u) \cdot \delta u + \partial_u b(u; p^u) \cdot \delta u$$
(3.40)

### J. Monnier

In the case of few components e.g.  $u = (u_1, u_2)$ , then one has:

$$j'(u) \equiv \left(\frac{\partial j}{\partial u_1}(u), \frac{\partial j}{\partial u_2}(u)\right)^T$$
 and  $j'(u) \cdot \delta u = \frac{\partial j}{\partial u_1}(u) \cdot \delta u_1 + \frac{\partial j}{\partial u_2}(u) \cdot \delta u_2$ 

The relationship between the gradient  $\nabla j(u)$  and the differential  $j^\prime(u)$  is:

$$\langle \nabla j(u), \delta u \rangle_{\mathbb{R}^m} = j'(u) \cdot \delta u$$
 for all  $\delta u \in U_h \subset \mathbb{R}^m$  (3.41)

# 3.11 Another example: control of an elastic membrane deformation \*

This is section is a "supplementary material" particularly for those doing the practical programming. If it is your case, you may skip this section.

Let us consider the deformation of an elastic membrane, fixed at its boundaries, and distorted by an external force f. We denote by d the displacement of the membrane,  $d : \Omega \subset \mathbb{R}^N \to \mathbb{R}^N$ . It satisfies the following linear second order elliptic system:

$$\begin{cases} -\Delta d = f + u \text{ in } \Omega \\ d = 0 \text{ on } \partial \Omega \end{cases}$$
(3.42)

where u is the control of the system, applied on  $\omega$  a subset of  $\Omega$ .

The goal is to characterize u such that the displacement d is as close as possible to a target displacement  $d_0$  given. Thus, we define the cost function:

$$j(u) = \int_{\Omega} |d - d_0|^2 dx + \alpha \int_{\omega} |u|^2 dx$$

where d is the solution of (3.42) given u. The second term of the cost has two consequences. In a mechanical point of view, it implies that we seek a control as small as possible (in  $L^2$  norm). In a mathematical point of view, it regularizes the cost function: it is the quadratic term essential in the previous existence-uniqueness theorem. In a numerical point of view it "convexifies" the functional; it is what we call a Tikhonov's regularization (see later).

The weight coefficient  $\alpha$ ,  $\alpha > 0$  makes a balance between the two terms, it has to be tuned by hand empirically (multi-objective function to be minimized...).

We set  $U = (L^2(\Omega))^N$ . The set of admissible control K is as follows:

$$U = \{ u \in U, \ u = 0 \text{ in } \Omega/\omega \}$$

an the optimal control reads:

$$\min_{u \in K} j(u) \tag{3.43}$$

Existence and uniqueness of solution.

i) Let us point out that in vertu of Lax-Milgram's theorem, given u, the direct model has we and we weak solution d in  $V = (H_0^1(\Omega))^N$ .

ii) Furthermore, since the mapping  $\mathcal{M} : u \in U \mapsto d \in V$  is affine, the cost function j is strongly convex  $(j(u) \ge \alpha ||u||_{L^2}^2)$ .

In other respect, K is a closed convex set of U. Therefore, in vertu of Theorem 3.9, there exists an unique minimum u to the optimal control problem (3.43).

**Exercice 3.26.** Write the optimality system which characterizes the solution u of the optimal control problem (3.43); both in the weak and classical form.

# 3.12 How to assess your adjoint computational code & gradient \*

We describe below how the adjoint code can be assessed. In short, it is verified through the scalar product property that the code actually computes the adjoint of the tangent linear code. Next, it is verified that the code computes correctly the partial derivatives of the cost function by comparing them with the value computed by finite differences: this is the important gradient test.

This section may be skipped in his first course reading.

### 3.12.1 The scalar product test

This test aims at checking if the adjoint code is actually the adjoint of the tangent linear code. This test supposes to have developed both the adjoint code and the linear tangent code.

This test is equivalent to check the relation (B.1).

- Given an arbitrary  $d\mathbf{k} \in \mathcal{K}$ , we compute using the tangent linear code :  $dY = \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right) \cdot d\mathbf{k}$
- Given an arbitrary  $dY^* \in \mathcal{Y}$ , we compute using the adjoint code :

$$d\mathbf{k}^* = \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right)^* \cdot dY^*$$

• The two following scalar products are computed:

$$. sp_1 = \left\langle dY^*, dY \right\rangle_{\mathcal{Y}}$$

$$sp_2 = \langle d\mathbf{k}^*, d\mathbf{k} \rangle_{\mathcal{K}}$$

• Finally, following the adjoint operator definition  $\langle A\mathbf{k}, Y \rangle = \langle \mathbf{k}, A^*Y \rangle (A(\cdot) \text{ linear})$ , we check that:  $sp_1 = sp_2$ .

Figure 3.4 (b) shows a typical example of the scalar product test.

Figure 3.4: Adjoint code validation: scalar product test

### 3.12.2 The gradient test

The objective of this test is to check if the adjoint variables  $dX^*$  computed by the adjoint code actually correspond to the partial derivatives of the cost function. This is a crucial test. This test requires to perform the direct code a few times and the adjoint code.

The Taylor expansion of the cost function j at  $\mathbf{k}$  for a small perturbation  $\alpha \, \delta \mathbf{k}$  (where  $\alpha \in \mathbb{R}^+$ ) reads:

$$j(\mathbf{k} + \alpha \,\delta \mathbf{k}) = j(\mathbf{k}) + \alpha \,\frac{\partial j}{\partial \mathbf{k}}(\mathbf{k}) \cdot \delta \mathbf{k} + o\left(\alpha \|\delta \mathbf{k}\|\right) \,. \tag{3.44}$$

It follows the uncentered finite difference approximation (order 1) and the centered finite difference approximation order 2:

$$\frac{j(\mathbf{k}+\alpha\,\delta\mathbf{k})-j(\mathbf{k}-\alpha\,\delta\mathbf{k})}{2\alpha} = \frac{\partial j}{\partial \mathbf{k}}(\mathbf{k}).\delta\mathbf{k} + O\left(\alpha^2 \|\delta\mathbf{k}\|^2\right). \tag{3.45}$$

Then, we set either

$$I_{\alpha} = \frac{j(\mathbf{k} + \alpha \,\delta \mathbf{k}) - j(\mathbf{k} - \alpha \,\delta \mathbf{k})}{2\alpha \,\frac{\partial j}{\partial \mathbf{k}}(\mathbf{k}) \cdot \delta \mathbf{k}} \,. \tag{3.46}$$

or:

$$I_{\alpha} = \frac{j(\mathbf{k} + \alpha \,\delta \mathbf{k}) - j(\mathbf{k})}{\alpha \,\frac{\partial j}{\partial \mathbf{k}}(\mathbf{k}) \cdot \delta \mathbf{k}} \,. \tag{3.47}$$

According to (3.44), one must have:  $\lim_{\alpha \to 0} I_{\alpha} = 1$ .

The gradient test consists to check this property as follows.

- For an arbitrary **k**, compute  $\frac{\partial j}{\partial \mathbf{k}}(\mathbf{k})$  with the adjoint code.
- Using the direct code, compute  $j(\mathbf{k})$ .
- For n = 0, ..., N:
  - Compute  $\alpha_n = 2^{-n}$ ;
  - Using the direct code, compute  $j(\mathbf{k} + \alpha_n \, \delta \mathbf{k})$ ;
  - Compute  $I_{\alpha_n}$ ;
- Check if  $\lim_{\alpha \to 0} I_{\alpha_n} = 1$  or not.

Figure 3.5 shows two results of the gradient test: at order 2 and at order 1.  $|I_{\alpha} - 1|$  is plotted vs  $\alpha$  in logarithmic scale.

The convergence is good until  $\alpha > 10^{-7}$ .

However, observe the difference of accuracy between the 1st order and 2nd order approximation.

If  $\alpha$  is smaller than  $\approx 10^{-7}$ , the approximation of the partial derivatives is not reliable anymore due to truncation errors.

(To show this, one can add a fix term in the Taylor expansion and notice that it is divided by the perturbation therefore increasing).



Figure 3.5: The adjoint code validation. Gradient test at order 1 (a), at order 2 (b).

**Exercice 3.27.** *Perform the gradient test for your practical problem. Is your gradient computation valid ?* 

## Chapter 4

## Equivalences between VDA, the BLUE and Bayesian inferences

If the forward model is linear, if the observation operator C is linear and if the functional J(u; y) is quadratic (in its primal variables) that is *if the considered inverse problem is the LQ problem* then the Best Linear Unbiased Estimator (BLUE) is equivalent to a particular VDA solution  $u^*$ .

This result is demonstrated in the first section.

Moreover, if assuming Gaussian probability density functions (pdf) then the same VDA solution is nothing else than the most probable solution, that is the Maximum A-Posteriori (MAP) estimation provided by the Bayes theorem.

This result is demonstrated in the third section.

In the second section, it is shown how the introduction of covariance operators may improve the robustness and the efficiency of the optimisation process by defining "consistent-model" metrics. These covariances operators (matrices in discrete version) may be interpreted as prior probability model of the studied phenomena. The outline of this chapter is as follows.

### Contents

| 4.1 The   | Best Linear Unbiased Estimator (BLUE) and VDA  | 83 |
|---|--|----|
| 4.1.1   | The least-square solutions depend on the metrics $\ldots$ $\ldots$ $\ldots$ $\ldots$ | 83 |
| 4.1.2   | The BLUE, scalar case  | 84 |
| 4.1.3   | Equivalence with a VDA solution  | 85 |
| 4.1.4   | The BLUE (vectorial case) & resulting optimal metric $\ldots \ldots \ldots$          | 87 |
| 4.1.5   | Sequential filters formalism   | 88 |
| 4.2 The   | e Bayesian view  | 89 |
| 4.3 Conclusion: in the LQ case, VDA, the Maximum A-Posteriori |  |    |
| $(\mathbf{M})$  | AP) estimation and BLUE can be equivalent  | 93 |

4.4 Another example of BLUE and VDA solution ..... 95

### 4.1 The Best Linear Unbiased Estimator (BLUE) and VDA

The basic VDA approach which has been presented in the previous chapter is purely deterministic. It provided the optimal control denoted by  $u^*$ , solution of (3.10). Below it is shown that the Best Linear Unbiased Estimator (BLUE) i.e. the most simple statistic estimator provides a particular VDA solution. This result suggests to define the norms (metrics) measuring the discrepancy between the model output and data as the inverse of the covariance matrix of errors, that is, see (3.33):

$$J_{obs}(y) = \frac{1}{2} \left\| C(y) - z_d \right\|_O^2$$
(4.1)

with O the inverse of a covariance matrix of errors.

### 4.1.1 The least-square solutions depend on the metrics

Let us consider the following simple but instructive example. One consider two measurements of a scalar quantity u:  $z_1 = 1$  and  $z_2 = 2$ .

'Naturally', one would seek u minimizing the following cost function:

$$j(u) = (u - z_1)^2 + (u - z_2)^2$$

It is a "explicit" (standard) least-square problem since u does not have to satisfy an underlying model. The solution is:  $u^* = \frac{3}{2}$ .

Now, let us assume that the second data had measured the quantity 2u and not u. Then, the two measurements are:  $z_1 = 1$  and  $z_2 = 4$ .

Again, 'naturally', one seek u minimizing the following cost function:

$$g(u) = (u - z_1)^2 + (2u - z_2)^2$$

The solution is:  $u^* = \frac{9}{5}$ . This solution differs from the first one !

This simple example illustrates that the solution of the minimization problem (here a trivial linear least-square problem) depends on the considered metrics (actually norms) employed in the cost function.

The standard least-square solution is dependent of the measurement metric and/or accuracy !

In presence of errors (which is always the case in real-life problems), it does not take into account the confidence we may have on each measurement.

In order to take into account the measurement accuracies, one can *introduce a-priori information in the norm* measuring the discrepancy.

### 4.1.2 The BLUE, scalar case

BLUE means Best Linear Unbiased Estimator. We refer the reader to Appendix A.4 for the employed statistics notations.

**Definition 4.1.** The Best Linear Unbiased Estimator (BLUE) of an aleatory variable  $\hat{X}$  based on data Z is as follows:

- a) It a linear function of Z. (Comment: it is the most simple).
- b) It is unbiased:  $E(\hat{U}) = u$ . (Comment: it is desirable).
- c) It has the smallest variance  $Var(\hat{X})$  among all unbiased linear estimation. (Comment: optimal accuracy is expected).

The BLUE for the basic example Let us denote the two measurements including errors as follows:

$$z_i = u_i + \varepsilon_i, \ i = 1, 2$$

The errors of measurements  $\varepsilon_i$  are supposed to be:

- unbiased:  $E(\varepsilon_i) = 0$ .

(Comment: sensors are unbiased; mean of errors vanishes).

- with a variance given:  $Var(\varepsilon_i) = \sigma_i, i = 1, 2.$ 

(Comment: the accuracy of each sensor is known).

- uncorrelated :  $E(\varepsilon_1 \varepsilon_2) = 0.$ 

(Comment. Measurements are independent hence the covariance vanishes; in addition, means vanish hence the relation).

Let us calculate the Best Linear Unbiased estimation (BLUE) denoted by  $u^*$ . (Here "best" means giving the lowest mean-square error).

By definition, the BLUE satisfies:  $u^* = a_1 z_1 + a_2 z_2$  (linear relationship). The coefficients  $a_i$  have to be determined. The errors are unbiased:

$$E(u^* - u) = 0$$

We have:

$$E(u^*) = (a_1 + a_2)u + a_1 E(\varepsilon_1) + a_2 E(\varepsilon_2) = (a_1 + a_2)u$$

Therefore:  $a_2 = (1 - a_1)$ .

### J. Monnier

Next, let us calculate the variance of  $u^*$ :

$$Var(u^*) = E((u^* - u)^2) = E((a_1z_1 + a_2z_2)^2)$$
  
=  $a_1^2 E(\varepsilon_1^2) + 2a_1a_2 E(\varepsilon_1\varepsilon_2) + a_2^2 E(\varepsilon_2^2)$   
=  $a_1^2\sigma_1^2 + a_2^2\sigma_2^2$   
=  $a_1^2\sigma_1^2 + (1 - a_1)^2\sigma_2^2$ 

By definition, the BLUE  $u^*$  minimize the variance. The latter is minimal if its derivative with respect to  $a_1$  vanishes; this reads:  $a_1 = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$ .

Therefore, the BLUE reads:

$$u^* = \frac{1}{\left(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}\right)} \left(\frac{1}{\sigma_1^2} z_1 + \frac{1}{\sigma_2^2} z_2\right)$$
(4.2)

Moreover:  $Var(u^*) = \frac{\sigma_1^2 \sigma_2^2}{\sigma_1^2 + \sigma_2^2}$ . Therefore:

$$\frac{1}{Var(u^*)} = \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}$$

### 4.1.3 Equivalence with a VDA solution

It is straightforward to verify that the BLUE  $u^*$  defined by (4.2) minimizes the following quadratic cost function:

$$j(u) = \frac{1}{2} \frac{1}{\sigma_1^2} (u - z_1)^2 + \frac{1}{2} \frac{1}{\sigma_2^2} (u - z_2)^2$$
(4.3)

Indeed, we have:

$$j''(u) = \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} = \frac{1}{Var(u^*)}$$

Therefore, the Hessian of j (defining the "convexity rate" of j) equals the estimation accuracy.

This property between the Hessian and the analysis quality can be simply illustrated in 1D, see Fig. 4.1.

Let us remark that the BLUE calculated above (with unbiased measurements) minimizes the following cost function too:

$$j(u) = \frac{1}{2}(u - z_1, u - z_2) \begin{pmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 \\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}^{-1} \begin{pmatrix} u - z_1 \\ u - z_2 \end{pmatrix} = \frac{1}{2} \|u - z\|_N^2$$
(4.4)

In this cost function expression, the norm N considers that the observations errors are correlated. Indeed, the extra diagonal terms are non vanishing.



Figure 4.1: In the LQ case: convexity=reliability. In 1D, the Hessian=second derivative hence the "convexity rate". The latter measures the pdf sharpness of the analysis. Figure extracted from F. Bouttier, P. Courtier course, ECMWF (www.ecmwf.int).

**Extension to** *m* **observations** The extension to *m* observations is straightforward. We denote the *m* measurements with errors by:  $z_i = u + \varepsilon_i$ , i = 1, ..., m. We assume that the errors are:

- unbiased:  $E(\varepsilon_i) = 0, i = 1, ..., m.$ 

- with a variance given:  $Var(\varepsilon_i) = \sigma_i, i = 1, ..., m$ .

- uncorrelated :  $E(\varepsilon_i \varepsilon_j) = 0$  for all  $(i, j), 1 \le i, j \le m, i \ne j$ .

Then, the Best Linear Unbiased Estimation reads:

$$u^* = \frac{1}{\sum_{i=1}^m \frac{1}{\sigma_i^2}} \left( \sum_{i=1}^m \frac{1}{\sigma_i^2} z_i \right) \text{ with } \frac{1}{Var(u^*)} = \sum_{i=1}^m \frac{1}{\sigma_i^2}$$

Therefore the BLUE minimizes:

$$j(u) = \frac{1}{2} \sum_{i=1}^{m} \frac{1}{\sigma_i^2} (u - z_i)^2$$
(4.5)

Again, the Hessian (here a scalar second derivative) defines the "convexity rate" and mesures the "sharpness" of the analysis:

$$j''(u) = \frac{1}{Var(u^*)}$$

A second simple and illustrative example of BLUE and VDA solution is proposed in a next section.

### 4.1.4 The BLUE (vectorial case) & resulting optimal metric

The goal of this section is to calculate the BLUE in a vectorial case; then, deduce the required norm in the VDA cost function to obtain the BLUE.

Let us consider the previous basic example but with n variables and m observations. We seek to estimate the *vectorial* quantity  $u = (u_1...u_n)^T \in \mathbb{R}^n$ . The measurements-observations are:  $z^{obs} = (z_1^{obs}...z_m^{obs})^T \in \mathbb{R}^m$ . We assume that the *observation operator* C is linear: C is a  $m \times n$  matrix. We consider errors e on the observations:

$$z^{obs} = Cu + e \qquad x \in \mathbb{R}^n, \ z \in \mathbb{R}^m \tag{4.6}$$

 $e, e \in \mathbb{R}^m$ , denotes the observation-modelling error.

This error e is supposed to be an aleatory variable which satisfies:

- 1) E(e) = 0 i.e. the sensors and/or model are unbiased. (This assumption may be reasonable (or not...) in real-world problems).
- 2) Its covariance matrix Cov(e),  $Cov(e) = E(ee^T)$  is given. That is the covariances of observations-modeling errors are supposed to be known.

This last assumption is generally not satisfied in real-world problems; however as shown below it is the adequate matrix to consider...

If the observations-modeling errors are not independent then the matrix Cov(e) is non-diagonal. Recall that each diagonal term of Cov(e) represents the confidence - accuracy one may has on the measurement-modeling error.

Under these assumptions and by using the *Gauss-Markov theorem*, the following result can be proved.

**Proposition 4.2.** The BLUE is the (unique) solution of the minimization problem  $\min_x j(u)$  with:

$$j(u) = \frac{1}{2} \|Cu - z^{obs}\|_{O}^{2} \text{ with } O = (Cov(e))^{-1}.$$
(4.7)

Recall the present discrete estimation problem would correspond to a linear direct model with a linear observation operator and errors on the observations.

Then this result demonstrates that if the exact covariances of the observation-modelling errors Cov(e) are known (this is unfortunately not the case for real-world problems) then the matrix

to consider in the metric of j(u) to obtain the BLUE is O with  $O = (Cov(e))^{-1}$ .

The reader can find more details in one of the numerous excellent books treating of the subject, see e.g. [19].

**Exercice 4.3.** Write the linear observation operator C in the following case: Two observations  $z_1, z_2$  of  $u_1$  are available, while one observation  $z_3$  of  $u_2$  only is available.

### 4.1.5 Sequential filters formalism

Filters are stochastic algorithms which operate recursively on streams of noisy input data to produce a statistically optimal estimation of the underlying state.

**The scalar case** Let go back to the simple example with two observations. In this case, the BLUE reads:

$$u^* = \frac{\sigma_2^2 z_1 + \sigma_1^2 z_2}{\sigma_1^2 + \sigma_2^2} = z_1 + \left(\frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2}\right) (z_2 - z_1)$$

Let us consider that:

a)  $z_1$  is a first estimation. It is the so-called *background* or *first guess* value.

b)  $z_2$  is an independent observation.

We denote the first-guess value  $z_1$  by  $u_b : u_b \equiv z_1$ . We denote the next observation  $z_2$  by  $z : z \equiv z_2$ . Following this point of view, the BLUE  $u^*$  reads :

$$u^* = u_b + \left(\frac{\sigma_b^2}{\sigma_b^2 + \sigma_0^2}\right)(z - u_b)$$
(4.8)

Using the terminology of sequential data assimilation e.g. the Kalman filter, this equation reads as follows:

"The best estimation  $u^*$  equals the first guess + the gain times the innovation  $(z - u_b)$ ".

The gain is defined as the multiplicative factor of the innovation. In a filter algorithm eg. the Kalman filter, defining the gain matrix is the key point.

**On the Kalman filter.** R. E. Kalman, Hungarian-born American electrical engineer, born in 1930.

In short, the Kalman filter is not affordable for CPU-time consuming models (despite the "easy" quadratic problem solved). Moreover it is not optimal for nonlinear models.

For non-linear models and/or expensive models, one has to consider Ensemble Kalman filters (e.g. EnKF, ETKF), or ... VDA ("4D-Var"). Hybrid methods exist too. Concerning sequential data assimilation methods (e.g. the Kalman filter), the reader may consult e.g. [19].

From Wikipedia's page: "Kalman filtering, also known as Linear Quadratic Estimation (LQE), works in a two-step process. In the prediction step, the Kalman filter produces estimations of the current state variables, along with their uncertainties. Once the outcome of the next measurement (necessarily corrupted with some amount of error, including random noise) is observed, these estimations are updated using a weighted average, with more weight being given to estimations with higher certainty. Because of the algorithm's recursive nature, it can run in real time using only the present input measurements and the previously calculated state and its uncertainty matrix; no additional past information is required".

### 4.2 The Bayesian view

Rev. Thomas Bayes (1702-1761), English statistician and philosopher. In probability and statistics, Bayes' theorem describes the probability of an event, based on prior knowledge of conditions that might be related to the event. Independently of Bayes, Pierre-Simon Laplace in 1774 used conditional probability to formulate the relation of an updated posterior probability from a prior probability, given evidence. Source: Wikipedia.

One of the many applications of Bayes' theorem is Bayesian inference, a particular approach to statistical inference.

**The Baye's law.** Let u be the uncertain parameter to be estimated and  $z^{obs}$  an observation. Let p(u) be the prior distribution of u.

 $p(u|z^{obs})$  (the probability of  $z^{obs}$  given u) is the *posterior distribution*. The Baye's theorem states that the posterior distribution satisfies:

$$p(u|z^{obs}) = \frac{p(z^{obs}|u)}{p(z^{obs})}p(u)$$
(4.9)

where  $p(z^{obs}|u)$  is the likelihood function for u.

Note that one can obtain  $p(z^{obs})$  by integrating over all u:  $p(z^{obs}) = \int p(z^{obs}|u) p(u)dx$ .



Figure 4.2: The Bayes theorem for discrete events. Figure extracted from Wikipedia.

Seeking the most probable parameter u. Let us seek the most probable parameter u, given the observations  $z^{obs}$  (and potentially given the background  $u_b$  too):

$$u^* = \arg\max_u \left( p(u|z^{obs} \text{ and } u_b) \right)$$
(4.10)

This may be done by defining the following cost function:

$$j(u) = -\log(p(u|z^{obs} \text{ and } u_b)) + c$$

$$(4.11)$$

for any constant c.

This choice of function log is adequate since it is a *monotonic function*. Then, it follows:

$$u^* = \arg\min j(u) \tag{4.12}$$

Assuming that the observation errors and the background errors are uncorrelated (this seems

highly reasonable):

$$p(z^{obs} \text{ and } u_b|u) = p(z^{obs} \text{ and } u_b)p(u)$$

Using the Bayes law (4.9), it follows:

$$j(u) = -\log p(z^{obs}|u) - \log p(u) + c$$
 (4.13)

for any constant c.

This calculation can be found e.g. in [30].

If assuming Gaussian pdf... The maximum likelihood approach is applicable to any pdf (probability density function)  $p(z^{obs}|u)$  and p(u). However, let us consider the special case of Gaussian pdf. Note that the choice of log function in the definition of j is perfectly adapted to Gaussian functions...

- We assume that the prior distribution p(u) is Gaussian. More precisely

$$p(u) \sim \mathcal{N}(u_b, \sigma_b^2)$$

with: the mean  $u_b$ , the variance  $\sigma_b^2$ , the standard deviation  $\sigma_b = ||B||$ . Here *B* denotes the background error covariance matrix as previously suggested:

$$B = Cov(e_b) = E(e_b e_b^T) \tag{4.14}$$

We have:

$$p(u) = \frac{1}{(2\pi\sigma_b)^{1/2}} \left(-\frac{1}{2}\|u - u_b\|_{B^{-1}}^2\right)$$

- We assume the likelihood  $p(z^{obs}|u)$  to be Gaussian too.

Note that this assumption is satisfied if p(u) is Gaussian and the model is linear. Otherwise this assumption is a-priori not satisfied...

We set:

$$p(z^{obs}|u) \sim \mathcal{N}(u_o, \sigma_o^2)$$

with : mean  $z^{obs}$ , variance  $\sigma_o^2$ , standard deviation  $\sigma_o = ||O||$ , O the observation error covariance matrix defined by:

$$O = Cov(e_o) = E(e_o e_o^T) \tag{4.15}$$

We have:

$$p(z^{obs}|u) = \frac{1}{(2\pi\sigma_o)^{1/2}} \exp\left(-\frac{1}{2}\|Cu - z^{obs}\|_{O^{-1}}^2\right)$$
(4.16)

By taking the log in (4.13) (and choosing an adequate constant value), it follows:

$$j(u) = \frac{1}{2} \|Cu - z^{obs}\|_{O^{-1}}^2 + \frac{1}{2} \|u - u_b\|_{B^{-1}}^2$$
(4.17)

with the norms O and B defined as presented above.

In summary, computing the most probable posterior value involves maximizing a product. Next, by applying the log function to the posterior pdf in the Gaussian case leads to the minimization of the "usual" VDA formulation for the LQ case (quadratic terms) with metrics defined from the error covariance matrices.

**Limitations of this analysis** The equivalency shown above assumes a Gaussian prior distribution of u (this may be reasonable) and a Gaussian likelihood  $p(z^{obs}|u)$ .

As already mentioned, this last assumption is a-priori not satisfied for non-linear problems. Moreover for an a-priori Gaussian likelihood function  $p(z^{obs}|u)$ , the equivalent cost function j(u) is quadratic. This is a-priori not true for  $p(z^{obs}|u)$  non Gaussian.



Figure 4.3: Figure extracted from A. Lorenc course, Met Office (JCSDA Summer School, 2009).

**The resulting MAP estimation** Following the Bayes theorem, see (4.9), the posterior distribution is a Gaussian distribution which reads (see eg. Lorenc's course):

$$p(u|z^{obs}) = \frac{1}{(2\pi\sigma_p)^{1/2}} \exp\left(-\frac{1}{2}\|u - u_p\|_{P^{-1}}^2\right)$$
(4.18)

with: the mean  $u_p$ , the standard deviation  $\sigma_p = ||P||$ , P satisfying:

$$P^{-1} = B^{-1} + C^T O^{-1} C (4.19)$$

And

$$u_p = u_b + PC^T O^{-1} (z^{obs} - C u_b) \tag{4.20}$$

That is the Maximum A-Posteriori (MAP) estimation.

## 4.3 Conclusion: in the LQ case, VDA, the Maximum A-Posteriori (MAP) estimation and BLUE can be equivalent

In this section, it has been shown that in the LQ case (the model is linear, the observation operator is linear therefore the cost function is strictly convex and it exists an unique minimum):

- . The VDA solution equals the BLUE if the exact error covariances are known and define the employed metric in j(u).
- . The Hessian ("convexity rate") of j(u),  $D^2 j$ , may be interpreted as the estimation accuracy.

With the previous notations, see e.g. (3.33), one has:

$$D^{2}j = C^{T}O^{-1}C + B^{-1} = [Cov(C(u^{true}) - z_{d})]^{-1}$$
(4.21)

Recall that  $D^2 j$  is  $m \times n$  matrix.

. If assuming Gaussian pdf, the VDA solution based (on the metrics above) equals the most probable solution.

It is the Maximum A-Posteriori (MAP) estimation provided by the Bayes theorem.

In summary, in the LQ case and if considering the error covariance matrices as metrics, all the equations - approach provide the same solution (optimal paramter) !

### What about in non-linear cases ?

If the direct model is linear and the observation operator is linear than the cost function is strictly convex (LQ problem); it admits a unique (global) minimum.

On the contrary, if the direct model is non-linear and/or the observation operator is non-linear than the cost function is a-priori non convex.

In such a case, the minimization process may converge to a local minimum, depending on the priors.

The minimization problem may be more or less well conditioned depending on the Hessian properties of j (its "convexity rate")... See eg. Fig. 3.1.

Moreover in such a case, the equivalence between a particular VDA solution and the BLUE does not hold anymore; the Gaussian assumption in the Bayes framework does not hold anymore too.

However, the analysis for the LQ case show what optimal metric, norms or covariance operators should be chosen.

This suggest to consider "model-consistent" norms (the covariance operators) in the definition of  $j(\cdot)$ ...

*Filters in non-linear cases.* Filtering methods can be extended to non-linear systems for example by linearizing the model at each time step; this the basic idea of the Extended Kalman's Filter. However in this case, the filter is not optimal anymore. Many others extensions exist. The Ensemble Kalman filter is the most common one. Briefly, it consists to perform a Monte-Carlo algorithm to estimate the covariance errors matrices. Others filters have been developed in particular to reduce the required CPU time to compute the gain matrices, see e.g. the SEEK filter. On the filters approaches, the reader may consult e.g. [19].

### 4.4 Another example of BLUE and VDA solution

These slides have been extracted from Prof. Olivier Thual's course (University of Toulouse, Fluid Mech. Institute).



