



**HAL**  
open science

# Optimisation of the total population size with respect to the initial condition for semilinear parabolic equations: Two-scale expansions and symmetrisations

Idriss Mazari, Grégoire Nadin, Ana Isis Toledo Marrero

## ► To cite this version:

Idriss Mazari, Grégoire Nadin, Ana Isis Toledo Marrero. Optimisation of the total population size with respect to the initial condition for semilinear parabolic equations: Two-scale expansions and symmetrisations. *Nonlinearity*, 2021, 34 (11), pp.7510-7539. 10.1088/1361-6544/ac23b9 . hal-03846882v2

**HAL Id: hal-03846882**

**<https://hal.science/hal-03846882v2>**

Submitted on 12 Jul 2024

**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.



Distributed under a Creative Commons Attribution 4.0 International License



PAPER • OPEN ACCESS

# Optimisation of the total population size with respect to the initial condition for semilinear parabolic equations: two-scale expansions and symmetrisations

To cite this article: Idriss Mazari *et al* 2021 *Nonlinearity* **34** 7510

View the [article online](#) for updates and enhancements.

## You may also like

- [A one-population Amari model with periodic microstructure](#)  
Nils Svanstedt, John Wyller and Elena Maljutina
- [Second-order two-scale analysis and numerical algorithms for the hyperbolic–parabolic equations with rapidly oscillating coefficients](#)  
Hao Dong, , Yu-Feng Nie *et al.*
- [Homogenization of biomechanical models of plant tissues with randomly distributed cells](#)  
Andrey Piatnitski and Mariya Ptashnyk

# Optimisation of the total population size with respect to the initial condition for semilinear parabolic equations: two-scale expansions and symmetrisations

Idriss Mazari<sup>1,\*</sup>, Grégoire Nadin<sup>2</sup> and Ana Isis Toledo Marrero<sup>3</sup>

<sup>1</sup> CEREMADE, UMR CNRS 753, Université Paris-Dauphine, Université PSL, Place du Maréchal De Lattre De Tassigny, Paris, F-75775 Paris cedex 16, France

<sup>2</sup> CNRS, Sorbonne Universités, UPMC Univ Paris 06, UMR 7598, Laboratoire Jacques-Louis Lions, F-75005, Paris, France

<sup>3</sup> Sorbonne Universités, UPMC Univ Paris 06, UMR 7598, Laboratoire Jacques-Louis Lions, F-75005, Paris, France

E-mail: [mazari@ceremade.dauphine.fr](mailto:mazari@ceremade.dauphine.fr), [gregoire.nadin@sorbonne-universite.fr](mailto:gregoire.nadin@sorbonne-universite.fr) and [ana-isis.toledo\\_marrero@sorbonne-universite.fr](mailto:ana-isis.toledo_marrero@sorbonne-universite.fr)

Received 2 April 2021, revised 12 August 2021

Accepted for publication 3 September 2021

Published 27 September 2021



CrossMark

## Abstract

In this article, we propose in-depth analysis and characterisation of the optimisers of the following optimisation problem: how to choose the initial condition  $u_0$  in order to maximise the spatial integral at a given time of the solution of the semilinear equation  $u_t - \Delta u = f(u)$ , under  $L^\infty$  and  $L^1$  constraints on  $u_0$ ? Our contribution in the present paper is to give a characterisation of the behaviour of the optimiser  $\bar{u}_0$  when it does not saturate the  $L^\infty$  constraints, which is a key step in implementing efficient numerical algorithms. We give such a characterisation under mild regularity assumptions by proving that in that case  $\bar{u}_0$  can only take values in the ‘zone of concavity’ of  $f$ . This is done using two-scale asymptotic expansions. We then show how well-known isoperimetric inequalities yield a full characterisation of maximisers when  $f$  is convex. Finally, we provide several numerical simulations in one and two dimensions that illustrate and exemplify the fact that such characterisations significantly improve the computational time. All our theoretical results are in the one-dimensional case

\*Author to whom any correspondence should be addressed.

Recommended by Dr Susanna Terracini.



Original content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](https://creativecommons.org/licenses/by/3.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

and we offer several comments about possible generalisations to other contexts, or obstructions that may prohibit doing so.

Keywords: reaction equation, optimal control, shape optimisation, two-scale expansions

Mathematics Subject Classification numbers: 35B30, 35B65, 35K15, 35K57, 35Q80, 35Q80.

(Some figures may appear in colour only in the online journal)

## 1. Introduction

### 1.1. Scope of the article

In this article, we propose to establish several results concerning an optimal control problem for a class of semilinear parabolic equations. Some aspects of this problem have been initially addressed by two of the authors in [36]. In the setting under consideration, the control variable to optimise is the initial condition. As we will see throughout the statement of the results, the form (e.g. convex or concave) of the semilinearity plays a crucial role in the analysis and calls for a detailed study of second order optimality conditions, which is our main result, theorem 1. In the case of convex semilinearities, using rearrangement arguments, we can give a full characterisation of maximisers, see theorem 2. Using theorem 1, we can improve an algorithm initially developed in [36], and we display numerical results in section 4.

**Initial motivation of the paper.** The origin of this paper is the study of an optimal control problem that arises naturally in mathematical biology and that deals with bistable reaction–diffusion equations. Namely, for a semilinear equation, what is the best possible initial condition (‘best’ being understood as maximising the integral of the solution at a certain time horizon  $T$ )? The complicated behaviour of bistable nonlinearities, which are neither convex nor concave, makes the analysis of this query very intricate. The two aforementioned results, theorems 1 and 2, enable us to show how complicated the behaviour of maximisers can be for such nonlinearities. Bistable equations are of central importance in mathematical biology [35] and, very broadly speaking, model the evolution of a subgroup of a population. Among their many applications, one may mention chemical reactions [39], neurosciences [13], phase transition [27], linguistic dynamics [40] or the evolution of diseases [35]. The last interpretation is of particular relevance to us, given that this model is used to design optimal strategies in order to control the spread of several mosquito borne diseases such as the dengue [7]; this was the main motivation in [3]. The strategy is to release a certain amount of Wolbachia carrying mosquitoes (Wolbachia is a bacterium that inhibits the transmission of mosquito borne diseases that individuals inherit from their mother) in a population of wild mosquitoes that can potentially transmit the diseases, in order to maximise the proportion of Wolbachia carrying mosquitoes at the final time. In mathematical terms: given a time horizon  $T$ .

*How should we arrange the initial population in order to maximise the population size at  $T$ ?*

Even without having stated it formally, we can make two observations on this problem: the first one is that, since the variable of the equation is the proportion of a subgroup, we need to enforce pointwise ( $L^\infty$ ) constraints. The second one is that we naturally have to add an  $L^1$  constraint for modelling reasons. Both of these constraints can in practice be very complicated to handle.

**Optimisation problems in mathematical biology.** Let us briefly sketch how this problem fits in the literature devoted to such optimisation and control problems for mathematical biology. Optimisation problems for reaction–diffusion equations have by now gathered a lot of attention from the mathematical community. Most of these optimisation problems are set in a stationary setting, that is, assuming that the population has already reached an equilibrium, and the main problems that have been considered often deal with the optimisation of the spatial heterogeneity [12, 16, 17, 22, 23, 30, 31, 33, 37] (we also refer to the recent surveys [19, 32]); most of these works deal with monostable nonlinearities. We also point to the recent [9] for the study of an optimal control problem for parabolic monostable equations. On the other hand, optimisation problems for bistable equations, which are the other paradigmatic class of equations in mathematical biology [35], have received a less complete mathematical treatment, but are now the topic of an intense research activity from the control point of view, see [3, 36] and the references therein. Related optimal control problems are not yet fully understood. More generally, less attention has been devoted to optimisation problem with respect to the initial condition for such semilinear evolution equations.

1.2. *Mathematical setup and statement of the results*

1.2.1. *Statement of the problem.* We work in  $\Omega = (0; \pi)$ . We consider a  $\mathcal{C}^2$  function  $f : [0; 1] \rightarrow \mathbb{R}$ , and the associated parabolic equation

$$\begin{cases} \partial_t u - \Delta u = f(u) & \text{in } \mathbb{R}_+ \times \Omega, \\ u(0, x) = u_0(x) & \text{in } \Omega, \\ \frac{\partial u}{\partial \nu}(t, x) = 0 & \text{in } \mathbb{R}_+ \times \partial\Omega, \end{cases} \tag{1}$$

where  $u_0$  is an initial condition satisfying the constraint

$$0 \leq u_0 \leq 1.$$

Since our initial motivation, as explained in the first paragraph of this introduction, is to maximise the proportion of a subgroup of a population, such an  $L^\infty$  constraint is natural. At the mathematical level, it should be noted that we could carry out the same analysis with any  $L^\infty$  constraint of the form  $0 \leq u_0 \leq \kappa$  by a simple change of variable.

We define, for any  $T > 0$ , the functional

$$\mathcal{J}_T(u_0) := \int_{\Omega} u(T, x) dx. \tag{2}$$

The goal is to maximise  $\mathcal{J}_T$  with respect to  $u_0$ . Since we are then again wondering how to maximise the proportion of a subgroup by controlling its distribution at the initial time, it is natural to introduce a  $L^1$  constraint on  $u_0$ . This constraint is encoded by a parameter  $m \in (0; |\Omega|)$  which is henceforth fixed.

These considerations lead us to defining our admissible class as

$$\mathcal{A} := \left\{ u_0 \in L^\infty(\Omega), 0 \leq u_0 \leq 1 \text{ a.e.}, \int_{\Omega} u_0 = m \right\}, \tag{3}$$

and the variational problem under scrutiny throughout this paper is

$$\max_{u_0 \in \mathcal{A}} \mathcal{J}_T(u_0). \tag{P}_f$$

This problem  $(\mathbf{P}_f)$  was directly addressed by two of the authors in [36], where expressions for the first and second order optimality conditions were provided. We need to recall them, to motivate and contextualise our results: if we consider  $u_0 \in \mathcal{A}$  and an admissible perturbation  $h_0$  at  $u_0$  (by ‘admissible perturbation’ we refer to the fact that  $h_0$  belongs to the tangent cone to the set  $\mathcal{A}$  at  $u_0$ . This tangent cone is the set of functions  $h \in L^\infty(\Omega)$  such that, for any sequence of positive real numbers  $(\varepsilon_n)_{n \in \mathbb{N}}$  decreasing to 0, there exists a sequence of functions  $(h_n)_{n \in \mathbb{N}} \in L^\infty(\Omega)^{\mathbb{N}}$  converging to  $h$  as  $n \rightarrow +\infty$ , and  $u_0 + \varepsilon_n h_n \in \mathcal{A}$  for every  $n \in \mathbb{N}$ ) then the first order Gâteaux-derivative of  $\mathcal{J}_T$  at  $u_0$  in the direction  $h_0$  is

$$\langle \nabla \mathcal{J}_T(u_0), h_0 \rangle = \int_{\Omega} h_0(x)p(0, x)dx \tag{4}$$

where  $p$  solves the adjoint equation

$$\begin{cases} -\partial_t p - \Delta p = f'(u)p & \text{in } (0, T) \times \Omega, \\ p(T, x) = 1 & \text{in } \Omega, \\ \frac{\partial p}{\partial \nu}(t, x) = 0 & \text{for all } t \in (0, T), \text{ for all } x \in \partial\Omega. \end{cases} \tag{5}$$

Here  $u$  is the solution of (1) with initial condition  $u_0$ .

The main result in [36] states the following:

**Theorem [36].** *There exist a solution  $\bar{u}_0 \in \mathcal{A}$  of  $(\mathbf{P}_f)$ . Moreover, setting  $\bar{u}$  as the solution of (1) associated with this optimal initial data and  $\bar{p}$  as the unique solution of (5) for  $u = \bar{u}$ , there exists a non-negative real value  $\bar{c}$  such that*

- (a) *If  $0 < \bar{u}_0(x) < 1$  then  $\bar{p}(0, x) = \bar{c}$ ,*
- (b) *If  $\bar{p}(0, x) > \bar{c}$ , then  $\bar{u}_0(x) = 1$ ,*
- (c) *If  $\bar{p}(0, x) < \bar{c}$ , then  $\bar{u}_0(x) = 0$ .*

Finally, for almost every  $x \in \{\bar{p}(0, \cdot) = \bar{c}\}$ , one has

$$f'(\bar{u}_0(x)) = -\bar{p}_t(0, x)/\bar{p}(0, x) \tag{6}$$

and the left-hand side belongs to  $L^p_{loc}(\Omega)$ .

