



HAL
open science

Nonparametric estimation of a multivariate density under Kullback-Leibler loss with ISDE

Louis Pujol

► **To cite this version:**

Louis Pujol. Nonparametric estimation of a multivariate density under Kullback-Leibler loss with ISDE. 2022. hal-03660157

HAL Id: hal-03660157

<https://hal.science/hal-03660157>

Preprint submitted on 5 May 2022

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.

Nonparametric estimation of a multivariate density under Kullback-Leibler loss with ISDE

Louis Pujol*

Abstract

In this paper, we propose a theoretical analysis of the algorithm ISDE, introduced in previous work. From a dataset, ISDE learns a density written as a product of marginal density estimators over a partition of the features. We show that under some hypotheses, the Kullback-Leibler loss between the proper density and the output of ISDE is a bias term plus the sum of two terms which goes to zero as the number of samples goes to infinity. The rate of convergence indicates that ISDE tackles the curse of dimensionality by reducing the dimension from the one of the ambient space to the one of the biggest blocks in the partition. The constants reflect a combinatorial complexity reduction linked to the design of ISDE.

Keywords— Multivariate Density Estimation, Independence Structure, Non-parametric Density Estimation

*Université Paris-Saclay, CNRS, Inria, Laboratoire de Mathématiques d'Orsay, 91405, Orsay, France. louis.pujol@universite-paris-saclay.fr

1 NOTATIONS

Let f be a density function (a nonnegative real function whose integral is equal to 1) over \mathbb{R}^d . If we think of f from a statistical viewpoint, it is natural to refer to the indices $\{1, \dots, d\}$ as the features.

Let $S \subset \{1, \dots, d\}$, we denote by f_S the marginal density of f over S . For all $x = (x_1, \dots, x_d) \in \mathbb{R}^d$

$$f_S(x) = \int \prod_{i \notin S} dx_i f(x). \quad (1.1)$$

With a slight abuse of notation, to highlight the fact that $f_S(x)$ does not depend on $(x_i)_{i \notin S}$, we write $f_S(x_S)$ instead of $f_S(x)$.

Let k be a positive integer not greater than d . We denote by Set_d^k the set of all subsets of $\{1, \dots, d\}$ with cardinal not greater than k and by Part_d^k the collection of all partitions of $\{1, \dots, d\}$ constructed with blocks in Set_d^k . We also use the shortcuts $\text{Set}_d = \text{Set}_d^d$ and $\text{Part}_d = \text{Part}_d^d$.

2 INTRODUCTION

In a previous work ([8]), we have introduced ISDE (Independence Structure Density Estimation). ISDE estimates a density f from a set of iid realizations X_1, \dots, X_N considering the Independence Structure (IS) hypothesis. This paper is devoted to a theoretical analysis of this algorithm. In this introduction, we review existing theory about IS, introduce ISDE, and set the goals of the present work.

2.1 Curse of dimensionality and independence Structure

Minimax Risk Let X_1, \dots, X_N be iid realizations of a random variable in \mathbb{R}^d admitting a density f . The goal of density estimation is to construct an estimator \hat{f} of the density. We can measure the hardness of such an estimation task using

the minimax framework. Assume that the true density belongs to some known model \mathcal{F} and let D be a (pseudo)distance on \mathcal{F} , the minimax risk is defined as follows:

$$\mathcal{R}(D, \mathcal{F}) := \inf_{\hat{f}} \sup_{f \in \mathcal{F}} \mathbb{E} \left[D(f, \hat{f}) \right] \quad (2.1)$$

where the inf is taken over all measurable functions from the data to \mathcal{F} . More specifically, a great part of the literature on the topic deals with the asymptotic regime of $\mathcal{R}(D, \mathcal{F})$ with respect to N .

Hölder Balls Let \mathcal{U} be an open subset of \mathbb{R}^d and $g : \mathcal{U} \rightarrow \mathbb{R}$ a function. Let $\gamma = (\gamma_1, \dots, \gamma_d) \in \mathbb{N}^d$ be a multiindex and let $|\gamma| = \sum_{i=1}^d \gamma_i$ be its order. The partial differentiate operator D^γ is defined as follows

$$D^\gamma g = \frac{\partial^{|\gamma|} g}{\partial_1^{\gamma_1} \dots \partial_d^{\gamma_d}}. \quad (2.2)$$

For a positive number β , if we denote by s the larger integer strictly lower than β and let $\delta = \beta - s \in (0, 1]$, g belongs to the Hölder ball $\mathcal{H}(\beta, L)$ where L is a positive real number if both following conditions are fulfilled

$$\begin{cases} \max_{|\gamma| \leq s} \sup_{x \in \mathcal{U}} \|D^\gamma g(x)\| \leq L \\ \max_{|\gamma| = s} \sup_{x, y \in \mathcal{U}} |D^\gamma g(x) - D^\gamma g(y)| \leq L \|x - y\|^\delta. \end{cases} \quad (2.3)$$

If g is defined on a close subset \mathcal{C} of \mathbb{R}^d , we say that $g \in \mathcal{H}(\beta, L)$ if the restriction of g to the interior of \mathcal{C} .

Minimax Risk over Hölder Balls In [3], the minimax rate of this family of functions was studied considering L_p distances. In particular, the result with the squared L_2 distance is the following

$$\mathcal{R}(\|\cdot\|_2^2, \mathcal{H}^\beta(d, H)) \sim N^{-\frac{2\beta}{2\beta+d}}. \quad (2.4)$$

We can interpret this bound as a manifestation of the curse of dimensionality because of its dependence on d . A solution is to consider the IS model introduced in [5].

Independence Structure For $k \leq d$, we define a family of functions:

$$\mathcal{D}_d^k = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \exists \mathcal{P} \in \text{Part}_d^k : f(x) = \prod_{S \in \mathcal{P}} f_S(x_S) \right\}. \quad (2.5)$$

In probabilistic terms, a density f over \mathbb{R}^d belongs to \mathcal{D}_d^k if we can group these features into independent blocks. Another viewpoint is that the random variable characterized by f admits a graphical model, a collection of disjoint fully connected cliques of size not greater than k . It was showed in [9] that

$$\mathbb{R} \left(\|\cdot\|_2^2, \mathcal{H}^\beta(d, H) \cap \mathcal{D}_d^k \right) \sim N^{-\frac{2\beta}{2\beta+k}}. \quad (2.6)$$

The striking fact here is that the hardness of the estimation problem is no longer related to the ambient dimension but instead to the size of the biggest block of the partition on which the density function is decomposable.

2.2 ISDE

As explained in [8], the density estimation problem under squared L_2 loss does not lead to a feasible algorithm. This is why we change the loss function to the Kullback-Leibler (KL) divergence. If \hat{f} is an estimator of f , the KL loss between f and \hat{f} is defined as

$$\text{KL} \left(f \parallel \hat{f} \right) = \int \log \left(\frac{f}{\hat{f}} \right) f. \quad (2.7)$$

This formulation is well suited to IS as it involves log-densities, and a log of a product of marginal densities is a sum of log of marginal densities. From an algorithmic viewpoint, this formulation allowed us to implement an algorithm with reasonable running time and memory usage (see [8] for details). ISDE operates as follows.

1. Two independent datasets W_1, \dots, W_m and Z_1, \dots, Z_n are extracted from X_1, \dots, X_N .
2. W_1, \dots, W_m is used to compute marginal density estimators $(\hat{f}_S)_{S \in \text{Set}_d^k}$. Any multidimensional density estimation procedure can be used for this step.
3. Z_1, \dots, Z_n is used to compute $(\ell_n(S))_{S \in \text{Set}_d^k}$ where

$$\ell_n(S) = \frac{1}{n} \sum_{i=1}^n \log \hat{f}_S((Z_i)_S). \quad (2.8)$$

4. The optimization problem

$$\max_{\mathcal{P} \in \text{Part}_d^k} \sum_{S \in \mathcal{P}} \ell_n(S) \quad (2.9)$$

is solved with an exact integer programming optimization procedure (branch-and-bound).

The output is a partition $\hat{\mathcal{P}}$ and an estimator of f taking the form $\hat{f}_{\hat{\mathcal{P}}} = \prod_{S \in \hat{\mathcal{P}}} \hat{f}_S$. The formulation with log-densities leads to a combinatorial complexity reduction. At first glance, the problem of density estimation under IS necessitates manipulating data structures of size P_d^k while ISDE only lies on data structures of size S_d^k .

2.3 Goal of this work

This paper is intended to provide a theoretical analysis of ISDE by upper-bounding the quantity $\text{KL}(f \parallel \hat{f})$. In particular, we will show that the introduction of IS tackles the curse of dimensionality and that the constants in the upper-bound reflect the combinatorial complexity reduction implemented in ISDE.

2.4 Organization of the paper

In section 3 we establish a first decomposition on the risk involving oracle partitions. In section 4 we introduce the regularity conditions on the proper density

and establish that an upper-bound for the uniform loss between marginal densities of f and marginal estimators is sufficient to obtain a convergence result for ISDE. In section 5 we show that it is possible to obtain an upper-bound for uniform estimation of marginal densities for a particular estimator. Then in section 6 we state the desired upper-bound for the estimator outputted by ISDE.

3 KULLBACK-LEIBLER RISK DECOMPOSITION

In this section, we show that the KL loss between f and $\hat{f}_{\hat{\mathcal{P}}}$, the estimator outputted by ISDE, decomposes as the sum of three terms with a clear interpretation.

3.1 Oracles partitions

We denote by $\hat{\mathcal{P}}$ the partition outputted by ISDE. Let $P_n[\cdot]$ denotes the empirical measure associated with the sample Z_1, \dots, Z_n and $P[\cdot]$ the measure associated with the true density f . $\hat{\mathcal{P}}$ is solution of the following optimisation problem :

$$\hat{\mathcal{P}} \in \arg \min_{\mathcal{P} \in \text{Part}_d^k} P_n \left(-\log(\hat{f}_{\mathcal{P}}) \right) \quad (3.1)$$

The partition $\hat{\mathcal{P}}$ is random depending on both W and Z . Let us define two other meaningful partitions.

$$\tilde{\mathcal{P}} \in \arg \min_{\mathcal{P} \in \text{Part}_d^k} P \left(-\log(\hat{f}_{\mathcal{P}}) \right) = \arg \min_{\mathcal{P} \in \text{Part}_d^k} \text{KL} \left(f \parallel \hat{f}_{\mathcal{P}} \right) \quad (3.2)$$

and

$$\mathcal{P}_* \in \arg \min_{\mathcal{P} \in \text{Part}_d^k} P \left(-\log(f_{\mathcal{P}}) \right) = \arg \min_{\mathcal{P} \in \text{Part}_d^k} \text{KL} \left(f \parallel f_{\mathcal{P}} \right). \quad (3.3)$$

$\tilde{\mathcal{P}}$ is a random partition depending on W but not on Z . It is the best combination of the estimators $(\hat{f}_S)_{S \in \text{Set}_d^k}$ if we consider that the quantities $\left(P(-\log \hat{f}_S)\right)_{S \in \text{Set}_d^k}$ are known.

