

Optimal projection for parametric importance sampling in high dimension

Maxime El Masri  ONERA/DTIS, ISAE-SUPAERO, Université de Toulouse

Jérôme Morio  ONERA/DTIS, Université de Toulouse

Florian Simatos ISAE-SUPAERO, Université de Toulouse

Date published: 2023-09-06 Last modified:

Abstract

We propose a dimension-reduction strategy in order to improve the performance of importance sampling in high dimensions. The idea is to estimate variance terms in a small number of suitably chosen directions. We first prove that the optimal directions, i.e., the ones that minimize the Kullback–Leibler divergence with the optimal auxiliary density, are the eigenvectors associated with extreme (small or large) eigenvalues of the optimal covariance matrix. We then perform extensive numerical experiments showing that as dimension increases, these directions give estimations which are very close to optimal. Moreover, we demonstrate that the estimation remains accurate even when a simple empirical estimator of the covariance matrix is used to compute these directions. The theoretical and numerical results open the way for different generalizations, in particular the incorporation of such ideas in adaptive importance sampling schemes.

Keywords: Parameter estimation, Importance sampling, Dimension reduction, Kullback–Leibler divergence, Projection

Contents

1	Introduction	2
2	Importance Sampling	3
3	Efficient dimension reduction	4
3.1	Projecting onto a low-dimensional subspace	4
3.2	Definition of the function ℓ	5
3.3	Main result of the paper	5
3.4	Choice of the number of dimensions k	7
3.5	Theoretical result concerning the projection on \mathbf{m}^*	7
4	Computational framework	8
4.1	General framework	8
4.2	Choice of the auxiliary density g' for the Gaussian model	9
4.3	Choice of the auxiliary density g' for the von Mises–Fisher–Nakagami model . . .	10

5	Numerical results on five test cases	11
5.1	Test case 1: one-dimensional optimal projection	11
5.1.1	Evolution of the partial KL divergence and spectrum	11
5.1.2	Numerical results	12
5.2	Test case 2: projection in 2 directions	12
5.2.1	Evolution of the partial KL divergence and spectrum	13
5.2.2	Numerical results	13
5.3	Test case 3: banana shape distribution	14
5.4	Application 1: large portfolio losses	14
5.5	Application 2: discretized Asian payoff	15
6	Conclusion	16
	Acknowledgement	17
7	Proof of Theorem 3.1 and Theorem 3.2 {#sec-proof}	17
	References	19

1 Introduction

Importance Sampling (IS) is a stochastic method to estimate integrals of the form $\mathcal{E} = \int \phi f$ with a black-box function ϕ and a probability density function (pdf) f . It rests upon the choice of an auxiliary density which can significantly improve the estimation compared to the naive Monte Carlo (MC) method (Agapiou et al. 2017), (Owen and Zhou 2000). The theoretical optimal IS density, also called zero-variance density, is defined by $\phi f / \mathcal{E}$ when ϕ is a positive function. This density is not available in practice as it involves the unknown integral \mathcal{E} , but a classical strategy consists in searching for an optimal approximation in a parametric family of densities. By minimising a “distance” to the optimal IS density, such as the Kullback–Leibler divergence, one can find optimal parameters in this family to get an efficient sampling pdf. Adaptive Importance Sampling (AIS) algorithms, such as the Mixture Population Monte Carlo method (Cappé et al. 2008), the Adaptive Multiple Importance Sampling method (Cornuet et al. 2012), or the Cross Entropy method (Rubinstein and Kroese 2011), estimate the optimal parameters adaptively by updating at intermediate levels (Bugallo et al. 2017).

These techniques work very well, but only for moderate dimensions. In high dimensions, most of these techniques fail to give suitable parameters for two reasons: (1) The weight degeneracy problem, for which the self-normalized likelihood ratios (weights) in the IS densities degenerate in the sense that the largest one takes all the mass, while all other weights are negligible so that the final estimation essentially uses only one sample. See for instance (Bengtsson, Bickel, and Li 2008) for a theoretical analysis in the related context of particle filtering. But even without likelihood ratios, such techniques may fail if they need to estimate high-dimensional parameters such as covariance matrices, whose size increases quadratically in the dimension (Ashurbekova et al. 2020), (Ledoit and Wolf 2004). The conditions under which importance sampling is applicable in high dimensions are notably investigated in a reliability context in (Au and Beck 2003): it is remarked that the optimal covariance matrix should not deviate significantly from the identity matrix. (El-Laham, Elvira, and Bugallo 2019) tackle the weight degeneracy problem by applying a recursive shrinkage of the covariance matrix, which is constructed iteratively with a weighted sum of the sample covariance estimator and a biased, but more stable, estimator. (2) The estimation of parameters in high dimensions, for which dimension reduction techniques can be applied. The idea was recently put forth to reduce the effective dimension by only estimating these parameters (in particular the covariance matrix) in suitable directions (El Masri, Morio, and Simatos 2021), (Uribe et al. 2021). In this paper we

delve deeper into this idea. The main contribution of the present paper is to identify the optimal directions in the fundamental case when the parametric family is Gaussian, and perform numerical simulations in order to understand how they behave in practice. In particular, we propose directions which, in contrast to the recent paper (Uribe et al. 2021), do not require the objective function to be differentiable, and moreover optimizes the Kullback–Leibler distance with the optimal density instead of simply an upper bound on it, as in (Uribe et al. 2021). In Section 3.1 we elaborate in more details on the differences between the two approaches.