The characterisation of  $\bar{u}_0$  with the help of  $p$  is almost complete here, except on the singular arc  $\omega = \{0 < u_0 < 1\}$ . Note first that this singular arc might have a positive measure. It was even proved in [36] that, if  $f$  is concave, then  $\omega \equiv \Omega$ . If  $f'$  is monotonic, equation (6) admits a unique solution and thus fully characterizes  $\bar{u}_0$ . But for a bistable nonlinearity  $f_\theta(u) = u(1 - u)(u - \theta)$ , equation (6) might have two solutions, one belonging to  $[0; \eta]$  and the other to  $(\eta; 1]$ , where  $\eta = \eta(\theta) \in (0; 1)$  is the unique real number such that  $f_\theta$  is convex in  $[0; \eta)$  and concave in  $(\eta; 1]$ . It is then necessary to distinguish between these two possible roots in order to completely characterize  $\bar{u}_0$  with  $p$ .

From the numerical point of view, the characterisation given by this result naturally leads to a gradient descent algorithm, which is not well-posed if we are not able to characterise  $\bar{u}_0$  on  $\omega$ . Let us briefly describe this algorithm, which we detail further in section 4 of the present paper, to explain the core difficulty and how our theoretical results enable us to bypass it: starting from an initial configuration  $u_0^0$ , we seek to improve it to obtain a better admissible candidate  $u_1^0$ . We first compute the adjoint state  $p_0^0$  associated with  $u_0^0$ . The problem arises if  $p_0^0$  has ‘flat zones’, in other words if there exists  $c_0$  (necessarily unique) such that

$$|\{p_0^0 > c_0\}| < m, \quad |\{p_0^0 \geq c_0\}| > m, \quad |\{p_0^0 = c_0\}| > 0.$$

On the set  $\{p_0^0 > c_0\}$ , we replace  $u_0^0$  by 1, while on  $\{p_0^0 < c_0\}$  we replace  $u_0^0$  by zero. On  $\{p_0^0 = c_0\}$  we must replace  $u_0^0$  with a root of (6). (6) can have two roots  $\mu_1^0$  and  $\mu_2^0$ . These roots can be distinguished by the convexity of  $f$ : up to a relabelling,  $f''(\mu_1^0) > 0$  and  $f''(\mu_2^0) \leq 0$ . In [36], the two possibilities were explored successively, which led to high computational costs. This was the main limitation of the numerical approach of [36]. Theorem 1 of the present paper shows that one should choose  $\mu_2^0$ . This significantly improves the running time of our algorithm and we refer to section 4 for examples.

**1.2.2. Related works.** A related problem has first been addressed by Garnier et al in [15], where the authors consider a bistable reaction term  $f(u) := u(1 - u)(u - \theta)$ , with  $\theta \in (0, 1)$ , over the full line  $\Omega = \mathbb{R}$ . In this earlier paper, the authors did not investigate  $(\mathbf{P}_f)$ , but they tried to optimize the initial datum in order to ensure the convergence to  $u \equiv 1$  when  $t \rightarrow +\infty$ . They investigated numerically the particular case  $u_0 := \mathbb{1}_{(-\alpha - \frac{\theta}{2}, -\alpha)} + \mathbb{1}_{(\alpha, \alpha + \frac{\theta}{2})}$ , and proved that in some situations the initial datum associated with  $\alpha = 0$  might lead to extinction (that is,  $u(t, x) \rightarrow 0$  as  $t \rightarrow +\infty$ ), while a positive  $\alpha > 0$  might lead to persistence (that is,  $u(t, x) \rightarrow 1$  as  $t \rightarrow +\infty$ ). Also, numerics for more general classes of initial datum indicate that fragmentation might favor species persistence. Hence, even if the problem we consider here is a bit different, we expect the maximiser to be fragmented, that is, non-smooth, for bistable nonlinearities.

More recently, this problem was also addressed in [24], in a slightly more general form, and for the criterion  $\int_{\Omega} |1 - u(T, x)|^2 dx$ . These authors investigated in particular various conditions ensuring that the maximiser  $\bar{u}_0$  is constant with respect to  $x$ , and, reciprocally, that the constant initial datum is a local maximiser.

**1.3. Main results of the paper**

The main contributions of this paper are the following:

- When  $u_0$  does not saturate the  $L^\infty$  constraints (i.e. when the set  $\omega := \{0 < u_0 < 1\}$  has positive measure), we prove in theorem 1 that any maximiser  $\bar{u}_0$  must necessarily be in a zone of concavity of  $f$ : in  $\omega$ ,  $f''(\bar{u}_0) \leq 0$ .
- When  $f$  is convex, theorem 2 characterizes explicitly a global maximiser, using rearrangement techniques. Here, the presence of Neumann boundary conditions prohibits using in a straightforward manner the results of [6], and we need to adapt some points of the proof to this case.
- When  $f$  is a bistable nonlinearity, we improve the algorithm initially introduced in [36] and display several numerical simulations. One-dimensional simulations are displayed that exemplify the fact that theorem 1 significantly improves the computational time of optimisation algorithms. We also provide two-dimensional simulations.

**1.3.1. Characterisation of the singular arc.** Let us first recall the expression of the second order derivative [36]:

$$\langle \nabla^2 \mathcal{J}_T(u_0), h_0 \rangle = \int \int_{(0;T) \times \Omega} f''(u(t, x)) p(t, x) h^2(t, x) dx dt, \tag{7}$$

where  $p$  solves (5) and  $h$  solves

$$\begin{cases} \frac{\partial h}{\partial t} - \Delta h = f'(u)h & \text{in } (0; T) \times \Omega, \\ h(0, x) = h_0(x), \\ \frac{\partial h}{\partial \nu} = 0 & \text{on } (0; T) \times \partial\Omega. \end{cases} \tag{8}$$

We can now state our main result:

**Theorem 1.** *Assume  $\Omega = (0; \pi)$ . Let  $\bar{u}_0$  be a solution of  $(\mathbf{P}_f)$ . If the set  $\Omega_{\bar{c}} := \{x \in \Omega : 0 < \bar{u}_0(x) < 1\}$  has a positive measure then, for almost every interior point  $x$  of  $\Omega_{\bar{c}}$ , there holds*

$$f''(\bar{u}_0(x)) \leq 0. \tag{9}$$

**Remark 1.** The method we put forth is reminiscent of one that was used in [36] to study the case of a constant initial condition and to prove that such a constant  $u_0$  was always a maximiser in the case of monostable nonlinearities. Here, working with interior point greatly complexifies the situation and calls for two-scale asymptotic expansions.

The main drawback to our approach is that it cannot cover the case of singular (e.g. Cantor-like) singular arcs, and it is a very interesting question to prove that such a property holds for any point of the singular arc  $\Omega_{\bar{c}}$ . We comment on the main difficulties of this approach in the conclusion and only state, for the moment, that the main problem is related to the ubiquitous problem of separation of phases in homogenisation [2].

**1.3.2. Convex nonlinearities and rearrangements.** We present, in this section, a characterisation of maximisers when the nonlinearity  $f$  is convex, using rearrangement and symmetrisation techniques. It should be noted that, since we are working with Neumann boundary conditions, it is not possible to use directly well-known parabolic isoperimetric inequalities [5, 6, 34]. We refer to [6, 28, 41] and the references therein for an introduction to parabolic isoperimetric inequalities, and only underline here that the most precise results available in the literature only encompass the case of Dirichlet boundary conditions. For Neumann boundary conditions, a large literature [11, 14, 25] is devoted to such questions. Usually, it involves a comparison of the solution with the solution of a Dirichlet, or of a mixed Dirichlet–Neumann problem, and makes use of constants appearing in relative isoperimetric inequalities. It is unclear whether these comparison results could be used in our case. Since we are working in the one-dimensional case, a direct adaptation of the proof of [6] yields the required results.

In order to state our result, let us introduce the following notation: let  $\tilde{u}_0$  be defined as

$$\tilde{u}_0 := \mathbb{1}_{(0;m)} \in \mathcal{A}. \tag{10}$$

**Theorem 2.** *Assume  $f$  is a convex,  $\mathcal{C}^1$  function such that  $f(0) = 0$ . Then  $\tilde{u}_0$  is a solution of  $(\mathbf{P}_f)$ .*

**Remark 2.**

- It should be noted that  $\tilde{u}_0$  appears as the solution of many other optimisation problems in population dynamics, most specifically for the monostable case [8, 17, 30] with Neumann boundary conditions.



- The maximiser  $\mathbb{1}_{(0;m)}$  is clearly not unique, since  $\mathbb{1}_{(1-m;1)}$  is also a maximiser for example. This provides an example of non-uniqueness of the maximiser.
- As a corollary,  $\tilde{u}_0 := \mathbb{1}_{(0;m)}$  is a minimiser of  $\mathcal{J}_T$  over  $\mathcal{A}$  if  $f$  is concave.

## 2. Proof of theorem 1

### 2.1. Notations, plan of the proof and first simplification

**Second order optimality conditions.** We recall the expression of the second order derivative of  $\mathcal{J}$ : for an admissible perturbation  $h_0$ , we have

$$\langle \nabla^2 \mathcal{J}_T(u_0), h_0 \rangle = \int_0^T \int_{\Omega} f''(u(t, x)) p(t, x) h^2(t, x) dx dt,$$

where  $p$  solves (5) and  $h$  solves (8). Let  $\bar{u}_0$  be a solution of  $(\mathbf{P}_f)$ . We assume that the set  $\Omega_{\bar{c}} := \{0 < \bar{u}_0 < 1\}$  has a non-empty interior and we want to prove that  $f''(\bar{u}_0) \leq 0$  almost everywhere on the interior of this set. To do this, we need the following expression of second order optimality conditions:

**Lemma 1.** For every  $h_0 \in L^\infty(\Omega)$  supported in  $\Omega_{\bar{c}}$ , such that  $\int_{\Omega} h_0 = 0$ , there holds

$$\iint_{(0;T) \times \Omega} f''(u(t, x)) p(t, x) h^2(t, x) dx dt \leq 0 \tag{11}$$

where  $h$  is the solution of (8) associated with the initial condition  $h_0$ .

**Proof of lemma 1.** We first notice the following thing: let, for any  $n \in \mathbb{N}^*$ , the set  $F_n$  be defined as

$$F_n := \left\{ \frac{1}{n} < u_0 < 1 - \frac{1}{n} \right\}.$$

Then, for any  $L^\infty$  function  $h_0$  supported in  $F_n$  (in the sense that  $h_0 \mathbb{1}_{F_n} = h_0$ ) such that  $\int_{\Omega} h_0 = 0$ , if  $h$  is the solution of (8) associated with  $h_0$ , we have

$$\iint_{(0;T) \times \Omega} f''(u(t, x)) p(t, x) h^2(t, x) dx dt \leq 0.$$

This is a consequence of the fact that, for any  $\tau$  such that  $|\tau|$  is small enough,  $\bar{u}_0 + \tau h_0$  is an admissible initial condition.