# Basic recalls of Probability & Statistics

(Random variables, continuous & discrete)

The probability density function (pdf) in case of a Normal (or Gaussian) distribution:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} exp(-\frac{(x-\bar{x})^2}{2\sigma^2})$$

The mean (expected value):

$$\mu \equiv E(\hat{X}) = \int_{a}^{b} x f(x) dx \approx \frac{1}{N} \sum_{i=1}^{N} \hat{X}_{n}$$

The variance:

ance:  $Var(\hat{X}) = E[(\hat{X} - \mu)^2]$ 

The standard deviation:  $\ \ \sigma = (Var(\hat{X}))^2$ 





1) The BLUE (called the estimate or the « analysis ») satisfies a linear relationship:

$$T_k = (1-k)T_1 + kT_2$$

 $\sigma_a^2 = (1-k)^2 \sigma_1^2 + k^2 \sigma_2^2$ 

2) Errors are supposed unbiased and uncorrelated, hence the variance writes:

$$k_a = \frac{\sigma_2^{-2}}{\sigma_1^{-2} + \sigma_2^{-2}}$$

Finally, we obtain the BLUE / estimate / analysis:

$$\boxed{T_a = \frac{\alpha_1 T_1 + \alpha_2 T_2}{\alpha_1 + \alpha_2}} \qquad \text{with} \quad \alpha_i = \frac{1}{\sigma_i^2} \ , \ i = 1, 2. \quad \text{and} \quad \alpha_a = \frac{1}{\sigma_a^2} = \alpha_1 + \alpha_2$$



Equivalency between the BLUE and an optimal solution (the variational one)

\*) The BLUE: 
$$T_a = \frac{\alpha_1 T_1 + \alpha_2 T_2}{\alpha_1 + \alpha_2}$$
 minimizes the following cost function:  
$$\overline{j(T) = \frac{1}{2} \frac{(T - T_1)^2}{\sigma_1^2} + \frac{1}{2} \frac{(T - T_2)^2}{\sigma_2^2}}$$

\*) One remark that the BLUE (obtained under the assumption of uncorrelated errors), minimizes the following cost function too:

 $\sigma_1^2$ 

$$j(T) = \frac{1}{2}(T - T_1, T - T_2) \begin{pmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2\\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}^{-1} \begin{pmatrix} T - T_1\\ T - T_2 \end{pmatrix}$$

i.e. with a-priori correlated errors ! (see the correlation coefficients in the extra diagonal terms).

### \*) T1 and T2 are two measurements uncorrolated: the correlation coefficient vanishs.



\*) T1 and T2 are two measurements corrolated:

$$\rho(T_1, T_2) = \frac{cov(T_1, T_2)}{\sigma_1 \sigma_2} = 0$$

$$\rho(T_1, T_2) = \frac{cov(T_1, T_2)}{\sigma_1 \sigma_2} \neq 0$$



How to improve the estimate ? How to assimilate extra measurements ?

### Assumption.

**O**ne knowns the « accuracy » of each measurement i.e. each variance is known (or equivalently the standard deviation):

 $\sigma_i$ , i = 1, ..., m (for m measurements) given.

Next, each extra measurement may « enhance the accuracy « of the estimate in the sense:

$$\frac{1}{\sigma_a^2} = \frac{1}{\sigma_1^2} + \ldots + \frac{1}{\sigma_m^2}$$

In fact, each extra measurement makes decrease the dispersion from the mean, but not necessarily the accuracy of the estimate if the extra measurements are not accurate !...

The BLUEstimate writes:

$$T_a = \frac{1}{\sum_{i=1}^{m} \frac{1}{\sigma_i^2}} \left( \sum_{i=1}^{m} \frac{1}{\sigma_i^2} T_i \right) \text{ with } \frac{1}{Var(T_a)} = \sum_{i=1}^{m} \frac{1}{\sigma_i^2}$$

Finally, what time is it?

## Chapter 5

# VDA for time-dependent PDE systems (parabolic, hyperbolic)

This chapter aims at extending the VDA formulation to unsteady PDE systems. Formally, the latter can be (well-posed) parabolic or hyperbolic systems, non-linear or not. The calculations derived in the present chapter are formal only.

The VDA algorithm obtained for a time-dependent system is classically called 4D-var. (Actually it is still called "4D-Var" even if the model is not 3D in space...).

In this chapter, the classical 4D-var algorithm is presented. Next, some variants are presented too: 3D-var and incremental version.

The outline of this chapter is as follows<sup>1</sup>.

### Contents

| 5.1        | The   | direct model and cost function  |
|------------|-------|---|
|            | 5.1.1 | Classical PDE models  |
|            | 5.1.2 | The general (and formal) direct model   |
|            | 5.1.3 | The cost function $\ldots \ldots \ldots$        |
| 5.2        | The   | optimization problem & regularization terms 105   |
|            | 5.2.1 | The optimization problem $\ldots \ldots \ldots$ |
|            | 5.2.2 | Regularization terms  |
| 5.3        | Deri  | vation of the optimality equations (in continuous form) 106   |
|            | 5.3.1 | The Linear Tangent Model (LTM) and the resulting gradient $\ . \ . \ . \ . \ 106$   |
|            | 5.3.2 | The adjoint model and the gradient expression   |
| <b>5.4</b> | Deri  | vation of the optimality system in a discrete formalism 110   |
| 5.5        | Algo  | $\operatorname{prithm}(\mathrm{s})$   |

<sup>&</sup>lt;sup>1</sup>Recall that the sections indicated with a \* are "to go further sections". They can be skipped in a first reading or if the reader is not particularly interested in deeper mathematical basis, mathematical proofs.

|     | 5.5.1   | The 4D-var algorithm  |  |
|-----|---|---|--|
|     | 5.5.2   | A few remarks   |  |
| 5.6 | The   | fundamental equations at a glance   |  |
| 5.7 | 7 Examples in Earth Sciences *                                |   |  |
|     | 5.7.1   | Twin experiments concept  |  |
|     | 5.7.2   | Identification of the topography in a 2d shallow-water model $\ . \ . \ . \ . \ 116$          |  |
|     | 5.7.3   | Identification of inflow boundary conditions in a 2d shallow-water model. Flood flows context |  |
| 5.8 | 8 Exercises: optimality systems of classical PDE models * 120 |   |  |
|     | 5.8.1   | Viscous Burgers' equation   |  |
|     | 5.8.2   | Diffusion equation with non constant coefficients   |  |
# 5.1 The direct model and cost function

Let us present a few classical examples of unsteady PDE models which could be "inverted" by the VDA method. Next, the considered general direct model is presented.

# 5.1.1 Classical PDE models

### Linear parabolic: the diffusion equation

The typical and simplest linear parabolic model is the 'heat equation'. This equation models any diffusion phenomena in a homogeneous media. Recall that diffusion phenomena occur in almost all applications: fluids, structures, chemical, images, population, finance etc.

### Non-linear parabolic / hyperbolic equation: Burgers equations

The Burgers equation is a scalar non-linear advection equation (non-linear advection term); its viscous version considers a diffusive term.

The unknown is u(x,t) the fluid velocity at point x and time t.

The forward (direct) model reads as follows. Find u(x,t) which satisfies:

$$\begin{cases} \partial_t u(x,t) - \nu \partial_{xx}^2 u(x,t) + u \partial_x u(x,t) = f(x,t) & \text{in } ]0, L[\times]0, T[\\ u(x,0) = u_0(x) & \text{in } ]0, L[\end{cases}$$
(5.1)

with boundary conditions.

With  $\nu > 0$ , the equation is parabolic. With  $\nu = 0$ , it is hyperbolic.

The Burgers equation ( $\nu = 0$ ) is the most basic non-linear hyperbolic scalar equation, very instructive to understand hyperbolic system solution behaviors.

## Non-linear hyperbolic system: The Saint-Venant equations (shallow-water model)\*

This is a "to go further section"

A model widely used for geophysical shallow flows or industrial thin-film flows, is the shallow-water model (called Saint-Venant equations in its 1D version). It offers a very interesting "reduced dimension" model for fluid flows occurring in thin geometries or under the long wave assumption. For geometries or flows models presenting a ratio  $\varepsilon$  very small,  $\varepsilon$  being the ratio between a characteristics depth  $h^*$  and a characteristic length scale  $L^*$ ,  $\varepsilon = h^*/L^*$ , one can reduce the free-surface 3D Navier-Stokes equations to a fixed boundary BVP: the 2D shallow-water equations. In practice, to be numerically consistent, one has to consider  $\varepsilon$  lower than  $\approx 0.5/0.01$ .

Thin films or shallow geometries are common both in industrial flows (eg coating) and in geophysical flows (eg tsunamis in oceans, rivers / glaciers / lavas dynamics and atmosphere). For environmental flow modeling, we refer the reader to [40].

Shallow-water models are non-linear hyperbolic systems, almost identical to the Euler equations in compressible fluid mechanics.

A first derivation of the shallow equations has been done by Barre de Saint-Venant (1871) for river hydraulics. Note that researches are still active in the topic e.g. in order to include non-linear turbulence terms, complex rheology laws etc. Roughly, these equations are obtained as follows. The Navier-Stokes equations free-surface, incompressible, are depth-integrated. Next, one develop asymptotic expansions with respect to the small (geometric) parameter  $\varepsilon$ . In the end, instead of the 3D Navier-Stokes free-surface model (mobile computational domain), we obtain a 2D model posed in a fixed domain. Thus the reduction of dimension and the computational time saving.

The variable are classically  $(h, \mathbf{q})$  the depth and the discharge of the flow. With such an assumption the pressure is not a direct variable anymore. In the equations h plays the role of p. The pressure value is nearly the hydrostatic value  $\rho gh$ .

Let us present the 2D shallow water equations in their conservative formulation. The state variables are the water depth h and the local discharge  $\mathbf{q} = h\mathbf{u}$ , where  $\mathbf{u} = (u, v)$  is the depth-averaged velocity vector.

On a 2D domain  $\Omega$  and for a computational time interval [0, T], the equations read:

$$\begin{cases} \partial_t h + \operatorname{div}(\mathbf{q}) = 0 & \operatorname{in} \Omega \times ]0, T] \\ \partial_t \mathbf{q} + \operatorname{div}\left(\frac{1}{h} \mathbf{q} \otimes \mathbf{q}\right) + \frac{1}{2}g\nabla h^2 + gh\nabla z_b = S_f & \operatorname{in} \Omega \times ]0, T] \\ h(0) = h_0, \quad \mathbf{q}(0) = \mathbf{q}_0 \\ + \operatorname{boundary \ conditions} \end{cases}$$
(5.2)

where g is the magnitude of the gravity,  $z_b$  the bed elevation,  $h_0$  and  $q_0$  the initial conditions for the state variables.

The term c denotes the local wave celerity,  $c = \sqrt{gh}$ .

In river hydraulics, one introduces the *semi-empirical* source term  $S_f$  modeling the *friction* at bottom. This term reads:

$$S_f(h, \mathbf{q}) = g \frac{n^2 \|\mathbf{q}\|_2}{h^{7/3}}$$

It is the Manning-Strickler law. The scalar friction parameter n is key parameter; however, it has to be empirically estimated... It was historically tuned by trial-test experiments. However, the present VDA method enables to efficiently estimate it.

Without the friction term  $S_f$ , these equations constitute a non-linear hyperbolic system. They are very similar to the Euler equation (compressible flow, unviscous, in a fixed domain).

This system must be closed with boundary conditions. Typical B.C. are as follows:

- at inflow, inflow discharge is imposed.

- at outflow, either incoming characteristic are imposed or a-priori given laws  $\mathbf{q}(H)$  ("rating curves") or imposed depth such as the equilibrium value  $h_n$  (the "normal depth"), see e.g. [40]. - at obstacle ('lateral boundaries'), wall conditions are imposed:  $\mathbf{q} \cdot \mathbf{n} = 0$ .

Let us point out that the solution  $(h, \mathbf{q})$  of the shallow-water system can develop discontinuities and shocks in finite time, even with a regular initial condition. In this case, the solution is locally not differentiable with respect to the "parameters" such as the boundary condition. Keeping in mind this warning, we will however develop an adjoint approach on this system.



Figure 5.1: Example of application of the 2d shallow water model to river hydraulics (3d freesurface flow computed in a 2d fixed domain). The computational domain and its boundaries.

### 5.1.2 The general (and formal) direct model

Let V be a Hilbert space (e.g. the Sobolev space  $H_0^1(\Omega)$  for a scalar linear PDE). Let us assume that we have a first order in time PDE model e.g. the heat equation, the Navier-Stokes equations or the Saint-Venant equations. Then, we consider the following functional space W in time:

$$W_V(0,T) = \{ f \ s.t. \ f \in L^2(0,T;V), \ \partial_t f \in L^2(0,T;V') \}$$

As previously, the state of the system is denoted by y. We have:  $y \in W_V(0,T)$ .

Let  $\mathcal{U}$  be the space of control parameters. it is supposed to be a Hilbert space. Since parameters can depend on time, we introduce the space:  $\mathcal{U}_T = L^2(0, T; \mathcal{U})$ .

 $\mathcal{U}_T$  is a Hilbert space with the norm defined from the following scalar product:

$$\langle u_1, u_2 \rangle_{\mathcal{U}_T} = \int_0^T \langle u_1(t), u_2(t) \rangle_{\mathcal{U}} dt$$

The direct model reads:

$$(\mathcal{D}) \begin{cases} \text{Given } (y_0(x), u(x, t)) \in H \times \mathcal{U}_T, \text{ find } y(x, t) \in W_V(0, T) \text{ such that }: \\ \partial_t y(x, t) + A(u(x, t); y(x, t)) = 0 \quad \text{in } \Omega \times ]0, T[ \\ y(x, 0) = y_0(x) \text{ in } \Omega \end{cases}$$
(5.3)

where  $A : \mathcal{U} \times V \longrightarrow V'$  is a partial derivatives operator,  $y_0 \in H$  is the initial condition, and  $u(t) \in \mathcal{U}$  is the parameter at instant t.

Of course, the solution y(x,t) depends on the I.C.  $y_0(x)$  and on the parameter u(x,t). Let us consider the following control variable:

$$\mathbf{c}(x,t) = (y_0(x), u(x,t))$$
(5.4)

The control  $\mathbf{c}(x,t)$  belongs to the space:  $\mathcal{C} = H \times \mathcal{U}_T$ . In the following, depending on the context, the state is denoted as follows:  $y(\mathbf{c};t), y(t), y(\mathbf{c})$  or simply y.

We assume that the direct model is well posed in the following sense.

Assumption 5.1. Given  $\mathbf{c} \in \mathcal{C}$  and T > 0, it exists a unique function  $y \in W_V(0,T)$  solution of the problem ( $\mathcal{D}$ ). Furthermore, this unique solution depends continuously of  $\mathbf{c}$ .

In other words, we assume the operator model defined by:

$$\mathcal{M}: \mathcal{C} \to V : \mathbf{c}(x,t) = (y_0(x), u(x,t)) \mapsto y(\mathbf{c}(x,t); x,t)$$
(5.5)

is continuous for all  $t \in [0, T[$ .

Furthermore, let us assume the state is differentiable with respect to the control parameters.

Assumption 5.2. The operator model  $\mathcal{M}(\mathbf{c})$  is continuously differentiable for all  $t \in [0, T[$ .

Under the assumption above, formally one can write:

$$y(\mathbf{c} + \delta \mathbf{c}; t) = y(\mathbf{c}; t) + \frac{dy}{d\mathbf{c}}(\mathbf{c}; t) \cdot \delta \mathbf{c} + o_0(\|\delta \mathbf{c}\|_c)$$
(5.6)

This differentiability property is necessary to calculate the cost function gradient, but it can be not satisfied / valid for any "point"  $\mathbf{c}$  or increment  $\delta \mathbf{c}$ .

Typically, non-linear hyperbolic systems (e.g. the shallow-water system) can generate nondifferentiable, even non-continuous, solutions with respect to some input parameters...

Recall that a mathematical result enabling to prove the differentiability is the implicit function theorem. To be satisfied, this theorem requires to have the linearized model well-posed.

### 5.1.3 The cost function

### Misfit to time-dependent observations: averaging in time

The typical objective function J reads:

$$J(\mathbf{c}; y) = J_{obs}(y) + \alpha_{reg} J_{reg}(\mathbf{c})$$
(5.7)

with

$$J_{obs}(y) = \int_{0}^{T} \left\| C(y(t)) - z^{obs}(t) \right\|_{R}^{2} dt$$
(5.8)

and  $J_{reg}(\mathbf{c})$  a regularization term to be determined.

As already discussed in the steady-state case, the choice of its expression depends on the problem and the potential priors. The definition of J is similar to those defined in the steady-state case, excepted the integration in time.

The cost function j to be minimized is defined from J as usual:

$$j(\mathbf{c}) = J(\mathbf{c}; y^{\mathbf{c}})$$
  
where  $y^{\mathbf{c}}$  is the (unique) solution of the direct model given  $\mathbf{c}$ . (5.9)

We have:  $j : \mathcal{C} \to \mathbb{R}$ .

**Pointwise observations** In practice observations are generally point-wise, moreover very sparse. Indeed, it is common to have sensors at few locations only.

Moreover, if these observations are provided as time-series, the misfit observation function reads:

$$J_{obs}(y) = \sum_{n=1}^{N} \left\| C(y(t_n^{obs})) - z_n^{obs} \right\|_{R}^{2}$$

where  $C(\cdot)$  is the observation operator. It is a-priori non linear.  $z_n^{obs}$  is the measurement at instant  $t_n^{obs}$ ,  $n = 1, \ldots, N$ . The complete dataset is:  $z^{obs} = \{z_n^{obs}\}_{n=1,\ldots,N}$ .

# 5.2 The optimization problem & regularization terms

As already defined in the steady-state case, a Variational Data Assimilation (VDA) formulation is nothing else than an optimal control problem with the cost function measuring the discrepancy between model output(s) and observations (datasets).

## 5.2.1 The optimization problem

Given the typical observation function defined by (5.7) and the corresponding cost function j, the VDA problem reads:

$$\begin{cases} \text{Find } \mathbf{c}^* = (y_0^*, u^*) \in H \times \mathcal{U}_T \text{ such that:} \\ j(\mathbf{c}^*) = \min_{H \times \mathcal{U}_T} j(\mathbf{c}) \end{cases}$$
(5.10)

Note that in the present context, the control parameter vector  $\mathbf{c}$  is a-priori time-dependent. In its discrete form, it is a-priori an extremely large vector !

As a consequence, the gradient of the cost function has likely to be computed by the adjoint method.

### 5.2.2 Regularization terms

As already investigated in the steady-state case, the regularization term  $J_{reg}$  acts (very) differently depending on its definition.

Recall that:  $\mathbf{c}(x,t) = (y_0(x), u(x,t))$ . Then, we naturally consider a regularization term for each control component:

$$J_{reg}(\mathbf{c}) = J_{reg}^0(y_0) + J_{reg}^u(u)$$
(5.11)

Classical regularization terms are as follows.

- A background value  $y_b$  of the initial condition  $y_0$  to attract the minimization algorithm to this value ( $y_b$  hopefully relatively good...):

$$J_{reg}^{0}(y_{0}) = \frac{1}{2} \left\| y_{0} - y_{b} \right\|_{B}^{2}$$
(5.12)

with B the inverse of the covariance matrix of the background error.

- For a time-independent model parameter u(x),

$$J_{reg}^{u}(u) = \frac{1}{2} \left\| \nabla u \right\|_{H}^{2}$$
(5.13)

Recall that such Tikhonov regularization terms enhance the "convex property" of the cost function, refer to the discussions in the steady-state case.

The reader may consult e.g. the book [25] for more details on the regularization of optimization problems and Tikhonov's regularization terms in particular.

**Exercice 5.3.** Investigate the influence of different regularization terms during your Programming Practical.

# 5.3 Derivation of the optimality equations (in continuous form)

Let us recall that the (unique) state of the system  $y(\mathbf{c})$  is assumed to be differentiable with respect to  $\mathbf{c}$ .

# 5.3.1 The Linear Tangent Model (LTM) and the resulting gradient

Recall that given a perturbation  $\delta \mathbf{c} \in C$ , the derivative of the state  $y(\mathbf{c})$  in the direction  $\delta \mathbf{c}$  (Gateaux's derivative) satisfies:

$$dy^{\delta \mathbf{c}} \equiv \frac{dy}{d\mathbf{c}}(\mathbf{c}) \cdot \delta \mathbf{c} = \frac{\partial y}{\partial y_0}(\mathbf{c}) \cdot \delta y_0 + \frac{\partial y}{\partial u}(\mathbf{c}) \cdot \delta u$$
(5.14)

### The Linear Tangent Model (LTM)

Let us derive the direct model  $(\mathcal{D})$  with respect to **c** in a direction  $\delta \mathbf{c}$  (at the "point"  $y(\mathbf{c})$ ). We obtain the so-called tangent linear model:

 $(\mathcal{LT}) \begin{cases} \text{Given } \mathbf{c} = (y_0, u) \in H \times \mathcal{U}_T \text{ and } y^{\mathbf{c}} \in W_V(0, T) \text{ solution of the direct problem } (\mathcal{D}), \\ \text{given } \delta \mathbf{c} = (\delta y_0, \delta u) \in H \times \mathcal{U}_T, \text{ find } dy^{\delta \mathbf{c}} \in W_V(0, T) \text{ such that:} \\ \partial_t dy^{\delta \mathbf{c}}(t) + \frac{\partial A}{\partial y}(u(t); y(t)) \cdot dy^{\delta \mathbf{c}}(t) = -\frac{\partial A}{\partial u}(u(t); y(t)) \cdot \delta u(t) \quad \forall t \in ]0, T[ \\ dy^{\delta \mathbf{c}}(0) = \delta y_0 \end{cases}$  (5.15)

The LTM describes the dynamic of the derivative of the state (solution of  $(\mathcal{D})$ ) for a perturbation  $\delta \mathbf{c}$  given.

**Remark 5.4.** If the operator A is linear with respect to the state variable y, then we have:  $\frac{\partial A}{\partial y}(u(t), y(t)) \cdot dy(t) = A(u(t), dy(t))$  and obviously the linearized tangent model is the same as the direct model (modulo the initial condition and the right hand side).

In order to compute  $dy^{\delta \mathbf{c}}$ , we can solve the LTM ( $\mathcal{LT}$ ). But for a different perturbation  $\delta \mathbf{c}^*$ , the LTM must be solved again to obtain  $dy^{\delta \mathbf{c}^*}$  !

Like in the steady-state case that is the reason why the adjoint equations are introduced. The latter enable to circumvent this issue.

### Differential and gradient expressions

Recall that in order to solve numerically the minimization problem, we use a descent algorithm. Hence, we need to evaluate the differentiable of j at any point  $\mathbf{c}$  and in any direction  $\delta \mathbf{c}$  i.e.  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  for any  $\delta \mathbf{c}$ . This is what we classically call improperly "the gradient". Indeed, it is actually the differentiable of j. Strictly speaking, the terminology "gradient" should refer to a vector in a finite dimension space.

Since the control variables has two distinct components,  $\mathbf{c} = (y_0, u)$ , then the "gradient" of  $j(\mathbf{c})$  reads as follows (it is an element of  $\mathcal{L}(H \times \mathcal{U}_T; \mathbb{R})$ ):

$$\nabla j(\mathbf{c}) = (\frac{\partial j}{\partial y_0}(\mathbf{c}), \frac{\partial j}{\partial u}(\mathbf{c}))$$

while we have:

$$j'(\mathbf{c}) \cdot \delta \mathbf{c} = \frac{\partial j}{\partial y_0}(\mathbf{c}) \cdot \delta y_0 + \frac{\partial j}{\partial u}(\mathbf{c}) \cdot \delta u$$

### The gradient expression depending on the term $dy^{\delta c}$

A differentiation of the cost function expression (3.8) provides the relation:

$$j'(\mathbf{c}) \cdot \delta \mathbf{c} = \partial_{\mathbf{c}} J(\mathbf{c}; y^{\mathbf{c}}) \cdot \delta \mathbf{c} + \partial_{y} J(\mathbf{c}; y^{\mathbf{c}}) \cdot dy^{\delta \mathbf{c}}$$
(5.16)

If considering  $J(\mathbf{c}; y)$  defined from a linear observation operator C, and  $J_{reg}(\mathbf{c})$  defined as previously, see (5.11), we have:

$$j'(\mathbf{c}) \cdot \delta \mathbf{c} = \int_0^T \left\langle C^* \Lambda_{\mathcal{O}}(C y^{\mathbf{c}}(t) - z^{obs}(t)), \, dy^{\delta \mathbf{c}}(t) \right\rangle_{V' \times V} dt + \left\langle B(y_0 - y_b), \delta y_0 \right\rangle_H + \left\langle \nabla u, \nabla(\delta u) \right\rangle_H$$
(5.17)

where  $\langle \cdot, \cdot \rangle_{V' \times V}$  denotes the duality product.

Here the operator  $\Lambda_{\mathcal{O}}$  is simply the canonical isomorphism from  $\mathcal{O}$  into  $\mathcal{O}'$ . In discrete dimension, it is simply equal to Identity.

 $C^* \in \mathcal{L}(\mathcal{O}', V')$  denotes the adjoint operator of C. The latter is defined by:

$$\forall \eta \in \mathcal{O}', \quad \forall \xi \in V \qquad \left\langle C^* \eta, \xi \right\rangle_{V' \times V} = \left\langle \eta, C \xi \right\rangle_{\mathcal{O}' \times \mathcal{O}}$$

Like in the steady-state case, the present expression of  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  depends explicitly on  $dy^{\delta \mathbf{c}}$ . The introduction of the adjoint equation will enable to circumvent this issue.

### 5.3.2 The adjoint model and the gradient expression

### The adjoint model

aТ

Let us derive the adjoint model of the problem. Recall that it enables to obtain an expression of  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  linearly dependent of  $\delta \mathbf{c}$ , therefore an explicit gradient expression for *any* direction  $\delta \mathbf{c}$ .

In the unsteady case, the starting point to derive the adjoint equation consists to write:

$$\int_0^T \langle (\mathcal{LT}), p \rangle_{V' \times V} dt = 0 \quad \forall p \in V$$

By integrating by part in time, we get:

$$-\int_{0}^{T} \left\langle dy(t), \partial_{t} p(t) \right\rangle_{V' \times V} dt + \int_{0}^{T} \left\langle \frac{\partial A}{\partial y} \cdot dy(t), p(t) \right\rangle_{V' \times V} dt + \left\langle p(T), dy(T) \right\rangle_{H} - \left\langle p(0), \delta y_{0} \right\rangle_{H} + \int_{0}^{T} \left\langle \frac{\partial A}{\partial u} \cdot \delta u, p(t) \right\rangle_{V' \times V} dt = 0$$

By making appear the adjoint operator, it comes:

$$\int_{0}^{T} \left\langle -\partial_{t} p(t) + \left[\frac{\partial A}{\partial y}\right]^{*} p(t), dy(t) \right\rangle_{V' \times V} dt$$

$$= -\left\langle p(T), dy(T) \right\rangle_{H} + \left\langle p(0), \delta y_{0} \right\rangle_{H} - \int_{0}^{T} \left\langle \left[\frac{\partial A}{\partial u}\right]^{*} p(t), \delta u(t) \right\rangle_{\mathcal{U}' \times \mathcal{U}} dt$$
(5.18)

Considering the expression (5.16) of  $j'(\mathbf{c})$  and if we try to make vanish the two terms in dy(t) above, the following equations are introduced:

$$(\mathcal{A}) \begin{cases} \text{Given } \mathbf{c} = (y_0, u) \in \mathcal{C} \text{ and } y^{\mathbf{c}} \in W_V(0, T) \text{ be the unique solution of the direct problem } (\mathcal{D}), \\ \text{find } p \in W_V(0, T) \text{ such that:} \\ -\partial_t p(t) + \left[\frac{\partial A}{\partial y}(u(t); y^{\mathbf{c}}(t))\right]^* p(t) = \partial_y J(\mathbf{c}; y^{\mathbf{c}}(t)) \quad \forall t \in ]0, T[ \\ p(T) = 0 \end{cases}$$
(5.19)

### J. Monnier

This is the adjoint equation accompanied with its "initial/final" condition.

### The gradient expression

By combining the equations (5.17), (5.18) (5.19), we obtain:

$$j'(\mathbf{c}) \cdot \delta \mathbf{c} = \partial_{\mathbf{c}} J(\mathbf{c}; y(\mathbf{c})) \cdot \delta \mathbf{c} + \langle p(0), \delta y_0 \rangle_H - \int_0^T \left\langle \left[\frac{\partial A}{\partial u}\right]^* p(t), \delta u(t) \right\rangle_{\mathcal{U}' \times \mathcal{U}} dt \qquad (5.20)$$

with  $\mathbf{c}(x,t) = (y_0(c), u(x,t)).$ 

The expression of  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  above does not depend anymore on  $dy^{\delta \mathbf{c}}$  ! This was the goal.

Moreover, by identification, we obtain each gradient components. Indeed, since:

$$j'(\mathbf{c}) \cdot \delta \mathbf{c} = \frac{\partial j}{\partial y_0}(\mathbf{c}) \cdot \delta y_0 + \frac{\partial j}{\partial u}(\mathbf{c}) \cdot \delta u , \qquad (5.21)$$

we have:

$$\begin{cases} \partial_{y_0} j(\mathbf{c}) = p^{\mathbf{c}}(0) + \partial_{y_0} J(\mathbf{c}; y(\mathbf{c})) . \delta y_0 \\ \partial_u j(\mathbf{c}) = -\left[\frac{\partial A}{\partial u}(u(t); y^{\mathbf{c}}(t))\right]^* p^{\mathbf{c}}(t) + \partial_u J(\mathbf{c}; y(\mathbf{c})) . \delta u \end{cases}$$
(5.22)

By considering the expression of  $J_{reg}$  previously defined, see (5.11), we obtain:

$$\begin{cases} \partial_{y_0} j(\mathbf{c}) = p^{\mathbf{c}}(0) + B(y_0 - y_b) \\ \partial_u j(\mathbf{c}) = -\left[\frac{\partial A}{\partial u}(u(t); y^{\mathbf{c}}(t))\right]^* p^{\mathbf{c}}(t) + (J^u_{reg})'(u) \end{cases}$$
(5.23)

### A few remarks

- The adjoint model is retroacting in time. Its initial condition must be given at final time T.