\mathcal{P}_* is not random. It is only a function of k . $f_{\mathcal{P}_*}$ can be interpreted as the Kullback-Leibler projection of f on the model \mathcal{D}_d^k thanks to the following property.

Proposition 3.1

$$f_{\mathcal{P}_*} \in \arg \min_{g \in \mathcal{D}_d^k} \text{KL}(f \| g) \quad (3.4)$$

Proof. Let $g \in \mathcal{D}_d^k$ and denote by \mathcal{P}_g a partition such that $g = \prod_{S \in \mathcal{P}_g} g_S$. We have

$$\text{KL}(f \| g) = \int \log \left(\frac{f}{g} \right) f \quad (3.5)$$

$$= \int \log \left(\frac{f}{f_{\mathcal{P}_g}} \right) f + \int \log \left(\frac{f_{\mathcal{P}_g}}{g} \right) f \quad (3.6)$$

$$= \text{KL}(f \| f_{\mathcal{P}_g}) + \sum_{S \in \mathcal{P}_g} \text{KL}(f_S \| g_S) \quad (3.7)$$

$$\leq \text{KL}(f \| f_{\mathcal{P}_g}) \quad (3.8)$$

with equality if $g = f_{\mathcal{P}_g}$. Then

$$\arg \min_{g \in \mathcal{D}_d^k} \text{KL}(f \| g) = \arg \min_{\mathcal{P} \in \text{Part}_d^k} \min_{g \in \mathcal{D}_d^k} \text{KL}(f \| g) \quad (3.9)$$

$$= \arg \min_{\mathcal{P} \in \text{Part}_d^k} \text{KL}(f \| f_{\mathcal{P}}) \quad (3.10)$$

□

3.2 Kullback-Leibler risk upper-bound

We are now in a position to establish a control of the Kullback-Leibler risk for $\hat{f}_{\hat{\mathcal{P}}}$ involving the oracles partitions.

Lemma 3.2: Kullback-Leibler risk control

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) \leq \text{KL} \left(f \| f_{\mathcal{P}_*} \right) + \sum_{S_* \in \mathcal{P}_*} \text{KL} \left(f_{S_*} \| \hat{f}_{S_*} \right) + (P - P_n) (\log \hat{f}_{\hat{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}}) \quad (3.11)$$

Proof. We start by decomposing $\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right)$ as follows

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) = \text{KL} \left(f \| f_{\mathcal{P}_*} \right) \quad (3.12)$$

$$+ \text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right) - \text{KL} \left(f \| f_{\mathcal{P}_*} \right) \quad (3.13)$$

$$+ \text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right) \quad (3.14)$$

$$+ \text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right). \quad (3.15)$$

Then, as $\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) \leq \text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right)$,

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) \leq \text{KL} \left(f \| f_{\mathcal{P}_*} \right) \quad (3.16)$$

$$+ \text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right) - \text{KL} \left(f \| f_{\mathcal{P}_*} \right) \quad (\text{i}) \quad (3.17)$$

$$+ \text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right) \quad (\text{ii}). \quad (3.18)$$

Now, we rewrite (i)

$$\text{KL} \left(f \| \hat{f}_{\mathcal{P}_*} \right) - \text{KL} \left(f \| f_{\mathcal{P}_*} \right) = \int \log \left(\frac{f(x)}{\hat{f}_{\mathcal{P}_*}(x)} \right) f(x) dx - \int \log \left(\frac{f(x)}{f_{\mathcal{P}_*}(x)} \right) f(x) dx \quad (3.19)$$

$$= \int \log \left(\frac{f_{\mathcal{P}_*}(x)}{\hat{f}_{\mathcal{P}_*}(x)} \right) f(x) dx \quad (3.20)$$

$$= \sum_{S_* \in \mathcal{P}_*} \int \log \left(\frac{f_{S_*}(x)}{\hat{f}_{S_*}(x)} \right) f_{S_*}(x) dx \quad (3.21)$$

$$= \sum_{S_* \in \mathcal{P}_*} \text{KL} \left(f_{S_*} \| \hat{f}_{S_*} \right). \quad (3.22)$$

And we upper-bound (ii). As $P_n \left[\log \hat{f}_{\hat{\mathcal{P}}} - \log \hat{f}_{\tilde{\mathcal{P}}} \right] \geq 0$:

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \| \hat{f}_{\tilde{\mathcal{P}}} \right) = P \left[-\log \hat{f}_{\hat{\mathcal{P}}} \right] - P \left[-\log \hat{f}_{\tilde{\mathcal{P}}} \right] \quad (3.23)$$

$$\leq P \left[\log \hat{f}_{\tilde{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}} \right] + P_n \left[\log \hat{f}_{\hat{\mathcal{P}}} - \log \hat{f}_{\tilde{\mathcal{P}}} \right] \quad (3.24)$$

$$= (P - P_n)(\log \hat{f}_{\tilde{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}}). \quad (3.25)$$

□

Three terms appear in the upper bound, and they can be easily interpreted.

- $\text{KL} (f \| f_{\mathcal{P}_*})$ is a bias term. It is the intrinsic error of the model \mathcal{D}_d^k and can be thought of as a distance from f to \mathcal{D}_d^k thanks to proposition 3.1.
- $\sum_{S_* \in \mathcal{P}_*} \text{KL} \left(f_{S_*} \| \hat{f}_{S_*} \right)$ is an approximation term. It is a random quantity depending on the sample W and represents the error made when $f_{\mathcal{P}_*}$ is estimated with $\hat{f}_{\mathcal{P}_*}$. Conditionally to W , it depends on how the estimation of log-likelihoods made thanks to Z is accurate and quantifies our ability to output the optimal partition.

In the sequel of the paper, we will focus on upper-bounding the approximation and selection terms as they are the random quantities of interest in our problem. We treat the bias term as a structural error, and we focus on upper-bounding the quantity

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} (f \| f_{\mathcal{P}_*}). \quad (3.26)$$

A study of the bias in a multivariate Gaussian framework can be found in appendix A.

4 CONDITIONS AND OBJECTIVE

4.1 Regularity conditions

Bounding condition Density estimation under Kullback-Leibler loss is known to be a challenging problem. One work by [2] has studied the asymptotic convergence rates for kernel estimators in a one-dimensional setting. It was shown that the tails of the kernel must be chosen appropriately regarding the tails property of the proper density to have convergent estimators. In this work, we restrict our attention to densities that are lower and upper bounded by some positive quantities. This is done to avoid hardly tractable tail behavior issues. In the sequel, we consider that the following bounding condition is valid for all $S \in \text{Set}_d^k$

$$e^{-A|S|} \leq f_S \leq e^{A|S|} \quad \forall S \in \mathcal{P}_* \quad (\text{BC})$$

Note that if we impose a positive lower bound on the marginal densities, we must consider that f is compactly supported. In the sequel, we will suppose that the support of f is $[0, 1]^d$.

Hölder Regularity We will consider in the sequel that it exists $\beta \in (0, 2]$ and $L > 0$ such that $f_S \in \mathcal{H}(\beta, L)$ for all $S \in \text{Set}_d^k$. We will use the following approximation property for functions in Hölder balls.

Lemma 4.1: Approximation

Let us consider that $g \in \mathcal{H}(\beta, H)$ with $\beta \in (0, 2)$ and the domain of g is $\mathcal{U} \subset \mathbb{R}^d$, then for all $x \in \mathcal{U}$ and u such that $x + u \in \mathcal{U}$. If $\beta \in (0, 1]$ then

$$|g(x) - g(x + u)| \leq L \|u\|^\beta. \quad (4.1)$$

If $\beta \in (1, 2]$ then

$$\left| g(x) - g(x + u) - \sum_{k=1}^d \partial_k u_k g(x) \right| \leq L \|u\|^\beta. \quad (4.2)$$

4.2 Objective

Our goal is to propose an estimation procedure for the collection of marginal densities $(f_S)_{S \in \text{Set}_d^k}$. If we are able to ensure, simultaneously for all $S \in \text{Set}_d^k$ a uniform control

$$\|\hat{f}_S - f_S\|_\infty \leq \epsilon_S < e^{-A|S|}(1 - e^{-A|S|}). \quad (\text{UC})$$

Then we can upper-bound the approximation term and the selection term thanks to the following proposition.

Proposition 4.2: Consequences of uniform control

If the uniform control (UC) is satisfied and (BC) is true, then

1. A bounding condition is satisfied by all the estimators $(\hat{f}_S)_{S \in \text{Set}_d^k}$

$$e^{-2A|S|} \leq \hat{f}_S \leq e^{2A|S|}. \quad (\widehat{\text{BC}})$$

2. For all $S \in \text{Set}_d^k$, the Kullback-Leibler divergence between f_S and \hat{f}_S can be upper-bounded

$$\text{KL}(f_S \parallel \hat{f}_S) \leq e^{2A|S|} \epsilon_S \quad (4.3)$$

3. Conditionnaly on W , the selection term can be upper-bounded with high probability. More precisely if $\delta_n \in (0, 1)$ we have

$$\mathbb{P} \left[\left| (P - P_n)(\log \hat{f}_{\hat{\mathcal{P}}} - \log \hat{f}_{\mathcal{P}}) \right| \geq 2d \sqrt{\frac{2Ak}{n}} \sqrt{\log \left(\frac{2S_d^k}{\delta} \right)} \middle| W \right] \leq \delta_n. \quad (4.4)$$

Proof. Proof of 1: Under (BC) we have

$$e^{-A|S|} - \|\hat{f}_S - f_S\|_\infty \leq \hat{f}_S \leq e^{A|S|} + \|\hat{f}_S - f_S\|_\infty. \quad (4.5)$$