Let us now consider  $h_0 \in L^\infty(\Omega_{\bar{c}})$  satisfying  $\int_{\Omega} h_0 = 0$ . We define, for any  $n \in \mathbb{N}$ ,

$$h_{0,n} := \mathbb{1}_{F_n} \left( h_0 - \frac{\int_{F_n} h_0}{|F_n|} \right).$$

For every  $n \in \mathbb{N}$ ,  $h_{0,n}$  is supported in  $F_n$  and verifies  $\int_{\Omega} h_{0,n} = 0$ . Hence, defining, for any  $n \in \mathbb{N}$ ,  $h_n$  as the solution of (8) associated with the initial condition  $h_{0,n}$  we have

$$\iint_{(0;T) \times \Omega} f''(u(t, x)) p(t, x) h_n^2(t, x) dx dt \leq 0. \tag{12}$$

However, there holds

$$h_{0,n} \xrightarrow{n \rightarrow \infty} h_0 \quad \text{in } L^2(\Omega),$$

which, by standard parabolic estimates, entails

$$h_n \xrightarrow{n \rightarrow \infty} h \quad \text{in } L^2((0; T) \times \Omega)$$

where  $h$  is the solution of (8) associated with the initial condition  $h_0$ . Passing to the limit  $n \rightarrow \infty$  in (12) yields the conclusion.  $\square$

Since we want to retrieve, from the second order optimality conditions (11), an information of  $f''(\bar{u}_0)$ , we need to find a perturbation  $h_0$  such that the ensuing solution  $h$  satisfies, roughly speaking,

$$h^2(t, x) \approx C\delta_{t=t_0}g(x),$$

for a certain function  $g$ . As  $h$  is a solution of a parabolic equation, one possibility to obtain such a behaviour is to choose a highly oscillating initial condition, say  $h_0(x) = \cos(kx)$  with a large integer  $k$ . This would give  $h^2(t, x) \approx e^{-tk^2} \cos(kx)^2$ , which, thanks to the Laplace method, does concentrate around  $t = 0$  up to a proper rescaling. This however is not particularly convenient, as such a perturbation is not admissible: it is not supported in  $\Omega_{\bar{c}}$ . To overcome this difficulty, we need to truncate such highly oscillating perturbations, thus choosing a perturbation of the form  $\theta(x)\cos(kx)$ , with  $\theta$  a cut-off function, leading to two-scale asymptotic expansions. The objective is to pick the correct function  $\theta$ .

In order to summarise our approach, let us fix notations: we pick an optimiser  $\bar{u}_0$ , we define  $\Omega_{\bar{c}} := \{0 < \bar{u}_0 < 1\}$ , and we set  $\mathring{\Omega}_{\bar{c}}$  as the interior of  $\Omega_{\bar{c}}$ . To prove theorem 1 we argue by contradiction: assume that, for some  $\delta > 0$ ,

$$|\{f''(\bar{u}_0) \geq \delta\} \cap \mathring{\Omega}_{\bar{c}}| > 0. \tag{13}$$

Since  $\mathring{\Omega}_{\bar{c}}$  is an open set, we can write it as a union of intervals

$$\mathring{\Omega}_{\bar{c}} = \bigcup_{k=0}^{\infty} (a_k; b_k). \tag{14}$$

By (13), there exists  $n_0 \in \mathbb{N}$  such that

$$|\{f''(\bar{u}_0) \geq \delta\} \cap (a_{n_0}; b_{n_0})| > 0,$$

so that there exists  $\epsilon > 0$  such that, for the same  $n_0$ , we have

$$|\{f''(\bar{u}_0) \geq \delta\} \cap (a_{n_0} + \epsilon; b_{n_0} - \epsilon)| > 0. \tag{15}$$

We fix such an  $\epsilon > 0$ .

To alleviate notations, define  $E := \{f''(\bar{u}_0) \geq \delta\} \cap (a_{n_0} + \epsilon; b_{n_0} - \epsilon)$ . As  $p(0, \cdot) > 0$  by the parabolic maximum principle, (15) yields

$$\int_{\Omega} f''(\bar{u}_0(\cdot))p(0, \cdot)\mathbb{1}_E > 0. \tag{16}$$

We approximate in  $L^1(a_{n_0}; b_{n_0})$  the function  $\mathbb{1}_E$  by a sequence  $\{\psi_k\}_{k \in \mathbb{N}}$  of uniformly bounded, non-negative,  $\mathcal{C}^\infty$  functions that are compactly supported in  $(a_{n_0}; b_{n_0}) \subset \overset{\circ}{\Omega}_{\bar{c}}$ . In particular, the sequence  $\{\psi_k^2\}_{k \in \mathbb{N}}$  also converges to  $\mathbb{1}_E$  in  $L^1(\Omega)$ , so that (16) implies that for  $K$  large enough

$$\int_{\Omega} f''(\bar{u}_0(\cdot))p(0, \cdot)\psi_K^2 > 0. \tag{17}$$

We fix  $K$  large enough so that (17) holds and we set, for this index  $K$ ,

$$\theta := \psi_K \in \mathcal{C}^\infty(\Omega).$$

The sequence of truncated, highly oscillating initial conditions is

$$\bar{h}_{k,0} := \theta(\cdot) (\cos(k\cdot) + \alpha_k)$$

where

$$\alpha_k = -\frac{\int_{\Omega} \theta \cos(k\cdot)}{\int_{\Omega} \theta} \tag{18}$$

simply ensures that  $\int_{\Omega} \bar{h}_{k,0} = 0$ . This constant does not play a role in the upcoming analysis for the following reason:

- (a) First, by setting  $h_{k,0} := \theta(x)\cos(kx)$  and by defining  $h_k$  as the solution of (8) associated with the initial condition  $h_k^0$  we shall show that

$$\iint_{(0;T) \times \Omega} f''(u)ph_k^2 \underset{k \rightarrow \infty}{\sim} \frac{C}{k^2} \int_{\Omega} f''(\bar{u}_0(\cdot))p(0, \cdot)\theta^2(\cdot) \tag{19}$$

for some constant  $C$ . This is the core of the proof, and will take up the remainder of this section of the paper.

It is also immediate by parabolic regularity to obtain that the sequence  $\{h_k\}_{k \in \mathbb{N}}$  is uniformly bounded in  $L^2((0; T) \times \Omega)$ .

- (b) Second we observe that, as  $\theta \in \mathcal{C}^4$ , the Riemann–Lebesgue lemma in particular ensure that  $\alpha_k = \mathcal{O}\left(\frac{1}{k^4}\right)$ .
- (c) If we now set  $z$  as the solution of (8) associated with the initial condition  $\theta(\cdot)$ , the solution  $\bar{h}_k$  associated with the (admissible) initial condition  $\bar{h}_{k,0}$  is given by  $\bar{h}_k = h_k + \alpha_k z$ . Then the second order derivative in the admissible direction  $\bar{h}_{k,0}$  is given by

$$\begin{aligned} \langle \nabla^2 \mathcal{J}_T(u_0), \bar{h}_{k,0} \rangle &= \iint_{(0;T) \times \Omega} f''(u(t, x))p(t, x)(\bar{h}_k)^2(t, x)dx dt \\ &= \iint_{(0;T) \times \Omega} f''(u(t, x))p(t, x)h_k^2(t, x)dx dt \\ &\quad + 2\alpha_k \iint_{(0;T) \times \Omega} f''(u(t, x))p(t, x)z(t, x)h_k(t, x)dx dt \\ &\quad + \alpha_k^2 \iint_{(0;T) \times \Omega} f''(u(t, x))p(t, x)z^2(t, x)dx dt. \end{aligned}$$

Taking into account (19) and the fact that  $\alpha_k = \mathcal{O}\left(\frac{1}{k^4}\right)$  leads to

$$\langle \nabla^2 \mathcal{J}_T(u_0), \bar{h}_{k,0} \rangle \underset{k \rightarrow \infty}{\sim} \frac{C}{k^2} \int_{\Omega} f''(\bar{u}_0(\cdot))p(0, \cdot)\theta^2(\cdot) > 0,$$

a contradiction.

(d) As a consequence, the theorem is proved, provided we can prove (19), and we henceforth focus on this point.

2.2. Asymptotic expansion of  $h_k$

Let  $\theta$  be given as above. We consider the following sequence of equations: let, for any  $k \in \mathbb{IN}$ ,  $h_k$  be the solution of

$$\begin{cases} \frac{\partial h_k}{\partial t} - \frac{\partial^2 h_k}{\partial x^2} = f'(u)h_k, \\ \frac{\partial h_k}{\partial \nu} = 0, \\ h_k(0, x) = h_{k,0}(x) = \theta(x) \cos(kx). \end{cases} \tag{20}$$

In this context, it is natural [1] to look for a two-scale asymptotic expansion of  $h_k$  of the form

$$h_k(t, x) \approx h_k^0(k^2t, x, kx) + \frac{1}{k}h_k^1(k^2t, x, kx) + \dots \tag{21}$$

which, after a formal identification at the first and second order, gives the following equations on  $h_k^0$  and  $h_k^1$ :

$$\begin{cases} \frac{\partial h_k^0}{\partial s} - \frac{\partial^2 h_k^0}{\partial y^2} = 0, \\ \frac{\partial h_k^0}{\partial \nu} = 0, \\ h_k^0(0, x, y) = \theta(x) \cos(y). \end{cases} \tag{22}$$

and

$$\begin{cases} \frac{\partial h_k^1}{\partial s} - \frac{\partial^2 h_k^1}{\partial y^2} = 2\frac{\partial^2 h_k^0}{\partial x \partial y}, \\ \frac{\partial h_k^1}{\partial \nu} = 0, \\ h_k^1(0, x, y) = 0. \end{cases} \tag{23}$$

Equation (22) can be solved explicitly, giving

$$h_k^0(s, x, y) = \theta(x) \cos(y)e^{-s}. \tag{24}$$

This, in turn, allows to solve equation (23) as

$$h_k^1(s, x, y) = -2s e^{-s} \theta'(x) \sin(y). \tag{25}$$

**Proposition 1.** *The asymptotic expansion (21) is valid in  $L^2(\Omega)$  in the following sense: there exists  $M > 0$  that depends on the time horizon  $T$  such that, if we define*

$$R_k := h_k(t, x) - h_k^0(k^2t, x, kx) - \frac{1}{k}h_k^1(k^2t, x, kx)$$

then, for any  $t \in (0; T)$ ,

$$\|R_k(t, \cdot)\|_{L^2(\Omega)} \leq \frac{M}{k^2}. \tag{26}$$

In particular,

$$\iint_{(0;T) \times \Omega} R_k^2 \leq \frac{M^2}{k^4}, \quad \int_0^T \|R_k\|_{L^2(\Omega)} \leq \frac{MT}{k^2}. \tag{27}$$

**Proof of proposition 1.** To prove this proposition, we write down the equation satisfied by  $R_k$ . Straightforward computations show that  $R_k$  solves

$$\partial_t R_k - \Delta R_k - f'(u)R_k := f'(u) \left( h_k^0 + \frac{1}{k} h_k^1 \right) + \frac{\partial^2 h_k^0}{\partial x^2} + 2 \frac{\partial^2 h_k^1}{\partial x \partial y} + \frac{1}{k} \frac{\partial^2 h_k^1}{\partial x^2}, \tag{28}$$

and all the functions on the right-hand side are evaluated at  $(k^2 t, x, kx)$  (we dropped this for notational convenience). We now introduce the following notations:

$$\begin{cases} W_0 := f'(u), \\ V_{0,k}(t, x) := h_k^0(k^2 t, x, kx) + \frac{1}{k} h_k^1(k^2 t, x, kx), \\ V_{1,k} := -\frac{\partial^2 h_k^0}{\partial x^2}(k^2 t, x, kx), \\ V_{2,k} := -2 \frac{\partial^2 h_k^1}{\partial x \partial y} - \frac{1}{k} \frac{\partial^2 h_k^1}{\partial x^2}. \end{cases}$$

First of all, since  $0 \leq u \leq 1$  and  $f \in \mathcal{C}^1$ , there exists  $M_0 > 0$  such that

$$\|W_0\|_{L^\infty((0;T) \times \Omega)} \leq M_0. \tag{29}$$

We gather the main estimates on source terms in the following lemma:

**Lemma 2.** *There exists  $\tilde{M} > 0$  such that*

$$\int_0^T \|V_{0,k}(t, \cdot)\|_{L^2(\Omega)} \leq \frac{\tilde{M}}{k^2}, \tag{30}$$

$$\int_0^T \|V_{1,k}(t, \cdot)\|_{L^2(\Omega)} dt \leq \frac{\tilde{M}}{k^2}. \tag{31}$$

$$\int_0^T \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)} dt \leq \frac{\tilde{M}}{k^2}. \tag{32}$$

**Proof of lemma 2.** We prove the three estimates separately. Let us recall the following consequence of the Laplace method: for any integer  $m \in \mathbb{N}^*$ , one has

$$\int_0^T t^{m-1} e^{-k^2 t} dt \underset{k \rightarrow \infty}{\sim} \frac{(m-1)!}{k^{2m}}. \tag{I}_m$$

**Proof of (30)**

By the triangle inequality we get, for any  $t \in (0; T)$ ,

$$\|V_{0,k}(t, \cdot)\|_{L^2} \leq \|h_k^0(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)} + \frac{1}{k} \|h_k^1(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)}.$$

We first use the explicit expressions (24) and (25) for  $h_k^0$  and  $h_k^1$  to obtain, using  $\|\theta\|_{L^\infty} \leq 1$ ,

$$\|h_k^0(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)}^2 = \int_{\Omega} \theta(x)^2 \cos(kx)^2 e^{-2k^2t} dx \leq e^{-2k^2t} |\Omega|, \tag{33}$$

and integrating this inequality between 0 and  $T$  gives

$$\int_0^T \|h_k^0(k^2t, x, k \cdot)\|_{L^2(\Omega)} dt \leq \frac{\tilde{M}}{k^2}.$$

In the same way, we have, for any  $t \in (0; T)$ ,

$$\|h_k^1(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)}^2 = 4k^4 t^2 e^{-2k^2t} \int_{\Omega} (\theta'(x))^2 \sin(x)^2 dx \leq Ck^4 t^2 e^{-2k^2t} |\Omega| \cdot \|\theta'\|_{L^\infty}^2, \tag{34}$$

for some constant  $C$ . Taking the square root and integrating in time we get, for a constant  $C'$ ,

$$\frac{1}{k} \int_0^T \|h_k^1(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)} dt \leq C' |\Omega| \cdot \|\theta'\|_{L^\infty} k \int_0^T t e^{-k^2t} dt.$$

Using  $(\mathbf{I}_m)$  with  $m = 2$  gives

$$\int_0^T \frac{1}{k} \|h_k^1(k^2t, \cdot, k \cdot)\|_{L^2(\Omega)} dt \leq \frac{M_1}{k^3}$$

for some constant  $M_1$ .