- Like in the steady-state case, by construction the adjoint model is linear, whatever if the direct model is linear or not.

- The linearized problem was supposed to be well-posed therefore the operator  $\partial_y A(u; y^u)$  is an isomorphism from V into V', and its adjoint operator  $(\partial_y A)^*(u; y^u)$  is an isomorphism from V into V' too, see e.g. [11]. As a consequence, the adjoint equation above is well-posed too.

- With the present definition of J, we have:  $\partial_y J(\mathbf{c}; y^{\mathbf{c}}(t)) = \int_0^T C^* \Lambda_{\mathcal{O}} (C y^{\mathbf{c}}(t) - z^{obs}(t)) dt$ . That is the RHS of the adjoint model contains the misfit to the observations:

$$C^*\Lambda_{\mathcal{O}}(Cy(t) - z^{obs}(t))$$

This term measures the discrepancy between the computed state and the observations, in the observations space. It is somehow a "nudging term".

- If the weak form of the direct model equation reads  $a(u; y, z) = l(u; z) \forall z$ , then the adjoint equation reads:

$$\partial_y a(u; y, p) \cdot z = \partial_y J(u; y) \cdot z \ \forall z$$

That is:

$$< \left[\frac{\partial A}{\partial y}(u(t); y^{\mathbf{c}}(t))\right]^* p(t), z(t) >= \partial_y a\left(u(t); y^k, p(t)\right) \cdot z(t)$$
(5.24)

- The gradient with respect to the initial condition equals the adjoint state at initial time (modulo the background term).

Exercice 5.5. Consider the time-dependent case of your practical (linear case or not).
a) Write the adjoint equations; both in the weak and the classical forms.
b) Write the gradient expression based on the adjoint method.

### The 1st order optimality system

As previously (steady-state system in the previous chapter), we can define the optimality system. It is the set of equations characterizing the optimal solution:

- the state equation,
- the adjoint state equation,
- the first order necessary optimality condition (the gradient must vanishes).

# 5.4 Derivation of the optimality system in a discrete formalism

The derivation of the equations in a discrete form involves easier calculations compared to those required in the infinite dimension spaces.

For a sake of clarity, let us consider here a simple explicit time scheme. Then, the formal discrete form of the direct model reads:

$$(D_d) \begin{cases} \text{Given the initial condition } y_0 \in \mathbb{R}^n, \text{ find } \{y_n\}_n \in \mathbb{R}^n \text{ such that :} \\ y_{n+1} = \mathcal{M}_n(y_n) \quad n = 0, 1, ..., N \text{ (the time steps)} \\ y_0 \text{ given} \end{cases}$$
(5.25)

where  $\mathcal{M}_n(y_n)$  denotes the non-linear model at instant  $t^n$ .

For a sake of simplicity, in all this section, only the initial condition  $y_0$  is sough. The cost function reads:

$$j(y_0) = \frac{1}{2} \sum_{n=0}^{N} (\mathcal{H}(y_n) - z_n)^T R^{-1} (\mathcal{H}(y_n) - z_n) + \alpha_0 \frac{1}{2} (y_0 - y_b)^T B^{-1} (y_0 - y_b)$$
(5.26)

where:

-  $\mathcal{H}: \mathbb{R}^n \to \mathbb{R}^m$  is the observation operator; it is a-priori non-linear.

-  $z_n$  is the observation vector at instant  $t_n$ ;  $y_b$  is the background (i.e. the first guess).

- R is the matrix of the covariance observation errors,
- B is the matrix of the covariance background errors.

We denote by:

o  $M_n = \partial_y \mathcal{M}_n(y_n)$  the LTM at  $y_n$ . Therefore, the adjoint operator is  $M_n^T$ . o  $H_n = \partial_y \mathcal{H}(y_n)$  the linear tangent observation operator at  $y_n$ .

Then, the adjoint model reads:

$$(A_d) \begin{cases} \text{Given the states } \{y_k\}_k \in \mathbb{R}^n, \text{ find } \{p_k\}_k \in \mathbb{R}^n \text{ such that }: \\ -p_{n-1} = M_n^T(p_n) + H_n^T R^{-1}(\mathcal{H}(y_n) - z_n) \quad n = N, ..., 1 \text{ (the time steps)} \\ p_N = 0 \end{cases}$$
(5.27)

The gradient (with respect to the initial condition  $y_0$ ) reads:

$$\nabla j(y_0) = p_0 + \alpha_0 B^{-1}(y_0 - y_b) \tag{5.28}$$

# 5.5 Algorithm(s)

The so-called 4D-var algorithm denotes the optimal control algorithm for unsteady PDEs systems (a-priori in 3 space-dimensions plus time therefore the 4D terminology).

## 5.5.1 The 4D-var algorithm

The adjoint model is time dependent, reverse in time.

Since the minimisation is performed by using an iterative descent algorithm (eg the quasi-Newton method BFGS), the algorithm reads as indicated in Fig. (5.2).

Given a first guess  $\mathbf{c}^0$ , compute  $\mathbf{c}^m$  making diminish the cost function j. To do so, at each iteration:

- 1) compute the cost function  $j(\mathbf{c})$  by solving the direct model from 0 to T,
- 2) compute the gradient  $j'(\mathbf{c})$  by solving the adjoint model from T to 0,

3) given the current iteration  $\mathbf{c}^n$ , the cost function value  $j(\mathbf{c}^n)$  and the gradient value  $\nabla j(\mathbf{c}^n)$ , compute a new iteration  $\mathbf{c}^{n+1}$  such that:

$$j(\mathbf{c}^{n+1}) < j(\mathbf{c}^n)$$

\*) Iterations until convergence



Cost function and its gradient

Figure 5.2: Optimal control algorithm for a time-dependent PDE model: the 4D-var algorithm.

## 5.5.2 A few remarks

## A 4D-var like algorithm: what for ?

The 4D-var algorithm can be used for different goals.

A) To estimate the value of uncertain or unknown input parameters (time dependent or not), it is an *identification problem*.

B) To *calibrate* the model in order to perform better predictions; forecasting is the goal. In this case, the data assimilation proceeds by "*analysis cycles*". Figure 5.3 represents one cycle.

In this context,

-) the first stage is called the "analysis step".

Observations of present and past time are assimilated to obtain the analysis i.e. the optimal state value (that is the optimal model trajectory).

-) the second stage consists to run the model in time: it is the *"forecasting step"*. It is expected that if the model fits better the observations in the past, it will be more accurate

for future.

Next, the forecast is used in the next analysis cycle etc



Figure 5.3: The 4D-var algorithm here used to identify the I.C.  $y_0$ . All observation within the assimilation time window are assimilated in the global least-square formulation (analysis step). Next, the resulting calibrated model is runned for prediction (forecasting step).

### The local sensitivity analysis: a rich information

As mentioned above, the adjoint method enables the evaluation of  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  independently of the linearized state  $dy^{\delta \mathbf{c}}$ . This is an explicit expression of all the "gradient" components of  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$ .

Computing a gradient, even without performing a descent algorithm, presents interest in a modeling point of view. Indeed, the gradient  $j'(\mathbf{c}) \cdot \delta \mathbf{c}$  represents the local sensitivity of the model output j with respect the parameters  $\mathbf{c}$ .

The gradient helps to understand the parameter influences of the model and to improve the understanding of the system.

Nevertheless this gradient - sensitivity is local in the sense that it is regarding to the given point  $\mathbf{c}^0$  only, and regarding to the observations used too.

Let us remark that this idea of computing sensitivities in order to better understand both the model and the physics can be applied in a context without observations.

Typically, it can be applied to cost functions depending on the state of the system only in view to study the system stability.

For example, it can be interesting to quantify the sensitivity of the following model output:

$$j(k) = \frac{1}{2} \int_{\omega} \|\nabla y\|^2 dx$$

with  $\omega$  a subset of  $\Omega$ .

### A concluding remark on hybrid approaches

Data assimilation aims at fusing in an "optimal way" all available information: the "physics" of the phenomena (e.g. conservation laws), the parameters (generally empirical), the initial condition (e.g. for weather forecast), the in-situ measurements, the remote-sensed measurements (e.g. extracted from various satellite datasets), prior probabilities (covariance operators defining the norms).

Additional measurements may be used to improve the "analysis", especially if its confidence (accuracy) can be estimated by expertise or a statistical method.

Estimating uncertain parameters of physical-based models could be addressed by "non-physically informed" ("blind") Machine Learning methods e.g. deep Neural Networks, if datasets are large enough.

A typical use of ML in the present context would be to define the first guess by a NN. Next, the VDA process would act as a physically-informed filter.

In other respect, note that the present physically-informed approach enables to introduce statistics and prior probabilities in the formulation.

Moreover, the VDA approach enables to assess the obtained estimations through the "physical model reading".

# 5.6 The fundamental equations at a glance

The unsteady (non-linear) PDE model reads:

$$(\mathcal{D}) \begin{cases} \text{Given } (y_0(t), u(x, t)) \in H \times \mathcal{U}_T, \text{ find } y(x, t) \in W_V(0, T) \text{ such that }:\\ \partial_t y(x, t) + A(u(x, t); y(x, t)) = L(u(x, t)) \quad \text{in } \Omega \times ]0, T[ \\ y(0) = y_0 \text{ in } \Omega \end{cases}$$
(5.29)

It is supposed to be well-posed.

The considered control parameter is:  $\mathbf{c}(x,t) = (y_0(x), u(x,t))$ . Then, the *model operator*  $\mathcal{M}(\mathbf{c})$  is defined as:  $\mathcal{M}(\mathbf{c}) = y^{\mathbf{c}}(x,t)$ .  $\mathcal{M}(\cdot)$  is a-priori non linear.

The considered **objective function** reads:

$$J(\mathbf{c};y) = \frac{1}{2} \int_0^T \left\| C(y(t)) - z^{obs}(t) \right\|_R^2 dt + \alpha_0 \frac{1}{2} \left\langle B(y_0 - y_b), y_0 - y_b \right\rangle_H + \alpha_u J^u_{reg}(u) \quad (5.30)$$

In semi-discrete form, the observations term  $J_{obs}(\mathbf{c}; y)$  reads:  $\frac{1}{2} \sum_{n=1}^{N} \left\| C(y(t_n^{obs})) - z_n^{obs} \right\|_{R}^{2}$ .

The **cost function** is defined as:  $j(\mathbf{c}) = J(\mathbf{c}; y^{\mathbf{c}})$ , where  $y^{\mathbf{c}}$  is the (unique) solution of the direct model given  $\mathbf{c}$ .

The optimization problem reads:

$$\begin{cases} \text{Find } \mathbf{c}^*(x,t) = (y_0^*(x), u^*(x,t)) \in H \times \mathcal{U}_T \text{ such that:} \\ j(\mathbf{c}^*) = \min_{\mathbf{c} \in H \times \mathcal{U}_T} j(\mathbf{c}) \end{cases}$$
(5.31)

The adjoint model reads:

$$(\mathcal{A}) \begin{cases} \text{Given } \mathbf{c} = (y_0, u) \text{ and } y^{\mathbf{c}} \in W_V(0, T) \text{ the unique solution of the direct problem } (\mathcal{D}), \\ \text{find } p \in W_V(0, T) \text{ such that:} \\ -\partial_t p(x, t) + \left[ \partial_y A(u(x, t); y^{\mathbf{c}}(x, t)) \right]^* p(x, t) = \partial_y J(\mathbf{c}; y^{\mathbf{c}}(x, t)) \quad \text{in } \Omega \times ]0, T[ \\ p(x, T) = 0 \end{cases}$$
(5.32)

The partial derivatives of  $j(\cdot)$  with respect to  $y_0(x)$  and u(x,t), i.e. the **gradient components**, read:

$$\begin{cases} \partial_{y_0} j(\mathbf{c}(x,t)) &= p^{\mathbf{c}}(x,0) + \alpha_0 B(y_0 - y_b)(x) \\ \partial_u j(\mathbf{c}(x,t)) &= (-[\partial_u A(u; y^{\mathbf{c}})]^* + [L'(u)]^*) \cdot p^{\mathbf{c}}(x,t) + \alpha_u (J^u_{reg})'(u(x,t)) \end{cases}$$
(5.33)

# 5.7 Examples in Earth Sciences \*

The historical and extremely complex example is briefly presented in Appendix: weather forecasting. For this illustration, the information and images have been extracted from the literature, see the cited sources and references.

Moreover, we present a few numerical results extracted from (not very recent) research studies. The first two subsections present applications in river hydraulics (developed originally at INP Grenoble, then at INSA & Mathematics Institute of Toulouse by the author and his collaborators M. Honnorat, J. Marin, see their work in [24, 33]). Models are based on the 2d shallow-water equations. Such "toy test cases" demonstrate the capabilities of the VDA method applied to river hydraulics. The reference computational code in this topic is DassFlow (Data Assimilation for Free Surface Flows, [17]).

When addressing a DA problem, one first analyses what we call twin experiments.

## 5.7.1 Twin experiments concept

The principle of twin experiments is as follows.

First, we generate data by running the direct model. Next, we can add noise or not to these 'synthetic' data.

Second, we perform the optimal control process from a first guess different than the "true" one which has been used to generate the synthetic data.

Thus, the goal is to recover the set of parameters which gave the data.

Since the "true" optimal control value is known, thorough investigations can be led.

After the validation of the direct model, next those of the adjoint model and the gradient expression (see in a next section how to do), twin experiments constitute the next step to assess and investigate the VDA formulation.

## 5.7.2 Identification of the topography in a 2d shallow-water model

We perform here twin experiments.

The first experiment aims at identifying the topography in an academic case. The domain is  $30 \ m$  long and  $4 \ m$  large. Its topography is defined by:

$$z_b(x,y) = 0.9 \exp\left(-\frac{1}{4}(x-10)^2\right) \exp\left(-(y-1)^2\right) + 0.7 \exp\left(-\frac{1}{8}(x-20)^2\right) \exp\left(-2(y-3)^2\right)$$
(5.34)

The inflow boundary is at x = 0, the outflow boundary at x = 30. Boundaries y = 0 and y = 4 are walls. We use a rectangular structured mesh of dimension  $90 \times 20$ .

Bed roughness, defined by its Manning coefficient, is uniform (n = 0.025). (see the definition of the source term  $S_f$  in the shallow-water model). We impose a constant discharge  $q^{in} = 8 m^3/s$ 

at x = 0 and a constant water height  $h_{out} = 1.4 m$  at x = 30. We obtain a steady state solution after about 80 s of simulation.

Figure 5.4 shows the water height of this steady state solution and the topography.



Figure 5.4: Identification of the topography. Topography and steady state elevation. From [24, 34].

From this steady state solution, we extract the observations:  $h^{obs}$  and  $u^{obs}$  every 0.02 s during 20 s on each cell.

That is we fully observe the steady-state flow.

The objective of this test case is to recover the topography. The first guess is a flat bottom.

We run the data assimilation process with the following cost function:

$$j_1(z_b) = \frac{1}{2} \int_0^T \left( \left\| h(t) - h^{obs}(t) \right\|_{\Omega}^2 + \left\| \mathbf{q}(t) - \mathbf{q}^{obs}(t) \right\|_{\Omega}^2 \right) dt , \qquad (5.35)$$

Figure C.9 shows the cost function and the norm of its gradient (normalized values by the initial values) vs iterations (a) and the identified topography (b). We can notice that convergence is obtained and the reference topography is well recovered.

# 5.7.3 Identification of inflow boundary conditions in a 2d shallowwater model. Flood flows context.

We consider a toy test case which includes many features of real river flows. The computational domain contains a main channel (river) and floodplains, see figures 5.6 and 5.7).



Figure 5.5: Value of the cost function and of the norm of its gradient, normalized by their initial values (a) and the identified topography (b)

Again, the present test case is a twin experiment. At the inflow boundary, we set the inflow discharge shown if figure 5.8 (a) simulating a flood event.

Then we perform a forward run to generate observations at points 1 and 2 shown with black stars in figure 5.7(a).

Then, we suppose that the inflow discharge is constant  $(4.95 \, m^3 s^{-1})$ , and we try to retrieve its real value by assimilating observations.

We present in Figure 5.8 the identified inflow discharge for different experiments. In Fig. 5.8(a), observations are h and  $\mathbf{q}$  at each cell and each time step. In Fig. 5.8(b), observations are h at point 1 and  $(h, \mathbf{q})$  at point 2, both at each time step. In Fig. 5.8(c), observations are h at point 1 only, but at each time step.

We can notice that the identified inflow discharge is good even with the observation of h at point 1 only.

In a practical point of view, such a test case show the ability of the method to identify inflow discharge in a river ...

One notice that the end of the flood event is not well identified. This is the "blind period" phenomena: for example in case (c), the inflow discharge after 270 s can not be identified because the information from the inflow boundary did not reach yet the gauge station.



Figure 5.6: The toy test case mesh (a) and the bathymetry (b)



Figure 5.7: The toy test case domain with the measurement locations (a); data vs time at point 1 (b)



Figure 5.8: In green : reference inflow discharge. In red : identified inflow discharge. (a): Observation of  $(h, \mathbf{q})$  everywhere; (b): Observation of h at point 1 and  $(h, \mathbf{q})$  at point 2. (c): Observation of h at point 1 only. From [24, 34]

# 5.8 Exercises: optimality systems of classical PDE models \*

## 5.8.1 Viscous Burgers' equation

The viscous Burgers' equation is the 1d simplification of the Navier-Stokes momentum equation. It is a scalar non-linear advection diffusion equation (non-linear advection term). The unknown is u(x, t) the fluid velocity at point x and time t.

The control variables we consider in the present example are: the initial condition  $u_0$  and the velocity value at one boundary extremity; the latter is denoted by v.

The forward (direct) model reads as follows.

Given  $\mathbf{c} = (u_0, v)$ , find u which satisfies:

$$\begin{cases} \partial_t u(x,t) - \nu \partial_{xx}^2 u(x,t) + u \partial_x u(x,t) = f(x,t) & \text{in } ]0, L[\times]0, T[\\ u(x,0) = u_0(x) & \text{in } ]0, L[\\ u(0,t) = v(t) ; u(L,t) = 0 & \text{in } ]0, T[ \end{cases}$$
(5.36)

We assume we have m observations points of the flow, continuous in time. Then, we seek to minimize the following cost function:

$$j(\mathbf{c}) = \frac{1}{2} \int_0^T \sum_{i=1}^m |u(x_i) - u_i^{obs}|^2 dt$$

**Exercice 5.6.** Write the optimality system corresponding to this data assimilation problem.

### 5.8.2 Diffusion equation with non constant coefficients

We consider the diffusion equation (or heat equation) in an inhomogeneous media. Let u be the quantity diffused and  $\lambda(x)$  be the diffusivity coefficient, non constant. The forward model we consider is as follows. Given  $\lambda$  and the flux  $\varphi$ , find u which satisfies:

$$\begin{cases} \partial_t u(x,t) - \partial_x(\lambda(x)\partial_x u(x,t)) = f(x,t) & \text{in } \Omega \times ]0, T[\\ u(x,0) = u_0(x) & \text{in } \Omega \\ -(\lambda(x)\partial_n u(x,t) = \varphi & \text{in } \Gamma_1 \times ]0, T[\\ u(x,t) = 0 & \text{in } \Gamma_0 \times ]0, T[ \end{cases}$$
(5.37)

with  $\partial \Omega = \Gamma_0 \cup \Gamma_1$ .

We assume we have measurements of the quantity u at boundary  $\Gamma_1$ , continuously in time. Then, we seek to minimize the following cost function:

$$j(\mathbf{c}) = \frac{1}{2} \int_0^T \int_{\partial\Omega} |u(x) - u^{obs}|^2 ds dt$$

Exercice 5.7. Write the optimality system corresponding to this data assimilation problem.

# Chapter 6

# **Dictionary-based model learning**

The present method has been first proposed in [12][37]. The idea is as interesting as simple. It is here detailed for ODE systems. Next, it is briefly presented in the case of scalar PDE models e.g. the viscous Burgers equation.

Acknowledgments. This chapter has been written with the great help of S. Agdestein, INSA student, during his 4th year summer intership (2019). The corresponding Python notebook has been in great part written by P. Le Dily, S. Ma during their 4th year research project in 2020. T. Malou, PhD IMT-INSA / CLS 2018-21 has helped too.

The outline of this chapter is as follows<sup>1</sup>.

### Contents

| 6.1 Lea | rning ODE terms from a dictionary   |
|---------|---|
| 6.1.1   | Basic principle   |
| 6.1.2   | The inverse problem $\ldots \ldots 125$ |
| 6.1.3   | The optimisation problem $\ldots \ldots 129$   |
| 6.1.4   | The minimisation algorithms: LASSO, STRidge $\ldots \ldots \ldots \ldots \ldots 132$  |
| 6.1.5   | When the goal is to identify model parameters $*$   |
| 6.2 Nur | nerical experiments $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $134$   |
| 6.2.1   | The direct model: non-linear dynamics of a spring $\ldots \ldots \ldots \ldots 134$   |
| 6.2.2   | The inverse problem & strategy to determine $\lambda$   |
| 6.2.3   | Numerical results   |
| 6.2.4   | Programming practical   |

<sup>&</sup>lt;sup>1</sup>Recall that the sections indicated with a \* are to go further sections. They can be skipped in a first reading or if the reader is not particularly interested in deeper mathematical basis, mathematical proofs.

| 6.3 | Lear                | ning PDE model terms                              | -4     |
|-----|---------------------|---|--------|
|     | 6.3.1               | Basic formalism for a generic scalar PDE equation | 4      |
|     | 6.3.2               | Case of the viscous Burgers equation              | 15     |
|     |                     |   |        |
| 6.4 | App                 | $\operatorname{endix}$                            | 8      |
| 6.4 | <b>App</b><br>6.4.1 | endix   | 8<br>8 |

# 6.1 Learning ODE terms from a dictionary

### 6.1.1 Basic principle

Let us consider the general ODE system of unknown  $\boldsymbol{x}(t), \, \boldsymbol{x}(t) = [x_1(t), x_2(t), \cdots, x_N(t)]^T \in \mathbb{R}^N$ :

$$\frac{d\boldsymbol{x}}{dt}(t) = \boldsymbol{f}(\boldsymbol{x}(t)), \ t \in [0,T]$$
(6.1)

with the initial condition:

$$\boldsymbol{x}(0) = \boldsymbol{x}_0, \ \boldsymbol{x}_0 \in \mathbb{R}^N \tag{6.2}$$

We have:  $\boldsymbol{f}: \mathbb{R}^N \to \mathbb{R}^N$ .

Let us assume that one has measurements of  $\mathbf{x}(t)$  and its derivatives  $\dot{\mathbf{x}}(t)$  at time instants  $t_1, \dots, t_M$ .

The present model learning problem aims at *identifying* the unknown EDO expression, that is the unknown expression of f(x(t)).

Note that for higher-order ODE, measuring  $\dot{\mathbf{x}}(t)$  means measuring higher-order derivatives of  $\mathbf{x}(t)$ ...

#### Example #1: the spring -mass system dynamic

Let a mass m be attached to a spring of constant k and friction coefficient  $\nu$ . We denote by  $x(\cdot)$  the mass displacement with respect to its equilibrium location. The second Newton law of motion provides the mass dynamics equation.

**Small displacements case: linear model** For small displacements, a linearized version of the complete mechanical model is sufficiently accurate. The linearized model reads:

$$m\ddot{x}(t) = -\nu\dot{x}(t) - kx(t)$$

By writing this second order ODE as a first order system, we obtain:

$$\dot{\boldsymbol{x}}(t) = \begin{bmatrix} \dot{\boldsymbol{x}}(t) \\ \dot{\boldsymbol{y}}(t) \end{bmatrix} = \boldsymbol{f}(\boldsymbol{x}(t)) = \begin{bmatrix} f_1(\boldsymbol{x}(t)) \\ f_2(\boldsymbol{x}(t)) \end{bmatrix} = \begin{bmatrix} y(t) \\ -\frac{k}{m}\boldsymbol{x}(t) - \frac{\nu}{m}\boldsymbol{y}(t) \end{bmatrix}$$
(6.3)

with  $x_1 = x$  and  $y \equiv x_2 = \dot{x}$  the mass velocity.

In the present example, the model term  $f_1$  contains a single term;  $f_2$  contains two terms.

All terms of f are first-order polynomials.

**Large displacements: the complete non linear model** If the displacements are moderate or large, non-linear effects become non negligible: an additional third order term has to be introduced in the expression of  $f(\mathbf{x})$ . In this case, the equation writes as follows, see Fig. 6.2 for notations:

$$m\ddot{x}(t) = -\nu\dot{x}(t) - k_1 x(t) - k_3 x^3(t)$$
(6.4)

### Example #2: the pendulum dynamics

Let us consider a pendulum of mass m and length L.

Its dynamics equations reads:  $mL\ddot{\theta}(t) = -mg\sin(\theta)$ , with g the gravity constant.

This system re-writes as:

$$\dot{\boldsymbol{x}}(t) = \begin{bmatrix} \dot{\boldsymbol{x}}(t) \\ \dot{\boldsymbol{y}}(t) \end{bmatrix} = \boldsymbol{f}(\boldsymbol{x}(t)) = \begin{bmatrix} f_1(\boldsymbol{x}(t)) \\ f_2(\boldsymbol{x}(t)) \end{bmatrix} = \begin{bmatrix} y(t) \\ -\frac{g}{L}sin(\boldsymbol{x}(t)) \end{bmatrix}$$
(6.5)

with:  $x_1 = x = \theta$  the rotation angle with respect to vertical,  $y = x_2 = \omega$  the angular velocity.

The RHS term  $f_1$  contains a single term; same for  $f_2$ .

In the present case, the RHS does not contain polynomial terms only; it contains the sin term too.

## 6.1.2 The inverse problem

The two examples above are very basic. However, in real-life ODE models, the variety of terms constituting the RHS is quite limited too. In others words, a few terms only constitute the great variety of the employed ODE models: polynomials terms and combination between each others, trigonometric terms etc

The basic idea is to define an a-priori dictionary which may contain the relevant terms to represent the measured phenomena.

Next, using a large amount of measurements (the dataset) and the compress sensing property of  $L^1$ -norm based optimisations, one expects to identify the relevant terms among all terms present in the dictionary.

Notations. The matrix notations are detailed in Section 6.4.1.

In all the sequel, the employed vector size indices are as follows:

•  $m = 1, \cdots, M$ , the time instant index,

- $n = 1, \cdots, N$ , the state component index,
- $d = 1, \dots, D$ , the dictionary component index.

### The a-priori defined dictionary D

We build up a dictionary containing all a-priori relevant terms in the RHS f.

For  $\boldsymbol{x} \in \mathbb{R}^N$ , a typical dictionary may be:

$$\boldsymbol{D}(\boldsymbol{x}^{T}) = \begin{bmatrix} 1 \quad \boldsymbol{x}^{T} \quad (\boldsymbol{x}^{T})^{2} \quad (\boldsymbol{x}^{T})^{3} \quad \cdots \quad \sin(\boldsymbol{x}^{T}) \quad \cos(\boldsymbol{x}^{T}) \quad \sin(2\boldsymbol{x}^{T}) \quad \cdots \end{bmatrix} \in \mathbb{R}^{1 \times D}$$
(6.6)

By convention,  $\mathbf{D}$  is a line vector.

 $(\boldsymbol{x}^T)^2$  denotes the line vector containing all the quadratic possible terms:

$$(\boldsymbol{x}^T)^2 = \begin{bmatrix} x_1^2, & x_1x_2, & x_2^2, & x_1x_3, & \cdots & , x_N^2 \end{bmatrix}$$
 (6.7)

Similarly,  $(\boldsymbol{x}^T)^3$  denotes the line vector containing all the cubic possible terms; etc. Also we have:  $\sin(\boldsymbol{x}^T) = [\sin(x_1), \sin(x_2), \cdots, \sin(x_N)].$ 

### The ODE system re-formulated in function of D

Let us write the ODE system (6.1) from the dictionary **D** terms.

To do wo, we introduce the *coefficients matrix*  $\mathbf{K} = \begin{bmatrix} \mathbf{k}_1 & \mathbf{k}_2 & \cdots & \mathbf{k}_N \end{bmatrix} \in \mathbb{R}^{D \times N}$  such that:

$$\dot{\boldsymbol{x}}^{T}(t) = \boldsymbol{D}(\boldsymbol{x}^{T}(t))\boldsymbol{K}, \quad t \in [0,T]$$
(6.8)

Let us denote by  $\mathbf{k}_n$ ,  $n = 1, 2, \dots, N$ ,  $\mathbf{k}_n \in \mathbb{R}^D$ , the column vector corresponding to the derivative of  $x_n$  in the function set  $\mathbf{D}(\mathbf{x}^T)$ . We obtain:

$$\dot{x}_n(t) = \boldsymbol{D}(\boldsymbol{x}^T(t))\boldsymbol{k}_n, \quad t \in [0, T]$$
(6.9)

Eqn (6.8) is equivalent to Eqn (6.1) with:

$$\boldsymbol{f}(\boldsymbol{x}) = \boldsymbol{K}^T \left( \boldsymbol{D}(\boldsymbol{x}^T) \right)^T, \quad \boldsymbol{x} \in \mathbb{R}^N$$
(6.10)

Recall that:  $\boldsymbol{f} : \mathbb{R}^N \mapsto \mathbb{R}^N$ .

Of course, it is expected (and assumed) that all terms of  $\mathbf{f}$  are present in the dictionary  $\mathbf{D}$ .

**Remark 6.1.** The components of D must be linearly independent to well define the coefficient matrix K. They don't have to be orthogonal. However, if they are orthogonal, the results may be better (see later).

The use of the transpose operator T is necessary when introducing the (discrete) matrices, in particular to adopt the convention that columns represent the field component at time instant n.

### The discrete fields and system

**Time discretization** Let  $\mathbf{t}_M = [t_1, t_2, \cdots, t_M]^T$  be the discrete time instants set with:  $t_0 = 0$  and  $t_M = T$ .