Now, under (UC)

$$e^{-A|S|} - \|\hat{f}_S - f_S\|_\infty \geq e^{-A|S|} - e^{-A|S|}(1 - e^{-A|S|}) = e^{-2A|S|} \quad (4.6)$$

and

$$e^{A|S|} + \|\hat{f}_S - f_S\|_\infty \leq e^{A|S|} + e^{-A|S|}(1 - e^{-A|S|}) \quad (4.7)$$

$$\leq e^{A|S|} + e^{3A|S|} (e^{-A|S|}(1 - e^{-A|S|})) \quad (4.8)$$

$$= e^{A|S|} + e^{2A|S|}(1 - e^{-A|S|}) = e^{2A|S|}. \quad (4.9)$$

Proof of 2: Let us compute

$$\text{KL} \left(f_S \| \hat{f}_S \right) = \int \log \left(\frac{f_S}{\hat{f}_S} \right) f_S \quad (4.10)$$

$$\leq \int \left(\frac{f_S - \hat{f}_S}{\hat{f}_S} \right) f_S \quad (4.11)$$

Using $(\widehat{\text{BC}})$, $1 / \hat{f}_S \leq e^{2A|S|}$

$$\leq e^{2A|S|} \|f_S - \hat{f}_{S, h_m}\|_\infty \quad (4.12)$$

$$\leq e^{2A|S|} \epsilon_S \quad (4.13)$$

Proof of 3: Let $S \in \text{Set}_d^k$, under $(\widehat{\text{BC}})$ we have $\log \hat{f}_S \in [-2A|S|, 2A|S|]$. Using Hoeffding inequality, we obtain

$$\mathbb{P} \left[\left| (P - P_n) \log \hat{f}_S \right| \geq \sqrt{\frac{2A|S|}{n}} \sqrt{\log \frac{2S_d^k}{\delta} |W} \right] \leq \frac{\delta}{S_d^k}. \quad (4.14)$$

Now, by union bound :

$$\mathbb{P} \left[\sup_{S \in \text{Set}_d^k} \left| (P - P_n) \log \hat{f}_S \right| \geq \sqrt{\frac{2A|S|}{n}} \sqrt{\log \frac{2S_d^k}{\delta} |W} \right] \leq \delta. \quad (4.15)$$

This leads to :

$$\mathbb{P} \left[2d \sup_{S \in \text{Set}_d^k} \left| (P - P_n) \log \hat{f}_S \right| \leq 2d \sqrt{\frac{2Ak}{n}} \sqrt{\log \frac{2S_d^k}{\delta} |W} \right] \geq 1 - \delta. \quad (4.16)$$

Now, we remark that

$$|(P - P_n) \log \hat{f}_{\tilde{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}}| = \left| \sum_{S \in \tilde{\mathcal{P}}} (P - P_n) \log \hat{f}_S - \sum_{S \in \hat{\mathcal{P}}} (P - P_n) \log \hat{f}_S \right| \quad (4.17)$$

$$\leq \sum_{S \in \tilde{\mathcal{P}}} |(P - P_n) \log \hat{f}_S| + \sum_{S \in \hat{\mathcal{P}}} |(P - P_n) \log \hat{f}_S| \quad (4.18)$$

$$\leq 2d \sup_{S \in \text{Set}_d^k} |(P - P_n) \log \hat{f}_S|. \quad (4.19)$$

Then, we have

$$\mathbb{P} \left[|(P - P_n) \log \hat{f}_{\tilde{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}}| \geq 2d \sqrt{\frac{2Ak}{n}} \sqrt{\log \frac{2S_d^k}{\delta} |W|} \right] \leq \delta. \quad (4.20)$$

□

5 UNIFORM DENSITY ESTIMATION FOR MARGINAL DENSITIES

5.1 For a fixed S

In this subsection, we fix a subset of variables $S \in \text{Set}_d^k$ and we study the problem of constructing an estimator \hat{f}_S based on the sample W_1, \dots, W_m giving a control of $\|f_S - \hat{f}_S\|_\infty$ in order to verify (UC). We decompose the error as a sum of a bias and a variance term as follows

$$\|f_S - \hat{f}_S\|_\infty \leq \underbrace{\|f_S - \mathbb{E}[\hat{f}_S]\|_\infty}_{\text{Bias}} + \underbrace{\|\mathbb{E}[\hat{f}_S] - \hat{f}_S\|_\infty}_{\text{Variance}}. \quad (5.1)$$

5.1.1 Bias upper-bound

Choice of the kernel function In the following we will use density estimator based on an ancillary function K called kernel. K is a nonnegative integrable function on \mathbb{R} such that $\int K(x)dx = 1$, we consider the following assumptions on K :

$$\begin{cases} \forall x \in \mathbb{R}, K(-x) = K(x) \\ \text{Supp}(K) \in [-1, 1] \\ \int xK(x)dx = 0 \\ \|K\|_\infty < \infty \end{cases} \quad (\text{A.K})$$

We will also assume that, if $K_{h,x}^S : u \mapsto \frac{1}{h^{|S|}} \prod_{k \in S} K\left(\frac{x_k - u_k}{h}\right)$, the family of function

$$\mathcal{F}_S = \{K_{h,x}^S, h > 0, x \in \mathbb{R}^{|S|}\} \quad (5.2)$$

is a bounded VC class of functions. It means that it exists positive numbers A and ν such that for any probability measure P over $\mathbb{R}^{|S|}$ and any $\tau \in (0, 1)$ we have

$$\mathcal{N}(\mathcal{F}_S, L_2(P), \tau) \leq \left(\frac{A\|K\|_\infty}{\tau}\right)^\nu \quad (5.3)$$

where $\mathcal{N}(\mathcal{F}_S, L_2(P), \tau)$ is the τ -covering number of \mathcal{F}_S for the $L_2(P)$ distance. As proved in [1] this condition is met for almost all classical kernels. An example of kernel function K satisfying all the assumptions is the Epanechnikov kernel K_{Epa}

$$K_{\text{Epa}}(x) = \frac{3}{4}(1 - x^2)\mathbb{1}_{[-1,1]}(x). \quad (5.4)$$

Boundary issue One must be aware of the issue induced by the fact that f_S is supported on $[0, 1]^{|S|}$. We define the usual kernel density estimator (KDE) as follows. Let h be a positive real number. The KDE for the marginal density f_S associated with the kernel K , the bandwidth $h > 0$, and the sample W is defined as

$$\hat{f}_{h,S}^{\text{KDE}}(x) = \frac{1}{mh^{|S|}} \sum_{i=1}^m \prod_{k \in S} K\left(\frac{(W_i)_k - x_k}{h}\right) \quad (5.5)$$

We remark that even in the samples W_1, \dots, W_m belong to $[0, 1]^d$, there is no reason to have $\hat{f}_{h,S}^{\text{KDE}}$ supported in $[0, 1]^{|S|}$.

In this setting, the bias of the classical KDE does not go to zero as $h \rightarrow 0$, illustrating the boundary issue induced by estimating a compactly supported density.

Proposition 5.1: Boundary issue

Let $f_S \in \mathcal{H}(2, L)$, the bias

$$\left\| \mathbb{E} \left[\hat{f}_{h,S}^{\text{KDE}} \right] - f_S \right\|_{\infty} \quad (5.6)$$

does not tend to 0 as $h \rightarrow 0$.

Proof.

$$\mathbb{E} \left[\hat{f}_{h,S}^{\text{KDE}} \right] (0) - f_S(0) = \frac{1}{h^{|S|}} \int_{[0,1]^{|S|}} f_S(t) \prod_{k \in S} K \left(\frac{t_k}{h} \right) dt - f_S(0) \quad (5.7)$$

$$= \int_{[-1,1]^{|S|}} [f_S(hu) - f_S(0)] \prod_{k \in S} K(u_k) du_k \quad (5.8)$$

As $\int xK(x) = 0$

$$= \int_{[-1,1]^{|S|}} \left[f_S(hu) - f_S(0) - h \sum_{k \in S} u_k \partial_k f_S(0) \right] \prod_{k \in S} K(u_k) du_k \quad (5.9)$$

Now as $f_S(x) = 0$ for $x \notin [0, 1]^{|S|}$

$$\begin{aligned} &= \int_{[0,1]^{|S|}} \left[f_S(hu) - f_S(0) - h \sum_{k \in S} u_k \partial_k f_S(0) \right] \prod_{k \in S} K(u_k) du_k \\ &\quad - f_S(0) \int_{[-1,1]^{|S|} \setminus [0,1]^{|S|}} \prod_{k \in S} K(u_k) du_k \\ &\quad - h \sum_{k \in S} \partial_k f_S(0) \int_{[-1,1]^{|S|} \setminus [0,1]^{|S|}} u_k \prod_{k \in S} K(u_k) du_k \end{aligned} \quad (5.10)$$

The third term in the final sum tends to 0 with h . The same is true for the first term as

$$\begin{aligned} & \left| \int_{[0,1]^{|S|}} \left[f_S(hu) - f_S(0) - h \sum_{k \in S} u_k \partial_k f_S(0) \right] \prod_{k \in S} K(u_k) du_k \right| \\ & \leq \int_{[0,1]^{|S|}} \left| f_S(hu) - f_S(0) - h \sum_{k \in S} u_k \partial_k f_S(0) \right| \prod_{k \in S} K(u_k) du_k \end{aligned} \quad (5.11)$$

$$\leq Lh^2 \sum_{k \in S} \int_{[0,1]^{|S|}} u_k^2 \prod_{k \in S} K(u_k) du_k \quad (5.12)$$

$$\leq L\sigma_K^2 h^2. \quad (5.13)$$

Now, as $\int_{[-1,1]^{|S|} \setminus [0,1]^{|S|}} \prod_{k \in S} K(u_k) du_k = \frac{2^{|S|}-1}{2^{|S|}}$, we conclude that

$$\lim_{h \rightarrow 0} \mathbb{E} \left[\hat{f}_{h,S}^{\text{KDE}} \right] (0) = f_S(0) \frac{2^{|S|}-1}{2^{|S|}} \geq e^{-A|S|} \frac{2^{|S|}-1}{2^{|S|}} > 0. \quad (5.14)$$

□

Mirror-Image KDE To correct the boundary bias previously introduced, a solution is to add a correction to the estimator $\hat{f}_{h,S}^{\text{KDE}}$ near the boundary of the domain of definition. Let us define three mirroring operations for a number $x \in [0, 1]$

$$M^{-1}(x) = -x; \quad M^0(x) = x; \quad M^1(x) = 2 - x. \quad (5.15)$$

We define the mirror-image KDE as a KDE constructed over the sample W augmented with mirror reflections of each point over all axis. An illustration of this operation in dimension 2 is given by fig. 1.