Summing these two contributions gives (30).

**Proof of (31).** This follows from the same arguments, by simply observing that

$$V_{1,k}(t, x) = -e^{-k^2t} \theta''(x) \cos(kx).$$

**Proof of (32).** We once again split the expression and estimate separately

$$\int_0^T \left\| \frac{\partial^2 h_k^1(k^2t, \cdot, k \cdot)}{\partial x \partial y} \right\|_{L^2(\Omega)} dt \quad \text{and} \quad \frac{1}{k} \int_0^T \left\| \frac{\partial^2 h_k^1(k^2t, \cdot, k \cdot)}{\partial x^2} \right\|_{L^2(\Omega)} dt.$$

We first observe that for any  $t \in (0; T)$ , we have

$$\frac{\partial^2 h_k^1(k^2t, \cdot, k \cdot)}{\partial x \partial y} = -4k^2 t e^{-k^2t} \theta''(x) \cos(y).$$

In particular, for any  $t \in (0; T)$

$$\left\| \frac{\partial^2 h_k^1(k^2t, \cdot, k \cdot)}{\partial x \partial y} \right\|_{L^2(\Omega)} \leq 4k^2 t e^{-k^2t} |\Omega| \cdot \|\theta''\|_{L^\infty}$$

so that the Laplace method  $(\mathbf{I}_m)$  with  $\lambda = 2$  gives the bound

$$\int_0^T \left\| \frac{\partial^2 h_k^1(k^2t, \cdot, k \cdot)}{\partial x \partial y} \right\|_{L^2(\Omega)} dt \leq \frac{M_2}{k^2}$$

for some constant  $M_2$ . The proof of the control of the second term follows along exactly the same lines.  $\square$

Let us now prove estimate (26). The equation on  $R_k$  rewrites

$$\partial_t R_k - \Delta R_k - W_0 R_k = W_0 V_{0,k} + V_{1,k} + V_{2,k}. \tag{35}$$

Multiplying the equation by  $R_k$  and integrating by parts in space gives

$$\begin{aligned} \frac{1}{2} \partial_t \int_{\Omega} R_k^2 + \int_{\Omega} |\nabla R_k|^2 - \int_{\Omega} W_0 R_k^2 &\leq \|R_k\|_{L^2(\Omega)} \\ &\times (M_0 \|V_{0,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)}). \end{aligned}$$

In other words, bounding  $W_0$  by  $M_0$  and defining  $g(t) := \|R_k(t, \cdot)\|_{L^2(\Omega)}^2$  we obtain

$$\frac{1}{2} g'(t) \leq M_0 g(t) + \sqrt{g(t)} (M_0 \|V_{0,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)}).$$

Furthermore,  $R_k(0, \cdot) = 0$ . We thus obtain, by the Gronwall lemma, for any  $t \in (0; T)$ ,

$$\sqrt{g(t)} e^{-M_0 t} \leq \int_0^t e^{-M_0 s} (M_0 \|V_{0,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)}) ds.$$

Hence, by lemma 2 we get for some constant  $N_0$  and any  $t \in (0; T)$ ,

$$\|R_k(t, \cdot)\|_{L^2(\Omega)} \leq N_0 \int_0^T (\|V_{0,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)} + \|V_{2,k}(t, \cdot)\|_{L^2(\Omega)}) dt \leq \frac{N_0 \tilde{M}}{k^2}.$$

$\square$

### 2.3. Back to the proof

We turn back to the proof of theorem 1 and, more precisely, to the proof of (19).

**Proof of theorem 1.** We use the same  $\theta$  as above and the same notation  $\alpha_k$  as in the introduction of the proof (equation (18)). Let us now consider the initial perturbation  $\bar{h}_{k,0} := \theta(x)(\cos(kx) + \alpha_k)$ . We recall that  $z$  is the solution of (8) with initial condition  $\theta$  and that  $h_k$  is the solution of (8) with initial condition  $\theta(\cdot)\cos(k\cdot)$ . Hence, since the equation is linear,

$$\bar{h}_k = h_k + \alpha_k z.$$

By parabolic regularity,

$$\sup \left( \sup_{k \in \mathbb{N}} \|h_k\|_{L^2((0;T) \times \Omega)}, \|z\|_{L^2((0;T) \times \Omega)} \right) < \infty. \tag{36}$$

Since  $\theta \in \mathcal{C}^4$ , the Riemann–Lebesgue lemma implies

$$\alpha_k = \mathcal{O}_{k \rightarrow \infty} \left( \frac{1}{k^4} \right). \tag{37}$$

We define

$$F(t, x) := f''(\bar{u}(t, x))\bar{p}(t, x)$$

so that the second order derivative of  $\mathcal{J}_T$  in  $\bar{u}_0$  rewrites as

$$\begin{aligned} \langle \nabla^2 \mathcal{J}_T(u_0), \bar{h}_{k,0} \rangle &= \iint_{(0;T) \times \Omega} F \bar{h}_k^2 \\ &= \iint_{(0;T) \times \Omega} F h_k^2 + 2\alpha_k \iint_{(0;T) \times \Omega} F z h_k + \alpha_k^2 \iint_{(0;T) \times \Omega} F z^2. \end{aligned} \tag{38}$$

We focus on the first term:

$$\begin{aligned} \iint_{(0;T) \times \Omega} F h_k^2 &= \iint_{(0;T) \times \Omega} F(t, x) (R_k(t, x) + V_{0,k}(t, x))^2 \, dx \, dt, \\ &= \iint_{(0;T) \times \Omega} F(t, x) (R_k(t, x)^2 + 2R_k(x, t)V_{0,k}(t, x) + V_{0,k}(t, x)^2) \, dx \, dt. \end{aligned} \tag{39}$$

From the assumptions on  $f$  and the estimates on  $u$  and  $p$ , it easy to see that

$$\|F\|_{L^\infty((0;T) \times \Omega)} \leq M_3. \tag{40}$$

Gathering (40) and (27) it follows that

$$\iint_{(0;T) \times \Omega} F(t, x) R_k(t, x)^2 \, dx \, dt = \mathcal{O}(k^{-4}) \tag{41}$$

and similarly, gathering (30) and (26) we obtain

$$\iint_{(0;T) \times \Omega} F(t, x) R_k(t, x) V_{0,k}(t, x) \, dx \, dt = \mathcal{O}(k^{-4}) \tag{42}$$

Let us now study the term

$$\begin{aligned} \iint_{(0;T) \times \Omega} F(t, x) V_{0,k}(t, x)^2 \, dx \, dt &= \iint_{(0;T) \times \Omega} F(t, x) \left( h_k^0(t, x) + \frac{1}{k} h_k^1(t, x) \right)^2 \\ &= \iint_{(0;T) \times \Omega} F(t, x) \left( h_k^0(t, x)^2 + 2\frac{1}{k} h_k^0(t, x) h_k^1(t, x) \right. \\ &\quad \left. + \frac{1}{k^2} h_k^1(t, x)^2 \right) \, dx \, dt \end{aligned}$$

Once again we split the expression. Applying the Cauchy–Swchartz inequality, and using estimates (33) and (34) it follows that the second term verifies

$$\begin{aligned} \left| 2\frac{1}{k} \iint_{(0;T) \times \Omega} F(t, x) h_k^0(t, x) h_k^1(t, x) \, dx \, dt \right| &\leq 2\frac{M_3}{k} \int_0^T \|h_k^0\|_{L^2(\Omega)} \|h_k^1\|_{L^2(\Omega)} \, dt \\ &\leq 2M_4 k \int_0^T t e^{-2k^2 t} \, dt, \\ &= \mathcal{O}\left(\frac{1}{k^3}\right). \end{aligned} \tag{43}$$



The last step in the above expression follows directly from  $(\mathbf{I}_m)$  with  $m = 2$ .

We obtain in a similar way the following estimate on the third term:

$$\begin{aligned} \left| \frac{1}{k^2} \iint_{(0;T) \times \Omega} F(t, x) h_k^1(t, x)^2 \, dx \, dt \right| &\leq \frac{M_3}{k^2} \int_0^T \|h_k^1\|_{L^2(\Omega)}^2 \, dt \leq M_4 k^2 \int_0^T t^2 e^{-2k^2 t} \, dt, \\ &= \mathcal{O}\left(\frac{1}{k^4}\right). \end{aligned} \tag{44}$$

In this case we applied  $(\mathbf{I}_m)$  for  $m = 3$ .

Finally, let us study the first term which can be written as

$$\iint_{(0;T) \times \Omega} F(t, x) h_k^0(t, x)^2 \, dx \, dt = \int_0^T e^{-2k^2 t} G(t) \, dt \tag{45}$$

where  $G(t) := \int_{\Omega} F(t, x) \theta(x)^2 \cos(kx)^2 \, dx$  is a continuous function of time as a consequence of parabolic regularity.

However,  $\theta$  was chosen so that

$$\int_{\Omega} f''(\bar{u}_0(\cdot)) p(0, \cdot) \theta^2 > 0.$$

As  $\cos(k \cdot)^2 = \frac{1}{2} (1 + \cos(2 \cdot)) \xrightarrow[k \rightarrow \infty]{} \frac{1}{2}$ , it follows that for any  $k$  large enough  $G(0) > 0$ . Furthermore, from the Laplace method, when  $k \rightarrow \infty$ , one has

$$\iint_{(0;T) \times \Omega} F(t, x) h_k^0(t, x)^2 \, dx \, dt \sim \frac{1}{2k^2} G(0). \tag{46}$$

Gathering the estimates in (41), (42), (43), (44), (46) and plugging them into the second derivative of the functional  $\mathcal{J}_T$  given by (38) it follows that

$$\iint_{(0;T) \times \Omega} F h_k^2 \underset{k \rightarrow \infty}{\sim} \frac{1}{2k^2} G(0). \tag{47}$$

We go back to (38). By (36), (37) and by (47) we have

$$\langle \nabla^2 \mathcal{J}_T(u_0), \bar{h}_{k,0} \rangle = \iint_{(0;T) \times \Omega} F h_k^2 + \mathcal{O}\left(\frac{1}{k^4}\right) \underset{k \rightarrow \infty}{\sim} \frac{1}{2k^2} G(0). \tag{48}$$

This means that for  $k$  sufficiently large,

$$k^2 \langle \nabla^2 \mathcal{J}_T(u_0), \bar{h}_{k,0} \rangle > 0$$

which contradicts the fact that  $\bar{u}_0$  is a maximiser of  $\mathcal{J}_T$ . The proof of the theorem is complete. □

### 3. Proof of theorem 2

The proof follows essentially from the same arguments as in [6]. We thus only present the main steps that are in order so as to apply the methods of [6]. We define  $g(u) := f(u) + cu$ , with  $c > \|f'\|_{L^\infty}$ , so that  $g$  is increasing.