The discrete state matrix X We set:

$$\boldsymbol{X} = \boldsymbol{x}^{T}(\boldsymbol{t}_{M}) = \begin{bmatrix} \boldsymbol{x}^{T}(t_{1}) \\ \boldsymbol{x}^{T}(t_{2}) \\ \vdots \\ \boldsymbol{x}^{T}(t_{M}) \end{bmatrix} = \begin{bmatrix} x_{1}(t_{1}) & x_{2}(t_{1}) & \cdots & x_{N}(t_{1}) \\ x_{1}(t_{2}) & x_{2}(t_{2}) & \cdots & x_{N}(t_{2}) \\ \vdots & \vdots & \ddots & \vdots \\ x_{1}(t_{M}) & x_{2}(t_{M}) & \cdots & x_{N}(t_{M}) \end{bmatrix} \in \mathbb{R}^{M \times N}$$
(6.11)

The lines of X correspond to the time instant. The columns correspond to the different state components  $x_1, x_2, \dots, x_N$ .

The discrete derivatives matrix  $\dot{X}$  We set:

$$\dot{\boldsymbol{X}} = \dot{\boldsymbol{x}}^{T}(\boldsymbol{t}_{M}) = \begin{bmatrix} \dot{\boldsymbol{x}}^{T}(t_{1}) \\ \dot{\boldsymbol{x}}^{T}(t_{2}) \\ \vdots \\ \dot{\boldsymbol{x}}^{T}(t_{M}) \end{bmatrix} = \begin{bmatrix} \dot{x}_{1}(t_{1}) & \dot{x}_{2}(t_{1}) & \cdots & \dot{x}_{N}(t_{1}) \\ \dot{x}_{1}(t_{2}) & \dot{x}_{2}(t_{2}) & \cdots & \dot{x}_{N}(t_{2}) \\ \vdots & \vdots & \ddots & \vdots \\ \dot{x}_{1}(t_{M}) & \dot{x}_{2}(t_{M}) & \cdots & \dot{x}_{N}(t_{M}) \end{bmatrix} \in \mathbb{R}^{M \times N}$$
(6.12)

The discrete dictionary Let us build up the dictionary at all instants  $t_1, t_2, \dots, t_M$  as follows:

$$\boldsymbol{D}(\boldsymbol{X}) = \begin{bmatrix} \boldsymbol{D}\left(\boldsymbol{x}^{T}(t_{1})\right) \\ \vdots \\ \boldsymbol{D}\left(\boldsymbol{x}^{T}(t_{M})\right) \end{bmatrix} = \begin{bmatrix} | & | & | & | & | \\ 1 & \boldsymbol{X} & \boldsymbol{X}^{2} & \boldsymbol{X}^{3} & \cdots & \sin(\boldsymbol{X}) & \cdots \\ | & | & | & | & | & | & \cdots \end{bmatrix} \in \mathbb{R}^{M \times D}$$
(6.13)

D(X) is a matrix  $M \times D$ .

For example, the third column reads:

$$\boldsymbol{X}^{2} = \begin{bmatrix} (\boldsymbol{x}^{T}(t_{1}))^{2} \\ (\boldsymbol{x}^{T}(t_{2}))^{2} \\ \vdots \\ (\boldsymbol{x}^{T}(t_{M}))^{2} \end{bmatrix} = \begin{bmatrix} x_{1}^{2}(t_{1}) & x_{1}(t_{1})x_{2}(t_{1}) & x_{2}^{2}(t_{1}) & \cdots & x_{N}^{2}(t_{1}) \\ x_{1}^{2}(t_{2}) & x_{1}(t_{2})x_{2}(t_{2}) & x_{2}^{2}(t_{2}) & \cdots & x_{N}^{2}(t_{2}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{1}^{2}(t_{M}) & x_{1}(t_{M})x_{2}(t_{M}) & x_{2}^{2}(t_{M}) & \cdots & x_{N}^{2}(t_{M}) \end{bmatrix}$$

**The discrete system** The discrete system corresponding to (6.8) follows:

$$\dot{x}_n(t_m) = \boldsymbol{D}(\boldsymbol{x}^T(t_m))\boldsymbol{k}_n \quad \text{for } m = 1, \cdots, M, \text{ and } n = 1, \cdots, N$$
(6.14)

In matrix form, it reads:

$$\dot{\boldsymbol{X}} = \boldsymbol{D}(\boldsymbol{X})\boldsymbol{K}$$
 in  $\mathbb{R}^{M \times N}$  (6.15)

### The measurement / observations matrices

Let us introduce the measurements matrices.

**Observation of the state**  $\mathbf{x}(t)$  As already mentioned, it is assumed that the state X has been measured at all instants  $t_1, t_2, \dots, t_M$ . Then, we define  $\mathbf{X}^{\text{obs}}$  as:

$$\boldsymbol{X}^{\text{obs}} = \begin{bmatrix} \boldsymbol{x}_{1}^{\text{obs}} & \boldsymbol{x}_{2}^{\text{obs}} & \cdots & \boldsymbol{x}_{N}^{\text{obs}} \end{bmatrix} = \begin{bmatrix} x_{11}^{\text{obs}} & x_{12}^{\text{obs}} & \cdots & x_{1N}^{\text{obs}} \\ x_{21}^{\text{obs}} & x_{22}^{\text{obs}} & \cdots & x_{2N}^{\text{obs}} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1}^{\text{obs}} & x_{M2}^{\text{obs}} & \cdots & x_{MN}^{\text{obs}} \end{bmatrix} \in \mathbb{R}^{M \times N}$$
(6.16)

The scalar value  $x_{mn}^{\text{obs}}$  denotes the observation (measurement) of the state value  $x_n(t_m)$ , for  $m = 1, \dots, M$  and  $n = 1, \dots, N$ .

The vectorial value  $\boldsymbol{x}_n^{\text{obs}}$  (and matrix  $\boldsymbol{X}^{\text{obs}}$ ) correspond to the discrete version of the exact vectorial value  $x_n(\boldsymbol{t}_M)$  (and matrix  $\boldsymbol{X}$  respectively).

**Observation of the state derivatives**  $\dot{\mathbf{x}}(t)$  Similarly, we define:

$$\dot{\boldsymbol{X}}^{\text{obs}} = \begin{bmatrix} \dot{\boldsymbol{x}}_{1}^{\text{obs}} & \dot{\boldsymbol{x}}_{2}^{\text{obs}} & \cdots & \dot{\boldsymbol{x}}_{N}^{\text{obs}} \end{bmatrix} = \begin{bmatrix} \dot{x}_{11}^{\text{obs}} & \dot{x}_{12}^{\text{obs}} & \cdots & \dot{x}_{1N}^{\text{obs}} \\ \dot{x}_{21}^{\text{obs}} & \dot{x}_{22}^{\text{obs}} & \cdots & \dot{x}_{2N}^{\text{obs}} \\ \vdots & \vdots & \ddots & \vdots \\ \dot{x}_{M1}^{\text{obs}} & \dot{x}_{M2}^{\text{obs}} & \cdots & \dot{x}_{MN}^{\text{obs}} \end{bmatrix} \in \mathbb{R}^{M \times N}$$
(6.17)

The scalar value  $\dot{x}_{mn}^{\text{obs}}$  denotes the observation of the derivative value  $\dot{x}_n(t_m)$ , for  $m = 1, \dots, M$ and  $n = 1, \dots, N$ .

### Observations frequency vs model dynamics frequency \*

We have assumed that one has measured the quantity  $x_n(t)$  at the time instants  $t_1, \dots, t_M$ .

This implicitly implies that one has measured  $x_n(t)$  and its derivatives at the aforementionned time scale. This may be a serious limitating feature of the method, even if one can (roughly)

### J. Monnier

estimate  $\dot{x}_{mn}^{\text{obs}}$  by Finite Differences from the values  $x_{mn}^{\text{obs}}$ .

Moreover, one has assumed above that the time sampling  $\{t_1, \dots, t_M\}$  correspond to the time numerical scheme too.

This implicitly means that one has assumed that one are able to observe the system at the numerical time grid scale.

For real-life problems, the time grid of measurements are greatly larger than the required computational grid.

Indeed, one generally have sparse, relatively low frequent and noisy measurements of the modelled phenomena.

The accuracy of measurements for the derivatives may be very rough and are relatively accurate at large scale only. The latter are generally sensitive to the measurements errors.

In summary, the considered time sampling  $\{t_1, \dots, t_M\}$  represent the observed time scale. And in the present description, the latter has to be fine enough to provide stable and accurate numerical scheme too...

## 6.1.3 The optimisation problem

### The non-differentiable convex optimization problem

Given the observations matrix  $X^{obs}$  containing the state measurements and its derivatives at all instants, the model learning problem is written as follows. Solve the N following (non differentiable) optimization problems:

$$(P_n): \min_{\boldsymbol{k}_n \in \mathbb{R}^D} \left\{ \|\boldsymbol{D}(\boldsymbol{X}^{\text{obs}})\boldsymbol{k}_n - \dot{\boldsymbol{x}}_n^{\text{obs}}\|_2^2 + \lambda \|\boldsymbol{k}_n\|_1 \right\} \quad \text{for } n = 1, \cdots, N$$
(6.18)

with  $\lambda > 0$  a weight parameter.

Recall that  $[\mathbf{k}_n]_{n=1,\dots,N} = \mathbf{K}$  denotes the coefficients of the derivatives  $[\dot{x}_n]_{n=1,\dots,N} = \dot{\mathbf{x}}^T$  in the dictionary  $\mathbf{D}$ .

 $\forall n = 1, \cdots, N,$ 

$$\boldsymbol{k}_n \in \mathbb{R}^D, \boldsymbol{X}^{\text{obs}} = \begin{bmatrix} \boldsymbol{x}_1^{\text{obs}} & \cdots & \boldsymbol{x}_N^{\text{obs}} \end{bmatrix} \in \mathbb{R}^{M \times N} \text{ and } \boldsymbol{D}(\boldsymbol{X}^{\text{obs}}) \in \mathbb{R}^{M \times D}$$
 (6.19)

with N the number of states, M the number of instants, D the number of terms in the dictionary D.

Finally, once the dictionary has been defined, this "model learning method" aims at solving a standard linear optimization problem, involving a  $L^1$ -regularization term.

In matrix form In matrix form, the N problems  $(P_n)$  can be reformulated as follows. Find the matrix K, solution of:

$$(P): \min_{\boldsymbol{K} \in \mathbb{R}^{D \times N}} \left\{ \|\boldsymbol{D}(\boldsymbol{X}^{\text{obs}})\boldsymbol{K} - \dot{\boldsymbol{X}}^{\text{obs}}\|_{2}^{2} + (\lambda N) \|\boldsymbol{K}\|_{1} \right\}$$
(6.20)

To obtain the implication between the N problems  $(P_n)$  and Problem (P), the weight parameter has to be equal to  $(\lambda N)$ .

The matrix norm definitions are recalled in Section 6.4.1. Recall that the  $L^1$ -norm is the "closest convex relaxation" of the pseudo-norm  $l^0$ ,

$$\|\boldsymbol{k}\|_{0} = \#\{k_{d}: k_{d} \neq 0, \ d = 1, \cdots, D\} = \sum_{d=1}^{D} \mathbf{1}_{k_{d} \neq 0}$$
(6.21)

 $L^2$ -norm vs  $L^1$ -norm. Note that if setting the minimisation problem as a simple least-square problem i.e. in the  $L^2$ -norm only, then the optimal solution **K** would be optimal in the least-square sense given *entire* dictionary **D**. This would not enable to identify the correct model terms. The  $L^1$ -norm with its compress sensing properties does.

The present model learning method is illustrated in Fig. 6.1.



Figure 6.1: The model learning method as presented in [12]. Illustration based on the (chaotic) Lorenz dynamical system.

### An identifiability condition

It has been demonstrated in the 2000's that a "sparse" signal k (that is containing a great majority of vanishing coefficients) can be identified from an amount of measurements much lower than for a "complete" signal, see e.g. in [13].

### J. Monnier

Let us apply the estimation presented in [13] in the present context.

Let us consider  $\mathbf{k} \in \mathbb{R}^D$ , k is supposed to be sparse.

Let  $D \in \mathbb{R}^{M \times D}$  be an *orthonormal* matrix. We consider the measurements  $y^{obs} = Dk$ ,  $y \in \mathbb{R}^{M}$ .

The following identifiability condition holds, see [13]:

$$M \ge \operatorname{cst} \mu^2(\boldsymbol{D}) \, \|\boldsymbol{k}\|_{\#0} \, \log(D) \tag{6.22}$$

with:  $cst \approx 1$ ,  $\mu(\mathbf{D}) = \sqrt{P} \max_{m,d}(|D_{md}|)$  a conditioning measure of  $\mathbf{D}$ ,  $\mu(\mathbf{D}) \approx 1$  if  $\mathbf{D}$  orthonormal.

Therefore, the number of measurements M has to be at least as  $\mathcal{O}(D\log(D))$ .

### In the case the measurements are perturbed, inaccurate

In real-life problems, even if the measurements of the state  $\mathbf{X}$  are very accurate, the measurements of its derivatives  $\dot{\mathbf{X}}$  are likely not, see the discussion above. These derivatives may be derived from the measurements of the state.

Therefore in practice one has to consider the following perturbed system:

$$\dot{\boldsymbol{X}}^{\text{obs}} = \boldsymbol{D}(\boldsymbol{X}^{\text{obs}})\boldsymbol{K} + \eta \boldsymbol{Z}$$
(6.23)

with  $\boldsymbol{Z}$  an aleatory variable following e.g. the normal distribution, and  $\eta$  the standard deviation value.

This uncertainty introduces an additional (aleatory) RHS in the ODE system.

### This Model Learning method vs VDA

Let us wonder how the present Model Learning method "degenerates" if the state derivatives cannot be measured.

The method previously presented requires to have measured both the state X and its derivative  $\dot{X}$  at all time steps, equivalently at a given time scale, cf the previous discussion. The state values  $X^{\text{obs}}$  are observed at instants  $t_M = [t_1, \cdots, t_M]^T$ .

If one cannot estimate the derivatives  $\dot{X}$  (at least at the same time scale), the method may re-write as follows.

Given a dictionary D, given the initial state  $x_0 \in \mathbb{R}^N$ , given the observations instants  $t_M = [t_1, \cdots, t_M]^T \in \mathbb{R}^M$ , given the observations matrix  $X^{\text{obs}} \in \mathbb{R}^{M \times N}$ , find the matrix coefficients

**K** satisfying:

$$\min_{\boldsymbol{K} \in \mathbb{R}^{M \times D}} \left\{ \|\boldsymbol{X}_{\boldsymbol{K}} - \boldsymbol{X}^{\text{obs}}\|_{2}^{2} + \lambda \|\boldsymbol{K}\|_{1} \right\}$$
(6.24)

under the contraints 
$$\begin{cases} \dot{\boldsymbol{x}}_{\boldsymbol{K}}^{T}(t) = \boldsymbol{D}(\boldsymbol{x}_{\boldsymbol{K}}^{T}(t))\boldsymbol{K}, \ t \in [0,T] \\ \boldsymbol{x}_{\boldsymbol{K}}(0) = \boldsymbol{x}_{0} \\ \boldsymbol{X}_{\boldsymbol{K}} = \boldsymbol{x}_{\boldsymbol{K}}^{T}(\boldsymbol{t}_{M}) \end{cases}$$
(6.25)

Therefore in this case, the model learning method turns into an optimal control problem where K is the control variable.

**Remark 6.2.** In the present case, the "problem components" are not independent of each others. Indeed, changing  $\mathbf{k}_n \in \mathbf{K}$ , implies a change of all others components through the model. This is not the case in Problem (6.18). Indeed in this last case, each vector  $\mathbf{k}_n$ ,  $n = 1, \dots, N$  can be computed independently of each other.

# 6.1.4 The minimisation algorithms: LASSO, STRidge

To solve non-differentiable but convex optimization problems such as (6.18), the historical algorithm is the Least Absolute Shrinkage and Selection Operator (LASSO) algorithm, [41]. This algorithm is recalled in Section 6.4.2.

Among other possible algorithms, let us mention the STRidge (Sequentially Thresholded Ridge regression) algorithm and the ElasticNet algorithm, [45], see Section 6.4.2.

A Ridge regression consists to consider the  $L^2$ -norm in (6.18) instead of  $L^1$ . The obtained Least Square solution present small coefficients (but a-priori non vanishing) for the less important terms present in the dictionary.

The STRidge (Sequentially Thresholded Ridge regression) version consists to set to 0 the values lower than a threshold, therefore imposing parsimony.

ElasticNet combines these two approaches by considering the two norms  $L^2$  and  $L^1$ .

### A few remarks

• As already discussed in the optimal control sections, the weight parameter  $\lambda$  of the regularization term, a  $L^1$ -norm term in LASSO has two roles: "convexifying" the cost function and tuning the regularization.

In the LASSO algorithm case, this defines the importance of parsimony.

• System (6.8) is a linear problem in K. The identification of an optimal value of  $\lambda$  in the optimisation problem (6.18) may be done by a cross-validation strategy.

# 6.1.5 When the goal is to identify model parameters \*

Let us consider the scalar ODE model:

$$\frac{d\boldsymbol{x}}{dt}(t) = \boldsymbol{f}(\boldsymbol{\nu}(t); \boldsymbol{x}(t)), \ t \in [0, T]$$
(6.26)

with  $\boldsymbol{\nu}(t) = [\nu_1(t), \nu_2(t), \cdots]^T$  an unknown, uncertain parameter to be identified.

### **Basic** principle

Following the principle of the model learning method above, one can introduce in the dictionary combinations of terms with  $\nu$ .

For example:  $\nu(t)x_1(t)x_2(t)$ , etc.

Obviously this make highly increase the dictionary dimension...

As a consequence, this approach seems a-priori possible for a very few number L of parameters,  $\nu(t) = (\nu_1(t), \dots, \nu_L(t))$ , only.

### Examples of reduction of a model parameter $\nu(t)$

In real-life modeling problems, one may have good a-priori on the parameters form. In this case, one may be able to reduce a time or space dependent parameter to a very few number of unknown parameters. For examples:

- A sinusoidal form with unknown frequency, phase and amplitude:  $\nu(t) = \nu_0 \sin(\omega_0 t + \varphi_0)$ . This reduces the time dependent parameter  $\nu(t)$  to three parameters only.
- An exponentially decreasing parameter :  $\nu(t) = \nu_0 e^{-a_0 t}$ . Two parameters only in this case.

## A general parametrization of $\nu(t)$

More generally, given a function basis  $(\psi_n)_{n=1,\dots,N_{\psi}}$  with  $N_{\psi} \ll M$ , one may set:

$$\nu(t) = \sum_{n=1}^{N_{\psi}} \nu_n \psi_n(t)$$
(6.27)

In this case, one has  $N_{\psi}$  parameters jonly;

However, in the dictionary, these  $N_{\psi}$  additionnal unknown parameters may multiply the total number of entries in **D** by  $N_{\psi}$ .

It would provide  $\boldsymbol{D} \in \mathbb{R}^{M \times N_{\psi}D}$ . This may be prohibitive...

# 6.2 Numerical experiments

We illustrate the model learning method in the case of the simple non-linear oscillator previously described.

# 6.2.1 The direct model: non-linear dynamics of a spring

The dimensionless equations



Figure 6.2: The spring system and notations. (Image extracted from S. Agdestein's INSA internship report).

The considered model is the non linear spring dynamic equation (6.4).

The state of the system in *physical dimensions* is denoted by  $\tilde{x}$  [m]; it represents the mass displacement with respect to its equilibrium position, see Fig. 6.2. The final time is denoted by  $\tilde{T}$  [s].

The equation in physical dimensions reads:

$$m\ddot{\tilde{x}}(\tilde{t}) = -\tilde{\nu}\dot{\tilde{x}}(\tilde{t}) - \tilde{k}_1\tilde{x}(\tilde{t}) - \tilde{k}_3\tilde{x}^3(\tilde{t}), \ \tilde{t} \in [0,\tilde{T}]$$

$$(6.28)$$

with initial conditions.

The parameter m is the spring mass,  $\tilde{\nu}\left[\frac{Ns}{m}\right]$  is a friction coefficient,  $\tilde{k}_1\left[\frac{N}{m}\right]$  is a resistivity related coefficient,  $\tilde{k}_3\left[\frac{N}{m^3}\right]$  is the actual resistivity coefficient.

The dimensionless system, in 1st order form Let us normalize and dimensionless the equations.

L denotes the spring length.  $T_0$  is the oscillation duration in the case of small displacements:  $T_0 = 2\pi \sqrt{\frac{m}{k_1}} = \frac{2\pi}{\omega_0}.$ 

The dimensionless (and normalized) variables are:

$$x = \frac{\tilde{x}}{L}$$
 and  $t = \frac{\tilde{t}}{T_0}$ 

In the 1st order form, the dimensionless equation reads:

$$\dot{\boldsymbol{x}}(t) = \begin{bmatrix} \dot{\boldsymbol{x}}(t) \\ \dot{\boldsymbol{y}}(t) \end{bmatrix} = \begin{bmatrix} y(t) \\ -\nu y(t) - k_1 x(t) - k_3 x^3(t) \end{bmatrix} \equiv \boldsymbol{f}(\boldsymbol{x}(t)) = \begin{bmatrix} f_1(\boldsymbol{x}(t)) \\ f_2(\boldsymbol{x}(t)) \end{bmatrix}, \ t \in [0, T]$$
(6.29)

with  $y = \dot{x}$  and  $T = \frac{\tilde{T}}{T_0}$ ,  $\nu = \frac{T_0}{m}\tilde{\nu} = \frac{2\pi\tilde{\nu}}{\sqrt{m\tilde{k}_1}}$ ,  $k_1 = 4\pi^2$ ,  $k_3 = \frac{L^2T_0^2}{m}\tilde{k}_3 = 4\pi^2 L^2 \frac{\tilde{k}_3}{\tilde{k}_1}$ .

Note that smaller displacements are, lower the non linear terms are too.

Using the previous notations, we have:

$$\boldsymbol{f}(\boldsymbol{x}(t)) = \begin{bmatrix} f_1(\boldsymbol{x}(t)) \\ f_2(\boldsymbol{x}(t)) \end{bmatrix} = \begin{bmatrix} y(t) \\ -\nu y(t) - k_1 x(t) - k_3 x^3(t) \end{bmatrix}$$

### Preliminary direct simulations

The considered inverse problem is the following: given measurements, do we manage to identify the dynamics equation terms ?

Let us recall first the following evidence: before solving any inverse problem, one must perfectly understand the phenomena and the direct model.

Following this remark, we perform few numerical direct simulations. Fig. 6.3 illustrates different direct simulations.



Figure 6.3: Direct simulation in the case of very small displacements (Case 1), for two values of  $k_3$ , see Tab. 6.1). (L) The state x(t) vs time t. (R) The phase portrait.

#### Setting up the dictionary D

After having observed the system dynamics, typical questions a modeler may wonder are: is there non linear effects ? Is there dissipation effects (due to friction) ? Recall that a friction classically correspond to a third order term. Then, one selects the following a-priori model terms:

$$\boldsymbol{D}(\boldsymbol{x}^{T}) = \begin{bmatrix} 1 & x & y & x^{2} & xy & y^{2} & x^{3} & x^{2}y & xy^{2} & y^{3} \end{bmatrix} \in \mathbb{R}^{1 \times 10}$$
(6.30)

The direct model (6.29) reads in function of the dictionary **D** as follows:

$$\dot{\boldsymbol{x}}^{T}(t) = \boldsymbol{D}(\boldsymbol{x}^{T}(t))\boldsymbol{K}^{\text{exact}}, \ t \in [0,T]$$
(6.31)

with  $\boldsymbol{K}^{\mathrm{exact}}$  defined by:

$$\boldsymbol{K}^{\text{exact}} = \begin{bmatrix} 0 & 0 \\ 0 & -k_1 \\ 1 & -\nu \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & -k_3 \\ 0 & 0 \\ 0 & -k_3 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} 1 \\ (x) \\ (x) \\ (x^2) \\ (x^2) \\ (x^2y) \\ (xy^2) \\ (xy^2) \\ (y^3) \end{pmatrix}$$
(6.32)

Therefore, the exact model reads:

$$\begin{bmatrix} \dot{x} & \dot{y} \end{bmatrix} = \begin{bmatrix} 1 & x & y & x^2 & xy & y^2 & x^3 & x^2y & xy^2 & y^3 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & -k_1 \\ 1 & -\nu \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$
(6.33)

### The time discrete system

The time sampling is fine:  $\mathbf{t}_M = [t_m]_{m=1,\dots,200}$ , with  $t_m = (m-1)\Delta t$  and  $\Delta t = 0.05$  (s). We have:  $T = t_M = 10.0$  (s). The dimension tabs are: N = 2, M = 200 and D = 10.

The time discrete system reads:  $\dot{X} = D(X)K$ , with

$$\dot{\mathbf{X}} = \begin{bmatrix} \dot{x}(t_1) & \dot{y}(t_1) \\ \dot{x}(t_2) & \dot{y}(t_2) \\ \vdots & \vdots \\ \dot{x}(t_{200}) & \dot{y}(t_{200}) \end{bmatrix}$$
(6.34)
$$\boldsymbol{D}(\boldsymbol{X}) = \begin{bmatrix} 1 & x(t_1) & y(t_1) & x^2(t_1) & xy(t_1) & \cdots & y^3(t_1) \\ 1 & x(t_2) & y(t_2) & x^2(t_2) & xy(t_2) & \cdots & y^3(t_2) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x(t_{200}) & y(t_{200}) & x^2(t_{200}) & xy(t_{200}) & \cdots & y^3(t_{200}) \end{bmatrix}$$
(6.35)

### Synthetic measurements: twin experiments

**Generation of the measurements** To generate the measurements, the direct model is performed with some parameters values. This provides the measurements for the model term identification phase. Then the inverse problem consists to recover the employed parameters values.

Such a procedure based on synthetic measurements generated by the model is called "twin experiments".

The employed parameter dimensionless values are indicated in Tab. 6.1.

Two set of measurements are generated. A first one with an initial displacement  $x_0$  equal to 3% of L. A second one with  $x_0$  equal to 15% of L.

| Parameters values |                         | · - | Initial value $x_0$ |              |
|-------------------|-------------------------|-----|---------------------|--------------|
| $k_1$             | $4\pi^2 \approx 39.478$ |     | Case 1              |              |
| ν                 | 0.562                   |     | Small displacements | $x_0 = 0.03$ |
| $k_3$             | 1.421                   |     | Case 2              |              |
| $T_0$             | 7.118                   |     | Large displacements | $x_0 = 0.15$ |

Table 6.1: (L) The (dimensionless) parameter values used to generate the two measurements sets. Recall that L = 1.

(R) The two different initial displacements corresponds to the two cases.

The measurements of the state x(t) are the computed values of x at all time instants  $t_1, \dots, t_{200}$ . The measurements of its derivatives  $y(t) = \dot{x}(t)$  are computed as the centered Finite Difference approximations at these same time instants. The values of x and  $\dot{x}$  constitutes the observation matrix  $X^{\text{obs}}$ .

To define the observation matrix  $\dot{X}^{\text{obs}}$ , one approximates the second derivatives of x by Finite Differences again.

**Remark.** Measuring second derivatives of a field x in real-world may be highly critical, in particular at "small" scale. It may be easier at larger scale than at the scale of  $x^{obs}$ . This point constitutes a limitation of the present model learning method.

Note that if performing the dimensionless dynamic system then the observations matrices have to be dimensionless too. **Learning set & testing set** To evaluate the method, the measurements set are split in two subsets:

- a first one containing the measurements from 0 to T/2, employed for the learning phase,

- a second one to compare the prediction capabilities (after the learning phase) vs the true model dynamics, from T/2 to T.

Then, the matrices are split as follows:

$$\boldsymbol{D}(\boldsymbol{X}^{\text{obs}}) = \begin{bmatrix} \boldsymbol{D}(\boldsymbol{X}^{\text{obs.learn}}) \\ \boldsymbol{D}(\boldsymbol{X}^{\text{obs.test}}) \end{bmatrix} \begin{array}{c} 50\% \\ 50\% \end{array}$$
(6.36)

$$\dot{\boldsymbol{X}}^{\text{obs}} = \begin{bmatrix} \dot{\boldsymbol{X}}^{\text{obs.learn}} \\ \dot{\boldsymbol{X}}^{\text{obs.test}} \end{bmatrix} \begin{bmatrix} 50\% \\ 50\% \end{bmatrix}$$
(6.37)

This decomposition is illustrated in Fig. 6.4.



Figure 6.4: The train set [0, T/2] and the test set [T/2, T]: (L) x(t) vs time t. (R) The phase portrait.

Training data are in blue (learning interval). "Testing data" are in orange (validation interval).

### 6.2.2 The inverse problem & strategy to determine $\lambda$

We solve here the following inverse problem: identifying the model terms with their corresponding coefficients  $k_{\Box}$  and  $\nu$ , given the observations previously defined. This reformulates as follows.

Given  $(\mathbf{X}^{\text{obs.learn}}, \dot{\mathbf{X}}^{\text{obs.learn}})$  and the dictionary **D**, identify the dominating model terms among

those present in  $\mathbf{D}$ .

### Recall of the optimization formulation

To solve this model learning problem, we solve the N convex non differentiable optimization problems 6.18 by using the LASSO algorithm, see Section 6.4.2 for details.

### J. Monnier

Recall the optimization problem to be solved:

$$\boldsymbol{K}^{*}(\lambda) = \underset{\boldsymbol{K} \in \mathbb{R}^{D \times N}}{\operatorname{arg\,min}} \left\{ \|\boldsymbol{D} \left( \boldsymbol{X}^{\operatorname{obs.learn}} \right) \boldsymbol{K} - \dot{\boldsymbol{X}}^{\operatorname{obs.learn}} \|_{2}^{2} + \lambda \|\boldsymbol{K}\|_{1} \right\}$$
(6.38)

with  $\lambda > 0$  a weight parameter to be a-priori set.

### Computing strategy in terms of $\lambda$

When solving this bi-objective optimization problem (6.38), the setting of the weight parameter  $\lambda$  is critical.

Larger  $\lambda$  is, more the compress sensing feature (i.e.  $||K||_0$ ) is important i.e. less model terms is selected.

On the contrary, smaller  $\lambda$  is, less the compress sensing feature is important.

In practice, it is interesting to follow the following strategy.