This estimator is an extension to every dimension of the one proposed in [6]. The formal definition is

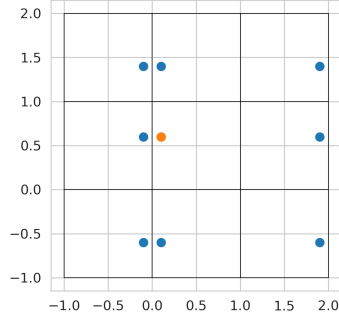
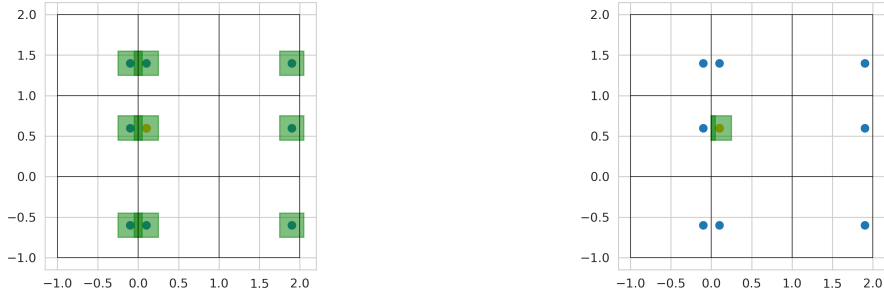


Figure 1: A datapoint in $[0, 1]^d$ (in orange) and his 8 mirror-images (in blue)



(a) A kernel is fitted over all points of the augmented dataset

(b) Restriction to the unit hypercube

Figure 2: Construction of the mirror-image KDE

$$\hat{f}_{m,S}^{\text{MI}}(x) = \mathbb{1}_{[0,1]^{|S|}}(x) \frac{1}{mh^{|S|}} \sum_{i=1}^m \sum_{a \in \{-1,0,1\}^{|S|}} \prod_{k \in S} K\left(\frac{M^{a_k}(W_i)_k - x_k}{h}\right). \quad (5.16)$$

It consists in summing multidimensional kernel over the points of the augmented samples and restrict the domain of the obtained function to $[0, 1]^{|S|}$ as illustrated in fig. 2. Roughly speaking it consists in flipping the part of $\hat{f}_{h,S}^{\text{KDE}}$ that fall outside $[0, 1]^{|S|}$ inside it. $\hat{f}_{m,S}^{\text{MI}}$ is supported on $[0, 1]^{|S|}$ and $\int_{[0,1]^{|S|}} \hat{f}_{m,S}^{\text{MI}}(x) dx = 1$.

Bias for mirror-image KDE Under an ad-hoc condition on the partial derivatives of f_S at the boundary of $[0, 1]^{|S|}$ it is possible to bound the bias for the mirror-image KDE. Our result is an extension of the lemma 3.1 in [6] to every dimension and every $\beta \in (0, 2]$ while the analysis in the original paper was restricted to bi-dimensional densities and $\beta = 2$. With our proof strategy, we find a better constant in the upper-bound for $|S| = 2$ and $\beta = 2$.

Proposition 5.2: Bias for mirror-image KDE

Let us assume that for all sequence $(x_n)_{n \in \mathbb{N}}$ in $[0, 1]^{|S|}$, if x_n converges to a boundary point of $[0, 1]^{|S|}$, then for all $k \in S$ $\lim_{n \rightarrow \infty} \partial_k f_S(x_n) = 0$. Then

$$\left\| f_S - \mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] \right\|_{\infty} \leq C_1 h^{\beta} \quad (5.17)$$

where $C_1 = L|S|^{\beta/2} (2\|K\|_{\infty})^{|S|}$ if $\beta < 2$ and $C_1 = L|S|$ if $\beta = 2$

Proof. We define f_S^{MI} as the function defined over $[-1, 2]^{|S|}$ such that for all $x \in [0, 1]^{|S|}$ and $a \in \{-1, 0, 1\}^{|S|}$

$$f_S^{\text{MI}}(M^a(x)) = f_S(X) \quad (5.18)$$

where $M^a(x) = (M^{a_k}(x_k))_{k \in S}$. The property that the partial derivatives of f_S vanish near the boundary of $[0, 1]^{|S|}$ ensures that $\partial_k f_S^{\text{MI}}$ is continuous on $(-1, 2)^{|S|}$ and so $f_S^{\text{MI}} \in \mathcal{H}(2, L)$.

Let $x \in [0, 1]^{|S|}$, we want to bound $\left| f_S(x) - \mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}}(x) \right] \right|$. Assume first that $x \in [0, 1/2]^{|S|}$ and denote by \mathcal{A} the set $\{k \in S : x_k < h\}$. We start by considering the situation where $|\mathcal{A}| \geq 1$. For all $k \in \mathcal{A}$ and all $t \in [0, 1]$, $K\left(\frac{t-(2-x_k)}{h}\right) = 0$ because the support of K is $[-1, 1]$, $h \leq 1/2$ and $x_k < h$. For all $k \in S \setminus \mathcal{A}$ and all $t \in [0, 1]$, $K\left(\frac{t-(2-x_k)}{h}\right) = 0$ and $K\left(\frac{t-(-x_k)}{h}\right) = 0$. Then the expected value of $\hat{f}_{m,S}^{\text{MI}}$ at the point x can be written as

$$\mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] (x) = \sum_{\mathcal{B} \subset \mathcal{A}} \frac{1}{h^{|S|}} \int_{[0,1]^{|S|}} \prod_{k \in \mathcal{B}} K\left(\frac{t_k + x_k}{h}\right) \prod_{k \in S \setminus \mathcal{B}} K\left(\frac{t_k - x_k}{h}\right) f_S(t) dt \quad (5.19)$$

Now, for $\mathcal{B} \in \mathcal{A}$, we denote $x_{\mathcal{B}}$ the vector such that $(x_{\mathcal{B}})_k = x_k$ if $k \notin \mathcal{B}$ and $(x_{\mathcal{B}})_k = -x_k$ if $k \in \mathcal{B}$. We have

$$\mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] (x) = \sum_{\mathcal{B} \subset \mathcal{A}} \int_{\chi_{\mathcal{B}}^S} \prod_{k \in S} K(u_k) f_S(x_{\mathcal{B}} + uh) du \quad (5.20)$$

where $\chi_{\mathcal{B}}^S = \{u \in [-1, 1]^{|S|} : x_{\mathcal{B}} + uh \in [0, 1]^{|S|}\}$. We see that $\chi_{\mathcal{B}}^S = \prod_{k \in S} [\underline{u}_k, \bar{u}_k]$ where $\underline{u}_k = -x_k/h$ if $k \in \mathcal{B}$, -1 otherwise and $\bar{u}_k = -x_k/h$ if $k \in \mathcal{A} \setminus \mathcal{B}$, 1 otherwise. What is more, as $f_S^{\text{MI}} = f_S$ on $[0, 1]^{|S|}$ we have

$$\mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] (x) = \sum_{\mathcal{B} \subset \mathcal{A}} \int_{\chi_{\mathcal{B}}^S} \prod_{k \in S} K(u_k) f_S^{\text{MI}}(x_{\mathcal{B}} + uh) du. \quad (5.21)$$

Now, as $(\chi_{\mathcal{B}}^S)_{\mathcal{B} \subset \mathcal{A}}$ forms a partition of $[-1, 1]^{|S|}$, we have

$$f_S(x) = \sum_{\mathcal{B} \subset \mathcal{A}} \int_{\chi_{\mathcal{B}}^S} \prod_{k \in S} K(u_k) f_S(x) du \quad (5.22)$$

$$= \sum_{\mathcal{B} \subset \mathcal{A}} \int_{\chi_{\mathcal{B}}^S} \prod_{k \in S} K(u_k) f_S^{\text{MI}}(x_{\mathcal{B}}) du \quad (5.23)$$

We denote by $\delta_{\mathcal{B}}(u, \beta)$ the quantity

$$\begin{cases} f_S^{\text{MI}}(x_{\mathcal{B}} + uh) - f_S^{\text{MI}}(x_{\mathcal{B}}) & \text{if } \beta \in (0, 1] \\ f_S^{\text{MI}}(x_{\mathcal{B}} + uh) - f_S^{\text{MI}}(x_{\mathcal{B}}) - h \sum_{k \in S} u_k \partial_k f_S^{\text{MI}}(x_{\mathcal{B}}) & \text{if } \beta \in (1, 2] \end{cases} \quad (5.24)$$

From lemma 4.1 we have

$$|\delta_{\mathcal{B}}(u, \beta)| \leq Lh^\beta \|u\|^\beta \quad (5.25)$$

And, as $\int xK(x) = 0$, we have

$$\left| f_S(x) - \mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] (x) \right| = \left| \sum_{\mathcal{B} \subset \mathcal{A}} \int_{\mathcal{X}_{\mathcal{B}}^S} \prod_{k \in S} K(u_k) \delta_{\mathcal{B}}(u) du \right| \quad (5.26)$$

$$\leq Lh^\beta \int_{[-1,1]^{|S|}} \prod_{k \in S} K(u_k) \|u\|^\beta du \quad (5.27)$$

If $\beta = 2$

$$\int_{[-1,1]^{|S|}} \prod_{k \in S} K(u_k) \|u\|^\beta du = \int_{[-1,1]^{|S|}} \prod_{k \in S} K(u_k) \sum_{k \in S} u_k^2 du \quad (5.28)$$

$$= \sum_{k \in S} \int_{-1}^1 K(u) u^2 du \quad (5.29)$$

$$\leq \sum_{k \in S} \int_{-1}^1 K(u) du = |S|. \quad (5.30)$$

If $\beta < 2$

$$\int_{[-1,1]^{|S|}} \prod_{k \in S} K(u_k) \|u\|^\beta du \leq \|K\|_\infty^{|S|} \int_{[-1,1]^{|S|}} \|u\|^\beta du \quad (5.31)$$

$$\leq \|K\|_\infty^{|S|} \sqrt{|S|}^\beta \int_{[-1,1]^{|S|}} du \quad (5.32)$$

$$= |S|^{\beta/2} (2\|K\|_\infty)^{|S|}. \quad (5.33)$$

Then $\sup_{x \in [0, 1/2]^{|S|}} \left| f_S(x) - \mathbb{E} \left[\hat{f}_{m,S}^{\text{MI}} \right] (x) \right| \leq C_1 h^\beta$. By symmetry the same inequality is true when the sup is taken over $[0, 1]^{|S|}$.

□

Then considering the mirror-image KDE leads to a correction of the boundary issue previously mentioned.