**Reduction to a bang-bang maximiser.** We recall that bang-bang functions are defined as characteristic functions of subsets of  $\Omega$ , that is, functions only taking values 0 and 1. First, as  $f$  is convex, it follows from the same arguments as proposition 6 of [36] that  $\mathcal{J}_T$  is convex. Hence, one can restrict to maximisers among the extremal points of  $\mathcal{A}$ . These are exactly bang-bang function:  $u_0$  satisfies  $u_0 = 0$  or 1 almost everywhere on  $(0, \pi)$ .

**Reduction to a periodic problem.** Next, for all  $t \in [0, T]$ , we extend  $u(t, \cdot)$  to  $(-\pi; \pi)$  by symmetrisation with respect to 0, and we then extend it to  $\mathbb{R}$  by  $2\pi$ -periodicity. The Neumann boundary conditions at  $x = 0$  and  $x = \pi$  ensure that the extended function is of class  $C^1$ , and it thus satisfies the equation on the torus:

$$\begin{cases} \frac{\partial u}{\partial t} - \Delta u + cu = g(u) & \text{in } (0; T) \times \mathbb{T}, \\ u(0, \cdot) = u_0. \end{cases} \tag{49}$$

Let us also recall some basic facts about rearrangements.

**Periodic rearrangements.** We recall the definition of the periodic rearrangement: for any periodic function  $u : \mathbb{T} \rightarrow \mathbb{R}_+$  if we identify  $\mathbb{T}$  with  $[-\pi; \pi]$  there exists a unique symmetric (with respect to 0) non-increasing function  $u^* : \mathbb{T} \rightarrow \mathbb{R}$  that has the same distribution function as  $u$ .  $u^*$  is called the periodic rearrangement of  $u$ . We recall that the distribution function of  $u$  is

$$\mu_u(t) := \text{Vol}(\{u \geq t\})$$

and that  $u^*$ , the periodic rearrangement of  $u$ , is the left inverse of  $\mu_u$ .

**Proposition 2.** *Let  $v$  the solution of (49) associated with the initial datum  $u_0^*$ . Then*

$$\forall t \in (0; T), \quad \forall r \in (0; \pi), \quad \int_{-r}^r v(t, x)dx \geq \int_{-r}^r u^*(t, x)dx.$$

*In particular, taking  $t = T$  and  $r = \pi$ :*

$$\int_{-\pi}^{\pi} v(T, x)dx \geq \int_{-\pi}^{\pi} u^*(T, x)dx = \int_{-\pi}^{\pi} u(T, x)dx. \tag{50}$$

As explained in the introduction, proposition 2 follows simply by adapting minor points in the proofs of [6], and so we will simply indicate the main steps. The core idea is the following: the comparison results and Talenti-type inequalities one finds in the rearrangement literature rely on integrating the solution of the equation we are working with on its level sets and using the isoperimetric inequality. Thus, these proofs generally work only in the case of Dirichlet boundary conditions (for a recent, analogous result in the case of Robin boundary conditions we refer to [4]), as these conditions guarantee that, if the solution is non-negative, none of its level sets intersects the boundary of the domain. However, in our case, since the Neumann boundary conditions allow to symmetrise the solution  $u$  and to obtain a solution on the torus  $\mathbb{T}$ , the boundary of the domain is empty and the isoperimetric inequality holds, so that the proofs are identical.

**Proof of proposition 2.** For the sake of readability, we break down the main steps in establishing the inequality

$$\forall t \in (0; T), \quad \forall r \in (0; \pi), \quad \int_{-r}^r v(t, x)dx \geq \int_{-r}^r u^*(t, x)dx.$$

- Comparison result for elliptic equations: the first step is to compare the solutions of two elliptic problems. Let  $\varepsilon > 0$ , let  $\varphi \in L^2(\mathbb{T})$ ,  $\varphi \geq 0$ , let  $\psi \in L^2(\mathbb{T})$ ,  $\psi \geq 0$  satisfying

$$\forall r \in (0; \pi), \quad \int_{-r}^r \psi^* \geq \int_{-r}^r \varphi^*,$$

and let  $w_\varphi, z_\psi$  be the solutions to

$$\begin{cases} -\Delta w_\varphi + \varepsilon w_\varphi = \varphi & \text{in } \mathbb{T}, \\ w_\varphi \in W^{1,2}(\mathbb{T}), \end{cases} \tag{51}$$

and

$$\begin{cases} -\Delta z_\psi + \varepsilon z_\psi = \psi^* & \text{in } \mathbb{T}, \\ z_\psi \in W^{1,2}(\mathbb{T}). \end{cases} \tag{52}$$

Then there holds:

$$\forall r \in (0; \pi), \quad \int_{-r}^r z_\psi \geq \int_{-r}^r w_\varphi. \tag{53}$$

To obtain (53), we may follow the standard steps of [42]: we assume that the level sets of (51) have measure zero (to cover the case of level sets of positive measure, one can argue as in [42], to which we refer for the sake of brevity). Let  $\tau > 0$  be a real number. Integrating (51) on  $\{w_\varphi \geq \tau\}$  yields

$$-\int_{\{w_\varphi=\tau\}} \frac{\partial w_\varphi}{\partial \nu} = \int_{\{w_\varphi \geq \tau\}} \varphi - \varepsilon \int_{\{w_\varphi \geq \tau\}} w_\varphi \leq \int_0^{\mu_{w_\varphi}(\tau)} \varphi^* - \varepsilon \int_0^{\mu_{w_\varphi}(\tau)} w_\varphi^*,$$

where the last inequality comes from the Hardy–Littlewood inequality. We recall that from the co-area formula

$$\text{for a.e. } \tau, \mu'_{w_\varphi}(\tau) = - \int_{\{w_\varphi=\tau\}} \frac{1}{|\nabla w_\varphi|}. \tag{54}$$

From the Cauchy–Schwarz inequality and the isoperimetric inequality we obtain

$$\begin{aligned} 4 \leq \text{Per}(\{w_\varphi = \tau\}) &\leq \int_{\{w_\varphi=\tau\}} \frac{1}{|\nabla w_\varphi|} \int_{\{w_\varphi=\tau\}} |\nabla w_\varphi| \\ &\leq -\mu'_{w_\varphi}(\tau) \int_{\{w_\varphi=\tau\}} |\nabla w_\varphi| \leq -\mu'_{w_\varphi}(\tau) \left( \int_0^{\mu_{w_\varphi}(\tau)} \varphi^* - \varepsilon \int_0^{\mu_{w_\varphi}(\tau)} w_\varphi^* \right). \end{aligned}$$

Since  $w_\varphi^*$  is the left inverse of  $\mu_{w_\varphi}$ , standard arguments [5] imply

$$\forall \xi \in (0; \pi), \quad -4(w_\varphi^*)'(\xi) \leq \int_0^\xi \varphi^* - \varepsilon \int_0^\xi w_\varphi^*. \tag{55}$$

It should be noted that if we work with (52) instead of (51) every inequality becomes an equality since  $z_\psi = z_\psi^*$ , and thus  $z_\psi$  satisfies

$$-4(z_\psi^*)'(\xi) = \int_0^\xi \psi^* - \varepsilon \int_0^\xi z_\psi^*. \tag{56}$$

Defining  $Z_\varphi := \int_0^\xi (z_\varphi - w_\varphi^*)$  we hence have

$$-Z_\varphi'' + \frac{\varepsilon}{4}Z_\varphi \geq 0, \quad Z_\varphi(0) = 0.$$

Furthermore, integrating (51) and (52) on the torus we obtain, by the equimeasurability of a function and its rearrangement,

$$\int_0^\pi w_\varphi^* = \frac{1}{2} \int_{\mathbb{T}} w_\varphi = \frac{1}{2\varepsilon} \int_{\mathbb{T}} \varphi \leq \frac{1}{2\varepsilon} \int_0^\pi \psi = \int_0^\pi z_\psi^* \tag{57}$$

so that

$$Z_\varphi(\pi) \geq 0.$$

By the maximum principle,  $Z_\varphi \geq 0$ , which concludes the proof.

- Comparison result for parabolic equations: for this second step, we follow the strategy of [6], which relies on a Picard iteration scheme. Namely, we discretise the parabolic problem in time: let  $N \in \mathbb{N}^*$  be a discretisation step. For  $u_0 \in \mathcal{A}$ , we define the sequences  $\{u_{0,k}, v_k\}_{k=0,\dots,N}$  as the solutions to

$$u_{0,0} = u_0, \quad \forall k \in \{0, \dots, N_1\}, \begin{cases} -\Delta u_{0,k+1} + \frac{1}{N}u_{0,k+1} = \frac{1}{N}u_{0,k} + g(u_{0,k}) & \text{in } \mathbb{T}, \\ u_{0,k} \in W^{1,2}(\mathbb{T}), \end{cases} \tag{58}$$

and

$$v_{0,0} = u_0^*, \quad \forall k \in \{0, \dots, N_1\}, \begin{cases} -\Delta v_{0,k+1} + \frac{1}{N}v_{0,k+1} = \frac{1}{N}v_{0,k} + g(v_{0,k}) & \text{in } \mathbb{T}, \\ v_{0,k} \in W^{1,2}(\mathbb{T}) \end{cases} \tag{59}$$

respectively. As  $g$  is convex and increasing, we have  $g(v)^* = g(v^*)$ . Thus, we can prove inductively that for every  $N$  and for every  $k \in \{0, \dots, N\}$  there holds

$$\forall r \in (0; \pi), \quad \int_{-r}^r v_{0,k} \geq \int_{-r}^r u_{0,k}^*$$

and it remains to pass to the limit  $N \rightarrow \infty$  to recover the result. □

**Conclusion.** Assume that  $u_0$  is a bang-bang maximiser of problem  $(\mathbf{P}_f)$ . Symmetrise it and extend it by periodicity. Consider the symmetric decreasing rearrangement  $u_0^*$  of its extension. Then  $\int_{-\pi}^\pi v(T, x)dx \geq \int_{-\pi}^\pi u(T, x)dx$  by proposition 2, where  $v$  is the solution of the periodic equation (49) associated with the initial datum  $u_0^*$ . Clearly,  $v(t, \cdot)$  and  $u(t, \cdot)$  are symmetric with respect to  $x = 0$  for all time  $t > 0$ . Hence,  $\int_0^\pi v(T, x)dx \geq \int_0^\pi u(T, x)dx$ . Also, one easily remarks that  $v$  restricted to  $(0, \pi)$  is the solution of the parabolic equation with Neumann boundary conditions (1), associated with the initial datum  $u_0^*$  restricted to  $(0, \pi)$ . On the other hand, as  $u_0$  is bang-bang, one has  $u_0^* = \mathbb{1}_{(-m,m)}$ . Hence,  $\mathbb{1}_{(0,m)}$  increases the criterion in  $(\mathbf{P}_f)$ . Thus, it is a solution of  $(\mathbf{P}_f)$ .

#### 4. Numerical analysis in the bistable framework

As we explained in the introduction, the behaviour of optimisers vary wildly depending on the shape of the reaction term  $f$ . To exemplify this phenomenon, we use the bistable nonlinearity that motivated [36], namely,  $f(u) := u(1 - u)(u - \theta)$ , with  $\theta \in (0, 1)$ .

When considering the optimisation problem  $(\mathbf{P}_f)$ , the fact that the set  $\{p = c\}$  may have a positive measure or, in other words, that a solution may not be the characteristic function of a set, leads to several difficulties in terms of numerical methods, because standard gradient methods or fixed-point algorithms fail to capture what this so-called ‘singular arc’ should be replaced with.