- One starts from the minimal large value  $\lambda$  that provides 1 model term only. The term  $\|\mathbf{K}\|_1$  in the cost function is so dominant that the algorithm provides 1 component optimal solution only.
- Next, the optimal solution is performed with  $\lambda$  divided e.g. by a factor 2 at each run. This is done until a small value of  $\lambda$  which provides the total number of potential terms, see Fig. 6.5 and Fig. 6.6(Down)à.

In this last case, the misfit term  $\|D(X^{\text{obs.learn}})K - \dot{X}^{\text{obs.learn}}\|_2^2$  in the cost function is so dominant that the algorithm does not compress the components number.



Figure 6.5: The number of model term numbers increases when  $\lambda$  decreases, from 0 to the total number of terms

### Determining an "optimal value" of $\lambda$

**Cross-validation method** As already mentioned, the identification of an optimal value of  $\lambda$  in the optimisation problem (6.18) may be done by the classical cross-validation strategy, here based e.g. on various dynamic trajectories datasets.

The optimal criteria to determine an optimal value of  $\lambda$  is then the misfit term. The goal is to

fit at best data without over-fitting them.

Another "performance criteria" approach Another approach to analyse the optimization process with respect to  $\lambda$  can be as follows.

- The optimization problem (6.38) for numerous values of λ as previously described. Starting from a large enough value of λ, it is next decreased until a minimal value makes appear all terms, see Fig. 6.5. A sequence of optimal solutions {**K**<sup>\*</sup><sub>1</sub>(λ),..., **K**<sup>\*</sup><sub>nr</sub>(λ)} is obtained.
- For each optimal solution  $\mathbf{K}^*(\lambda)$ , the following *performance criteria* is computed.

$$\mathcal{C}^{perf}(K^*(\lambda),\eta) = \|\boldsymbol{D}(\boldsymbol{X}^{\text{obs.test}})K^*(\lambda) - \dot{\boldsymbol{X}}^{\text{obs.test}}\|_2 + \eta \|K^*(\lambda)\|_0$$
(6.39)

The coefficient  $\eta$  is empirically set. A few values can be considered.

• Finally, "optimal" values of  $\lambda$  may be defined as follows:

$$\lambda^*(\eta) = \operatorname*{arg\,min}_{\lambda>0} \left\{ \mathcal{C}^{perf}(K^*(\lambda), \eta) \right\}$$
(6.40)

This strategy is illustrated in the numerical results below, see e.g. Fig. 6.6. This method helps to analyse the optimization processes.

### 6.2.3 Numerical results

### Analysis of the performance criteria

The figures 6.6 (Up) and 6.8 (Up) present the values of the performance criteria  $\mathcal{C}^{perf}(\mathbf{K}^*(\lambda),\eta)$  vs  $\lambda$ .

In Fig. 6.6 (Top),  $C^{perf}(\mathbf{K}^*(\lambda), \eta)$  is plotted for three different values of  $\eta$ :  $\eta = 10^{-4}\kappa(\mathbf{D})$ ,  $10^{-3}\kappa(\mathbf{D})$  and  $10^{-2}\kappa(\mathbf{D})$ ,  $\kappa(\mathbf{D})$  the conditioning number of  $\mathbf{D}$ .

Figures 6.6 and 6.8 differ from the initial displacements values  $x_0$  only, see Tab. 6.1.



Figure 6.6: Case 1:  $x_0 = 3\%$  of L (small displacements case). (Up) The performance criteria  $C^{perf}(\mathbf{K}^*(\lambda), \eta)$  (indicated as  $J_{\eta}^{\text{test}}(\mathbf{K}^*(\lambda))$  vs  $\lambda$  for three different values of  $\eta$  ( $\eta = 10^{-4}\kappa(\mathbf{D})$ ,  $10^{-3}\kappa(\mathbf{D})$ ,  $10^{-2}\kappa(\mathbf{D})$ ). These curves suggest to consider  $\lambda^* \approx 0.08$ .

(Down) For  $\lambda^* \approx 0.08$ , we recover the 3 non vanishing components of K. Recall that here the non-linear term is too small to be influential.

### Recall of the model to be recovered

Let us recall the (dimensionless) coefficient matrix  $\mathbf{K}$  value which has to be recovered:

$$\boldsymbol{K}^{\text{exact}} = \begin{bmatrix} 0 & 0 & (1) \\ 0 & -39.478 & (x) \\ 1 & -0.562 & (y) \\ 0 & 0 & (x^2) \\ 0 & 0 & (xy) \\ 0 & 0 & (y^2) \\ 0 & -1.421 & (x^3) \\ 0 & 0 & (x^2y) \\ 0 & 0 & (xy^2) \\ 0 & 0 & (xy^2) \\ 0 & 0 & (y^3) \end{bmatrix}$$
(6.41)

Therefore, one has to identify 4 non vanishing coefficients: 1 dominating (the linear term coefficient  $k_1$ ) and 2 a-priori less influent terms (the non linear terms coefficients  $k_3$  and  $\nu$ ).

#### Identified model terms & results analysis

Case 1): small displacements The first numerical experiment is first performed with  $\lambda = 1$ .

In this case, the algorithm makes vanish all coefficient excepted  $(-k_1)$ , see Fig. 6.6 for  $\lambda = 1$ .

Its value is estimated at -31.211. The corresponding model does not enable to accurately describe the spring dynamics.

For  $\lambda = 0.5$ , two non vanishing coefficients are identified, see Fig. 6.6 for  $\lambda = 0.5$ . The obtained estimations are:  $1 \approx 0.435$ ,  $(-k_1) \approx (-35.038)$ .

These values correspond to the model:  $\dot{x} = 0.435y$  and  $\dot{y} = -35.038x$ . By taking into account the I.C., its solution reads:  $x(t) = 0.03 \cos(\sqrt{0.435 \times 35.038}t)$ . This cannot be the correct model since its solution x(t) does not damp as the observations do, see Fig. 6.4.

The performance criteria seems minimal for  $\lambda^* \equiv \lambda \approx 0.08$ , see Fig. 6.6, with 3 identified terms.

For  $\lambda = 0.117$ , that is for a greater value than  $\lambda^*$ , the optimization algorithm provides the following estimations:

$$\boldsymbol{K}^{*}(\lambda = 0.117) = \begin{bmatrix} 0 & 0 \\ 0 & -39.479 \\ 0.868 & -0.562 \\ 0 & 0$$

The linear terms coefficients only are detected. The second column terms are accurately identified. The estimation of the coefficient 1 is still inaccurate.

For the half of the previous value,  $\lambda = 0.055$  ( $\lambda < \lambda^* \approx 0.08$ ), all the terms including the non

linear ones are accurately identified, see Fig. 6.6:

$$\boldsymbol{K}^{*}(\lambda = 0.055) = \begin{bmatrix} 0 & 0 \\ 0 & -39.478 \\ 0.938 & -0.562 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & -1.420 \\ 0 & 0 \\$$

For lower values of  $\lambda$ , more non vanishing coefficients are detected. However, in these cases the performance criteria is increasing, see Fig. 6.6. This is a typical overfitting behavior.

It is to worth to notice that each dictionary term appears in its importance order when  $\lambda$  decreases, see Fig. 6.7.

This remarkable property enables to identify the model terms by their order of importance.

In summary, one has easily detected the two dominant terms: the coefficients 1 and (-35.038). Moreover, as expected in the small displacements case, it is more tricky to detect the non linear term(s).

This result is consistent with small displacement dynamics because in such case the non linear terms are negligible. Non linear terms are almost useless in the model for small displacements. Despite this low influence, the non linear terms appear when the value of  $\lambda$  is decreased.

Case 2): larger displacements In this case and following the same procedure, one easily detect all coefficients, see Fig. 6.8.

Moreover, the estimated values of all coefficients are accurate, see (6.41) for the exact values. The detailed results are not shown here.

### 6.2.4 Programming practical

Please consult the INSA Moodle page of the course. Upload the provided computational code (written in Python). Next, the goal is to make your own numerical experiments by e.g.:

- changing numerical parameters,
- changing the density of observations and by perturbing their values (unperfect measurements),



Figure 6.7: The model terms appear by their importance order when  $\lambda$  decreases. Here a test case with "small displacements". (L) The first column of the matrix K. (R) The second column of the matrix K.

Read the increasing terms number from 1 to the maximum value, when  $\lambda$  decreases. This test case is not the one analyzed ! Fig. to be updated...

• changing the optimisation formulation and employed algorithm.

### 6.3 Learning PDE model terms

Formally, the model learning method previously described can be extended to PDE models. This has first been done in [37].

### 6.3.1 Basic formalism for a generic scalar PDE equation

Let us consider the formal generic 2D scalar PDE equation:

$$\partial_t u + A\left((\nu_1, \cdots, \nu_P); (c_1, c_2), \partial_x u, \partial_y u, \partial_{xx}^2 u, \partial_{yy}^2 u, u \partial_x u, u \partial_y u, \dots\right) = f(x, t) \text{ in } \Omega \times (0, T)$$

$$(6.44)$$

with u(x,t) the unknown and  $\nu = (\nu_1, \cdots, \nu_P)$  a multi-components parameter.

### Examples



Figure 6.8: Same as Fig. 6.6 but in Case 2 i.e. for a larger initial displacement:  $x_0 = 15\%$  of L. (Up) In this case and if considering  $C^{perf}(K^*(\lambda), \eta)$  as performance criteria, we obtain  $\lambda^* \approx 0.269$ . (Down) The minimum of the performance criteria correspond to the 4 non vanishing components of K. Here, the non-linear term is easily detected and correctly quantified.

The advection-diffusion equation In this formalism, the 2D linear advective-diffusion equation reads:

$$\partial_t u + A((\nu, c_1, c_2); \partial_x u, \partial_y u, \partial_{xx}^2 u, \partial_{yy}^2 u) = f(x, t)$$
(6.45)

with:

$$A = -\partial_x(\nu(x)\partial_x u(x,t)) - \partial_y(\nu(x)\partial_y u(x,t)) + c_1\partial_x u(x,t) + c_2\partial_y u(x,t)$$
(6.46)

with  $c = (c_1, c_2)$  a given velocity field.

The viscous Burgers equation The 2D viscous Burgers's equation corresponds to:

$$A = -\partial_x(\nu(x)\partial_x u(x,t)) - \partial_y(\nu(x)\partial_y u(x,t)) + u_1\partial_x u(x,t) + u_2\partial_y u(x,t)$$
(6.47)

### 6.3.2 Case of the viscous Burgers equation

#### **Dictionary definition**

In the viscous 1D Burgers equation case, the a-priori dictionary may be defined for example as:

$$\boldsymbol{D}(u;\nu) = \begin{bmatrix} 1, & u, & u^2, & c_1 u \partial_x u & \nu_1 \partial_{xx}^2 u & \cdots & , & \nu_P \partial_{xx}^2 u \end{bmatrix}$$
(6.48)

Recall that if considering third order terms, waves dispersion can be captured.

By adopting the same formalism as in the ODE case, one has:

$$u_t = \boldsymbol{D}(u, \nu) \mathbf{K} \tag{6.49}$$

with  $\mathbf{K}$  the matrix coefficient.

### The discrete fields and system

On has to consider a time-space grid. Let us consider for simplicity the following basic grid:  $[(x_m, t_m)]_{m=1,\dots,M}$ .

Using this basic space-time grid", the discrete fields and matrices read:

$$\boldsymbol{U}_{M} = \begin{bmatrix} u(x_{1}, t_{1}) \\ \vdots \\ u(x_{M}, t_{M}) \end{bmatrix}$$
$$\boldsymbol{D}(\boldsymbol{U}_{M}, \boldsymbol{\nu}_{P}) = \begin{bmatrix} \boldsymbol{D}(u(x_{1}, t_{1}), \boldsymbol{\nu}_{1..P}) \\ \vdots \\ \boldsymbol{D}(u(x_{M}, t_{M}), \boldsymbol{\nu}_{1..P}) \end{bmatrix}$$

The discrete system reads:

$$\boldsymbol{U}_t^M = \boldsymbol{D}(\boldsymbol{U}_M, \boldsymbol{\nu}_P) \mathbf{K}$$
(6.50)



Figure 6.9: Illustration of the PDE model learning method from [37] with an application to the viscous Burger's equation.

Given  $U^{\text{obs}}$  the measurements of the field u, the method steps are as follows:

### J. Monnier

- 1. Compute by Finite Differences the derivatives of  $U^{obs}$ ,
- 2. Build up the dictionary  $\boldsymbol{D}(\boldsymbol{U}, \partial_x \boldsymbol{U}_x, \partial_{xx}^2 \boldsymbol{U}, \cdots) \in \mathbb{R}^{M \times D}$ ,
- 3. Solve the non differentiable optimization problem:

$$oldsymbol{K} = rgmin_{oldsymbol{K} \in \mathbb{R}^{D imes N}} \left\{ \|oldsymbol{D}(\mathbf{U}, \cdots)oldsymbol{K} - oldsymbol{U}_t\|_2^2 + \lambda \|oldsymbol{K}\|_1 
ight\}$$

The present learning methods have been implemented in the Python PDE\_FIND package, [12, 37].

The latter includes various basic PDE examples. Also, some codes are available on the INSA Moodle page of your course.

### 6.4 Appendix

### 6.4.1 Notations matricielles & normes

On définit par  $\mathbb{R}$  l'espace des scalaires (nombres réels), par  $\mathbb{R}^N$  l'espace de vecteurs de dimension N et par  $\mathbb{R}^{M \times N}$  l'espace de matrices réelles de dimension  $M \times N$ . La notation matricielle consiste à traiter les scalaires, les vecteurs et les matrices de la même manière : comme des matrices. Ainsi, on a – par convention –, les égalités suivantes :  $\mathbb{R} = \mathbb{R}^{1 \times 1}$  et  $\mathbb{R}^N = \mathbb{R}^{N \times 1}$ . Pour un scalaire x on a donc l'égalité x = [x].

### Fonctions iscalaires;;

So t une fonction scalaire F :

$$\begin{array}{rccc} F: & \mathbb{R} & \to & \mathbb{R} \\ & x & \mapsto & F(x) \end{array}$$

Alors, pour un vecteur  $\boldsymbol{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$ , on note  $F(\boldsymbol{x}) = \begin{bmatrix} F(x_1) \\ \vdots \\ F(x_N) \end{bmatrix} \in \mathbb{R}^N = \mathbb{R}^{N \times 1}$ , tandis que  $F(\boldsymbol{x}^T) = \begin{bmatrix} F(x_1) & \cdots & F(x_N) \end{bmatrix} \in (\mathbb{R}^N)^T = \mathbb{R}^{1 \times N}$ .

De même, pour une matrice  $\boldsymbol{X} = \begin{bmatrix} x_{11} & \cdots & x_{1N} \\ \vdots & \ddots & \vdots \\ x_{M1} & \cdots & x_{NN} \end{bmatrix} \in \mathbb{R}^{M \times N}$ , on désigne par  $F(\boldsymbol{X})$  la matrice

$$\begin{bmatrix} F(x_{11}) & \cdots & F(x_{1N}) \\ \vdots & \ddots & \vdots \\ F(x_{M1}) & \cdots & F(x_{NN}) \end{bmatrix} \in \mathbb{R}^{M \times N}.$$

### Fonctions ;;vectorielles;;

Soit F une fonction à valeur vectorielle :

$$F: \mathbb{R} \to \mathbb{R}^{P}$$
$$x \mapsto F(x) = \begin{bmatrix} F_{1}(x) \\ \vdots \\ F_{P}(x) \end{bmatrix}$$

Les deux définitions suivantes sont équivalentes

$$\boldsymbol{F}(\boldsymbol{x}^{T}) = \begin{bmatrix} F_{1}(\boldsymbol{x}^{T}) \\ \vdots \\ F_{P}(\boldsymbol{x}^{T}) \end{bmatrix} = \begin{bmatrix} F_{1}(x_{1}) & \cdots & F_{1}(x_{N}) \\ \vdots & \vdots \\ [F_{P}(x_{1}) & \cdots & F_{P}(x_{N}) \end{bmatrix} = \begin{bmatrix} F_{1}(x_{1}) & \cdots & F_{1}(x_{N}) \\ \vdots & \ddots & \vdots \\ F_{P}(x_{1}) & \cdots & F_{P}(x_{N}) \end{bmatrix} \in \mathbb{R}^{P \times N}$$

$$\boldsymbol{F}(\boldsymbol{x}^{T}) = \begin{bmatrix} \boldsymbol{F}(x_{1}) & \cdots & \boldsymbol{F}(x_{N}) \end{bmatrix} = \begin{bmatrix} F_{1}(x_{1}) \\ \vdots \\ F_{P}(x_{1}) \end{bmatrix} & \cdots & \begin{bmatrix} F_{1}(x_{N}) \\ \vdots \\ F_{P}(x_{N}) \end{bmatrix} \end{bmatrix} = \begin{bmatrix} F_{1}(x_{1}) & \cdots & F_{1}(x_{N}) \\ \vdots & \ddots & \vdots \\ F_{P}(x_{1}) & \cdots & F_{P}(x_{N}) \end{bmatrix} \in \mathbb{R}^{P \times N}$$

Cette notation donne donc des libertés en considérant des matrices par blocs.

### Normes vectorielles et matricielles

Soit  $\boldsymbol{x} \in \mathbb{R}^N$ . On définit les normes- $\ell_1$  et - $\ell_2$  :

$$\begin{cases} \|\boldsymbol{x}\|_{1} = \sum_{n=1}^{N} |x_{n}| \\ \|\boldsymbol{x}\|_{2} = \sqrt{\sum_{n=1}^{N} x_{n}^{2}} \end{cases}$$

Soit  $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1 & \cdots & \boldsymbol{x}_N \end{bmatrix} = \begin{bmatrix} \boldsymbol{x}_{11} & \cdots & \boldsymbol{x}_{1N} \\ \vdots & \ddots & \vdots \\ \boldsymbol{x}_{M1} & \cdots & \boldsymbol{x}_{MN} \end{bmatrix} \in \mathbb{R}^{M \times N}$ . On définit les normes- $\ell_1$  et - $\ell_2$  de la

matrice  $\boldsymbol{X}$ :

$$\begin{cases} \|\boldsymbol{X}\|_{1} = \sup_{\boldsymbol{y} \in \mathbb{R}^{N}} \frac{\|\boldsymbol{X}\boldsymbol{y}\|_{1}}{\|\boldsymbol{y}\|_{1}} = \max_{n=1,\cdots,N} \|\boldsymbol{x}_{n}\|_{1} = \max_{n=1,\cdots,N} \sum_{m=1}^{M} |x_{mn}| \\ \\ \|\boldsymbol{X}\|_{2} = \sup_{\boldsymbol{y} \in \mathbb{R}^{N}} \frac{\|\boldsymbol{X}\boldsymbol{y}\|_{2}}{\|\boldsymbol{y}\|_{2}} = \sqrt{\lambda_{\max}(\boldsymbol{X}^{T}\boldsymbol{X})} \\ \\ \|\boldsymbol{X}\|_{F} = \sqrt{\sum_{n=1}^{N} \|\boldsymbol{x}_{n}\|_{2}^{2}} = \sqrt{\sum_{n=1}^{N} \sum_{m=1}^{M} x_{mn}^{2}} \end{cases}$$

avec  $\sqrt{\lambda_{\max}(\boldsymbol{X}^T\boldsymbol{X})}$  la plus grande valeur singulière de  $\boldsymbol{X}$ , et  $\|\cdot\|_F$  la norme de Frobenius.

#### Algorithmes d'optimisation 6.4.2

Soit le problème

$$\min_{\boldsymbol{k}\in\mathbb{R}^{D}} \{\|\boldsymbol{k}\|_{0}\} \text{ sous la contrainte } \boldsymbol{y} = \boldsymbol{D}\boldsymbol{k}$$
(6.51)

avec  $\boldsymbol{y} \in \mathbb{R}^M$  et  $\boldsymbol{D} \in \mathbb{R}^{M \times D}$  donnés. Dans cette section, on va énoncer quelques méthodes pour trouver des solutions à ce problème.

### LASSO

Least Absolute Shrinkage and Selection Operator (LASSO) [41] consiste à résoudre le problème

$$\min_{\boldsymbol{k} \in \mathbb{R}^D} \left\{ \|\boldsymbol{D}\boldsymbol{k} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{k}\|_1 \right\}$$
(6.52)

remplaçant la norme- $\ell_0$  par une norme- $\ell_1$ , rendant le problème convexe (mais non-différentiable). Si  $\boldsymbol{D}$  est orthonormée, i.e.  $\boldsymbol{D}^T \boldsymbol{D} = \boldsymbol{I}$ , alors la solution du problème de minimalisation est donnée par  $\boldsymbol{k} = [k_p]_{p=1,\dots,P}^T$ , avec  $k_p = \hat{k}_p \max(0, |\hat{k}_p| - \lambda)$  et  $\hat{\boldsymbol{k}} = \underset{\boldsymbol{k}' \in \mathbb{R}^P}{\operatorname{arg\,min}} \{ \|\boldsymbol{D}\boldsymbol{k}' - \boldsymbol{y}\|_2^2 \} = (\boldsymbol{D}^T \boldsymbol{D})^{-1} \boldsymbol{D}^T \boldsymbol{y}$  la solution du problème des moindres carrées pûr. Le LASSO est implémenté dans l'algorithme ?? avec une méthode de gradients proximaux rapide.

### Algorithm \* Normaliser D

\*  $L \leftarrow \left\| \boldsymbol{D}^T \boldsymbol{D} \right\|_2$  {Constante de Lipschitz}

 $^{*} oldsymbol{k} \leftarrow \mathbf{0}_{\mathbb{R}^{P}}$ 

$$^{*} oldsymbol{k}_{0} \leftarrow oldsymbol{0}_{\mathbb{R}^{P}}$$

for iters  $= 1, \cdots, \text{maxiters}$  {Fast proximal gradient method (FISTA)}

\* 
$$\boldsymbol{z} \leftarrow \boldsymbol{k} + \frac{\text{iter}}{\text{iter}+1} (\boldsymbol{k} - \boldsymbol{k}_0)$$
  
\*  $\boldsymbol{k}_0 \leftarrow \boldsymbol{k}$   
\*  $\boldsymbol{z} \leftarrow \boldsymbol{k} - \boldsymbol{D}^T (\boldsymbol{D}\boldsymbol{z} - \boldsymbol{y}) / L$   
for  $p = 1, \cdots, P$   
\*  $k_p \leftarrow \text{sign}(z_p) \max(0, |z_i| - \frac{\lambda}{L})$   
endfor  
endfor

\* grands\_coeffs =  $[p: k_p \neq 0]$  {Moindres carrées sur les coefficients non nuls}

\* 
$$\boldsymbol{k}$$
 [grands\_coeffs] =  $\underset{\boldsymbol{k}' \in \mathbb{R}^{|\text{grands_coeffs}|}}{\operatorname{arg min}} \|\boldsymbol{D}$  [:, grands\_coeffs]  $\boldsymbol{k}' - \boldsymbol{y}\|_2^2$ 

\* Dénormaliser  $\boldsymbol{k}$ 

\* Return  $\boldsymbol{k}$ 

LASSO( $\boldsymbol{D}, \boldsymbol{y}, \lambda$ , maxiters

### Ridge

Ridge consiste à résoudre le problème

$$\min_{\boldsymbol{k}\in\mathbb{R}^D}\left\{\|\boldsymbol{D}\boldsymbol{k}-\boldsymbol{y}\|_2^2+\lambda\|\boldsymbol{k}\|_2^2\right\}$$
(6.53)

remplaçant la norme- $\ell_0$  par une norme- $\ell_2$ , rendant le problème convexe (et différentiable). Ridge ne garantit pas que les coefficients de  $\boldsymbol{k}$  soient nuls, mais les incitent à être petits. On peut remarquer que la solution du problème de minimisation est  $\boldsymbol{k} = (\boldsymbol{D}^T \boldsymbol{D} + \lambda \boldsymbol{I}_D)^{-1} \boldsymbol{D}^T \boldsymbol{y}$ , où  $\boldsymbol{I}_D$ est la matrice d'identité de dimension D.

### STRidge

Sequential Threshold Rigde regression (STRidge) consiste à étendre le principe Ridge en tronquant les petits coefficients. Soit une tolérence et un nombre d'itérations à effectuer. On définit alors l'algorithme **??**. Le choix de la tolérance peut être automatisé, en utilisant des techniques d'apprentissage.

Algorithm \* Trouver  $\boldsymbol{k}$  solution de  $\min_{\boldsymbol{k} \in \mathbb{R}^P} \|\boldsymbol{D}\boldsymbol{k} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{k}\|_2^2$ , c.-à-d.  $\boldsymbol{k} = (\boldsymbol{D}^T \boldsymbol{D} + \lambda \boldsymbol{I}_P)^{-1} \boldsymbol{D}^T \boldsymbol{y}$ 

if iters  $\neq 0$ 

\* grands\_coeffs  $\leftarrow [p : |k_p| \ge \text{tol}]$ 

\*  $\boldsymbol{k} [\sim \text{grands\_coeffs}] \leftarrow 0$ 

\* k [grands\_coeffs]  $\leftarrow$  STRidge (D [:, grands\_coeffs], y, tol, iters -1) {appel récursif}

endif

\* Return  $\boldsymbol{k}$ 

 $STRidge(\boldsymbol{D}, \boldsymbol{y}, tol, iters)$ 

### ElasticNet

L'algorithme ElasticNet [45] consiste à combiner les approcher de LASSO et de Ridge, en pénalisant et une norme- $\ell_1$  et une norme- $\ell_2$ , pour des coefficients  $\lambda_1$  et  $\lambda_2$  respectivement donnés. On résout donc le problème

$$\min_{\boldsymbol{k}\in\mathbb{R}^D}\left\{\|\boldsymbol{D}\boldsymbol{k}-\boldsymbol{y}\|_2^2+\lambda_1\|\boldsymbol{k}\|_1+\lambda_2\|\boldsymbol{k}\|_2^2\right\}$$
(6.54)

### Appendices

### Appendix A

# A few mathematical and optimization recalls

The outline of this chapter is as follows.

In this chapter, few basic mathematical and numerical concepts are shortly recalled. The reader is invited to consult his previous courses or courses available on-line for more details. We suggest the reading of the following excellent books: [1, ?] for recalls in Sobolev spaces, numerical analysis of PDE equations and differentiable optimization; also [11] for any recall or improvement in functional analysis.

### A.1 Differential calculus, functional analysis

### A.1.1 Differentials

First, let us recall the differentiability and differential definitions. In all the sequel, by default we denote by: X a Banach space, V a Hilbert space, (.,.) its scalar product and  $\|.\|$  its norm. The functional j is defined from X (or V) into  $\mathbb{R}$ .  $\Omega$  denotes a bounded subset of  $\mathbb{R}^d$ ;  $\Omega$  is a Lipschitz domain (roughly, its boundary  $\partial\Omega$  can be locally described as a Lipschitz graph).

The Fréchet-differential is defined as follows.

**Definition A.1.** Let  $u \in X$ . The functional j is Fréchet-differentiable at point u, if there exists a linear form continuous  $L, L \in \mathcal{L}(X, \mathbb{R}) = X'$ , such that:

$$j(u+\delta u) = j(u) + L(\delta u) + o(\delta u) \quad \text{with } \lim_{\delta u \to 0} \frac{|o(\delta u)|}{\|\delta u\|} = 0$$

We call L the (Fréchet-)differential of j at point u, and we denote:  $L = \frac{dj}{du}(u)$ .

In all this manuscript, the differential  $\frac{dj}{du}(u)$  can be denoted as follows too: dj(u), j'(u), depending on the context...

Next, we define the Gâteaux differential.

**Definition A.2.** Let u and  $\delta u \in X$ . The functional j is Gâteaux-differentiable at point u in the direction  $\delta u$ , if there exists a linear form continuous  $L, L \in \mathcal{L}(X, \mathbb{R}) = X'$ , such that:

$$\lim_{\varepsilon \to 0^+} \frac{1}{\varepsilon} [j(u + \varepsilon \delta u) - j(u)] = L(\delta u)$$

The Gâteaux-differentiable concept is weaker than the Fréchet-differential one ("Fréchet  $\implies$  Gâteaux" but the reciprocal is not true). Nevertheless its definition can be more convenient. In all the manuscript, when "differentiable" is mentioned, it means Fréchet-differentiable (hence Gâteaux-differentiable in any direction  $\delta u$ ).

The definition extension to the second derivative is natural. Next the following Taylor expansion order 2 holds:

**Proposition A.3.** If the functional j is twice differentiable then:

$$j(u+\delta u) = j(u) + j'(u)(\delta u) + \frac{1}{2}j''(u)(\delta u, \delta u) + o(\|\delta u\|^2) \quad \text{with } \lim_{\delta u \to 0} \frac{o(\|\delta u\|^2)}{\|\delta u\|^2} = 0$$

where j''(u) is a bilinear form defined from  $X \times X$  into  $\mathbb{R}$ .

**Exercise A.4.** Let X be a Banach space and j be a functional:  $j : X \to \mathbb{R}$ . - Calculate the differential expressions at the "point"  $u_0$  in the direction  $\delta u$ , in the case: 1)  $j(u) = \int_{\Omega} u^p dx$ , with  $X = L^p(\Omega), 1 \le p < \infty$ .

- Calculate the 1st and 2nd order differential at  $u_0$  in the direction  $\delta u$  and  $(\delta u, \delta u)$  respectively, in the following cases:

2)  $j(u) = \nabla u$ , with  $X = H^1(\Omega)$ . 3)  $j(u) = \int_{\Omega} \|\nabla u\|^2 dx$ , X a subspace of  $H^1(\Omega)$ .

**Exercise A.5.** Let V be a subspace of  $H^1(\Omega)$  and let a(.,.) be a bilinear form defined from  $V \times V$  into  $\mathbb{R}$ . Give the differential expression  $\partial_u a(u_0, v) \cdot \delta u$ , in the case: 1)  $a(u, v) = \int_{\Omega} u^p v dx$ ,  $1 \leq p < \infty$ . 2)  $a(u, v) = \int_{\Omega} \nabla u \nabla v dx$ .