5.1.2 Variance upper-bound

To upper-bound the variance of the mirror-image KDE, we will use corollary 15 of [4]. Our setting is not the same as in this paper as we deal with mirror-image KDE. Then in order to obtain the same result, we must ensure that the family of functions

$$\mathcal{F}_S^{\text{MI}} = \{K_{x,h}^{\text{MI}} | x \in [0, 1]^{|S|}, h \in (0, 1/2)\} \quad (5.34)$$

where

$$K_{x,h}^{\text{MI}} : u \mapsto \frac{1}{h^{|S|}} \mathbb{1}_{[0,1]^{|S|}} \sum_{a \in \{-1,0,1\}^{|S|}} \prod_{k \in S} K\left(\frac{M^{a_k}(u_k) - x_k}{h}\right) \quad (5.35)$$

is a bounded VC class of function. We know that \mathcal{F}_S is a bounded VC class of function. The results of section 2.6 of [11] indicate that a family of functions is a bounded VC class if and only if the associated collection of sublevels is a VC class of sets. Now, we remark that the sublevels of functions in $\mathcal{F}_S^{\text{MI}}$ can be written as intersections of sublevels of functions in \mathcal{F}_S intersected with $[0, 1]^{|S|}$. Then, as intersections preserve the VC class property for collection of sets (see [10]), $\mathcal{F}_S^{\text{MI}}$ is a bounded VC class of functions, and the corollary 15 of [4] applies, leading to the following result.

Proposition 5.3: Variance

Let $h_{m,S}$ be a bandwidth in $(0, 1/2)$ and $\delta_m \in (0, 1)$. With probability $1 - \delta_m$

$$\left\| \hat{f}_{S,h_{m,S}} - \mathbb{E} \left[\hat{f}_{m,h_{m,S}} \right] \right\|_{\infty} \leq C_2 \sqrt{\frac{\log(1/h_{m,S}) + \log(2/\delta_m)}{mh_{m,S}^{|S|}}}. \quad (5.36)$$

The constant C_2 depends on $|S|$, on $\|K\|_{\infty}$ and on $\|K'\|_{\infty}$.

5.1.3 Conclusion

Now, as we have for a control of the bias and the variance term for every bandwidth $h_{m,S} \in (0, 1/2)$, by choosing appropriately $h_{m,S}$ it is possible to bound $\|f_S - \hat{f}_{S,h_{m,S}}\|_\infty$.

Proposition 5.4: Convergence

Choosing $h_{m,S} \asymp (1/m)^{\frac{1}{2\beta+|S|}}$, it exists a constant C_S such that with probability at least $1 - \delta_m$

$$\|f_S - \hat{f}_{S,h_{m,S}}\|_\infty \leq C_S \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{\beta}{2\beta+|S|}}. \quad (5.37)$$

5.2 Uniformity over Set_d^k

We have just established a control in high probability for the quantity $\|f_S - \hat{f}_{S,h_{m,S}}\|_\infty$ for a given S . Our objective is to have such a control uniformly over Set_d^k . Applying a union-bound, we obtain the following result.

Proposition 5.5: Uniform control of uniform error over all subsets

Let us denote $C_k = \max_{S \in \text{Set}_d^k} C_S$. We have, with probability at least $1 - S_d^k \delta_m$

$$\sup_{S \in \text{Set}_d^k} \left\| \hat{f}_{S,h_m^S} - f_S \right\|_\infty \leq C_k \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{2}{4+k}}. \quad (5.38)$$

In particular with the choice $\delta_m = \frac{2}{m S_d^k}$, we have that with probability $1 - \frac{2}{m}$

$$\sup_{S \in \text{Set}_d^k} \left\| \hat{f}_{S,h_m^S} - f_S \right\|_\infty \leq C_k \sqrt{3 \log m + 2 \log S_d^k} \left(\frac{1}{m}\right)^{\frac{2}{4+k}}. \quad (5.39)$$

6 THEOREM

Let now m_0 be the smallest integer m such that

$$C_k \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{2}{4+k}} \leq e^{-A|S|} (1 - e^{-A|S|}). \quad (6.1)$$

If $m \geq m_0$, we know that on an event \mathcal{A}_m^k of probability at least $1 - S_d^k \delta_m$

$$\|f_S - \hat{f}_{S, h_m}\|_\infty \leq C_k \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{\beta}{2\beta+k}}. \quad (6.2)$$

Then, on \mathcal{A}_m^k (UC) is satisfied with $\epsilon_S = C_k \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{\beta}{2\beta+k}}$ for all $S \in \text{Set}_d^k$. As a consequence, using proposition 4.2, we have

$$\sum_{S \in \mathcal{P}_*} \text{KL} \left(f_S \| \hat{f}_{S, h_m^s} \right) \leq e^{2Ak} |\mathcal{P}_*| C_k \sqrt{\log m + 2 \log(2/\delta_m)} \left(\frac{1}{m}\right)^{\frac{2}{4+k}}. \quad (6.3)$$

And, on \mathcal{A}_m^k , for $\delta_n \in (0, 1/S_d^k)$ with probability $1 - S_d^k \delta_n$

$$(P - P_n)(\log \hat{f}_{\hat{\mathcal{P}}} - \log \hat{f}_{\hat{\mathcal{P}}}) \leq 2d \sqrt{\frac{Ak}{n}} \sqrt{\log(2/\delta_n)}. \quad (6.4)$$

Now, with the choice $\delta_m = 1/(2S_d^k m)$ and $\delta_n = 1/(2S_d^k n)$ we obtain the following result

Theorem 6.1: Final bound

If for all $S \in \text{Set}_d^k$, $f \in \mathcal{H}(2, L)$, and for all sequence $(x_n)_{n \in \mathbb{N}}$ such that $\lim_{n \rightarrow \infty} x_n = x^*$ where x^* belongs to the boundary of $[0, 1]^{|S|}$ and for all $k \in S$ $\lim_{n \rightarrow \infty} \partial_k f_S(x_n) = 0$. With the choice $\hat{f}_S = \hat{f}_{h_{m,S}, S}^{\text{MI}}$ where $h_{m,S} \approx \left(\frac{1}{m}\right)^{\frac{1}{2+|S|}}$, we have with probability at least $(1 - 1/m)(1 - 1/n)$

$$\begin{aligned} \text{KL} \left(f \parallel \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \parallel f_{\mathcal{P}_*} \right) &\leq e^{2Ak} \sqrt{2} |\mathcal{P}_*| C_k \sqrt{\log m + \log (S_d^k)} \left(\frac{1}{m} \right)^{\frac{\beta}{2\beta+k}} \\ &\quad + 2d \sqrt{\log n + \log (S_d^k)} \sqrt{\frac{Ak}{n}} \end{aligned} \quad (6.5)$$

Ignoring logarithmic factors, the rate of convergence of the approximation term is $\left(\frac{1}{m}\right)^{\frac{\beta}{2\beta+k}}$. The dependence of this quantity in k illustrates that ISDE tackles the curse of dimensionality for the density estimation problem under KL loss in the same spirit that [9] showed that his estimator does for the squared L_2 loss.

Ignoring logarithmic factors again, the rate of convergence of the selection term is $\frac{1}{\sqrt{n}}$. This is a classical rate of convergence for hold-out procedures with bounded loss (see corollary 8.8 in [7]).

The term $\log(S_d^k)$ in the upper-bound illustrates the combinatorial complexity reduction operated by ISDE. The presence of the log of the number of hypotheses is classical for hold-out procedures with bounded loss (see again corollary 8.8 in [7]). In our context, we have reduced the combinatorial complexity from the number of partitions P_d^k to the number of subsets S_d^k .

7 CONCLUSION

In this paper, we have studied the convergence properties of ISDE. In particular, we have shown that under suitable assumptions on the true density and for the mirror-image KDE as marginal density estimator, we can provide an upper-bound

valid with high-probability of the quantity

$$\text{KL} \left(f \| \hat{f}_{\hat{\mathcal{P}}} \right) - \text{KL} \left(f \| f_{\mathcal{P}_*} \right). \quad (7.1)$$

This bound highlights how ISDE tackles the curse of dimensionality and reduces the combinatorial complexity of the density estimation problem under IS compared to a brute-force approach. These results offer a theoretical validation of the empirical observations presented in [8].

To complete the study, it let to study how the bias term $\text{KL} \left(f \| f_{\mathcal{P}_*} \right)$ behaves. It is hard to give a precise statement on this quantity in a general setting. One simple situation is when $f \in \mathcal{D}_k^d$. In this case $\text{KL} \left(f \| f_{\mathcal{P}_*} \right) = 0$. This bias can also be explicitly evaluated in some multivariate Gaussian frameworks, see appendix A.

Acknowledgement This work was supported by the program Paris Region Ph.D. of DIM Mathinnov and was partly supported by the French ANR Chair in Artificial Intelligence TopAI - ANR-19-CHIA-0001. The author is thankful to Marc Glisse and Pascal Massart for their constructive remarks on this work.

References

- [1] Evarist Giné and Armelle Guillou. Rates of strong uniform consistency for multivariate kernel density estimators. In *Annales de l'Institut Henri Poincaré (B) Probability and Statistics*, volume 38, pages 907–921. Elsevier, 2002.
- [2] Peter Hall. On kullback-leibler loss and density estimation. *The Annals of Statistics*, pages 1491–1519, 1987.
- [3] Rafael Hasminskii, Ildar Ibragimov, et al. On density estimation in the view of kolmogorov's ideas in approximation theory. *The Annals of Statistics*, 18(3):999–1010, 1990.
- [4] Jisu Kim, Jaehyeok Shin, Alessandro Rinaldo, and Larry Wasserman. Uniform convergence rate of the kernel density estimator adaptive to intrinsic volume

- dimension. In *International Conference on Machine Learning*, pages 3398–3407. PMLR, 2019.
- [5] Oleg Lepski. Multivariate density estimation under sup-norm loss: oracle approach, adaptation and independence structure. *Annals of Statistics*, 41(2):1005–1034, 2013.
- [6] Han Liu, Larry Wasserman, and John Lafferty. Exponential concentration for mutual information estimation with application to forests. *Advances in Neural Information Processing Systems*, 25, 2012.
- [7] Pascal Massart. Concentration inequalities and model selection. 2007.
- [8] Louis Pujol. Isde: Independence structure density estimation. *arXiv preprint arXiv:2203.09783*, 2022.
- [9] Gilles Rebelles et al. Lp adaptive estimation of an anisotropic density under independence hypothesis. *Electronic journal of statistics*, 9(1):106–134, 2015.
- [10] Aad Van Der Vaart and Jon A Wellner. A note on bounds for vc dimensions. *Institute of Mathematical Statistics collections*, 5:103, 2009.
- [11] AW van der Vaart, A.W. van der Vaart, A. van der Vaart, and J. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, 1996.