Let us first recall the main principles of the numerical algorithm introduced in [36] and explain the difficulty related to  $\{p = c\}$  further. Given the initial condition at the  $n$ -step  $u_0^n$ , we construct  $u_0^{n+1} = u_0^n + h_0^n$ , where  $h_0^n$  maximises (4) and is an admissible perturbation. Since the adjoint at the  $n$ th step  $p_0^n$  may have level sets of positive measure, one cannot directly apply the bathtub principle and choose  $h_0^n$  as the difference of characteristic functions of two level sets of  $p_0^n$ ; we must thus describe what happens on the singular arc, that is, on the level set  $\{p_0^n = c^n\}$  where  $c^n$  is chosen so that

$$|\{p_0^n > c^n\}| < m, |\{p_0^n \geq c^n\}| > m, |\{p_0^n = c^n\}| > 0. \tag{60}$$

We first define, in this case,  $u_0^{n+1} = 1$  on  $\{p_0^n > c^n\}$ ,  $u_0^{n+1} = 0$  on  $\{p_0^n < c^n\}$ , and it remains to fix the value of  $u_0^{n+1}$  on  $\omega_n$ . Defining  $\omega_n := \{p_0^n = c^n\}$  and discretising equation (5) on  $\omega_n$  we obtain, with an explicit finite difference scheme

$$-\left(\frac{p_0^n(dt, x) - c^n}{dt}\right) = f'(u_0^{n+1})c^n \tag{61}$$

and the value on  $u_0^{n+1}$  on  $\omega_n$  must be a root of (61). However, for bistable nonlinearities, this equation may have two roots, say  $\mu_1^n$  and  $\mu_2^n$ . In this case, these two roots can be distinguished through the convexity of  $f$ . In other words, if we have two roots, up to relabelling,

$$f''(\mu_1^n) > 0, \quad f''(\mu_2^n) < 0. \tag{62}$$

In [36] this difficulty is overcome by examining the two different possibilities and choosing the best one, which significantly lessens the performance of the algorithm, but theorem 1 allows to overcome this difficulty by choosing directly the root  $\mu_2^n$ , which is in the ‘concavity’ zone of  $f$ .

##### 4.1. Comparison of different numerical methods in the one-dimensional case

In this section, we want to study an example in order to compare the performances of our numerical algorithm with other well known optimisation algorithms to solve general non-linear problems under constraints. More precisely, we will consider the following numerical methods:

- *Method 1*: the numerical algorithm introduced in [36], which we improve using theorem 1, and that will be referred to as *our algorithm*.
- *Method 2*: the *interior-point* method, which is used to solve optimisation problems with linear equality and inequality constraints by applying the Newton method to a sequence of equality constrained problems. For a more detailed description of this method see for instance [10].

- *Method 3*: the *sequential quadratic programming* (SQP), which solves a sequence of optimisation sub-problems, each of which optimizes a quadratic model of the objective function subject to a linearisation of the constraints, see for instance [38].
- *Method 4*: the *simulated annealing* method, which is a probabilistic technique used to approximate global optimisation in a large search space. See for instance [21] for more details on this technique.

We used the MATLAB platform to perform the simulations. Methods 2 and 3 are already coded in the MATLAB function 'fmincon' while methods 1 and 4 were coded for the experiment.

**Setting the data.** Let us consider  $\Omega = (-50; 50)$ , and  $m = 13$ ; thus the admissible set is defined as follows:

$$\mathcal{A}_{13} = \left\{ u_0 \in L^1(\Omega) : 0 \leq u_0(x) \leq 1, \quad \text{and} \quad \int_{\Omega} u_0(x) dx = 13 \right\}.$$

Note that this set is defined by two inequalities and an equality constraint. We aim at maximising the quantity  $\mathcal{J}_T(u_0) := \int_{\Omega} u(T, x) dx$  for  $T = 25$  and we use a bistable reaction term  $f(u) := u(1 - u)(u - 0.25)$ .

In order to compare the performance of the four algorithms under the same conditions, we consider the same discretisation of  $\Omega$ . Moreover, the solution of the equation is systematically computed by the Crank–Nicolson method, and, for the initialisation we consider the same  $u_0^0$  given by a single block of mass 13. The value of the objective function at each iteration is numerically approximated by the rectangle rule. In particular, for the initialisation we have  $\mathcal{J}_{25}(u_0^0) = 29.42$ .

The results of the simulations are shown in figure 1 and table 1. For this example, our algorithm turns out to be faster than other well-known algorithms. Moreover, the evaluation of the objective function differs in less than 1% with respect to the best result obtained with the SQP method which takes more than twice the run-time of our algorithm.

Though the solution given by the SQP method is clearly more regular than the others, the profile of the local optimisers found by simulated annealing and by our algorithm do not seem to be far from this profile. Indeed, the solutions obtained through methods 1, 3 and 4 are qualitatively similar. On the other hand, the interior-point method gives a significantly different optimum, which seems to point out the good performance of our algorithm. It is important, however, to highlight that since uniqueness is not guaranteed in general, one cannot ensure that the algorithms have converged to a global maximiser but only to a local one.

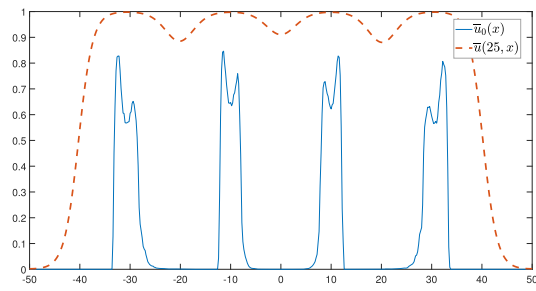
#### 4.2. Numerical simulations in the two-dimensional case

We only considered in the present paper the one-dimensional case. We now display some numerical results obtained in dimension 2, for which new patterns might arise.

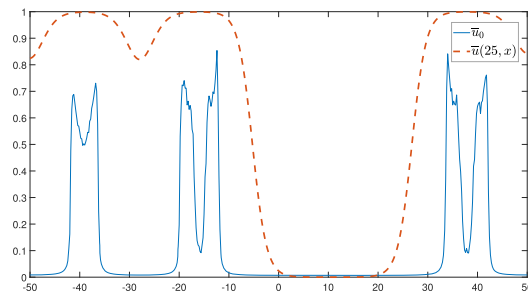
To solve the reaction–diffusion equation in the two-dimensional case, we consider the alternating direction implicit method which is a classical method to solve parabolic problems in two or three dimensions. As in the one-dimensional case, the algorithm and routines were coded in MATLAB.

We consider a square domain  $\Omega = (-10; 10) \times (-10; 10)$ , discretised uniformly by squares of side  $dx = 0.22$ . We fix  $T = 30$  for all subsequent simulations. We first tackle the case of a bistable reaction term

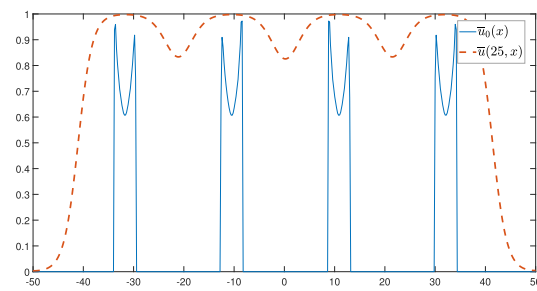
$$f(u) = u(1 - u)(u - 0.25).$$



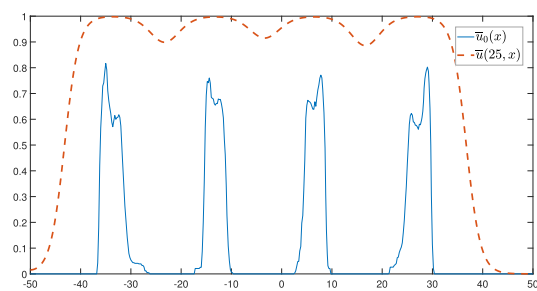
(a) Our algorithm



(b) Interior-point method



(c) Sequential quadratic programming

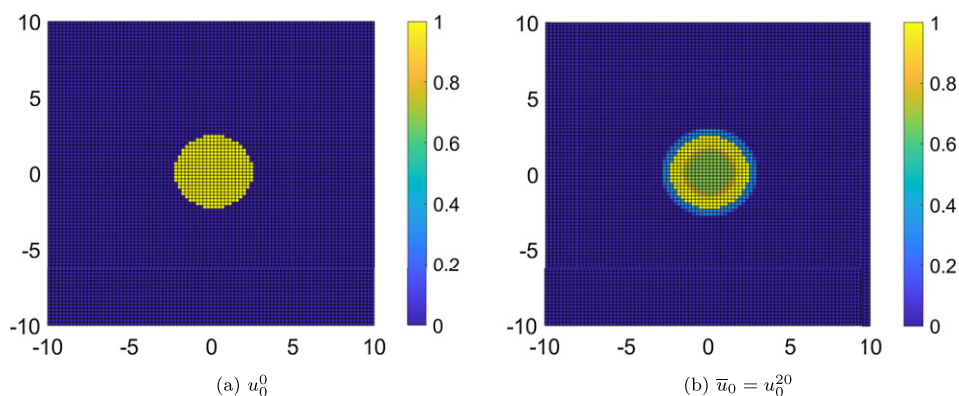


(d) Simulated annealing

**Figure 1.** Optimum found by means of the four different numerical algorithms.

**Table 1.** Comparing algorithms.

| Algorithm           | Objective function $\mathcal{J}_{25}(\bar{u}_0)$ | Run-time (in seconds) |
|---------------------|--|-----------------------|
| Our algorithm       | 77.9864  | 452                   |
| Interior point      | 65.6175  | 676                   |
| SQP                 | 78.7672  | 1342                  |
| Simulated annealing | 77.6238  | 4148                  |



**Figure 2.** In the left-hand side, we show the input of the algorithm, given by the ball of radius  $r = \sqrt{5.8}$  centered at the origin. In the right-hand side, we display the local optimum found by the numerical algorithm after 20 iterations, which looks radial, but is no longer a bang-bang distribution: it does not only take values 1 and 0.

In a second paragraph, we study the monostable case

$$f(u) = (u + 0.25)u(1 - u).$$

The justification for this second case is that this is a non-concave monostable nonlinearity. It is hence not covered by the theoretical results of [36].

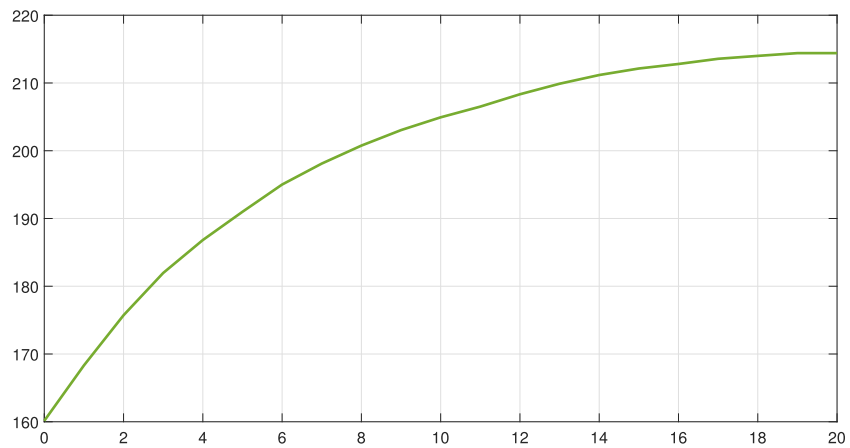
4.2.1. The bistable case.

**Example 1.** The algorithm is initialised with a ball of full density located in the middle of the domain  $\Omega$ . The mass is fixed to  $m = 5.8\pi$ , see figure 2(a). After 20 iterations, the algorithm converges to the local optimum showed in figure 2(b). The evolution of the objective function  $\mathcal{J}_{30}$  through iterations is showed in figure 3.

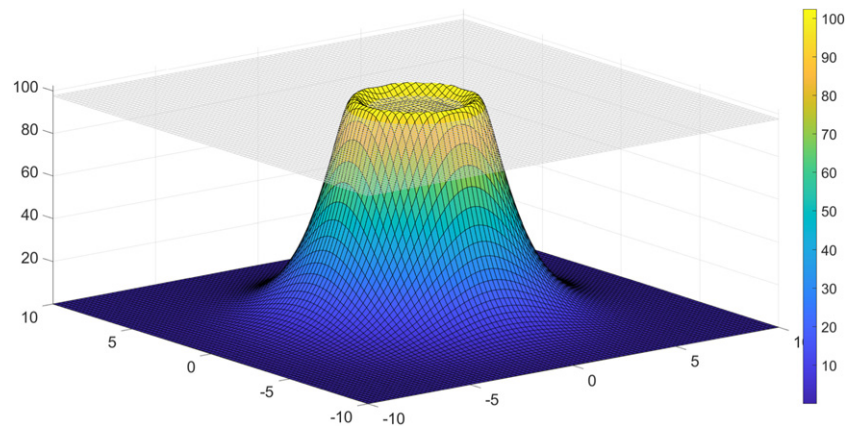
One might see that the local optimum found by the numerical algorithm is no longer a bang-bang function but a circular ball with less mass in the middle and a slightly bigger ratio. Looking at the adjoint state defined as the solution of equation (5), associated to this initial data, one might see that the area in the middle of the circle corresponds to a set where the adjoint state remains constant, see figure 4.