**Exercice A.6.** Let V be a Hilbert space, let a(.,.) be a bilinear form defined from  $V \times V$  into  $\mathbb{R}$ , l a linear form from V into  $\mathbb{R}$ . Let j be the functional defined by:  $j(u) = \frac{1}{2}a(u, u) - l(u)$ . Prove that j is twice differentiable in V, and

$$j'(u)(v) \equiv j'(u) \cdot v = a(u, v) - l(v) \quad \forall u \in V, \ \forall v \in V$$
$$j''(u)(v, w) = a(v, w) \quad \forall u, v, w \in V$$

**Exercise A.7.** Let  $V = L^2(\Omega)$ , a(.,.) the bilinear symmetric form defined on  $V \times V$  by:  $a(u,v) = \int_{\Omega} uvdx$  and l(.) the linear form defined by:  $l(v) = \int_{\Omega} fvdx$ , with  $f \in L^2(\Omega)$ . We set:  $j(u) = \frac{1}{2}a(u,u) - l(u)$ . By identifying V and V', prove that: j'(u) = u - f.

### A.1.2 Green formulas

The formulas are key tools in the derivation of weak formulations of the model equations. Let  $\Omega$  be a bounded subset of  $\mathbb{R}^d$ ,  $\Omega$  regular in the sense  $C^1$  (roughly, its boundary  $\partial\Omega$  is locally the graph of a  $C^1$  function). Let  $n = (n_i)_{1 \le i \le d}$  be the external unit normal of the boundary. In  $\Omega$ , dx denotes the volumic measure (Lebesgue's measure in dimension d), and ds the surface measure (Lebesgue's measure in dimension (d-1) on the subvariety  $\partial\Omega$ ).

We have the following important equalities.

**Theorem A.8.** (Green formulae). Let  $\Omega$  be a bounded subset of  $\mathbb{R}^d$ ,  $C^1$ -regular, let v be a function in  $C^1(\overline{\Omega})$  (its support is bounded in  $\overline{\Omega}$ ). We have:

$$\int_{\Omega} \frac{\partial v}{\partial x_i}(x) dx = \int_{\partial \Omega} v(x) n_i ds \quad 1 \le i \le d$$

**Corollary A.9.** (Integration by part / Green formula). Let  $\Omega$  be a bounded subset of  $\mathbb{R}^d$ ,  $C^1$ -regular, let u and v be two functions in  $C^1(\overline{\Omega})$  (with their support bounded in  $\overline{\Omega}$ ). We have:

$$\int_{\Omega} u(x) \frac{\partial v}{\partial x_i}(x) dx = -\int_{\Omega} \frac{\partial u}{\partial x_i}(x) v(x) dx + \int_{\partial \Omega} u(x) v(x) n_i ds \quad 1 \le i \le d$$

**Corollary A.10.** Let  $\Omega$  be a bounded subset of  $\mathbb{R}^d$ ,  $C^1$ -regular, let u be a function in  $C^2(\overline{\Omega})$ , and v be a function in  $C^1(\overline{\Omega})$  (with their support bounded in  $\overline{\Omega}$ ). We have:

$$\int_{\Omega} \Delta u(x)v(x)dx = -\int_{\Omega} \nabla u(x)\nabla v(x)dx + \int_{\partial\Omega} \frac{\partial u}{\partial n}(x)v(x)ds \quad 1 \le i \le d$$

with:  $\frac{\partial u}{\partial n} = \nabla u \cdot n.$ 

### A.1.3 The Riesz-Fréchet theorem

The Riesz-Fréchet representation theorem makes a connection between a Hilbert space and its dual space. Its result is useful to derive rigorously the adjoint operator including non trivial boundary conditions.

**Theorem A.11.** (*Riesz-Fréchet*). Let V be a Hilbert space and L be an element of V'. It exists an unique element  $l \in V$  such that:

$$L(v) \equiv < L, v >_{V' \times V} = (l, v) \ \forall v \in V$$

Furthermore:  $||l||_V = ||L||_{V'}$ .

In other words, the dual space V' can be represented by elements of V and the scalar product (., .) of V. Furthermore this defines an isomorphism between V' and V.

**Exercice A.12.** Let V be a subspace of  $H^1(\Omega)$ , let l(.) be a bilinear form defined from V into  $\mathbb{R}$ , and let a(.,.) be a bilinear form defined from  $V \times V$  into  $\mathbb{R}$ . Apply the Riesz-Fréchet representation theorem to the following cases: 1)  $l(v) = \int_{\Omega} fv dx, f \in L^2(\Omega)$ . 2)  $a(u, v) = \int_{\Omega} \nabla u \nabla v dx$ . 3)  $a(u, v) = \int_{\Omega} \lambda \nabla u \nabla v dx + \int_{\Omega} cu^3 v dx$ , with  $\lambda$  and u given in  $L^{\infty}(\Omega)$ .

### A.1.4 The adjoint operator

Let E and F be two Banach spaces, E' and F' their dual space respectively. Let A be a *linear operator*, unbounded, defined from  $D(A) \subset E$  into F; D(A) dense into E (see e.g. [11] for recalls of these concepts). Then its *adjoint operator*  $A^*$  satisfies:

$$\langle v, Au \rangle_{F',F} = \langle A^*v, u \rangle_{E',E} \quad \forall u \in D(A) \subset E, \ \forall v \in D(A^*) \subset F'$$

In the case E = V Hilbert space (hence a reflexive Banach space),  $A : V \to V'$ , <u>A linear</u>, we have  $A^* : V \to V'$ , linear, such that:

$$\langle v, Au \rangle_{V,V'} = \langle A^*v, u \rangle_{V',V} \quad \forall u \in V, \ \forall v \in V$$
(A.1)

**Exercice A.13.** Let V be a subspace of  $H^1(\Omega)$ , let A be the operator defined from V into V' as follows:

$$A(u) = -div(\lambda \nabla u) + cu$$

with  $\lambda$  and u given in  $L^{\infty}(\Omega)$ .

Give an expression of  $A^*$  (potentially in a weak form), in the following cases: 1)  $V = H_0^1(\Omega)$ . 2)  $V = H_{\Gamma_0}^1(\Omega) = \{v \in H^1(\Omega), v = 0 \text{ on } \Gamma_0\}$ , with  $\partial \Omega = \Gamma_0 \cup \Gamma_1$ .

**Exercice A.14.** Let V be a subspace of  $H^1(\Omega)$ , let A be the operator defined from V into V' as follows:

$$A(u) = -div(\lambda \nabla u) + w \cdot \nabla u$$

with  $\lambda$  given in  $L^{\infty}(\Omega)$ , and w a vector field given in  $(L^{\infty}(\Omega))^d$ .

Give an expression of  $A^*$  (potentially in a weak form), in the following cases: 1)  $V = H_0^1(\Omega)$ . 2)  $V = H_{\Gamma_0}^1(\Omega) = \{ v \in H^1(\Omega), v = 0 \text{ on } \Gamma_0 \}$ , with  $\partial \Omega = \Gamma_0 \cup \Gamma_1$ ,  $w \cdot n = 0 \text{ on } \Gamma_0$ ,  $w \cdot n < 0 \text{ on } \Gamma_1$ .

### A.2 Differentiable optimization

This section is a short summary of the optimization chapter of [1]; we refer to this book for the proofs and more details.

### A.2.1 Convexity

**Proposition A.15.** Let V be a Hilbert space, let j be a differentiable functional defined from V into  $\mathbb{R}$ . The following assertions are equivalent:

i) j is convex in V,

ii)

$$j(v) \ge j(u) + (j'(u), v - u) \quad \forall u, v \in V$$

iii)

$$(j'(u) - j'(v), u - v) \ge 0 \quad \forall u, v \in V$$

**Exercice A.16.** Prove that if the functional j is twice-differentiable then j is convex if and only if:

$$j"(u)(w,w) \ge 0$$

### A.2.2 Optimality conditions

Let us recall the first order necessary conditions of optimality.

Let j be a functional defined from V to  $\mathbb{R}$ . The optimization problem to be solved reads:

$$\inf_{u \in K} j(u)$$

with K a closed subset of V.

The optimality conditions are more simple in the case K convex.

### Case A: K convex.

Indeed, the proof of first order optimality conditions tests the optimality of a candidate u in the *admissible direction* (v - u) as follows: w = u + t(v - u),  $w \in K$  with  $t \in [0, 1]$ .

**Theorem A.17.** (Euler's inequality). Let j be differentiable, and let K be convex. i) If u is a local minimum of j in K, then:

$$(j'(u), v-u) \ge 0 \quad \forall v \in K$$

ii) Furthermore, if j convex, then u is a global minimum.

The Euler's inequality is a *necessary* condition of optimality. If in addition j convex, then it is a *necessary and sufficient* condition (result ii)).

**Theorem A.18.** (Euler's equality). Let j be differentiable. If K = V or if the local minimum u is in the interior of K, then the previous inequality reads:

$$(j'(u), v) = 0 \quad \forall v \in K$$

**Exercice A.19.** Let A be  $p \times n$ -matrix, and  $b \in \mathbb{R}^p$ . Prove that the following least-square problem: problem:

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|^2$$

has one and only solution; and give the corresponding Euler's equality condition.

Case B: equality constraints, K affine subspace.

We consider the set of constraints K defined by:

$$K = \{v \in V, F(v) = 0\}$$

with  $F(v) = (F_1(v), ..., F_m(v))$  a mapping defined from V into  $\mathbb{R}^m$ . Then, the necessary optimality conditions make introduce the Lagrange multipliers and read as follows.

**Theorem A.20.** (Lagrange multipliers). Let j be differentiable and let F be  $C^1$  in K. The vectors  $(F'_i(v))_{1 \le i \le m}$  are supposed to be linearly independent. Then, if u is a local minimum of j in K, it exists m reals  $\lambda_1, ..., \lambda_m$ , called Lagrange multipliers, such that:

$$j'(u) + \sum_{i=1}^{m} \lambda_i F'_i(u) = 0$$

### A.2.3 Lagrangian and saddle-point

Let us define the lagrangian  $\mathcal{L}$  defined from  $V \times \mathbb{R}^m$  into  $\mathbb{R}$  by:

$$\mathcal{L}(v,\mu) = j(v) + \sum_{i=1}^{m} \mu_i F_i(v) = j(v) + \mu F(v)$$

If u is a local minimum of j in K, Theorem A.20 states that it exists  $\lambda \in \mathbb{R}^m$  such that  $(u, \lambda)$  is a stationnary point of  $\mathcal{L}$ , i.e.

$$\partial_v \mathcal{L}(u,\lambda) = 0$$
 and  $\partial_\mu \mathcal{L}(u,\lambda) = 0$ 

Indeed, we have:

$$\partial_{\mu}\mathcal{L}(u,\lambda) = F(u) = 0 \text{ if } u \in K$$

and

$$\partial_v \mathcal{L}(u,\lambda) = j'(u) + \lambda F'(u) = 0$$

The pair  $(u, \lambda)$ , local solution of the minimization problem with constraint equality, can be interpreted a *saddle point* of the Lagrangian  $\mathcal{L}$ , i.e. a point  $(u, \lambda)$  which satisfies locally in  $K \times \mathbb{R}^m$ :

$$\forall \mu \quad \mathcal{L}(u,\mu) \leq \mathcal{L}(u,\lambda) \leq \mathcal{L}(v,\lambda) \quad \forall v$$

The duality theorem, see [1], implies that under the assumption of Theorem A.20 and if j strictly convex, then  $(u, \lambda)$  satisfies:

$$\min_{v} \max_{\mu} \mathcal{L}(v,\mu) = \max_{\mu} \min_{v} \mathcal{L}(v,\mu)$$

Figure A.1: A 2d saddle point (or a mountain pass) represented for a pair of scalar variables  $(v, \mu)$ 

For more details, we refer to [1].

Few extra crucial and classical exercices are proposed; all of them can be read in [1].

**Exercice A.21.** (Quadratic optimization with linear constraints) Let A be a matrix order n, symmetric positive define. Let B be a rectangular matrix  $m \times n$ , let  $b \in \mathbb{R}^n$ . We seek to solve the optimization problem:

$$\inf_{x \in Ker(B)} j(x) \text{ with } j(x) = \frac{1}{2}(Ax, x) - (b, x)$$

Prove that it exists a solution if A positive; and the solution is unique if A positive define.
 Prove that the unique solution x\* satisfies:

$$Ax^* - b = B^*p \text{ with } p \in \mathbb{R}^m$$

**Exercice A.22.** Let  $f \in L^2(\Omega)$ , we consider the functional  $j: V = H_0^1(\Omega) \to \mathbb{R}$  defined by:

$$j(v) = \frac{1}{2} \int_{\Omega} \|\nabla v\|^2 dx - \int_{\Omega} f v dx$$



This functional represents for example the energy of an elastic membrane under the external force f, or the energy of any purely diffusive phenomena (eg a temperature, molecule diffusion etc).

We seek to solve:

$$\inf_{u \in V} j(u)$$

In others words, we seek the solution of the system at equilibrium (minimal energy).

1) Prove that the Euler equation gives the variational formulation:

$$\int_{\Omega} \nabla u \nabla v dx = \int_{\Omega} f v dx \quad \forall v \in H_0^1(\Omega)$$

(We remark that H<sup>1</sup>(Ω) is the natural energy space for second order linear elliptic models).
2) Write the classical form of this equation.
3)

a) Prove that the energy functional is strictly convex in V.

b) After discretization (using FEM for example), what numerical algorithms would you use to solve this minimization problem ?

c) After discretization (using FEM), what algorithm would you use to solve the linear system corresponding to the variational formulation ?

### A.3 Descent algorithms and Newton-like methods

For more details we refer to any excellent course on-line, or [1, 9].

The variational data assimilation method leads to the numerical minimization of a functional  $j; j : V \to \mathbb{R}, V$  Hilbert space. Thus, let us consider the following *discrete optimization problem*:

Find  $x^* \in \mathbb{R}^m$  such that:  $j(x^*) = \min_{x \in \mathbb{R}^m} j(x)$ ,

where the (discrete) cost function  $j : \mathbb{R}^m \to \mathbb{R}$  is supposed to be  $C^1$  (i.e. continuously differentiable).

If m is large (classical case in data assimilation), an affordable approach is to solve this minimization problem by using a *descent algorithm*.

Then the minimization is local only since if j is not convex, it gives at best *a local minimum*, on contrary to global minimization methods such as genetic algorithms for example.

### A.3.1 Gradient type methods

A descent type algorithm is based on the *information of the gradient* of j; that is why these algorithms require j to be  $C^1$ .

All gradient type methods can be formulated as follows:

$$\begin{cases} x^0 & \text{given} \\ x^{k+1} = x^k + \alpha^k d^k & \text{for } k = 0, \dots \end{cases}$$

where  $d^k$  is the direction of descent,  $d^k \in \mathbb{R}^m$ , and  $\alpha^k$  is the step of descent,  $\alpha_k \in \mathbb{R}^+$ .

By definition,  $d^k$  is a descent direction if and only if it satisfies:  $\langle d^k, \nabla j(x^k) \rangle < 0$ . Thus, it exists a scalar value  $\alpha^k$  such that:  $j(x^{k+1}) < j(x^k)$ .

Exercice A.23. Prove this assertion.

Hint. Write a second order Taylor expansion:  $j(x^{k+1}) = j(x^k) + \alpha^k < d^k, \nabla j(x^k) > + (\alpha^k)^2 < H_j^k d^k, H_j^k d^k > + O(\alpha^k ||d^k||).$ 

Then, defining a descent algorithm consists to define the direction  $d^k$  and the descent step  $\alpha^k$  (linear search).

The steepest descent method. The most simple gradient method consists to "follow" the local steepest descent at each step i.e.:  $d^k = -\nabla j(x^k)$ .

This algorithm seems "natural", it is simple and it is interesting in view to illustrate descent algorithms in general. Nevertheless, it should never be used because absolutely inefficient.

It can be proved that if j is quadratic, the steepest descent method, associated with the socalled optimal step descent, converges. Nevertheless, at each iterate, one has:  $\langle d^k, d^{k+1} \rangle = 0$ i.e. a zig-zag behavior, see Fig. A.2.

The Conjugate-Gradient method. An improved version of the gradient method is the conjugategradient method. In the case j quadratic, it can be proved that CG algorithm converges with at most m iterates. Hence in theory, the CG algorithm is a direct method and not a iterative one. Nevertheless in practice this algorithm converges in only few iterations compared to m; that the reason why it is considered as an efficient and iterative method.

For a short but illustrative presentation for the CG algorithm, we refer to "An Introduction to the Conjugate Gradient Method Without the Agonizing pain, by Shewchuk (1994)" available on-line.

**Remark A.24.** If j is strictly convex then CG algorithm is The right method to employ (typically j defined from a symmetrical positive define matrix).

But if j is not convex (classical case in non-linear data assimilation problems), a quasi-Newton method should be prefered instead (see below).

Initial point. Let us point out that the choice of the initial point  $x^0$  is very important, especially if j is not convex. As a matter of fact, it will determine the local minimum obtained. In the 2d case (m = 2), you can mind to a small area in mountains; the algorithm is driving



Figure A.2: The steepest descent algorithm with optimal step at each iterate, go down following a zig-zag route ! Courtesy of .

you to the nearest lake or valley bottom. Is that point you are looking for ?

### A.3.2 The linear search (step descent)

Given  $d^k$ , computing the descent step is called the linear search (it is a 1d problem). An efficient method is based on the Wolf's conditions. The idea is to compute an 'acceptable' step length  $\alpha^k$ that reduces the objective function 'sufficiently' (compared to its minimization in the direction  $d^k$  given). The Wolf's conditions (also called Armijo's rule) writes a follows. Determine  $\alpha^k$ such that:

$$j(x^{k} + \alpha^{k}d^{k}) \leq j(x^{k}) + \omega_{1} \alpha^{k} \langle \nabla j(x^{k}), d^{k} \rangle$$
$$\langle \nabla j(x^{k} + \alpha^{k}d^{k}), d^{k} \rangle \geq \omega_{2} \langle \nabla j(x^{k}), d^{k} \rangle$$

where  $\omega_1$  and  $\omega_2$  are two scalar values such that:  $0 < \omega_1 < \frac{1}{2}$  and  $\omega_1 < \omega_2 < 1$ . Practical values can be for example:  $\omega_1 = 10^{-4}$ ,  $\omega_2 = 0.99$ .

### A.3.3 Newton method (2nd order)

Let us consider the first order necessary condition of minimum:  $\nabla j(x) = 0$ . Next, let us write the Newton's method (also called the Newton-Raphson method) to solve this non-linear system of equations in x. We set:  $G(x) = \nabla j(x)$ . We obtain the Newton's iterates:

$$x^{k+1} = x^k - [DG(x^k)]^{-1} \cdot G(x^k)$$

Thus, the Newton's method can be viewed as a descent method with:

$$d^{k} = -\left[\nabla^{2} j(x^{k})\right]^{-1} \cdot \nabla j(x^{k})$$

The Newton's method is very efficient if converging (since its quadratic convergence; recall j must be  $C^2$ ).

The (potentially huge) drawback is the computation of the inverse of the Hessian  $\left[\nabla^2 j(x^k)\right]^{-1}$  at reasonable cost...

For large dimensional problems (m large), the computation of the Hessian matrix is often not affordable. Then, two options are : Gauss-Newton method or Quasi-Newton method.

### A.3.4 Gauss-Newton method

A first option is to solve the necessary condition by a Gauss-Newton algorithm (potentially incomplete). The Gauss-Newton method, used to solve non-linear least squares problems, can be seen as a modification of the Newton method to find a minimum. Unlike Newton's method, the Gauss-Newton algorithm has the advantage that second derivatives are not required. At the end, the algorithm consists to approximate:  $\left[\nabla^2 j(x^k)\right]$  by  $\left[\nabla j(x^k)\right]^T$   $\left[\nabla j(x^k)\right]$ .

Its convergence is not guaranteed, not even locally like it is in the Newton method case. The rate of convergence of the Gauss-Newton algorithm may "approach order 2". As usual, its convergence will depend on the choice of the initial guess  $x^0$  and the matrix condition number.

### A.3.5 Quasi-Newton method

A second option (the most classical) is to use a Quasi-Newton method. The basic idea of a Quasi-Newton algorithm is to approximate the inverse of the Hessian. Briefly, the descent direction  $d^k$  is computed as follows:

$$d^{k} = -W^{k} \cdot \nabla j(x^{k}) \tag{A.2}$$

where  $W^k$  is a matrix symmetric positive define which *approximates* the inverse Hessian  $[\nabla^2 j(x^k)]^{-1}$ , and using the gradient information only.

Thus, these methods remain first order methods and the computation of the gradient only is required (instead of the Hessian for Newton's method).

In some Quasi-Newton methods (not the one described below),  $d^k$  is computed by solving the linear system:  $Z^k d^k = -\nabla j(x^k)$ , with  $Z^k$  which approximates the Hessian  $[\nabla^2 j(x^k)]$ .

The most classical and efficient Quasi-Newton method is the so-called BFGS (for the names of its authors Broyden-Fletcher-Goldfarb-Shanno). The version so-called 'L-BFGS' (for Limited memory) is particularly efficient for m large. The BFGS algorithms have been implemented in all good computational librairies and software.

In the BFGS method, the matrices  $W^k$  are updated step-by-step using the formula BFGS inverse which requires the gradient only.

The matrix  $W^k$  is not stored in memory, but only couples of vectors which allow to rebuilt it. For more details, we refer the reader to [9].

An excellent implementation of L-BFGS algorithm in Fortran is the routine M1QN3 by C. Lemarechal, J.C. Gilbert from INRIA, France (see website dedicated, free program available on line). It is used for example by Meteo-France to compute its forecasts (4 times a day). It has been implemented into Scilab too.

**Stopping criteria** Few criteria can (must ?) be apply to stop the minimization process. First, one can test the (normalized) norm of the gradient:

$$\frac{\left\|\nabla j(x^k)\right\|}{\left\|\nabla j(x^0)\right\|} \leq \varepsilon \; .$$

with  $\varepsilon \approx 10^{-5} - 10^{-7}$  typically.

Furthermore, since data assimilation can be employed only with a good understanding of the physical system modeled (it cannot be used reasonably as a black box...), a possible criteria is the stationarity either of the control variable and/or the cost function (both interpreted in a physical point of view). That is:

$$\frac{\|x^k - x^{k-1}\|}{\|x^k\|} \le \varepsilon_1 \text{ and } / \text{ or } \frac{|j(x^k) - j(x^{k-1})|}{|j(x^k)|} \le \varepsilon_2$$

with  $\varepsilon_i$  defined by the end-user (with its own understanding of the physical system).

In addition, of course a maximum of iterate must be defined, in the case no 'convergence' is reached in a reasonable number of iterates.

### A.4 Basic notations of statistics

Below very basic definitions (and notations) of probability and statistics are recalled.

Let  $\hat{X}$  be an aleatory (random) variable.

Its average, also called the expected value, expectation or mean, is defined by:

$$\mu \equiv E[\hat{X}] = \int_{a}^{b} x f(x) dx$$

It is the probability-weighted average over all samples, f(x) being the probability density function (pdf).

The covariance between two random variables  $\hat{X}$  and  $\hat{Y}$  is denoted by  $cov(\hat{X}, \hat{Y})$  or  $\sigma_{\hat{X}\hat{Y}}$ . It measures how the two random variables change together.

If the two variables  $\hat{X}$  and  $\hat{Y}$  are independent then  $cov(\hat{X}, \hat{Y}) = 0$ . The reciprocal is not true. Random variables whose covariance is zero are called uncorrelated. The definition is:

$$\sigma_{\hat{X}\hat{Y}} \equiv cov(\hat{X}, \hat{Y}) = E[(\hat{X} - E[\hat{X}])(\hat{Y} - E[\hat{Y}])]$$

We have the property:

$$cov(\hat{X}, \hat{Y}) = E[\hat{X}\hat{Y}] - E[\hat{X}]E[\hat{Y}]$$

The form cov(.,.) is bilinear symmetrical positive define. The associated quadratic form defines the variance. Its normalized version is the correlation coefficient  $\rho(\hat{X}, \hat{Y})$ :

$$\rho(\hat{X}, \hat{Y}) = \frac{\sigma_{\hat{X}\hat{Y}}}{\sigma_{\hat{X}}\sigma_{\hat{Y}}}$$

 $\rho$  measures the strength of the linear relation between  $\hat{X}$  and  $\hat{Y}$ . It is a dimensionless number. If  $\rho = 0$  the two variables are uncorrolated; if  $\rho = \pm 1$ , the two are linearly dependent.

The variance of a random variable or distribution, denoted by  $Var(\hat{X})$  (or  $\sigma_{\hat{X}}^2$ ), is a measure of how far the numbers lie from the average.

In other words, the variance is a measure of the amount of variation of values of the variable, taking account of all possible values and their probabilities or weightings (not just the extremes which give the range).

It is the expected value of the squared difference between the variable realisation and the variable mean:

$$\sigma_{\hat{X}}^2 \equiv Var(\hat{X}) = cov(\hat{X}, \hat{X})$$
  
=  $E[(\hat{X} - \mu)((\hat{X} - \mu)]$   
=  $E[(\hat{X} - \mu)^2]$ 

The standard deviation ("ecart-type" in French) is the square root of the variance. It is denoted by  $\sigma_{\hat{X}}$ . It has the same dimension as the random variable  $\hat{X}$ .

Equivalently, it is the square root of the average value of  $(\hat{X} - E(\hat{X})^2)$ .

Roughly, the standard deviation measures the variation (dispersion) from the average. A low standard deviation indicates that data tend to be very close to the mean; high standard deviation indicates that data are spread out over a large range of values, Fig. A.3.

In statistics, an estimator is a rule to calculate the estimation of a given quantity from observed data.



Figure A.3: A plot of a normal (Gaussian) distribution. Each band has a width of 1 standard deviation  $\sigma$ . Image from Wikipedia webpage.

The Mean Square Error (MSE) of an estimator quantifies the difference between values provided by the estimator and the true values.

The mean-square error measures how far the collection of values are in average from a estimated parameter x. It is defined by:

$$MSE(\hat{X}) = E((\hat{X} - x)^2)$$

We have:

$$MSE(\hat{X}) = Var(\hat{X}) + (Bias(\hat{X}, x))^2$$

Hence, if the random variable  $\hat{X}$  is unbiased, the variance equals the MSE.

### Appendix B

### Algorithmic - Automatic Differentiation \*

The outline of this chapter is as follows.

### Contents

| A.1 Diffe | erential calculus, functional analysis $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $155$                       |
|-----------|---|
| A.1.1     | Differentials   |
| A.1.2     | Green formulas  |
| A.1.3     | The Riesz-Fréchet theorem   |
| A.1.4     | The adjoint operator  |
| A.2 Diffe | $erentiable \ optimization \ \ldots \ $ |
| A.2.1     | Convexity   |
| A.2.2     | Optimality conditions   |
| A.2.3     | Lagrangian and saddle-point   |
| A.3 Desc  | cent algorithms and Newton-like methods   |
| A.3.1     | Gradient type methods   |
| A.3.2     | The linear search (step descent)  |
| A.3.3     | Newton method (2nd order)   |
| A.3.4     | Gauss-Newton method   |
| A.3.5     | Quasi-Newton method   |
| A.4 Basi  | c notations of statistics   |

At the present stage, the reader known why it can be interesting to consider an adjoint model, what for. In practice, the strategy to write the adjoint model may be different according to the forward code structure. One of the three possibilities to derive an adjoint code consists to derive the forward *source code*. In C or Fortran it may be possible as a "source-to-source"

transformation" process; in C++ by operators over-loading techniques.

In the present chapter, the forward source code is supposed to be written in Fortran (like many efficient computational codes).. In next section, the advantages and drawbacks of algorithmic-"automatic" differentiation is briefly presented. Next, we present the link between differential calculus (as employed in the course to derive the adjoint equation and the cost gradient) and a source code differentiation. The presentation is done in the framework of the automatic differentiation software Tapenade [23]. Moreover the conceptual way to re-employ your direct linear solver (eg based on your favourite linear algebra library) in the adjoint code; indeed the latter does not have to be differentiated. Finally, are presented some useful tricks for MPI instructions (see "how to adjointize the forward MPI instructions ?"), and how to save a bit of your memory. Indeed the memory will likely be your greatest issue (vs CPU time) while performing your "backward" adjoint code.

Finally let us remark thats the algorithmic differentiation described below to obtain the adjoint code (therefore the gradient) is similar to the so-called "backpropagation gradient" in the neural newtwork - deep learning community. More precisely, the "backpropagation gradient" is a particular case of the algorithmic differentiation (retropropagating by following the chain rule and saving in memory the intermediate states).

### B.1 What adjoint code ?

In practice, there exists three approaches to compute the adjoint state variable, then the gradient of the cost function.

- 1. We write the continuous adjoint equations, we discretize them using an appropriate numerical scheme (a-priori the same numerical method than the direct model), we discretize the expression of the (continuous) gradient. We obtain the so-called *discretized continuous gradient*
- 2. We discretize the direct model. We derive these discrete equations in order to obtain the adjoint discrete equations. It is calculations in finite dimension spaces, that is why the present approach is the most famous one in engineering. We obtain the so-called *discrete gradient*.
- 3. The *computational gradient* is obtained from the differentiation of the forward code directly !

The three approaches are possible. The choice should be done upon the human time development required. The three gradients obtained are not exactly the same since propagation of all different errors is not the same for each approach. Very few numerical analysis results exits on the topic (e.g. in case of a linear PDE discretized using conforming finite element, the first two approaches are the same up to the scheme errors).

For non-linear and very large scale problems, some end-users claim that the tiny differences between the three gradients can justify differences of local convergence behaviours observed...

The algorithmic differentiation approach (last approach) may present two advantages. First, it may ensure a better consistency between the computed cost function, which is the output of

the direct code, and its gradient since it is the computed cost function which is differentiated (consistency including all types of errors: schemes, rounding errors, iterative algorithms etc). Second, a large part of this extensive task can be automated using algorithmic differentiation, see e.g. [20], if the direct code has been designed to it !... In case of a direct code written in Fortran, and initially if designed to be differentiated by an algorithmic process, the adjoint code can be almost automatically derived using an automatic differentiation software. One of the most efficient automatic differentiation tool, source-to-source, is Tapenade, see [23, 22].

The computational software DassFlow presented previously [17] has been initially design to be differentiated by Tapenade; that is why one can obtain very fast the adjoint code up-to-date. As a matter of fact, let us point out that as soon as the direct model is modified (extra terms, numerical approximation changed etc), the corresponding adjoint code must be modified in consequence, whatever the approach chosen.