A Bias term in the multivariate Gaussian framework

In this appendix, we study the bias $\text{KL}(f \| f_{\mathcal{P}^*})$ in a multivariate Gaussian framework where exact computations are possible.

A.1 Model and notations

Model If Σ denotes a $d \times d$ definite positive matrix, we denote by f_{Σ} the density of a multivariate centered Gaussian random variable with covariance Σ and by $\Sigma_{\mathcal{P}}$ the matrix defined as follows.

$$\Sigma_{\mathcal{P}}(i, j) = \begin{cases} \Sigma(i, j) & \text{if } i \text{ and } j \text{ belongs to the same bloc in } \mathcal{P} \\ 0 & \text{else.} \end{cases} \quad (\text{A.1})$$

If S_1 and S_2 are subsets of $\{1, \dots, d\}$, we denote by $\Sigma(S_1, S_2)$ the $|S_1| \times |S_2|$ submatrix matrix of Σ where we keep the intersection of rows in S_1 and columns in S_2 , to keep notations compact, we write $\Sigma(S)$ instead of $\Sigma(S, S)$

For a multivariate Gaussian random variable with covariance Σ , $f_{\Sigma} \in \mathcal{D}_d^k$ is equivalent to the fact that it exists a permutation σ of $\{1, \dots, d\}$ such that $P_{\sigma} \Sigma P_{\sigma}^{-1}$ is block-diagonal with blocks of size smaller than $k \times k$. For clarity, in what follows, we will always consider that this property is met with $\sigma = \text{id}$, meaning that we restrict ourselves to partitions in which each block is made of consecutive features. This does not imply a loss of generality.

We now consider that $f = f_{\Sigma}$ and

$$\Sigma = \Sigma_{\mathcal{P}} + \epsilon \quad (\text{A.2})$$

where $\Sigma_{\mathcal{P}}$ is a block-diagonal covariance matrix corresponding to the independence structure \mathcal{P} and ϵ is a “small” (in a sense to be defined later) definite positive matrix. The question is how this perturbation influences the bias term. In order to answer it, we must control $\text{KL}(f \| f_{\mathcal{P}})$ for all \mathcal{P} in Part_d^k .

A.2 Some useful lemmas

Computation of KL losses The first useful result is an explicit computation of $\text{KL}(f_\Sigma \| f_{\Sigma_{\mathcal{P}}})$ for any \mathcal{P} in Part_d^k .

Lemma A.1: Exact computation of KL between two centered multivariate Gaussian

For every $\mathcal{P} \in \text{Part}_d^k$

$$\text{KL}(f_\Sigma \| f_{\Sigma_{\mathcal{P}}}) = \frac{1}{2} \left(\sum_{S \in \mathcal{P}} \log \det \Sigma(S) - \log \det \Sigma \right) \quad (\text{A.3})$$

Or if $\lambda_1 \leq \dots \leq \lambda_d$ are the eigenvalues of Σ and $\lambda_1^{\mathcal{P}} \leq \dots \leq \lambda_d^{\mathcal{P}}$ the eigenvalues of $\Sigma_{\mathcal{P}}$

$$\text{KL}(f_\Sigma \| f_{\Sigma_{\mathcal{P}}}) = \frac{1}{2} \sum_{i=1}^d \log \left(\frac{\lambda_i^{\mathcal{P}}}{\lambda_i} \right) \quad (\text{A.4})$$

Proof. The density f_Σ has the following expression.

$$f_\Sigma(x) = \frac{1}{(2\pi)^{d/2} (\det \Sigma)^{1/2}} \exp \left(-\frac{1}{2} x^T \Sigma^{-1} x \right). \quad (\text{A.5})$$

We compute the KL divergence between f_Σ and $f_{\Sigma_{\mathcal{P}}}$

$$\text{KL}(f_\Sigma \| f_{\Sigma_{\mathcal{P}}}) = \int \log \left(\frac{f_\Sigma(x)}{f_{\Sigma_{\mathcal{P}}}(x)} \right) f_\Sigma(x) dx \quad (\text{A.6})$$

$$\begin{aligned} &= \frac{1}{2} \log \frac{\det \Sigma_{\mathcal{P}}}{\det \Sigma} \underbrace{\int f_\Sigma(x) dx}_{=1} \\ &\quad + \frac{1}{2} \underbrace{\int x^T \Sigma_{\mathcal{P}}^{-1} x f_\Sigma(x) dx}_{=\text{Tr}(\Sigma_{\mathcal{P}}^{-1} \Sigma)} \\ &\quad + \frac{1}{2} \underbrace{\int x^T \Sigma^{-1} x f_\Sigma(x) dx}_{=\text{Tr}(\Sigma^{-1} \Sigma) = d} \end{aligned} \quad (\text{A.7})$$

$$= \frac{1}{2} (\log \det \Sigma_{\mathcal{P}} - \log \det \Sigma + \text{Tr}(\Sigma_{\mathcal{P}}^{-1} \Sigma) - d) \quad (\text{A.8})$$

Now, we define a permutation $\sigma_{\mathcal{P}}$ of $\{1, \dots, d\}$ such that :

$$\Sigma_{\mathcal{P}} = P_{\sigma_{\mathcal{P}}} \begin{pmatrix} \Sigma(S_1) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \Sigma(S_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \Sigma(S_M) \end{pmatrix} P_{\sigma_{\mathcal{P}}}^{-1} \quad (\text{A.9})$$

where $\{S_1, \dots, S_M\}$ denotes the blocks of \mathcal{P} . It is then clear that $\log \det \Sigma_{\mathcal{P}} = \sum_{S \in \mathcal{P}} \log \det \Sigma(S)$. We also have

$$\Sigma = P_{\sigma_{\mathcal{P}}} \begin{pmatrix} \Sigma(S_1) & \Sigma(S_1, S_2) & \dots & \Sigma(S_1, S_M) \\ \Sigma(S_2, S_1) & \Sigma(S_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & \Sigma(S_{M-1}, S_M) \\ \Sigma(S_M, S_1) & \dots & \Sigma(S_M, S_{M-1}) & \Sigma(S_M) \end{pmatrix} P_{\sigma_{\mathcal{P}}}^{-1}. \quad (\text{A.10})$$

Then

$$\Sigma_{\mathcal{P}}^{-1} \Sigma = P_{\sigma_{\mathcal{P}}} \begin{pmatrix} I_{|S_1|} & X_{1,2} & \dots & X_{1,M} \\ X_{2,1} & I_{|S_2|} & \ddots & \vdots \\ \vdots & \ddots & \ddots & X_{M-1,M} \\ X_{M,1} & \dots & \Sigma(S_M, S_{M-1}) & I_{|S_M|} \end{pmatrix} P_{\sigma_{\mathcal{P}}}^{-1} \quad (\text{A.11})$$

where for $i \neq j$, $X_{i,j}$ is a $|S_i| \times |S_j|$ matrix. Then $\text{Tr}(\Sigma_{\mathcal{P}}^{-1} \Sigma) = d$.

The formulation of the result involving the eigenvalues comes from the fact that $\det \Sigma = \prod_{i=1}^d \lambda_i$ and $\det \Sigma_{\mathcal{P}} = \prod_{i=1}^d \lambda_i^{\mathcal{P}}$. \square

Some computation of determinants We define the $p \times p$ matrix

$$A_{\sigma}^p = \begin{pmatrix} 1 & \sigma & \dots & \sigma \\ \sigma & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \sigma \\ \sigma & \dots & \sigma & 1 \end{pmatrix}. \quad (\text{A.12})$$

If k divides d we define

$$\Sigma_{\sigma}^{(d,k)} = \left(\begin{array}{c|c|c|c} A_{\sigma}^k & \mathbf{0} & \dots & \mathbf{0} \\ \hline \mathbf{0} & \ddots & \ddots & \vdots \\ \hline \vdots & \ddots & \ddots & \mathbf{0} \\ \hline \mathbf{0} & \dots & \mathbf{0} & A_{\sigma}^k \end{array} \right) \quad (\text{A.13})$$

For $\epsilon > 0$ we define

$$\Sigma_{\sigma, \epsilon}^{(d,k)} = \left(\begin{array}{c|c|c|c} A_{\sigma}^k & \epsilon & \dots & \epsilon \\ \hline \epsilon & \ddots & \ddots & \vdots \\ \hline \vdots & \ddots & \ddots & \epsilon \\ \hline \epsilon & \dots & \epsilon & A_{\sigma}^k \end{array} \right) \quad (\text{A.14})$$

Lemma A.2: Property for block matrices

- i $\det A_{\sigma}^p = [1 - \sigma]^{p-1} [1 + (p-1)\sigma]$
- ii $\det \Sigma_{\sigma}^{(d,k)} = [1 - \sigma]^{\frac{d}{k}(k-1)} [1 + (k-1)\sigma]^{\frac{d}{k}}$
- iii $\det \Sigma_{\sigma, \epsilon}^{(d,k)} = [1 - \sigma]^{\frac{d}{k}(k-1)} [1 + (k-1)\sigma + (d-k)\epsilon] [1 + (k-1)\sigma - k\epsilon]^{\frac{d}{k}-1}$

Proof. (i) We start by computing the eigenvalues of A_σ^p . We remark that

$$A - (1 - \sigma)I = \begin{pmatrix} \sigma & \dots & \sigma \\ \vdots & & \vdots \\ \sigma & \dots & \sigma \end{pmatrix}. \quad (\text{A.15})$$

Then it is clear that $x \in \mathbb{R}^p \in \ker(A_\sigma^p - (1 - \sigma)I) \Leftrightarrow x \in \{y \in \mathbb{R}^p : \sum_{i=1}^p y_i = 0\}$, which is a linear subspace of \mathbb{R}^p of dimension $p - 1$.

Then we remark that if $\mathbb{1}_p = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ then

$$A_\sigma^p \mathbb{1}_p = (1 + (p - 1)\sigma) \mathbb{1}_p \quad (\text{A.16})$$

Then $1 + (p - 1)\sigma$ is an eigenvalue of A and it could not be of multiplicity greater than 1 as we have just proven that $(1 - \sigma)$ has a multiplicity of $p - 1$. Using the fact that the determinant is the product of the eigenvalues, we obtain

$$\det A_\sigma^p = (1 - \sigma)^{p-1} (1 + (p - 1)\sigma) \quad (\text{A.17})$$

The proof of (ii) follows immediately as the determinant of a block-diagonal matrix is the product of the derminants of the diagonal blocks.