The paper is organised as follows: in Section 2 we recall the foundations of IS. In Section 3, we state our main theoretical result and we compare it with the current state-of-the-art. **@sec-proof** presents the proof of our theoretical result; Section 4 introduces the numerical framework that we have adopted, and Section 5 presents the numerical results obtained on five different test cases to assess the efficiency of the directions that we propose. We conclude in Section 6 with a summary and research perspectives.

2 Importance Sampling

We consider the problem of estimating the following integral:

$$\mathcal{E} = \mathbb{E}_f(\phi(\mathbf{X})) = \int \phi(\mathbf{x})f(\mathbf{x})d\mathbf{x},$$

where \mathbf{X} is a random vector in \mathbb{R}^n with standard Gaussian pdf f , and $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_+$ is a real-valued, non-negative function. The function ϕ is considered as a black-box function which is potentially expensive to evaluate, and this means that the number of calls to ϕ should be limited.

IS is an approach used to reduce the variance of the classical Monte Carlo estimator of \mathcal{E} . The idea of IS is to generate a random sample $\mathbf{X}_1, \dots, \mathbf{X}_N$ from an auxiliary density g , instead of f , and to compute the following estimator:

$$\widehat{\mathcal{E}}_N = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{X}_i)L(\mathbf{X}_i), \quad (1)$$

with $L = f/g$ the likelihood ratio, or importance weight, and the auxiliary density g , called importance sampling density, is such that $g(\mathbf{x}) = 0$ implies $\phi(\mathbf{x})f(\mathbf{x}) = 0$ for every \mathbf{x} (which makes the product ϕL well-defined). This estimator is consistent and unbiased but its accuracy strongly depends on the choice of the auxiliary density g . It is well known that the optimal choice for g is (Bucklew 2013)

$$g^*(\mathbf{x}) = \frac{\phi(\mathbf{x})f(\mathbf{x})}{\mathcal{E}}, \quad \mathbf{x} \in \mathbb{R}^n.$$

Indeed, for this choice we have $\phi L = \mathcal{E}$ and so $\widehat{\mathcal{E}}_N$ is actually the deterministic estimator \mathcal{E} . For this reason, g^* is sometimes called zero-variance density, a terminology that we will adopt here. Of course, g^* is only of theoretical interest as it depends on the unknown integral \mathcal{E} . However, it gives an idea of good choices for the auxiliary density g , and we will seek to approximate g^* by an auxiliary density that minimizes a distance between g^* and a given parametric family of densities.

In this paper, the parametric family of densities is the Gaussian family $\{g_{\mathbf{m}} : \mathbf{m} \in \mathbb{R}^n, \in \mathcal{S}_n^+\}$, where $g_{\mathbf{m}}$ denotes the Gaussian density with mean $\mathbf{m} \in \mathbb{R}^n$ and covariance matrix $\in \mathcal{S}_n^+$ with $\mathcal{S}_n^+ \subset \mathbb{R}^{n \times n}$ the set of symmetric, positive-definite matrices:

$$g_{\mathbf{m}}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}||^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^\top{}^{-1}(\mathbf{x} - \mathbf{m})\right), \quad \mathbf{x} \in \mathbb{R}^n.$$

with $\|\cdot\|$ the determinant of \cdot . Moreover, we will consider the Kullback–Leibler (KL) divergence to measure a “distance” between g^* and $g_{\mathbf{m},*}$. Recall that for two densities f and h , with f absolutely continuous with respect to h , the KL divergence $D(f, h)$ between f and h is defined by:

$$D(f, h) = \mathbb{E}_f \left[\log \left(\frac{f(\mathbf{X})}{h(\mathbf{X})} \right) \right] = \int \log \left(\frac{f(\mathbf{x})}{h(\mathbf{x})} \right) f(\mathbf{x}) d\mathbf{x}.$$

Thus, our goal is to approximate g^* by $g_{\mathbf{m}^*,*}$ with the optimal mean vector \mathbf{m}^* and the optimal covariance matrix Σ^* given by:

$$(\mathbf{m}^*, \Sigma^*) = \arg \min \{ D(g^*, g_{\mathbf{m},*}) : \mathbf{m} \in \mathbb{R}^n, \Sigma \in \mathcal{S}_n^+ \}. \quad (2)$$

In the Gaussian case, it is well-known that \mathbf{m}^* and Σ^* are simply the mean and variance of the zero-variance density (section 3.3.) (Rubinstein and Kroese 2011), (section 5.7.) (Rubinstein and Kroese 2017):

$$\mathbf{m}^* = \mathbb{E}_{g^*}(\mathbf{X}) \quad \text{and} \quad \Sigma^* = \text{Var}_{g^*}(\mathbf{X}). \quad (3)$$