## B.2 From mathematical differentiation to source code differentiation

We describe how to define the direct code, then what is the response of the adjoint code automatically generated and finally how to use it. We follow the presentation done in [24], see [34] too.

Let  $\mathcal{K}$  be the space of control variables and  $\mathcal{Y}$  the space of the forward code response. In the case of DassFlow, we have :

$$\mathbf{k} = (y_0, q_{in}, z_{out}, n, z_b)^T$$
 and  $Y = (y, j)^T$ 

Let us point out that we include both the state and the cost function in the response of the forward code.

We can represent the direct code as the operator  $\mathcal{M} : \mathcal{K} \longrightarrow \mathcal{Y}$ , see figure B.1.



Figure B.1: Representation of the direct model.

The tangent model becomes  $\frac{\partial \mathcal{M}}{\partial \mathbf{k}}(\mathbf{k}) : \mathcal{K} \longrightarrow \mathcal{Y}$ . It takes as input variable a perturbation of the control vector  $d\mathbf{k} \in \mathcal{K}$ , then it gives the variation  $dY \in \mathcal{Y}$  as output variable, see figure B.2:

$$dY = \frac{\partial \mathcal{M}}{\partial \mathbf{k}}(\mathbf{k}) \cdot d\mathbf{k}$$



Figure B.2: Representation of the tangent model.

The adjoint model is defined as the adjoint operator of the tangent model. This can be represented as follows:  $\left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}(\mathbf{k})\right)^* : \mathcal{Y}' \longrightarrow \mathcal{K}'$ . It takes  $dY^* \in \mathcal{Y}'$  an input variable and provides the adjoint variable  $d\mathbf{k}^* \in \mathcal{K}'$  at output, see figure B.3:

$$d\mathbf{k}^* = \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}(\mathbf{k})\right)^* \cdot dY^*$$



Figure B.3: Representation of the adjoint model.

Now, let us make the link between the adjoint code and the "computational gradient".

By definition of the adjoint, we have :

$$\left\langle \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right)^* \cdot dY^*, \, d\mathbf{k} \right\rangle_{\mathcal{K}' \times \mathcal{K}} = \left\langle dY^*, \, \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right) \cdot d\mathbf{k} \right\rangle_{\mathcal{Y}' \times \mathcal{Y}}$$
(B.1)

or, using the relations presented above:

$$\langle d\mathbf{k}^*, d\mathbf{k} \rangle_{\mathcal{K}' \times \mathcal{K}} = \langle dY^*, dY \rangle_{\mathcal{Y}' \times \mathcal{Y}}.$$
 (B.2)

If we set  $dY^* = (0, 1)^T$  and by denoting the perturbation vector  $d\mathbf{k} = (\delta y_0, \delta n, \delta z_b, \delta q^{in})^T$ , we obtain:

$$\left\langle \left(\begin{array}{c} 0\\1\end{array}\right), \left(\begin{array}{c} dy^*\\ dJ^*\end{array}\right) \right\rangle_{\mathcal{Y}'\times\mathcal{Y}} = \left\langle \left(\begin{array}{c} \delta y_0\\ \delta n^*\\ \delta z_b^*\\ \delta q^{in*}\\ \delta z_{out}^*\end{array}\right), \left(\begin{array}{c} \delta y_0\\ \delta n\\ \delta z_b\\ \delta q^{in}\\ \delta z_{out}\end{array}\right) \right\rangle_{\mathcal{K}'\times\mathcal{K}}$$
Moreover, we have by definition:

$$dj = \frac{\partial J}{\partial y_0}(\mathbf{k}) \cdot \delta y_0 + \frac{\partial J}{\partial n}(\mathbf{k}) \cdot \delta n + \frac{\partial J}{\partial z_b}(\mathbf{k}) \cdot \delta z_b + \frac{\partial J}{\partial q_{in}}(\mathbf{k}) \cdot \delta q_{in} + \frac{\partial J}{\partial z_{out}}(\mathbf{k}) \cdot \delta z_{out}$$

Therefore, the adjoint variable  $d\mathbf{k}^*$  (output of the adjoint code with  $dY^* = (0, 1)^T$ ) corresponds to the partial derivatives of the cost function j:

In summary, in order to compute the "computational gradient" (partial derivatives of the cost function j using differentiation of the forward code), first, one runs the direct code then one runs the adjoint code with  $dY^* = (0, 1)^T$  as input.

## **B.3** Automatic differentiation in short

Automatic differentiation softwares Automatic Differentiation (AD) (also called *algorithmic differentiation*) aims at computing the derivative of an program output. one the possible approach is the source-to-source one; it means the AD software (kind of a "superover-compiler"...) transforms the direct source code into a tangent linear source code and/or an adjoint source code.

Researches in this field are very active; nowadays few reliable AD tools are available, both for Fortran and C languages. Let us cite only OpenAD and Tapenade. There use is still quite technical and tricky. It requires a good knowledge of what is the derivative of the output of the model, next how it is translated in terms of a program instructions.

In Fortran and C languages, the only remaining limitation is the pointer management; nevertheless some easy rules allow to circumvent this limitation.

We present below the basic principles of a AD tool, and we will focus on the Tapenade. Tapenade [23, 22] is an automatic differentiation tool for Fortran and C programs. It is developed by the Tropics team at INRIA Sophia-Antipolis. Tapenade works using source code transformation : it builds the tangent and/or adjoint code automatically from the direct source code. We refer to its rich webpage and the FAQ section.

Concerning AD in general, we refer also to [20, 35].

**Some basic rules to respect for the direct code.** In order to be able to differentiate a Fortran source code using Tapenade, one must respect some basic rules of code structure. Typically, the code structure must be as follows:

. Preprocessing.

Reading mesh, parameters, data etc AND allocate all arrays.

. Computational kernel.

At the second level of the program tree, one gives the head routine to be differentiated to Tapenade. From this head routine, neither allocate nor pointer may appear.

. Post-processing.

Writing in files, visualization etc.



## **Cost function:** It measures the discrepancy between the simulated state and the observations

It is a function of the control variable k (Manning, topo., B.C., I.C.):

$$j(k) = \frac{1}{2} \int_0^T \|C y(k;t) - y^{obs}(t)\|_*^2 dt + \text{ regularization terms}$$

where C: operator of observations

 $\|.\|_*$ : norm in the observation space.

Optimization problem:

 $\begin{cases} Find k^* \text{ such that} \\ j(k^*) = \min_{K} j(k) \end{cases}$ 

where  $j(k) = J(k, y^k)$  and  $y^k$ 

is the solution of the model (state of the system)

## The state equation (forward model)

The control variable:  $k = (y_0, v(t))^T$ 

 $y_0$ : I.C.; v(t): parameters

State equation (nonlinear): Given k, find y(t) s.t.

$$\begin{cases} \partial_t y(t) + A(v(t); y(t)) = 0 \quad \forall t \in ]0, T[\\ y(0) = y_0 \end{cases}$$

This gives  $y^k(t)$ : the state,  $j(k) = J(k; y^k)$ : the cost function

## Since the cost computation is time-consuming: local minimization algorithm (quasi-Newton e.g. BFGS) Since many variables of control: introduction of the adjoint model to compute the gradient

Ref. JL Lions '68

Adjoint State equation Find  $\tilde{y}(t)$  s.t.

$$\begin{array}{ll} \partial_t \, \widetilde{y}(t) - \left[ \frac{\partial A}{\partial y} \big( y(t), v(t) \big) \right]^* \widetilde{y}(t) = C^* \Lambda_{\mathcal{O}} \left( C \, y(t) - y^{obs}(t) \right) & \forall t \in \left] 0, T \right[ \\ \widetilde{y}(T) &= 0 \end{array}$$

Note. a. Reverse time b.Contains the observations

Gradient obtained:

$$\frac{\partial j}{\partial y_0}(\mathbf{k}) = -\widetilde{y}(0) + N(y_0 - y_b)$$
$$\frac{\partial j}{\partial v}(\mathbf{k}) = \left[\frac{\partial A}{\partial v}\left(y(\mathbf{k}), v\right)\right]^* \widetilde{y}(\mathbf{k})$$

## The Optimality System is

State equation Ajoint State equation (contains observations) Ref. e Le Di

$$\nabla j(k) = 0$$

Ref. e.g. [Le Dimet - Talagrand '86] [Talagrand-Courtier'87]

Variationnal Data Assimilation (4D-var): optimal control process





#### **Remarks**

- The gradient gives local information.
   One run of the forward+adjoint models (without optimization) gives a local sensitivity information
- One can control the boundary conditions similarly

## What Adjoint Code? Differentiate / Discretize / Implement : in what order !?

## Continuous gradient: 1. Differenciate mathematically; 2. discretize; 3. implement

- Discrete gradient:
  - 1. discretize; 2. differenciate; 3. implement

#### Make figure

## • Code differenciation:

1. Implement; 2. Differentiate the code (Automatic Differentiation tool e.g. Tapenade)

If the forward code well suited, prefer Automatic Differentiation (more re-usable). NB. A.D. produces the gradient of the computed cost function



## From Equations to Automatic Differentiation (continued)

Definition of the adjoint model: For all  $\delta k \in \mathcal{K}, \tilde{Y} \in \mathcal{Y}$ ,

$$\left\langle \widetilde{Y}, \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right) \cdot \delta \mathbf{k} \right\rangle_{\mathcal{Y}} = \left\langle \left(\frac{\partial \mathcal{M}}{\partial \mathbf{k}}\right)^* \cdot \widetilde{Y}, \delta \mathbf{k} \right\rangle_{\mathcal{K}}$$

Thus,

$$\langle \tilde{Y}, dY \rangle_{\mathcal{Y}} = \langle \tilde{\mathbf{k}}, \delta \mathbf{k} \rangle_{\mathcal{K}}$$

If we set  $\widetilde{Y} = (0,1)$  we get

$$\left\langle \left(\begin{array}{c} 0\\1\end{array}\right), \left(\begin{array}{c} \widehat{dy}\\\widehat{dj}\end{array}\right) \right\rangle_{\mathcal{Y}} = \left\langle \left(\begin{array}{c} \widetilde{y_0}\\\widetilde{v}\end{array}\right), \left(\begin{array}{c} \delta y_0\\\delta v\end{array}\right) \right\rangle_{\mathcal{X}}$$

This means:

$$\widehat{dj}(\mathbf{k}, \delta \mathbf{k}) = \frac{\partial j}{\partial y_0}(\mathbf{k}) \cdot \delta y_0 + \frac{\partial j}{\partial v}(\mathbf{k}) \cdot \delta v = \langle \widetilde{y_0}, \delta y_0 \rangle_H + \langle \widetilde{v}, \delta v \rangle_U$$

Therefore the response of the adjoint code is the gradient

$$\frac{\partial j}{\partial y_0}(\mathbf{k}) = \widetilde{y_0} \qquad \qquad \frac{\partial j}{\partial v}(\mathbf{k}) = \widetilde{v}$$

## Principles of « Automatic Differentiation » Tapenade: source-to-source transformation

The « A.D. » (Algorithm Differentiation in fact) consists to compute derivatives of a given function M. The derivatives are exact modulo rounding errors. Ref. Hascoet et al. Tapenade tool INRIA Tropics

Given the control variable  $(\mathbf{k}_{\#} \subset X_0)$ , we have the forward code :

$$Y_{\#} = \mathcal{M}_{\#}(\mathbf{k}_{\#}) = m_p(X_{p-1}) \circ m_{p-1}(X_{p-2}) \circ \cdots \circ m_1(X_0)$$

Then in a direction  $\delta \mathbf{k}_{\#}$ 

$$\frac{\partial \mathcal{M}_{\#}}{\partial \mathbf{k}_{\#}}(\mathbf{k}_{\#}) \cdot \delta \mathbf{k}_{\#} = m'_{p}(X_{p-1}) \times m'_{p-1}(X_{p-2}) \times \cdots \times m'_{1}(X_{0}) \cdot \delta \mathbf{k}_{\#}$$

where m'\_p is the jacobian of the elementary operation m\_p.

The tangent code value is computed from right to left (matrix-vector products).

This can be inserted into the original source code.

This is the « forward mode » or « tangent mode » of the A.D. tool

## Principles of A.D. (continued)

Given the adjoint input variable  $\widetilde{Y}_{\#}$  , we have the adjoint code :

$$\left(\frac{\partial \mathcal{M}_{\#}}{\partial \mathbf{k}_{\#}}(\mathbf{k}_{\#})\right)^{T} \times \widetilde{Y}_{\#} = m_{1}^{\prime T}(X_{0}) \times m_{2}^{\prime T}(X_{1}) \times \cdots \times m_{p}^{\prime T}(X_{p-1}) \cdot \widetilde{Y}_{\#}$$

This can be computed from right to left (matrix-vector products). Nevertheless.

- the code is reverse way
- all the variables have to be either stored or all recomputed
- $\rightarrow$  A balance between storing & recomputing is done (« checkpointing »)

This cannot be inserted into the original source code !

This is the « adjoint mode » or « reverse mode »of the A.D. tool

₽3

Source-to-source transformation: an extra source code is (automatically) generated



Tapenade On-line Differentiation engine. Current version: all Fortran 95 excepted pointers & dynamic allocates

## An example. The « tangent » mode

Instruction élémentaire du code initial (code direct):

$$a_i = x \, b_i + \cos(a_i).$$

Ref. Hascoet et al. Tapenade tool **INRIA** Tropics

Ce qui peut s'écrire sous la forme:

$$m_k : \mathbb{R}^3 \longrightarrow \mathbb{R}^3$$
$$X_{k-1} = \begin{pmatrix} a(i) \\ b(j) \\ x \end{pmatrix} \longmapsto X_k = \begin{pmatrix} a(i) \\ b(j) \\ x \end{pmatrix}$$

On pose:

$$\delta X_{k} = m'_{k}(X_{k-1}) \cdot \delta X_{k-1} \qquad m'_{k}(X_{k-1}) = \begin{pmatrix} -\operatorname{SIN}(a(i)) & x & b(j) \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

alors: 
$$\begin{pmatrix} ad(i) \\ bd(j) \\ xd \end{pmatrix} = \begin{pmatrix} -SIN(a(i)) & x & b(j) \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} ad(i) \\ bd(j) \\ xd \end{pmatrix}$$

code tangent : ad(i) = -ad(i)\*SIN(a(i)) + x\*bd(j) + xd\*b(j)

Tangent code:

•Additional derivative instructions are inserted just before each instruction of the original code. This leads to a source code quite easy to read.

Differentiated variables retain the type and dimension of the original variable

•The code can be used either to compute the derivative in ONE direction given,

or to compute for ex. the jacobian of the flux (Newton algorithm)

An example. The « adjoint » mode

 $a_i = x \, b_j + \cos(a_i).$ 

On a:

$$\widetilde{X}_{k-1} = m'_k^T(X_k) \cdot \widetilde{X}_k$$

D'où

$$\begin{pmatrix} ab(i) \\ bb(j) \\ xb \end{pmatrix} = \begin{pmatrix} -SIN(a(i)) & 0 & 0 \\ x & 1 & 0 \\ b(j) & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} ab(i) \\ bb(j) \\ xb \end{pmatrix}$$

| code adjoint : | xb    | = | xb + b(j)*ab(i)  |
|----------------|-------|---|------------------|
|                | bb(j) | = | bb(j) + x*ab(i)  |
|                | ab(i) | = | -SIN(a(i))*ab(i) |

Instructions are computed in reverse order. Adjoint code: much more complicated to generate and the source code is extremely difficult to read. CPU time ~2-3 times higher

Principles of D.A. in C++: operator overloading. Tool of ref.: Adol-C. Nevertheless, very hard in « adjoint» mode...

## If and For statements

This includes functions like min, |.| etc

Direct mode A; if (bool) then B else C; D; 1st branch: A; B; D . This gives: A'; A; B'; B; D'; D 2<sup>nd</sup> branch: A; C; D. This gives: A'; A; C'; C; D'; D Hence: A'; A; if (bool) then {B'; B} else {C'; C} end if; D'; D;

Similarly A; for i=1:N, do B(i); D;

This gives: A'; A; for i=1:N, do {B'(i); B(i)}; D'; D;

#### Differentiable or not?

If the var. bool and/or the bounds of the for statement depend on the control var: not differentiable ! In practice, programs are piecewise differentiable; and non-differentiable points (rare !?) should be identified by users ...

In adjoint mode:

- branches of IF statements are treated separately,
  - then the results are assembled after the derivation.
- we have to follow the graph of the program in reverse order
  - => mind to have an equilibrated original call graph !
- all « intermediate solutions » (for all time step) are a-priori stored
   > very large memory required

## => Snapshot and checkpointing strategy

## Snapshot and checkpointing in Tapenade

Store All (SA) strategy. Structure of the adjoint code:

- 1. A forward sweep, identical to the original program augmented with instructions that memorize the intermediate values before they are overwritten, and instructions that memorize the control flow in order to reproduce the control flow in reverse during the **backward sweep**.
- 2. **Start the backward sweep**, that actually computes the derivatives, and uses the memorized values both to compute the derivatives of each instruction and to choose its flow of control.



The isolated push are memorizations of intermediate values and control.

The group of 3 push before the subroutine call is a snapshot because the call will be checkpointed.

## Checkpointing and snapshot in Tapenade

Differentiated var retain the type and dimension of the original variables

#### **Forward Sweep**

... push(v2) v2 = 2\*sin(v1) + 5 if (v2>0.0) then push(v4) ; v4 = v2 + p1\*v3/v2 push(p1) ; p1 = p1 + 0.5 push(true) else push(false) endif push(v2) ; push(v4) ; push(p1) call sub1(v2,v4,p1) ...

#### **Backward Sweep**

```
...

pop(p1) ; pop(v4) ; pop(v2)

call sub1_b(v2,v2b,v4,v4b,p1)

pop(control)

if (control) then

pop(p1) pop(v4)

v3b = v3b + v4b*p1/v2

v2b = v2b + v4b*(1-p1*v3/(v2*v2))

v4b = 0.0

endif

pop(v2)

v1b = v1b + v2b*2*cos(v1)

v2b = 0.0 ...
```

#### Remarks

- The adjoint of one original instruction is very often a sequence of instructions.
- sub1\_b is the result of checkpointing on sub1,
- and therefore itself consists of a forward sweep followed by a backward sweep.
- Intrinsic functions are directly differentiated : no call to a sin\_b procedure but to cos(v1).

## **B.4** Exercices

Let us consider a computational Fortran code solving a very simple dynamic system: the scalar linear ODE order one: y(t) = f(y(t)u(t)) with f(y(t), u(t) = ay(t) + u(t), with initial condition. y(t) is the state of the system, while u(t) is the control.

The goal here is to differentiate first by hand the source code (algorithmic differentiation), next using Tapenade software (automatic differentiation).

The direct source code is the following.

```
! Basic program designed to be differentiated using Tapenade
! The head of the tree to be differentiated by Tapenade is
! the first routine of computational part
program dynamic_system
 implicit none
 integer i
 real*8 :: dt, T0, TN, y0
! all allocate must be in the main program
! (out of the computational part which will be diffrentiated)
 real*8, dimension(:), allocatable :: u, y
 integer :: nstep
 T0 = 0.d0
 TN = 1.d0
  !time step (should be defined by the stability condition)
 dt = 0.001d0
 nstep = int((TN-T0)/dt)
 !allocate instructions must be out of the part differentiated by Tapenade
 allocate(y(nstep))
 allocate(u(nstep))
  !initial state and initial control
  !here they are constant (u could be a time dependent vector)
 y0 = 3.4d0
 u = 2.6d0
  !head routine of the computational part
  !=>the top tree routine of the code thar will be differentiated by Tapenade
 call solve_model( y0, u, y, dt, nstep)
end program dynamic_system
!The direct model is:
!y'(t) = f(y(t)u(t)) + I.C.
!Scheme: Euler explicit
subroutine solve_model( y0, u, y, dt, nstep )
```

```
implicit none
 integer, intent(in) :: nstep
           intent(in) :: dt, y0
 real*8,
 real*8,
          dimension(nstep), intent(in) :: u
          dimension(nstep), intent(out) :: y
 real*8,
 real*8 :: f
 integer :: i
  ! initialize the state variable
 y(1) = y0
  ! time loop
 do i = 1, nstep -1
    y(i+1) = y(i) + dt * f(y(i),u(i))
 end do
 return
end subroutine solve_model
! Right-hand side f
! f(y(t),u(t)) = a y(t) + u(t)
real*8 function f(y,v)
 implicit none
 real*8, intent(in) :: y, v
 real*8, parameter :: a = 1
 f = a * y + v
end function f
```

**Exercice B.1.** Algorithmic differentiation by hand.

1) Write the tangent linear codes of the subroutine solve\_model and function f.

2) Write the corresponding adjoint codes.

**Exercice B.2.** Algorithmic differentiation using Tapenade.

1) Add in the direct code the computation of a cost function j.

2) Use Tapenade software (web interface on-line) in order to obtain:

a) the tangent linear codes of the subroutine solve\_model and function f.

b) the corresponding adjoint codes.

3) Read carefully and analyze the two routines automatically generated.

## B.5 Recycle your linear solver

Let us recall that, roughly, the adjoint of a linear system is the same linear system but transposed. (Notice that the latter feature brings some quite heavy difficulties...). Thus, one should not derive the linear solver instructions one-by-one using an automatic differentiation tool (source-to-source) but instead recycle it ! Here, we describe how to generate the adjoint of a routine containing a call to a linear solver.

Next, we present a small Fortran program including the call to a linear solver; we apply the method on the codes automatically generated by Tapenade.

#### B.5.1 How to adjoint-ize your linear solver ?

#### The direct routine

Let us consider two input parameters c and d (hence *active variables* in Tapenade terminology). Then, we have A matrix and b vector defined as follows:

$$\begin{pmatrix} A \\ b \end{pmatrix} = f(c,d) = \begin{pmatrix} f_1(c,d) \\ f_2(c,d) \end{pmatrix}$$

Let x be the vector solution of the linear system: Ax = b. Next, we compute the cost function j.

We present in Fig. B.4 the direct routines dependencies.



Figure B.4: Direct routines representation

#### The linear tangent routine generated

The linear tangent code derived by Tapenade is defined as follows:

- $\dot{c}$  and d are the corresponding tangent variables of the parameters, they are input parameters.
- A and b satisfy:

$$\left(\begin{array}{c}\dot{A}\\\dot{b}\end{array}\right) = df(c,d) \cdot \left(\begin{array}{c}\dot{c}\\\dot{d}\end{array}\right)$$

Next, we have to treat the linear solver. If we differentiate the linear system:

$$Ax = b \tag{B.3}$$

we find:  $A\dot{x} + \dot{A}x = \dot{b}$  or, after rewrite:

$$A\dot{x} = \dot{b} - \dot{A}x \tag{B.4}$$

At this stage, we already know A and x by running direct routine. The variables  $\dot{A}$ ,  $\dot{b}$  are output variables of the Tapenade tangent linear routines. Next, we call the linear solver in order to obtain:  $\dot{x}$ . The latter represents the derivative value of x at point (c, d) in the direction  $(\dot{c}, \dot{d})$  given.

Thus, we are able to compute  $\dot{x}$  without deriving the linear solver instructions. We can call directly the same linear solver (we recycle it).

Next, we obtain  $\dot{j}$  the gradient value at point (c, d) in the direction  $(\dot{c}, \dot{d})$  given (it is a scalar value like j).

We present in Fig. B.5 the scheme representing the linear tangent routines.



Figure B.5: Linear tangent routines representation

#### The adjoint routines generated

Let us recall that the adjoint code is deduced from the tangent linear code, in the reverse order. The output variable of the adjoint routine are  $\bar{c}$  and  $\bar{d}$ . In a Tapenade point of view, they are the adjoint variables of (c, d) (same type); while in a mathematical point of view, they are the *gradient* with respect to c and d.

The input variable of the adjoint cost computation routine is  $\overline{j}$ . Let us recall that  $\overline{j}$  must be set to 1, see Section B.2.

The input variable of the adjoint linear system routine is  $\bar{x}$ .

The adjoint code can be splitted in three steps, see Fig. B.5:

- 1) From  $\overline{j}$  to  $\overline{x}$  (let us assume it is a separated routine: cost computation, adjoint)
- 2) From  $\bar{x}$ , find  $\bar{A}$  and  $\bar{b}$
- **3)** From  $\overline{A}$  and  $\overline{b}$ , find  $\overline{c}$  and  $\overline{d}$ .

The first and third steps can be obtained directly by running the Tapenade adjoint code. The adjoint of the third step writes:

$$\begin{pmatrix} \bar{c} \\ \bar{d} \end{pmatrix} = df^*(c,d) \cdot \begin{pmatrix} \bar{A} \\ \bar{b} \end{pmatrix}$$
(B.5)

#### The adjoint of the linear system

The second step involves a call to the linear solver. Thus, let us detail how to compute  $\overline{A}$  and  $\overline{b}$  with the use of a linear solver as a black box. In other words, we detail below the adjoint of Step 2) only.

Input variable is  $\bar{x}$ , while output variables are  $\bar{A}$ ,  $\bar{b}$ .

In the linear tangent code, we have :  $A\dot{x} = \dot{b} - \dot{A}x$  or equivalently:

$$\begin{array}{rcl} 2a) & \dot{b'} &=& \dot{b} - \dot{A}x \\ 2b) & A\dot{x} &=& \dot{b'} \end{array}$$

The adjoint is reverse. Thus let us consider first the instruction 2b)  $A\dot{x} = \dot{b}'$  only. It can be written as follows (see e.g. [22]):

$$\begin{pmatrix} \dot{x} \\ \dot{b}' \end{pmatrix} = \begin{pmatrix} 0 & A^{-1} \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} \dot{x} \\ \dot{b}' \end{pmatrix}$$
(B.6)

#### Two fundamental remarks.

a) Let us point out that by convention, at left-hand sides are the variables at output, while at right-hand sides are the variables at input (of the routine).

b) Let us point out that input variables in the linear tangent routine become output variables in the adjoint routine. Furthermore, by convention, output variables are initialized at 0.

Let us write the adjoint of the instruction 2b). Since  $\dot{b}$  is the input variable of the instruction, its adjoint  $\bar{b}$  will be the output one (and set at 0 when entering into the routine). Similarly, since  $\dot{x}$  is the output variable, its adjoint  $\bar{x}$  will be the input one. The adjoint instruction of (B.6) writes:

$$\left(\begin{array}{c} \bar{x}\\ \bar{b}' \end{array}\right) = \left(\begin{array}{c} 0 & 0\\ A^{-T} & 1 \end{array}\right) \times \left(\begin{array}{c} \bar{x}\\ \bar{b}' \end{array}\right)$$

Hence:

$$\begin{cases} \bar{x} = 0\\ \bar{b}' = A^{-T}\bar{x} + \bar{b}' \end{cases}$$

Thus the system rewrites:

$$\begin{cases} A^T \bar{b}' = \bar{x} \\ \bar{x} = 0 \end{cases}$$

And one has to solve the linear system:  $A^T \bar{b}' = \bar{x}$ . It can be done using the linear solver called as a black box (we recycle it !). It gives  $\bar{b}'$ .

Now, let us consider the adjoint of the instruction 2a):

$$2a) \quad \dot{b'} = \dot{b} - \dot{A}x \tag{B.7}$$

This linear instruction writes as follows:

$$(\dot{b}', \dot{b}, \dot{A}) = (\dot{b}', \dot{b}, \dot{A}) \times \begin{pmatrix} 0 & 0 & 0 \\ 1 & 1 & 0 \\ -x & 0 & 1 \end{pmatrix}$$

The corresponding adjoint instruction writes:

$$(\bar{b}', \bar{b}, \bar{A}) = (\bar{b}', \bar{b}, \bar{A}) \times \begin{pmatrix} 0 & 1 & -x^T \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Therefore, the adjoint instruction of (B.7) writes:

$$\begin{cases} \bar{b}' = 0\\ \bar{b} = \bar{b}' + \bar{b}\\ \bar{A} = -\bar{b}'x^T + \bar{A} \end{cases}$$

The variables  $\bar{A}$ ,  $\bar{b}$  are output variables, hence set at 0 when entering into the adjoint routine.

Therefore, in summary the adjoint of the tangent linear instructions  $A\dot{x} = \dot{b} - \dot{A}x$  (ie Step 2)) can be re-write:

$$\begin{cases} A^T b = \bar{x} \\ \bar{A} = -\bar{b}x^T \\ \bar{x} = 0 \end{cases}$$

The first instruction can be solved using the same linear solver than the direct routine ("recycle your linear solver"). The second instruction writes too:  $\bar{A}_{ij} = -\bar{b}_i x_j$ .

#### Remark B.3.

- Let us notice that the matrix  $\overline{A}$  is the same type of A with the same sparse profile (even  $if -\overline{b}x^T$  is a-priori a full matrix). As a matter of fact, one needs only the adjoint values of coefficients  $A_{i,j}$  (and one do not need the others coefficients).

- Let us recall that  $(\bar{c}, d)$  are the components of the gradient with respect to (c, d) respectively; they are obtained by (B.5).

- Since the direct model is Ax = b, the state of the system is x, the adjoint state is  $\overline{b}$  (the solution of the transposed linear system) and  $\overline{x}$  is the right-hand side which contains the observations and misfit terms (it equals to  $(x - x^{obs})$  if the observation operator equals identity).

We summarize all steps in Fig. B.6.



Figure B.6: Adjoint routine representation

#### B.5.2 A simple example

We illustrate the previous procedure of the linear solver recycling on a very simple example. In all the following, one can replace the routine *sys\_solver* by any linear system solver routine. This example has been written by R. Madec, research engineer, IMT Toulouse.