(iii) We determine the eigenvalues of $\Sigma_{\sigma,\epsilon}^{(d,k)}$. To this end we will find a set of d linearly independent eigenvectors. We remark that

$$\Sigma_{\sigma,\epsilon}^{(d,k)} \mathbb{1}_d = (1 + (k - 1)\sigma + (d - k)\epsilon) \mathbb{1}_d. \quad (\text{A.18})$$

Then $1 + (k - 1)\sigma + (d - k)\epsilon$ is an eigenvalue of $\Sigma_{\sigma,\epsilon}^{(d,k)}$ with multiplicity at least 1. Now, we remark that for all integer $0 \leq i \leq \frac{d}{k} - 1$ and $2 \leq j \leq k$ we have

$$\Sigma_{\sigma,\epsilon}^{(d,k)} (e_{ik+1} - e_{ik+j}) = (1 - \sigma)(e_{ik+1} - e_{ik+j}). \quad (\text{A.19})$$

Then $(1 - \sigma)$ is an eigenvalue of $\Sigma_{\sigma,\epsilon}^{(d,k)}$ with multiplicity at least $\frac{d}{k}(k - 1)$. Finally, if for $i < j$ we denote by $\mathbb{1}_j^i = \sum_{k=i}^j e_k$, we remark that for all integer $1 \leq i \leq \frac{d}{k} - 1$

$$\Sigma_{\sigma,\epsilon}^{(d,k)} \left(\mathbb{1}_k^0 - \mathbb{1}_{(i+1)k}^{ik+1} \right) = (1 + (k - 1)\sigma - k\epsilon) \left(\mathbb{1}_k^0 - \mathbb{1}_{(i+1)k}^{ik+1} \right). \quad (\text{A.20})$$

Then $1 + (k - 1)\sigma + (d - k)\epsilon$ is an eigenvalue of $\Sigma_{\sigma,\epsilon}^{(d,k)}$ with multiplicity at least $\frac{d}{k} - 1$.

As $1 + \frac{d}{k}(k - 1) + (\frac{d}{k} - 1) = d$, we now that the eigenvalues of $\Sigma_{\sigma,\epsilon}^{(d,k)}$ are $1 + (k - 1)\sigma + (d - k)\epsilon$, $1 - \sigma$ and $1 + (k - 1)\sigma + (d - k)\epsilon$ with multiplicity 1, $\frac{d}{k}(k - 1)$ and $\frac{d}{k} - 1$.

□

A.3 Control of the bias

Almost independence structure the following property precise the KL loss between $f_{\Sigma_{\sigma,\epsilon}^{(d,k)}}$ and $f_{\Sigma_{\sigma}^{(d,k)}}$, with a particular look at the situation where $\epsilon \rightarrow 0$.

Proposition A.3: Almost independence

$$\begin{aligned} \text{KL} \left(f_{\Sigma_{\sigma,\epsilon}^{(d,k)}} \| f_{\Sigma_{\sigma}^{(d,k)}} \right) &= -\frac{1}{2} \log \left(1 + \frac{d - k}{1 + (k - 1)\sigma} \epsilon \right) \\ &\quad - \frac{1}{2} \left(\frac{d}{k} - 1 \right) \log \left(1 - \frac{k}{1 + (k - 1)\sigma} \epsilon \right) \end{aligned} \quad (\text{A.21})$$

At the limit $\epsilon \rightarrow 0$

$$\text{KL} \left(f_{\Sigma_{\sigma,\epsilon}^{(d,k)}} \| f_{\Sigma_{\sigma}^{(d,k)}} \right) \underset{\epsilon \rightarrow 0}{=} \frac{d(d - k)}{4(1 + (k - 1)\sigma)^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.22})$$

Proof. As $\Sigma_{\sigma}^{(d,k)}$ is a block-diagonal submatrix of $\Sigma_{\sigma,\epsilon}^{(d,k)}$, using lemma A.1 we have

$$\text{KL} \left(\Sigma_{\sigma,\epsilon}^{(d,k)} \| \Sigma_{\sigma}^{(d,k)} \right) = \frac{1}{2} \log \left(\frac{\det \Sigma_{\sigma}^{(d,k)}}{\det \Sigma_{\sigma,\epsilon}^{(d,k)}} \right) \quad (\text{A.23})$$

Let $\beta = 1 + (k - 1)\sigma$ Now, using lemma A.2, we have

$$\text{KL}(\Sigma_{\sigma, \epsilon}^{(d,k)} \parallel \Sigma_{\sigma}^{(d,k)}) = \frac{1}{2} \left[\log \left(\frac{\beta^{\frac{d}{k}}}{[\beta + (d - k)\epsilon][\beta - k\epsilon]^{\frac{d}{k} - 1}} \right) \right] \quad (\text{A.24})$$

$$= \frac{1}{2} \left[\frac{d}{k} \log \beta - \log(\beta + (d - k)\epsilon) - \left(\frac{d}{k} - 1 \right) \log(\beta - k\epsilon) \right] \quad (\text{A.25})$$

$$= \frac{d}{2k} \log \beta - \frac{\log \beta}{2} - \frac{1}{2} \log \left(1 + \frac{d - k}{\beta} \epsilon \right) - \frac{1}{2} \left(\frac{d}{k} - 1 \right) \log \beta - \frac{1}{2} \left(\frac{d}{k} - 1 \right) \log \left(1 - \frac{k}{\beta} \epsilon \right) \quad (\text{A.26})$$

$$= -\frac{1}{2} \log \left(1 + \frac{d - k}{\beta} \epsilon \right) - \frac{1}{2} \left(\frac{d}{k} - 1 \right) \log \left(1 - \frac{k}{\beta} \epsilon \right) \quad (\text{A.27})$$

Now, we use that

$$\log \left(1 + \frac{d - k}{\beta} \epsilon \right) \underset{\epsilon \rightarrow 0}{=} \frac{d - k}{\beta} \epsilon - \frac{(d - k)^2}{2\beta^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.28})$$

and

$$\log \left(1 - \frac{k}{\beta} \epsilon \right) \underset{\epsilon \rightarrow 0}{=} -\frac{k}{\beta} \epsilon - \frac{k^2}{2\beta^2} \epsilon^2 + o(\epsilon^2). \quad (\text{A.29})$$

And we obtain

$$\begin{aligned} \text{KL}(\Sigma_{\sigma, \epsilon}^{(d,k)} \parallel \Sigma_{\sigma}^{(d,k)}) &\underset{\epsilon \rightarrow 0}{=} -\frac{d - k}{2\beta} \epsilon + \frac{(d - k)^2}{4\beta^2} \epsilon^2 \\ &\quad + \frac{\left(\frac{d}{k} - 1\right) k}{2\beta} \epsilon + \frac{\left(\frac{d}{k} - 1\right) k^2}{4\beta^2} \epsilon^2 + o(\epsilon^2) \end{aligned} \quad (\text{A.30})$$

$$\underset{\epsilon \rightarrow 0}{=} \frac{(d - k)^2 + kd - k^2}{4\beta^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.31})$$

$$\underset{\epsilon \rightarrow 0}{=} \frac{d^2 - 2kd + k^2 + kd - k^2}{4\beta^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.32})$$

$$\underset{\epsilon \rightarrow 0}{=} \frac{d(d - k)}{4\beta^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.33})$$

□

Optimal structure for small k The following proposition establish that if $\Sigma = A_\sigma^d$ and $k < d$, \mathcal{P}_* is composed of a maximum number of blocks of size k .

Proposition A.4: Optimal Structure

Suppose that $\Sigma = A_\sigma^d$. A structure $s = (s_i)_{i=1}^M$ is a list of positive integer with $\sum_{i=1}^M s_i = d$. To a structure is associated a partition with blocks of consecutive features with size s_1, \dots, s_M . For any structure s we have

$$\text{KL}(f_\Sigma \| f_{\Sigma_s}) = \frac{1}{2} \left(\sum_{i=1}^M \log \left(\frac{1 + (s_i - 1)\sigma}{1 - \sigma} \right) - \log \left(\frac{1 + (d - 1)\sigma}{1 - \sigma} \right) \right). \quad (\text{A.34})$$

If we denote by p and r the only integers such that $d = pk + r$ where $r < k$, we have

$$s^* = \underbrace{(k, \dots, k)}_{p \text{ times}}, r \quad (\text{A.35})$$

Proof. We combine lemma A.1 and lemma A.2 to obtain

$$\text{KL}(f_\Sigma \| f_{\Sigma_s}) = \frac{1}{2} \left(\sum_{S \in \mathcal{P}} \log \det A_\sigma^{s_i} - \log \det A_\sigma^d \right) \quad (\text{A.36})$$

$$= \frac{1}{2} \left[\sum_{S \in \mathcal{P}} (s_i - 1) \log(1 - \sigma) + \log(1 + (s_i - 1)\sigma) - (d - 1) \log(1 - \sigma) - \log(1 + (d - 1)\sigma) \right] \quad (\text{A.37})$$

$$= \frac{1}{2} \left(\underbrace{\sum_{i=1}^M (s_i - 1)}_{=d-M} - (d - 1) \right) \log(1 - \sigma) + \frac{1}{2} \left(\sum_{i=1}^M \log(1 + (s_i - 1)\sigma) - \log(1 + (d - 1)\sigma) \right) \quad (\text{A.38})$$

$$= \frac{1}{2} \left(\sum_{i=1}^M \log \left(\frac{1 + (s_i - 1)\sigma}{1 - \sigma} \right) - \log \left(\frac{1 + (d - 1)\sigma}{1 - \sigma} \right) \right) \quad (\text{A.39})$$

Now we want to prove that the structure minimizing $\text{KL}(f_\Sigma \| f_{\Sigma_s})$ is $(\underbrace{k, \dots, k}_p, r)$.