**Example 2.** In this case, we keep the same discretization and initial mass  $m$  of the previous example, but we consider an initial data which is a stripe of full density dividing our domain into two equal regions of zero density, see figure 5(a). The algorithm converges after 38 iterations





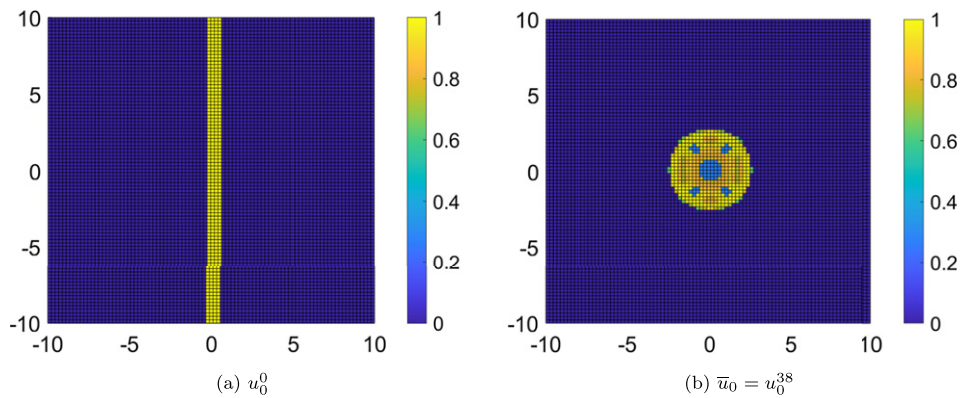
**Figure 3.** Evolution of the objective function from the initialisation  $\mathcal{J}_{30}(u_0^0) = 160.1$  to the last iteration  $\mathcal{J}_{30}(u_0^{20}) = 214.4$ .



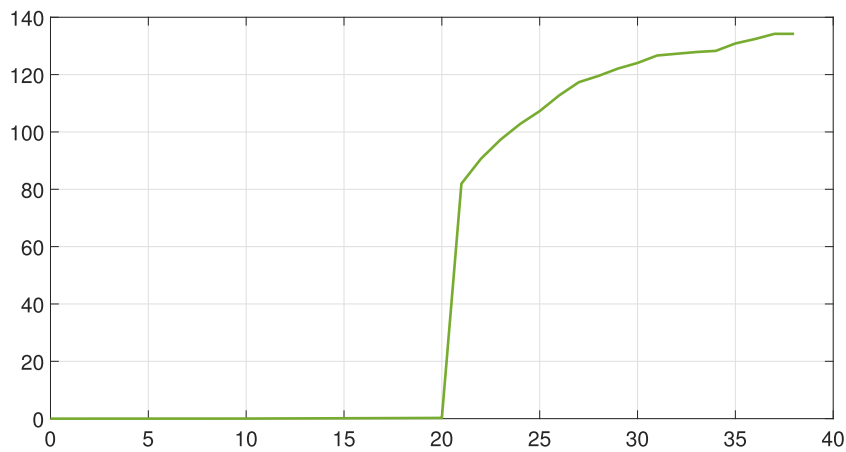
**Figure 4.** The figure shows the surface given by the solution  $\bar{p}(0, x)$  of the adjoint problem defined by the equation (5) associated to the initial data  $\bar{u}_0$  found by our algorithm. The plane colored in gray, is associated to the value  $\bar{c}$  described by (6) and thus for every  $x \in \Omega$  such that  $p_0(\mathbf{x}) = \bar{c}$ , one has  $0 < \bar{u}_0(x) < 1$ , see figure 2(b).

and the local optimum is displayed in figure 5(b). The corresponding variations of the objective function is showed in figure 6.

We observe that, in this case, the value of the objective function remains very low during the first 20 iterations. This fact, together with the radial geometry of the optimum found by the algorithm suggests that this stripe geometry is not optimal. It should also be pointed out that the geometry of the local optimum is interesting: indeed, it shows regions of zero density (i.e. the optimum  $\bar{u}_0$  found by the algorithm is equal to 0 in these regions) in the middle of regions of full density (i.e. where  $\bar{u}_0 = 1$ ), which exemplifies the phenomenon described in the one-dimensional case in [15].



**Figure 5.** In the left-hand side is showed the input of the algorithm, given by the stripe of width  $r = 0.91$  centred at the origin. In the right-hand side, the local optimum found by the numerical algorithm after 38 iterations.



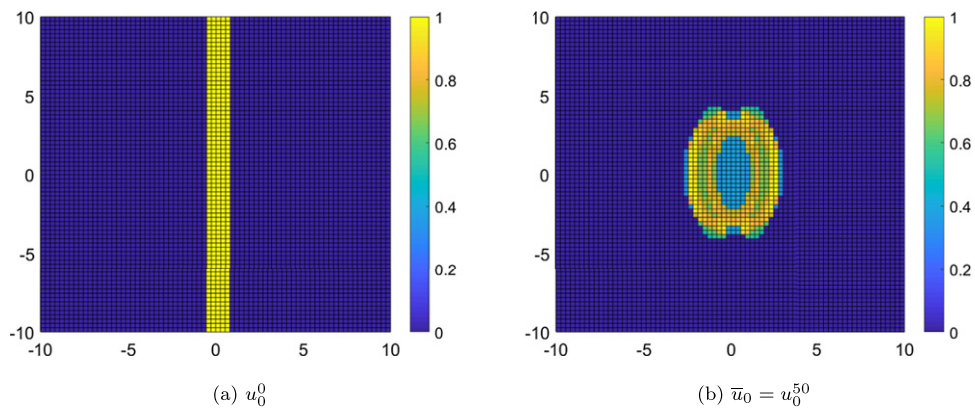
**Figure 6.** Evolution of the objective function from the initialisation  $\mathcal{J}_{30}(u_0^0) = 5.5 \times 10^{-6}$  to the last iteration  $\mathcal{J}_{30}(u_0^{38}) = 134.2$ .

Another relevant feature is that the optima found in the first and second examples are different, which indicates that our algorithm converge to local optima, and thus that the choice of the initial distribution  $u_0^0$  is crucial.

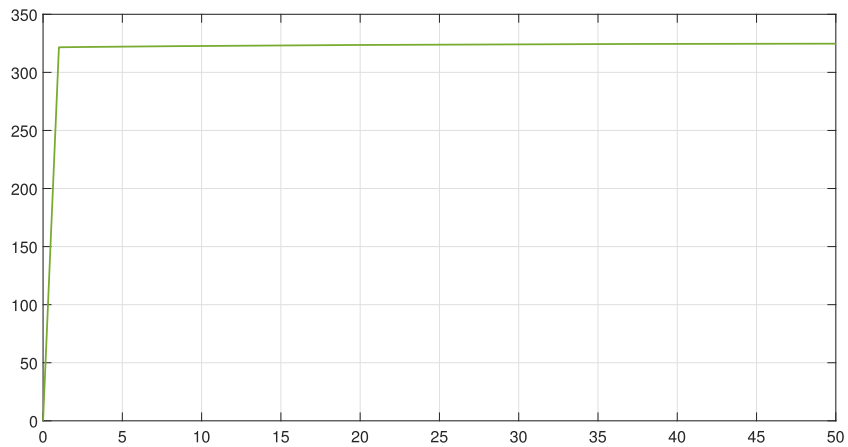
**Example 3.** For this example we keep the settings of the previous one, but we consider a higher initial mass  $m = 27$ . The geometry of the initial distribution is a stripe of full density dividing the domain into two regions of zero density, like in the example 2, see figure 7(a).

The corresponding local optimum found by the numerical algorithm is showed in figure 7(b). As in the previous case, the optimiser reflects a low density zone ringed by a high density region. This gap is clearly filled by diffusion as time evolves.

This example suggests once again the non optimality of stripe-like initial distributions. Note from figure 8 that, despite the considerable increase of the initial mass with respect to



**Figure 7.** In the left-hand side is showed the input of the algorithm, given by the stripe of width  $r = 1.6$  centered at the origin. In the right-hand side, the local optimum found by the numerical algorithm after 50 iterations.



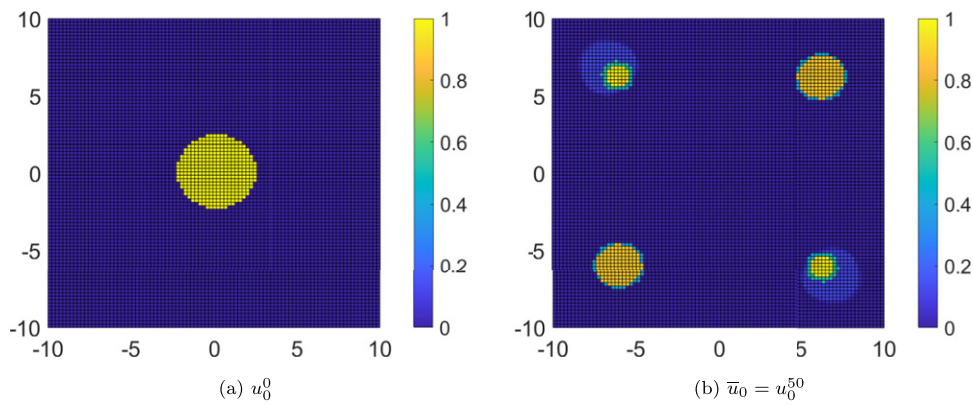
**Figure 8.** Evolution of the objective function from the initialisation  $\mathcal{J}_{30}(u_0^0) = 3 \times 10^{-5}$  to the last iteration  $\mathcal{J}_{30}(u_0^{50}) = 324$ .

example 2, the values of the objective function  $\mathcal{J}_{30}(u_0^0)$  associated to the stripe is of the order of  $10^{-5}$ , which is very low compared with the value associated to the final distribution.

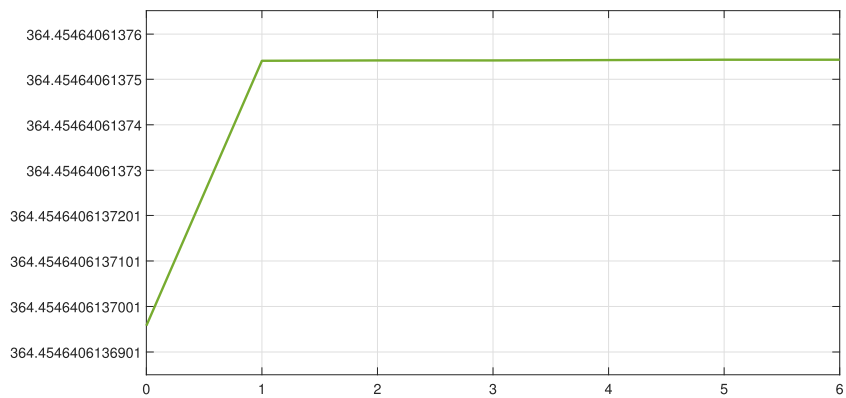
Finally, let us mention that a possible approximation of the maximiser was discussed in the appendix of [26]. Namely, in this thesis, the author replaced the maximiser  $\bar{u}_0$  by its mean on each of the connected components of  $\{\bar{u}_0 > 0\}$ . This gives pretty good results in several cases. It would be good to manage to quantify analytically the difference of criterion between this approximated initial datum and the global maximiser.

4.2.2. *The non-concave monostable case.* We now present some simulations in the case

$$f(u) = (u + 0.25)u(1 - u),$$



**Figure 9.** In the left-hand side we display the input of the algorithm, given by a centered ball. In the right-hand side, we show the local optimum found by the numerical algorithm after 50 iterations.