#### Direct code

We present the Fortran source code of the direct subroutine written in file *simple\_operation.f90*:

```
subroutine simple_operation(c,d,x)
  implicit none
 integer, parameter :: n=2
 double precision, intent(in) :: c
 double precision, intent(in) :: d
 double precision, dimension (n), intent(out):: x
 double precision, dimension (n,n) :: A
 double precision, dimension (n)
                                   :: b
!definition of matrix A
 A(1,1) = 3*c + 4*d
 A(1,2) = c + 18*d
 A(2,1) = 2.5*d
 A(2,2) = 3.4*c
!definition of vector b
 b(1) = 8*c*d
 b(2) = 2*d
!solver Ax=b
 call sys_solver(A,b,x,n)
```

## end subroutine simple\_operation

#### Linear tangent code

We generate the linear tangent code using Tapenade. To do so, we type in a terminal (or in a Makefile):

```
tapenade -forward -head simple_operation \
    -vars "c d" -outvars "x" \
```

```
-html -O $(PWD)/forward simple_operation.f90
```

Then, the file *simple\_operation\_d.f90* is created in the *forward* directory. This file is the following:

```
!
         Generated by TAPENADE
                                   (INRIA, Tropics team)
  Tapenade 3.4 (r3375) - 10 Feb 2010 15:08
!
ļ
!
 Differentiation of simple_operation in forward (tangent) mode:
                 of useful results: x
!
   variations
Į.
   with respect to varying inputs: c d
ļ
   RW status of diff variables: x:out c:in d:in
SUBROUTINE SIMPLE_OPERATION_D(c, cd, d, dd, x, xd)
 IMPLICIT NONE
 INTEGER, PARAMETER :: n=2
 DOUBLE PRECISION, INTENT(IN) :: c
 DOUBLE PRECISION, INTENT(IN) :: cd
 DOUBLE PRECISION, INTENT(IN) :: d
 DOUBLE PRECISION, INTENT(IN) :: dd
 DOUBLE PRECISION, DIMENSION(n), INTENT(OUT) :: x
 DOUBLE PRECISION, DIMENSION(n), INTENT(OUT) :: xd
 DOUBLE PRECISION, DIMENSION(n) :: b
 DOUBLE PRECISION, DIMENSION(n) :: bd
 DOUBLE PRECISION, DIMENSION(n, n) :: a
 DOUBLE PRECISION, DIMENSION(n, n) :: ad
 INTEGER :: i,j
 ad(1, 1) = 3*cd + 4*dd
 a(1, 1) = 3*c + 4*d
 ad(1, 2) = cd + 18*dd
 a(1, 2) = c + 18*d
 ad(2, 1) = 2.5*dd
 a(2, 1) = 2.5*d
 ad(2, 2) = 3.4*cd
 a(2, 2) = 3.4*c
 bd(1) = 8*(cd*d+c*dd)
 b(1) = 8*c*d
 bd(2) = 2*dd
 b(2) = 2*d
 CALL SYS_SOLVER_D(a, ad, b, bd, x, xd, n)
```

```
END SUBROUTINE SIMPLE_OPERATION_D
```

Using the method presented in subsection B.5.1, we replace the unknown linear tangent of the linear system solver:

```
CALL SYS_SOLVER_D(a, ad, b, bd, x, xd, n)
by
CALL SYS_SOLVER(a, b, x, n)
!matrix vector product : bd = bd - Ad x
do j=1,n
    do i=1,n
    bd(i)=bd(i) - ad(i,j)*x(j)
    end do
end do
CALL SYS_SOLVER(a, bd, xd, n)
```

#### Adjoint code

We generate the adjoint code using Tapenade. To do so, we type in a terminal (or in a Makefile):

```
tapenade -backward -head simple_operation \
    -vars "c d" -outvars "x" \
    -html -0 $(PWD)/backward simple_operation.f90
```

Then, the file *simple\_operation\_b.f90* is created in the *backward* directory. This file is the following:

```
(INRIA, Tropics team)
i
         Generated by TAPENADE
!
  Tapenade 3.4 (r3375) - 10 Feb 2010 15:08
I.
!
 Differentiation of simple_operation in reverse (adjoint) mode:
!
   gradient
                 of useful results: x
!
   with respect to varying inputs: x c d
   RW status of diff variables: x:in-zero c:out d:out
!
SUBROUTINE SIMPLE_OPERATION_B(c, cb, d, db, x, xb)
  IMPLICIT NONE
 INTEGER, PARAMETER :: n=2
 DOUBLE PRECISION, INTENT(IN) :: c
 DOUBLE PRECISION :: cb
 DOUBLE PRECISION, INTENT(IN) :: d
 DOUBLE PRECISION :: db
 DOUBLE PRECISION, DIMENSION(n) :: x
 DOUBLE PRECISION, DIMENSION(n) :: xb
 DOUBLE PRECISION, DIMENSION(n, n) :: a, at
 DOUBLE PRECISION, DIMENSION(n, n) :: ab
 DOUBLE PRECISION, DIMENSION(n) :: b
 DOUBLE PRECISION, DIMENSION(n) :: bb
 INTEGER :: i,j
 a(1, 1) = 3*c+ 4*d
```

```
a(1, 2) = c + 18*d
 a(2, 1) = 2.5*d
 a(2, 2) = 3.4*c
 b(1) = 8*c*d
 b(2) = 2*d
 CALL SYS_SOLVER_B(a, ab, b, bb, x, xb, n)
 db = 2*bb(2)
 bb(2) = 0.D0
 cb = 3.4*ab(2, 2) + 8*d*bb(1)
 ab(2, 2) = 0.D0
 db = db + 2.5*ab(2, 1) + 8*c*bb(1)
 ab(2, 1) = 0.D0
 cb = cb + ab(1, 2)
 db = db + 18*ab(1, 2)
 ab(1, 2) = 0.D0
 cb = cb + 3*ab(1, 1)
 db = db + 4*ab(1, 1)
 xb = 0.D0
END SUBROUTINE SIMPLE_OPERATION_B
```

Using the method presented in subsection B.5.1 we replace the unknown adjoint of the linear system solver:

CALL SYS\_SOLVER\_B(a, ab, b, bb, x, xb, n)

by the following instructions:

```
!at = transpose(a)
 at(1, 1) = 3*c + 4*d
 at(2, 1) = c + 18*d
 at(1, 2) = 2.5*d
 at(2, 2) = 3.4*c
!new solver
 call sys_solver(a,b,x,n)
 call sys_solver(at,xb,bb,n)
 ab=0.d0
 do j=1,n
    do i=1,n
 if (ab(i,j).ne.0.) ab(i,j)=-bb(i)*x(j)
     end do
 end do
!At = transpose(A)
 do j=1,n
```

```
do i=1,n
   At(i, j) = A(j,i)
   end do
end do
!new solver
   call linear_solver(A,b,x,n)
   call linear_solver(At,xb,bb,n)
   do j=1,n
      do i=1,n
   if (Ab(i,j).ne.0.) Ab(i,j)= - bb(i)*x(j)
      end do
end do
```

We have recycled the linear solver (which can be called as a black box).

## **B.6** Complement: MPI instructions

We present the procedure to obtain the adjoint code of a Fortran code calling MPI basic routines and using the automatic differentiation tool Tapenade. This part is skipped since no time enough.

## B.7 Complement: optimize your memory. Checkpointing.

When you use Tapenade to create an adjoint code, you have no mean to know the memory consumption that will take Tapenade. The latter is a direct consequence of the "PUSH" and "POP" instructions introduced by Tapenade in your code. The "PUSH" instruction is keeping in memory the variable until a "POP" of this variable is done. Every adjoint code is divided into two part: the first part, which is just a copy of your direct code with add of "PUSH" instructions filling the memory; the second part written by Tapenade and running a backward code with the adjoint variables, and freeing the memory with POP instructions.

Memory problems can appear if you have "PUSH" es of, for example, your state variable at every time steps on a simulation.

We present below a way to avoid memory faults by using the checkpointing strategy of Tapenade in a "smart" way.

**Checkpointing strategy** First thing is to anticipate and identify which part of the "PUSH" and "POP" instructions in your code is potentially greedy in memory. To narrowing it, it should appear in the biggest loop (i.e. with the most iterations) of the algorithm and occur where your state variables are.

An example of order of magnitude. In DassFlow software, based on the 2D shallow-water equations, the main "PUSH" and "POP" contributions are related to the state variable in the time step loop.

For a 25 000 cells mesh with 100 000 time steps, if we do nothing, as our state variable counts 3 fields, this would cost just for "pushing" this variable:  $3 \times 25000 \times 100000 \times 8 = 6 \times 10^{10} \text{o} = 60$  Go of RAM which is a lot !

(NB. The number 8 represents the octet size of a double precision number; it is 4 for a simple precision).

Most of the time, it is only one variable that make your system fall.

## B.8 Practical work Tapenade

As practical work, we generate the adjoint code of a Fortran source code using the software Tapenade. Next, validate the codes obtained and we perform sensitivities analysis.

The direct code solves the 1D heat equation  $\partial_t u = K \partial_x^2 u$  using finite differences, with a time scheme either explicit in time or implicit in time. In the latter case, we call a linear solver (as a black box).

In the practical work, we consider two cases (two direct programs):

an explicit FD solver with MPI instructions.

an implicit FD solver with a call to a linear solver as a black box.

Sorry, no time enough to consider this part.

## B.9 Example of a software architecture designed for Automatic Differentiation

See the toy DassFlow code (TP automatic differentiation using Tapenade software).

## Appendix C

# Weather forecasting & VDA algorithm variants \*

The outline of this chapter is as follows.

#### Contents

| B.1 What adjoint code ?  |  |  |  |  |
|--|--|--|--|--|
| B.2 From mathematical differentiation to source code differentiation . 171   |  |  |  |  |
| B.3 Automatic differentiation in short   |  |  |  |  |
| B.4 Exercices  |  |  |  |  |
| B.5 Recycle your linear solver 184   |  |  |  |  |
| B.5.1 How to adjoint-ize your linear solver ? $\dots \dots $ |  |  |  |  |
| B.5.2 A simple example $\dots \dots \dots$                   |  |  |  |  |
| B.6 Complement: MPI instructions   |  |  |  |  |
| B.7 Complement: optimize your memory. Checkpointing 193  |  |  |  |  |
| B.8 Practical work Tapenade  |  |  |  |  |
| B.9 Example of a software architecture designed for Automatic Dif-<br>ferentiation   |  |  |  |  |

## C.1 Data for weather forecasting

#### The historical and one of the most challenging DA problem: weather forecast

The Variational Data Assimilation (VDA) method has been introduced by Y. Sasaki in the mid 50's [38, 39] for applications in meteorology. An improvement of the method based on the optimal control theory and the adjoint equations has been introduced by G. Marchuk, F.-X. Le Dimet and O. Talagrand in the mid 80's, see e.g. [27] and references therein; see also the

pioneer studies [15, 32, 31].

Since the early 2000th, operational VDA systems are implemented in all large Weather Prediction Centers (WPC, eg. ECMWF, Meteo-France, Met Office, NOAA etc). The most complex and sophisticated data assimilation systems are probably elaborated in these WPC. The dynamic of the atmosphere is 3D, turbulent, multi-scale, multi-physics. The mathematical models are based on the conservation laws: mass, momentum, energy. The resulting equations are highly non-linear, fully coupled. Unknowns are the following 3D fields: temperature, pressure, humidity, velocities etc. The complexity of the (direct) model is huge: at each time step (typically 10 minutes for mid-range forecast), the discrete unknown contains billions of components.



Figure C.1: Weather prediction based on data assimilation processes optimally fusing all the available information. Image source: CNRM, Meteo-France.

Forecasts have been greatly improved during the last two decades in particular due the introduction of more and more sophisticated data assimilation processes. The latter enable to benefit of the exponentially exploding observations systems. Employed algorithms are variants of the presented VDA algorithm combined with the Ensemble Kalman Filter (EnKF) (not studied in the present course). Some information on weather forecasting are presented in Appendix 5), .

Observations are heterogeneous in nature, resolution and accuracy; they are distributed more or less regularly in space and time. The numbers of available observations is exploding; it

196



Figure C.2: Observations number explosion. Plot: numbers used in the global scale model ARPEGE (and classified by nature) vs months (period 2002-2010). Image source: CNRM, Meteo-France.

is about dozens of millions per day, see Fig. Satellite measurements (surface temperature, nebulosity, radiance etc) are difficult to compare to the model outputs. Measurements in-situ (temperature, pressure, humidity, wind) are heterogeneously distributed and very local values. In complex systems, both observations and models are affected by multi-source uncertainties. Furthermore, extra difficulties are specific to geophysical flows: it is non reproducible phenomena, observations are sparse and uncertain, phenomena are multi-scale and models of course incomplete; moreover initial states cannot be fully measured.

A great challenge of the WPC is to reconstruct the initial state of the atmosphere i.e. the atmosphere variables at present time ! Indeed even if the observations are more and more dense and accurate, they are not dense enough to accurately reconstruct the atmosphere state everywhere. Moreover the atmosphere dynamic depends strongly on its initial state; actually it is well known its behavior may be chaotic in medium-long time range.

The *purpose of data assimilation* is to reconstruct as accurate as possible the state of the atmosphere at present time (the initial condition) using all available information i.e.:

- . the physical and conservation laws governing the flow dynamic (mathematical model),
- . the observations (sparse and heterogeneous in nature),
- . the statistic errors of the observations and prior probabilistic behavior of the various fields.

The VDA algorithms aim at fusing as best as possible all these information in view to recover the atmosphere state at present time.

#### Orders of magnitude from the operational ECMWF forecast system (in 2009...)

Below are presented some order of magnitudes for Numerical Weather Prediction (NWP). All information and figures below come from documents written by Y. Tremolet and M. Fisher from the European Centre for Medium-Range Weather Forecasts (ECMWF) based at Reading, UK. ("ECMWF is an intergovernmental organization supported by 34 States. It provide operational medium- and extended-range forecasts and a state-of-the-art super-computing facility for scientific research. It pursues scientific and technical collaboration with satellite agencies and with the European Commission").

Forecast has improved the last twelve years in particular because of the introduction of data assimilation process into the complex dynamical atmosphere models, which allow to benefit more and more from the global various observing systems. Methods employed are 4d-var, 3d-var and ensemble methods.

Orders of magnitude presented below date from 2009 - 2010.

The unknowns of the (non-linear, fully coupled) mathematical models are 3D fields: temperature, pressure, humidity, velocities.

The dimension of the state (number of unknowns) is about  $10^9$ .

Time step is about 10 minutes. Discretization is about 16 kms in horizontal (triangles) and with about 90 layers in vertical (0 - 80 km meshed).

Observations are heterogeneous in nature and in space. They are in dozens of millions (see figures below).

Satellites measurements are: surface temperature, nebulosity etc. They are complex to interpret and to compare to the model outputs.

Measurements in-situ are: temperature, pressure, humidity, wind etc.

The configuration uses a 12h cycling window, with a 4D-var incremental method.

The outer loop (and forecast) resolution is 25 km. The inner loops resolutions are between 200 km and 80 km.

On average, 9 million observations are assimilated per 12h cycle. 96% of assimilated data is from satellites. On average, 4D-Var runs on 1536 CPUs in 1h10.

The incremental method with appropriate preconditioning allows the computational cost to be reduced to an acceptable level.



Figure C.3: Meteorology. Some observation coverage. Figure extracted from Tremolet - Fisher, ECMWF, Colloquium 2009.

## C.2 The incremental VDA (4D-Var) algorithm

For very large scale problems (e.g. oceans and atmosphere in geophysics), it can be unaffordable to perform the 4D-var process for the full model resolution (i.e. the direct model including the full physics, and solved on the fine grids). Furthermore, if the assimilation is required for prediction (e.g. weather forecast), the forecast needs to be performed faster than the real time... For smaller scale problems (e.g. based on reduced models such as the 2D shallow-water equations), even if full 4d-var is possible, it remains CPU time consuming (see the forthcoming examples in river hydraulics).

An option to reduce the computational cost is to elaborate an incremental 4D-var method.

The basic idea is to combine empirically in the minimization process low-resolution and/or linearized equation terms with the full physics model.



Figure C.4: Meteorology. Observation sources. Figure extracted from Tremolet - Fisher, ECMWF, Colloquium 2009.



Figure C.5: Meteorology. Observation Numbers (in milions per day). Source: ECMWF.

Keeping the discrete notations above, the so-called *innovation vector*  $d_n$  is defined by :

$$d_n = (z_n - \mathcal{H}(y_n)) \quad n \ge 0$$

The innovation vector  $d_n$  measures at time  $t_n$ , the discrepancy between the model output and the observed quantity, in the observation space.

The basic idea of the incremental 4D-var method is to modify the 4D-var algorithm as follows: 1) the iterative control variable corrections are performed using a low-resolution model (both

|                   | Screen      | ed     | Assimilated |        |
|-------------------|-------------|--------|-------------|--------|
| Synop:            | 450,000     | 0.3%   | 64,000      | 0.7%   |
| Aircraft:         | 434,000     | 0.3%   | 215,000     | 2.4%   |
| Dribu:            | 24,000      | 0.02%  | 7,000       | 0.1%   |
| Temp:             | 153,000     | 0.1%   | 76,000      | 0.8%   |
| Pilot:            | 86,000      | 0.1%   | 39,000      | 0.4%   |
| AMV's:            | 2,535,000   | 1.6%   | 125,000     | 1.4%   |
| Radiance data:    | 150,663,000 | 96.9%  | 8,207,000   | 91.0%  |
| Scat:             | 835,000     | 0.5%   | 149,000     | 1.7%   |
| GPS radio occult. | 271,000     | 0.2%   | 137,000     | 1.5%   |
| TOTAL:            | 155,448,000 | 100.0% | 9,018,000   | 100.0% |

#### Data count for one 12h 4D-Var cycle (0900-2100 UTC, 3 March 2008)

Figure C.6: Meteorology. Observation usage. Figure extracted from Tremolet - Fisher, ECMWF, Colloquium 2009. We are still far from using all available observations.

in physics and grids). It gives low-resolution inner-loops.

2) Once these low-resolution inner-loops have converged, the discrepancy between the original full model and the measurements (i.e. the innovative vector) is computed.\*) Repeat the process.

For a sake of simplicity, we define the operator model  $\mathcal{M}$  as follows :

$$y_n = \mathcal{M}(y_0)$$

If we set:  $y_b = y_0 + \delta y_0$ , the sough quantity becomes  $\delta y_0$  and we have:

$$j(y_0) = \frac{1}{2} \sum_{n=0}^{N} (\mathcal{H}(y_n) - z_n)^T R^{-1} (\mathcal{H}(y_n) - z_n) + \frac{1}{2} \delta y_0^T B^{-1} \delta y_0$$

Next, the "perturbation"  $\delta y_n$  corresponding to the perturbation  $\delta y_0$  is defined as follows:

$$y_n + \delta y_n = \mathcal{M}(y_0 + \delta y_0)$$

A formal linearization (Taylor's expansion order 1) gives :

$$\mathcal{M}(y_0 + \delta y_0) \approx \mathcal{M}(y_0) + M(y_0) \cdot \delta y_0$$
 hence  $y_n + \delta y_n \approx y_n + M(y_0) \cdot \delta y_0$ 

where M is the linear tangent model at time step  $t_n$ . Then :  $\delta y_n \approx M(y_0) \cdot \delta y_0$ .

Similarly:

$$\mathcal{H}(y_n + \delta y_n) \approx \mathcal{H}(y_n) + H_n \cdot \delta y_n$$

with  $H_n$  the linearized observation operator at time step  $t_n$ .

We have:

$$j(y_0 + \delta y_0) \approx g(\delta y_0) = \frac{1}{2} \sum_{n=0}^{N} (H_n \cdot \delta y_n - d_n)^T R^{-1} (H_n \cdot \delta y_n - d_n) + \frac{1}{2} \delta y_0^T B^{-1} \delta y_0$$

The basic idea of the incremental 4D-var method is to minimize the cost function with respect to the increment  $\delta y_0$ . Considering the cost function  $g(\delta y_0)$ , the inverse problem becomes *linear-quadratic optimal control problem* since  $\delta y_n \approx M(y_0) \cdot \delta y_0$ . Then the minimization can be performed using the Conjugate Gradient with preconditioning, hence very fast to compute.



Figure C.7: The incremental VDA (4D-Var) algorithm

In summary, the 4D-var incremental algorithm reads as follows (see Figure C.7).

\*) Inner loop.

The minimization process is performed with linearized operators (M and  $H_n$ ), furthermore potentially with coarser grids and simplified physics.

Since the cost function  $g(\delta y_0)$  is strictly convex (it is quadratic in  $\delta y_0$ ), the extremely efficient

Conjugate Gradient method with preconditioning can be used to minimize g.

\*) Outer loop.

Once the low-resolution minimization process has converged (the so-called "analysis step" is done), the original "high-resolution / full" model is performed (it is the so-called "prediction stage"). The prediction stage gives new innovation vectors (the discrepancy between the model and the observations).

Let us point out that the innovation vectors  $d_n$  are the crucial input of the assimilation method, thus they are computed using the high-resolution model.

\*) Update the increment and the reference state.

Thus, the incremental 4D-var method is a mix of calibration on linearized operators, full physic predictions and discrepancy measurements based on the fine innovation vector.

For more details, the reader should consult [10, 43], where the method is developed in a weather forecast context.

## C.3 Other VDA algorithm variants

**The sequential 3D-var** The 3D-var algorithm is a time-independent version of the 4D-var.

Historically, it has been used for time-dependent problems because the 4D-var was too much CPU-time consuming. The 3D-var algorithm can be described as follows. During a time interval, all observations are gathered, then the discrepancy between the model and the observations is computed at a single time step (eg the final time of the interval considered). In other words, all observations are gathered (and potentially averaged) at the assimilation time. Next, the minimization process is performed.

No time integration of the direct model nor of the adjoint model have to be performed; only steady-state assimilations are performed at a given instant.

Next the unsteady model is performed until the next assimilation time.

The 3D-var drawback applied to a time-dependent model is obvious: the model outputs and the observations are not compared at the right instants...

**Variant: 3D-FGAT** In the 3D-FGAT version (FGAT for First Guess at Appropriate Time), the linear tangent model and adjoint model are replaced by the identity operator. It is an empirical method between the 3D-var and the 4D-var. The only advantage compared to the 3D-var is the innovation vector is evaluated at the right observation times.



Figure C.8: Comparison of algorithms: 3D-var (sequential, time-independent) vs 4D-var (with time integration). Figure source: the ECMWF data assimilation course (www.ecmwf.int) by F. Bouttier, P. Courtier, [10].



Figure C.9: Left) Algo 3D-var: During a time interval, all observations are gathered. Right) Algo 3D-FGAT: each observation is assimilated at the right instant but the linear tangent and the adjoint models are replaced by the identity operator. Image source: A. Vidard's PhD manuscript [44].

#### Other deterministic data assimilation methods

Finally, let us cite other deterministic data assimilation algorithms: the "reduced 4D-var" method and the "4D-PSAS" method are extensions of the present formalism (adjoint method, optimality system, local minimization) while the third method "Back and Forth Nudging" is not based on the present optimal control formalism.

Reduced 4D-Var consisting to elaborate order reductions (based on Empirical Orthogonal Functions EOF) may be applied to the present control problem. This may lead to a drastic reduction of the control space dimension, see e.g. [36] and references therein. The 4D-PSAS (Physical space Statistical Analysis System), a *dual method*. It consists to set the minimization problem in the observation space instead of the state space, hence a lower-dimensional minimization problem, see [16, 4] and references therein.

In the case of a "perfect model" (no RHS error term is introduced into the direct model) then the 4D-Var and 4D-PSAS methods are equivalent. But in presence of a model error term (as a direct model source term), this term is naturally taken into account in the 4D-PSAS method. Furthermore, the dual approach performs minimization in the observation space, which is smaller than the state space.

Finally, let us cite the *Back and Forth Nudging* (BFN) method [5], which presents the advantage to circumvent the adjoint model derivation. From [5]: "The standard nudging technique consists in adding to the equations of the model a relaxation term that is supposed to force the observations to the model. The BFN algorithm consists in repeatedly performing forward and backward integrations of the model with relaxation (or nudging) terms, using opposite signs in the direct and inverse integrations, so as to make the backward evolution numerically stable". 

## Bibliography

- [1] G. Allaire. Numerical analysis and optimization. Oxford university press, 2007.
- [2] Mark Asch, Marc Bocquet, and Maëlle Nodet. Data assimilation: methods, algorithms, and applications. SIAM, 2016.
- [3] Richard C Aster, Brian Borchers, and Clifford H Thurber. *Parameter estimation and inverse problems*. Elsevier, 2018.
- [4] D. Auroux. Etude de differentes methodes d'assimilation de donnees pour l'environnement. *PhD thesis. Nice - Sophia-Antipolis university, France*, 2003.
- [5] Didier Auroux and Jacques Blum. A nudging-based data assimilation method: the back and forth nudging (bfn) algorithm. Nonlinear Processes in Geophysics, 15(2):305–319, 2008.
- [6] Combettes P.-L. Bauschke H. Convex Analysis and Monotone Operator Theory in Hilbert Spaces. Springer, 2011.
- [7] A. F. Bennett. Inverse modeling of the ocean and atmosphere. Cambridge University Press, 2005.
- [8] J. Blum, F.-X. Le Dimet, and Navon I.M. Data assimilation for geophysical fluids. In Temam and Tribbia, editors, *Handbook of Numerical Analysis*, volume 14. North-Holland, 2009.
- [9] J.F. Bonnans, J.C. Gilbert, C. Lemaréchal, and C.A. Sagastizábal. *Numerical optimization:* theoretical and practical aspects. Springer, 2006.
- [10] L. Bouttier and P. Courtier. Data assimilation concepts and methods. ECMWF Training course. www.ecmwf.int., 1999.
- [11] H. Brézis. Functional Analysis, Sobolev Spaces and Partial Differential Equations. Springer, 2010.
- [12] Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *PNAS*, 113(15), 04 2016.
- [13] Emmanuel Candès and Justin Romberg. Sparsity and incoherence in compressive sampling. Inverse problems, 23(3), 04 2007.

- [14] Alberto Carrassi, Marc Bocquet, Laurent Bertino, and Geir Evensen. Data assimilation in the geosciences: An overview of methods, issues, and perspectives. Wiley Interdisciplinary Reviews: Climate Change, 9(5):e535, 2018.
- [15] P. Courtier and O. Talagrand. Variational assimilation of meteorological observations with the direct adjoint shallow water equations. *Tellus*, 42(A):531–549, 1990.
- [16] Philippe Courtier. Dual formulation of four-dimensional variational assimilation. Quarterly Journal of the Royal Meteorological Society, 123(544):2449–2461, 1997.
- [17] DassFlow. "data assimilation for free surface flows". Open source computational software: http://www.math.univ-toulouse.fr/DassFlow.
- [18] Lawrence C. Evans. An Introduction to Mathematical Optimal Control Theory. University of California, Berkeley, 2014.
- [19] Geir Evensen. Data assimilation: the ensemble Kalman filter. Springer Science & Business Media, 2009.
- [20] A. Griewank. Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation. Number 19 in Frontiers in Appl. Math. SIAM, Philadelphia, PA, 2000.
- [21] Matthew CG Hall and Dan G Cacuci. Physical interpretation of the adjoint functions for sensitivity analysis of atmospheric models. *Journal of the atmospheric sciences*, 40(10):2537-2546, 1983.
- [22] L. Hascoet. Automatic Differentiation by Program Transformation. INRIA Sophia-Antipolis, TROPICS team. Available on-line., 2007.
- [23] L. Hascoet, V. Pascual, et al. Tapenade software. INRIA, project Tropics, wwwsop.inria.fr/tropics, 2012.
- [24] M. Honnorat. Assimilation de données lagrangiennes pour la simulation numérique en hydraulique fluviale. *PhD thesis. Grenoble IT, France*, 2007.
- [25] Barbara Kaltenbacher, Andreas Neubauer, and Otmar Scherzer. *Iterative regularization methods for nonlinear ill-posed problems*, volume 6. Walter de Gruyter, 2008.
- [26] Andreas Kirsch. An introduction to the mathematical theory of inverse problems, volume 120. Springer Science & Business Media, 2011.
- [27] F.X. Le Dimet and O. Talagrand. Variational algorithms for analysis assimilation of meteorological observations: theoretical aspects. *Tellus A*, 38:97–110, 1986.
- [28] J.L. Lions. Contrôle optimal de systèmes gouvernés par des équations aux dérivées partielles. Dunod, 1968.
- [29] J.L. Lions. Optimal control of systems governed by partial differential equations. Springer-Verlag, 1971.
- [30] S. Mallat. Cours collège de france "sciences des données". 2020.
- [31] GI Marchuk and VP Shutyaev. Iteration methods for solving a data assimilation problem. Russian Journal of Numerical Analysis and Mathematical Modelling, 9(3):265–280, 1994.
- [32] GI Marchuk and VB Zalesny. A numerical technique for geophysical data assimilation problems using pontryagin?s principle and splitting-up method. *Russian Journal of Numerical Analysis and Mathematical Modelling*, 8(4):311–326, 1993.
- [33] J. Monnier. Modèles numériques directs et inverses d'écoulements de fluides. Habilitation à Diriger des Recherches. Institut National Polytechnique de Grenoble (Grenoble IT), 2007.
- [34] J. Monnier, K. Larnier, P. Brisset, and F. et al. Couderc. Dassflow (data assimilation for free surface flows): numerical analysis report. Open source computational software: http://www.math.univ-toulouse.fr/DassFlow.
- [35] U. Naumann. The Art of Differentiating Computer Programs: An Introduction to Algorithmic Differentiation. SIAM, 2012.
- [36] Céline Robert, Eric Blayo, and Jacques Verron. Comparison of reduced-order, sequential and variational data assimilation methods in the tropical pacific ocean. Ocean Dynamics, 56(5-6):624–633, 2006.
- [37] Samuel H. Rudy, Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4), 04 2017.
- [38] Yoshikazu Sasaki. An objective analysis based on the variational method. J. Meteor. Soc. Japan, 36(3):77–88, 1958.
- [39] Yoshikazu Sasaki. Some basic formalisms in numerical variational analysis. *Monthly Weather Review*, 98(12):875–883, 1970.
- [40] O. Thual. *Hydrodynamique de l'environnement*. Editions de l'Ecole Polytechnique, 2010.
- [41] Robert Tibshirani. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society. Series B (Methodological), 58(1), 1996.
- [42] E. Trélat. Optimal control: theory and applications. Vuibert, Paris, 2008.
- [43] Y. Trémolet. Incremental 4d-var convergence study. Tellus A, 59(5):706–718, 2008.
- [44] A. Vidard. Vers une prise en compte de l'erreur modèle en assimilation de données 4dvariationnelle. *PhD thesis. Grenoble university J. Fourier, France*, 2001.
- [45] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. Journal of the Royal Statistical Society. Series B (Statistical Methodology), 67(2), 2005.