To do so we start by remarking that for any $s = (s_i)_{i=1}^M \neq (k, \dots, k, r)$ it exists $i \neq j$ such that $s_i \neq k$ and $s_j \neq k$. We will prove that, for our minimization problem it is always possible to find a better structure \tilde{s} with the following

- i if $s_i + s_j \leq k$, $\tilde{S} = (s_k)_{k \notin \{i,j\}} \cup (s_i + s_j)$
- ii if $\exists l > 0 : s_i + s_j = k + l$, $\tilde{S} = (s_k)_{k \notin \{i,j\}} \cup (k, l)$

To prove (i), we start from

$$2\text{KL}(f_\Sigma \| f_{\Sigma_s}) = \sum_{k=1}^M \log \left(\frac{1 + (s_k - 1)\sigma}{1 - \sigma} \right) - \log \left(\frac{1 + (d - 1)\sigma}{1 - \sigma} \right) \quad (\text{A.40})$$

$$\begin{aligned} 2\text{KL}(f_\Sigma \| f_{\Sigma_{\tilde{s}}}) &= \sum_{k=1, \dots, M, k \notin \{i,j\}} \log \left(\frac{1 + (s_k - 1)\sigma}{1 - \sigma} \right) + \log \left(\frac{1 + (s_i + s_j - 1)\sigma}{1 - \sigma} \right) \\ &\quad - \log \left(\frac{1 + (d - 1)\sigma}{1 - \sigma} \right). \end{aligned} \quad (\text{A.41})$$

Then to prove that $\text{KL}(f_\Sigma \| f_{\Sigma_s}) > \text{KL}(f_\Sigma \| f_{\Sigma_{\tilde{s}}})$ it is sufficient to prove that for all $a, b \geq 1$, $g(a, b) > 0$ where

$$g(a, b) = \log \left(\frac{1 + (a - 1)\sigma}{1 - \sigma} \right) + \log \left(\frac{1 + (b - 1)\sigma}{1 - \sigma} \right) - \log \left(\frac{1 + (a + b - 1)\sigma}{1 - \sigma} \right). \quad (\text{A.42})$$

Let us start by computing $\partial_1 g(a, b)$

$$\partial_1 g(a, b) = \frac{\sigma}{1 + (a - 1)\sigma} - \frac{\sigma}{1 + (a + b - 1)\sigma} \quad (\text{A.43})$$

$$= \frac{\sigma}{(1 + (a - 1)\sigma)(1 + (a + b - 1)\sigma)} (1 + (a + b - 1)\sigma - 1 - (a - 1)\sigma) \quad (\text{A.44})$$

$$= \frac{b\sigma^2}{(1 + (a - 1)\sigma)(1 + (a + b - 1)\sigma)} \geq 0. \quad (\text{A.45})$$

Then for any $b \geq 1$, $g(a, b)$ is nondecreasing in a . As a and b play similar roles in $g(a, b)$, we have that for any $a \geq 1$, $g(a, b)$ is nondecreasing in b . To prove that

$g(a, b) \geq 0$ it is sufficient to show that $g(1, 1) > 0$.

$$g(1, 1) = \log\left(\frac{1}{1-\sigma}\right) + \log\left(\frac{1}{1-\sigma}\right) - \log\left(\frac{1+\sigma}{1-\sigma}\right) \quad (\text{A.46})$$

$$= -\log((1-\sigma)(1+\sigma)). \quad (\text{A.47})$$

Now, as $\sigma \in (0, 1)$, $(1-\sigma)(1+\sigma) \in (0, 1)$. Then $\log((1-\sigma)(1+\sigma)) < 0$, implying $g(1, 1) > 0$.

To prove (ii) we start from

$$2\text{KL}(f_\Sigma \| f_{\Sigma_s}) = \sum_{k=1}^M \log\left(\frac{1+(s_k-1)\sigma}{1-\sigma}\right) - \log\left(\frac{1+(d-1)\sigma}{1-\sigma}\right) \quad (\text{A.48})$$

$$\begin{aligned} 2\text{KL}(f_\Sigma \| f_{\Sigma_s}) &= \sum_{k=1, \dots, M, k \notin \{i, j\}} \log\left(\frac{1+(s_k-1)\sigma}{1-\sigma}\right) + \log\left(\frac{1+(k-1)\sigma}{1-\sigma}\right) \\ &\quad + \log\left(\frac{1+(l-1)\sigma}{1-\sigma}\right) - \log\left(\frac{1+(d-1)\sigma}{1-\sigma}\right). \end{aligned} \quad (\text{A.49})$$

Then to prove that $\text{KL}(f_\Sigma \| f_{\Sigma_s}) > \text{KL}(f_\Sigma \| f_{\Sigma_s})$ it is sufficient to prove that for all $x \in [l, k]$, $h(x)$ attains its minimum at l or k where

$$h(x) = \log(1+(x-1)\sigma) + \log(1+((k+l)-x-1)\sigma). \quad (\text{A.50})$$

Let us start by computing $h'(x)$

$$h'(x) = \sigma \left(\frac{1}{1+(x-1)\sigma} - \frac{1}{1+((k+l)-x-1)\sigma} \right) \quad (\text{A.51})$$

$$= \frac{\sigma(1+((k+l)-x-1)\sigma - 1 - (x-1)\sigma)}{(1+(x-1)\sigma)(1+((k+l)-x-1)\sigma)} \quad (\text{A.52})$$

$$= \frac{\sigma^2}{(1+(x-1)\sigma)(1+((k+l)-x-1)\sigma)}((k+l)-2x). \quad (\text{A.53})$$

Then h increases from l to $(k+l)/2$ and decreases from $(k+l)/2$ to k and $h(l) = h(k)$ the minimum of h is attained on l and k .

□

Conclusion We finish this appendix by establishing a general upper bound of $\text{KL}(f_\Sigma \| f_{\mathcal{P}_*})$ where $\Sigma = \Sigma_{\sigma, \epsilon}^{(d, k^*)}$.

Theorem A.5: Upper-bound for the bias in a multivariate Gaussian framework

If $\Sigma = \Sigma_{\sigma, \epsilon}^{(d, k^*)}$ and if $k < k^*$. Let (p, r) be the unique couple of integer with $0 \leq r < k$ such that $k^* = pk + r$, we have

$$\begin{aligned} \text{KL}(f_{\Sigma} \| f_{\mathcal{P}_*}) &\leq \text{KL}\left(f_{\Sigma_{\sigma, \epsilon}^{(d, k^*)}} \| f_{\Sigma_{\sigma}^{(d, k^*)}}\right) + \frac{dp}{2k^*} \log\left(\frac{1 + (k-1)\sigma}{1-\sigma}\right) \\ &\quad + \frac{d}{2k^*} \log\left(\frac{1 + (r-1)\sigma}{1-\sigma}\right) - \frac{d}{2k^*} \log\left(\frac{1 + (k^*-1)\sigma}{1-\sigma}\right) \end{aligned} \quad (\text{A.54})$$

with

$$\text{KL}\left(f_{\Sigma_{\sigma, \epsilon}^{(d, k)}} \| f_{\Sigma_{\sigma}^{(d, k)}}\right) \underset{\epsilon \rightarrow 0}{=} \frac{d(d-k)}{4(1+(k-1)\sigma)^2} \epsilon^2 + o(\epsilon^2) \quad (\text{A.55})$$

Proof. Let us consider

- the structure $\tilde{s} = (\underbrace{k, \dots, k}_{p \text{ times}}, r)$, and $\mathcal{P}_{\tilde{s}}$ the associated partition of k^* features,
- the structure $s = (\underbrace{\tilde{s}, \dots, \tilde{s}}_{d/k^* \text{ times}})$, and \mathcal{P}_s the associated partition of d features,
- the structure $s_0 = (\underbrace{k^*, \dots, k^*}_{d/k^* \text{ times}})$ and \mathcal{P}_0 the associated partition of d features.

We can upper-bound the bias term $\text{KL}(f_{\Sigma} \| f_{\mathcal{P}_*})$ as follows

$$\text{KL}(f_{\Sigma} \| f_{\mathcal{P}_*}) \leq \text{KL}(f_{\Sigma} \| f_{\mathcal{P}_s}) \quad (\text{A.56})$$

$$= \int \log\left(\frac{f_{\Sigma}}{f_{\mathcal{P}_s}}\right) f_{\Sigma} \quad (\text{A.57})$$

$$= \int \log\left(\frac{f_{\Sigma}}{f_{\mathcal{P}_0}}\right) f_{\Sigma} + \int \log\left(\frac{f_{\mathcal{P}_0}}{f_{\mathcal{P}_s}}\right) f_{\Sigma} \quad (\text{A.58})$$

$$= \text{KL}(\Sigma_{\sigma, \epsilon}^{(d, k^*)} \| \Sigma_{\sigma}^{(d, k^*)}) + \int \log\left(\frac{f_{\mathcal{P}_0}}{f_{\mathcal{P}_s}}\right) f_{\Sigma} \quad (\text{A.59})$$

The blocks of the partition \mathcal{P}_s are subsets of blocks of the partition \mathcal{P}_0 , then

$$\int \log \left(\frac{f_{\mathcal{P}_0}}{f_{\mathcal{P}_s}} \right) f_{\Sigma} = \int \log \left(\frac{\prod_{S \in \mathcal{P}_0} f_S}{\prod_{S \in \mathcal{P}_0} (f_{\mathcal{P}_s})_S} \right) f_{\Sigma} \quad (\text{A.60})$$

$$= \sum_{S \in \mathcal{P}_0} \int \log \left(\frac{f_S}{(f_{\mathcal{P}_s})_S} \right) f_S \quad (\text{A.61})$$

$$= \sum_{S \in \mathcal{P}_0} \text{KL} \left(f_{A_{\sigma}^d} \| f_{(A_{\sigma}^d)_{\mathcal{P}_s}} \right) \quad (\text{A.62})$$

$$= \frac{d}{k^*} \text{KL} \left(f_{A_{\sigma}^d} \| f_{(A_{\sigma}^d)_{\mathcal{P}_s}} \right) \quad (\text{A.63})$$

Now, using proposition A.3

$$\begin{aligned} \text{KL} \left(f_{A_{\sigma}^d} \| f_{(A_{\sigma}^d)_{\mathcal{P}_s}} \right) &= \frac{1}{2} \left(p \log \left(\frac{1 + (k-1)\sigma}{1-\sigma} \right) + \log \left(\frac{1 + (r-1)\sigma}{1-\sigma} \right) \right. \\ &\quad \left. - \log \left(\frac{1 + (k^*-1)\sigma}{1-\sigma} \right) \right) \end{aligned} \quad (\text{A.64})$$

And, using proposition A.4, we have that

$$\text{KL} \left(f_{\Sigma_{\sigma, \epsilon}^{(d,k)}} \| f_{\Sigma_{\sigma}^{(d,k)}} \right) \underset{\epsilon \rightarrow 0}{=} \frac{d(d-k)}{4(1+(k-1)\sigma)^2} \epsilon^2 + o(\epsilon^2). \quad (\text{A.65})$$

Then we have proven the desired upper-bound.

□