**Figure 10.** Evolution of the objective function from the initialisation  $\mathcal{J}_{30}(u_0^0) = 364.4$  to the last iteration  $\mathcal{J}_{30}(u_0^{50}) = 364.5$ .

still working under the constraint that  $0 \leq u_0 \leq 1$ . The motivation behind this case is that this nonlinearity is monostable on  $(0; 1)$ , i.e. it only has one stable equilibrium, but it is not concave. As a consequence, the theoretical approach developed in [36] cannot guarantee that the optimiser  $\bar{u}_0$  is the characteristic function of a subset of  $\Omega$ .

The parameters of the simulations are still the same:  $T = 30$ , and the initial configuration is the same as in example 1: the initialisation is a ball of full density, with mass  $m = 5.8\pi$  located in the middle of the domain  $\Omega$ . We refer to figures 9 and 10. We should however point out that in this simulation, the value of the objective function remains almost constant, despite the fact that the final configuration is very different from the initial one. It is plausible that such monostable nonlinearities converge too quickly to the equilibrium.

### 5. Conclusion, open problems and possible extensions

We make, in this conclusion, several concluding remarks and comments about possible generalisations and extensions of the results presented in this paper. For each of them, we try to present the arguments that have led us to the conclusion that other approaches were necessary in general.

#### 5.1. Regarding the regularity of the singular arc

One of the main drawbacks of theorem 1 is the regularity assumption on the singular arc. Namely, we obtain, for a maximiser  $u_0$  of  $\mathcal{J}_T$ , the characterisation  $f''(u_0) \leq 0$  only on the interior of the singular arc  $\{0 < u_0 < 1\}$ . The method presented in this paper (two-scale expansions) strongly relies on the smoothness properties of the cut-off function  $\theta$ .

It may be tempting, in the one dimensional case, to overcome this difficulty arguing as in [29]: if we simply assume that the singular arc  $\omega$  is measurable, but has positive measure, one can show that for every  $K \in \mathbb{N}$  there exists  $h_K \in L^2(\Omega)$ , supported in  $\omega$ , that writes

$$h_K = \sum_{k \geq K} \alpha_{k,K} \cos(k \cdot), \|h_K\|_{L^2} = 1. \tag{63}$$

Using the fact that  $h_K$  only has high Fourier modes, one may hop for a two-scale expansion of the form

$$h_K \approx \sum_{k \geq K} \alpha_{k,K} \left( h_k^0(k^2 t, x, kx) + \frac{1}{k} h_k^1(k^2 t, x, kx) \right).$$

This is however *a priori* prohibited by the problem of separation of phase: to obtain such a description, one needs well-separated phases, in the sense of [2]. In the context of Fourier series, this would require, at the very least, that  $h_K$  should write as a lacunary Fourier series (typically,  $h_K = \sum_{j=0}^\infty \alpha_{j,K} \cos(K^j x)$ ). However, for such lacunary Fourier series, Zygmund’s theorem (see [18] for instance) prohibits that they have compact support, so that admissible perturbations cannot have this structure. This is a major drawback, and it is unclear whether or not one may be able to overcome this difficulty via a similar approach, or if an entirely new strategy needs to be devised.

Another approach would be to prove some regularity on  $\omega$  ensuring that almost every of its point lie in its interior. This is satisfied for example if  $\bar{u}_0$  is Riemann integrable (since, due to Lebesgue’s characterisation of Riemann integrable functions, almost every point is a continuity point of  $\bar{u}_0$ ). Riemann integrability is satisfied by BV functions. Unfortunately, we were not able to push the regularity further than  $L^\infty$ .

#### 5.2. Monostable nonlinearities

As seen in subsection 4.2.2 of this paper, the numerical approach we propose, based on theorem 1, works for general monostable nonlinearities. The theoretical tools are, however, not sufficient at this level to fully characterise optimisers. An interesting question would be to discuss whether or not optimisers in the monostable case are always bang-bang, are if some degeneracy zones can appear.

#### 5.3. The singular arc in higher dimensions

It may be plausible to adapt the methods of theorem 1 to obtain a characterisation of the singular arc analogous to that of theorem 1 in the case  $\Omega = \prod_{i=1}^N [0; a_i]$ ,  $a_i > 0$ . To do so, the

main difference with our proof would be to replace the initial perturbation  $\theta(x)\cos(kx)$  with  $\prod_{i=1}^N \theta(x) \cos(kx_i)$ .

#### 5.4. Rearrangement inequalities for other types of boundary conditions

In this work, we mostly dealt with the case of Neumann boundary conditions in the one-dimensional case. We ought to note two things: first, the proof of theorem 1 should hold in the case of Dirichlet or of Robin boundary conditions, provided the functions  $\cos(k\cdot)$ , in the proof, are replaced with the Dirichlet or Robin eigenfunctions of the Laplacian in the interval. Second, regarding theorem 2, the same type of results can be obtained in a straightforward manner for the case of Dirichlet boundary conditions, by applying directly [6]. The case of Robin boundary conditions may be encompassed by using the recent Talenti inequalities obtained in this case in [4]. Addressing the problem on the full line  $\mathbb{R}$  is more tricky, since in this case even the existence of a maximiser is unclear. We plan on investigating such matters in future works.

### Acknowledgments

The authors wish to warmly thank J Pertinand for scientific conversations. The authors also wish to thank the anonymous referee for her/his comments, which have helped us improve the quality of this paper. I Mazari was partially supported by the French ANR Project ANR-18-CE40-0013—SHAPO on Shape optimisation and by the Austrian Science Fund (FWF) through the Grant I4052-N32. I Mazari and G Nadin were partially supported by the Project ‘Analysis and simulation of optimal shapes—application to lifesciences’ of the Paris City Hall.

### References

- [1] Allaire G 1992 Homogenization and two-scale convergence *SIAM J. Math. Anal.* **23** 1482–518
- [2] Allaire G and Briane M 1996 Multiscale convergence and reiterated homogenisation *Proc. R. Soc. Edinburgh A* **126** 297–342
- [3] Almeida L, Privat Y, Strugarek M and Vauchelet N 2019 Optimal releases for population replacement strategies: application to Wolbachia *SIAM J. Math. Anal.* **51** 3170–94
- [4] Alvino A, Nitsch C and Trombetti C 2019 A Talenti comparison result for solutions to elliptic problems with Robin boundary conditions *Analysis of PDEs*
- [5] Alvino A, Trombetti G and Lions P-L 1990 Comparison results for elliptic and parabolic equations via Schwarz symmetrization *Ann. Inst. Henri Poincaré C* **7** 37–65
- [6] Bandle C 1980 *Isoperimetric Inequalities and Applications (Monographs and Studies in Mathematics)* (Boston, MA: Pitman)
- [7] Barton N H and Turelli M 2011 Spatial waves of advance with bistable dynamics: cytoplasmic and genetic analogues of allee effects *Am. Nat.* **178** E48–75
- [8] Berestycki H, Hamel F and Roques L 2005 Analysis of the periodically fragmented environment model: I. Species persistence *J. Math. Biol.* **51** 75–113
- [9] Bintz J and Lenhart S 2020 Optimal resource allocation for a diffusive population model *J. Biol. Syst.* **28** 945–76
- [10] Boyd S and Vandenberghe L 2004 *Convex Optimization* (Cambridge: Cambridge University Press)

- [11] Bramanti M 1991 Symmetrization in parabolic Neumann problems *Appl. Anal.* **40** 21–39
- [12] Caubet F, Deheuvels T and Privat Y 2017 Optimal location of resources for biased movement of species: the 1D case *SIAM J. Appl. Math.* **77** 1876–903
- [13] Evans J 1975 Nerve axon equations: iv the stable and unstable impulse *Indiana University Mathematics Journal* **24** 1169–90
- [14] Ferone V and Mercaldo A 2005 Neumann problems and Steiner symmetrization *Commun. PDE* **30** 1537–53
- [15] Garnier J, Roques L and Hamel F 2012 Success rate of a biological invasion in terms of the spatial distribution of the founding population *Bull. Math. Biol.* **74** 453–73
- [16] Inoue J and Kuto K 2017 On the unboundedness of the ratio of species and resources for the diffusive logistic equation *Discrete Continuous Dyn. Syst. - Ser. B* **26** 2441
- [17] Kao C-Y, Lou Y and Yanagida E 2008 Principal eigenvalue for an elliptic problem with indefinite weight on cylindrical domains *Math. Biosci. Eng.* **5** 315–35
- [18] Kovrizhkin O 2003 A version of the uncertainty principle for functions with lacunary Fourier transforms *J. Math. Anal. Appl.* **288** 606–33
- [19] Lam K-Y, Liu S and Lou Y 2020 Selected topics on reaction–diffusion–advection models from spatial ecology *Math. Appl. Sci. Eng.* **1** 150–80
- [20] Lamboley J, Laurain A, Nadin G and Privat Y 2016 Properties of optimizers of the principal eigenvalue with indefinite weight and Robin conditions *Calc. Var. Partial Differ. Equ.* **55** 144
- [21] Locatelli M 2000 Simulated annealing algorithms for continuous global optimization: convergence conditions *J. Optim. Theory Appl.* **104** 121–33
- [22] Lou Y 2008 *Some Challenging Mathematical Problems in Evolution of Dispersal and Population Dynamics* (Berlin: Springer) pp 171–205
- [23] Lou Y, Nagahara K and Yanagida E 2021 Maximizing the total population with logistic growth in a patchy environment *J. Math. Biol.* **82** 2
- [24] Duprez Y P M, Hélie R and Vauchelet N 2021 Optimization of spatial control strategies for population replacement, application to Wolbachia
- [25] Maderna C, Salsa S and Pucci C 1979 Symmetrization in Neumann problems *Appl. Anal.* **9** 247–56
- [26] Marrero J I T 2021 Reaction–diffusion equations and applications to biological control of dengue and inflammation PhD Thesis
- [27] Masi A D, Ferrari P A and Lebowitz J L 1986 Reaction–diffusion equations for interacting particle systems *J. Stat. Phys.* **44** 589–644
- [28] Mazari I 2021 Quantitative estimates for parabolic optimal control problems under  $l^\infty$  and  $l^1$  constraints in the ball
- [29] Mazari I, Nadin G and Privat Y 2021 in preparation
- [30] Mazari I, Nadin G and Privat Y 2020 Optimal location of resources maximizing the total population size in logistic models *Journal de Mathématiques Pures et Appliquées* **134** 1–35
- [31] Mazari I, Nadin G and Privat Y 2020 Shape optimization of a weighted two-phase Dirichlet eigenvalue
- [32] Mazari I, Nadin G and Privat Y 2020 Some challenging optimisation problems for logistic diffusive equations and numerical issues to appear in *Handbook of Numerical Analysis*
- [33] Mazari I and Ruiz-Balet D 2021 A fragmentation phenomenon for a non-energetic optimal control problem: optimisation of the total population size in logistic diffusive models *SIAM Journal on Applied Mathematics* **81** 153–72
- [34] Mossino J and Rakotoson J M 1986 Isoperimetric inequalities in parabolic equations *Ann. della Scuola Norm. Super. Pisa - Cl. Sci.* **4** 1351–73
- [35] Murray J D 1993 *Mathematical Biology* (Berlin: Springer)
- [36] Nadin G and Toledo Marrero A I 2020 On the maximization problem for solutions of reaction–diffusion equations with respect to their initial data accepted for publication in *Journal of Mathematical Modelling of Natural Phenomena* **15** 71
- [37] Nagahara K and Yanagida E 2018 Maximization of the total population in a reaction–diffusion model with logistic growth *Calc. Var.* **57** 80
- [38] Nocedal J and Wright S J 2006 *Numerical Optimization* 2nd edn (New York: Springer)
- [39] Perthame B 2015 *Parabolic Equations in Biology* (Berlin: Springer)

- [40] Prochazka K and Vogl G 2017 Quantifying the driving factors for language shift in a bilingual region *Proc. Natl Acad. Sci. USA* **114** 4365–9
- [41] Rakotoson J-M 2008 *Réarrangement Relatif* (Berlin: Springer)
- [42] Talenti G 1976 Elliptic equations and rearrangements *Ann. della Scuola Norm. Super. Pisa - Cl. Sci.* **4** 697–718