3 Efficient dimension reduction

3.1 Projecting onto a low-dimensional subspace

As g^* is unknown, the optimal parameters \mathbf{m}^* and Σ^* given by Equation 3 are not directly computable. However, we can sample from the optimal density as it is known up to a multiplicative constant. Therefore, usual estimation schemes start with estimating \mathbf{m}^* and Σ^* , say through $\hat{\mathbf{m}}^*$ and $\hat{\Sigma}^*$, respectively, and then use these approximations to estimate E through Equation 1 with the auxiliary density $g_{\hat{\mathbf{m}}^*, \hat{\Sigma}^*}$. Although the estimation of E with the auxiliary density $g_{\hat{\mathbf{m}}^*, \hat{\Sigma}^*}$ usually provides very good results, it is well-known that in high dimensions, the additional error induced by the estimations of \mathbf{m}^* and Σ^* severely degrades the accuracy of the final estimation (Papaioannou, Geyer, and Straub 2019), (Uribe et al. 2021). The main problem lies in the estimation of Σ^* which, in dimension n , involves the estimation of a quadratic (in the dimension) number of terms, namely $n(n+1)/2$. Recently, the idea to overcome this problem by only evaluating variance terms in a small number of influential directions was explored in (El Masri, Morio, and Simatos 2021) and (Uribe et al. 2021). In these two papers, the auxiliary covariance matrix is modeled in the form

$$= \sum_{i=1}^k (v_i - 1) \mathbf{d}_i \mathbf{d}_i^\top + I_n \quad (4)$$

where the \mathbf{d}_i 's are the k orthonormal directions which are deemed influential. It is easy to check that is the covariance matrix of the Gaussian vector

$$v_1^{1/2} Y_1 \mathbf{d}_1 + \dots + v_k^{1/2} Y_k \mathbf{d}_k + Y_{k+1} \mathbf{d}_{k+1} + \dots + Y_n \mathbf{d}_n$$

where the Y_i 's are i.i.d. standard normal random variables (one-dimensional), and the $n - k$ vectors $(\mathbf{d}_{k+1}, \dots, \mathbf{d}_n)$ complete $(\mathbf{d}_1, \dots, \mathbf{d}_k)$ into an orthonormal basis. In particular, v_i is the variance in the direction of \mathbf{d}_i , i.e., $v_i = \mathbf{d}_i^\top \Sigma^* \mathbf{d}_i$. In Equation 4, k can be considered as the effective dimension in which variance terms are estimated. In other words, in (El Masri, Morio, and Simatos 2021) and (Uribe et al. 2021), the optimal variance parameter is not sought in \mathcal{S}_n^+ as in Equation 2, but rather in the subset of matrices of the form

$$\mathcal{L}_{n,k} = \left\{ \sum_{i=1}^k (\alpha_i - 1) \frac{\mathbf{d}_i \mathbf{d}_i^\top}{\|\mathbf{d}_i\|^2} + I_n : \alpha_1, \dots, \alpha_k > 0 \text{ and the } \mathbf{d}_i \text{'s are orthogonal} \right\}.$$

The relevant minimization problem thus becomes

$$(\mathbf{m}_k^*, \mathbf{d}_k^*) = \arg \min \{D(g^*, g_{\mathbf{m}}) : \mathbf{m} \in \mathbb{R}^n, \mathbf{d} \in \mathcal{L}_{n,k}\} \quad (5)$$

instead of Equation 2, with the effective dimension k being allowed to be adjusted dynamically. By restricting the space in which the variance is looked up, one seeks to limit the number of variance terms to be estimated. The idea is that if the directions are suitably chosen, then the improvement of the accuracy due to the smaller error in estimating the variance terms will compensate the fact that we consider less candidates for the covariance matrix. In (El Masri, Morio, and Simatos 2021), the authors consider $k = 1$ and $\mathbf{d}_1 = \mathbf{m}^* / \|\mathbf{m}^*\|$. When f is Gaussian, this choice is motivated by the fact that, due to the light tail of the Gaussian random variable and the reliability context, the variance should vary significantly in the direction of \mathbf{m}^* and so estimating the variance in this direction can bring information. In Section 3.5, we use the techniques of the present paper to provide a stronger theoretical justification of this choice, see Theorem 3.2 and the discussion following it. The method in (Uribe et al. 2021) is more involved: k is adjusted dynamically, while the directions \mathbf{d}_i are the eigenvectors associated to the largest eigenvalues of a certain matrix. They span a low-dimensional subspace called Failure-Informed Subspace, and the authors in (Uribe et al. 2021) prove that this choice minimizes an upper bound on the minimal KL divergence. In practice, this algorithm yields very accurate results. However, we will not consider it further in the present paper for two reasons. First, this algorithm is tailored for the reliability case where $\phi = \mathbb{I}_{\{\phi \geq 0\}}$, with a function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$, whereas our method is more general and applies to the general problem of estimating an integral (see for instance our test case of Section 5.5). Second, the algorithm in (Uribe et al. 2021) requires the evaluation of the gradient of the function ϕ . However, this gradient is not always known and can be expensive to evaluate in high dimensions; in some cases, the function ϕ is even not differentiable, as will be the case in our numerical example in Section 5.4. In contrast, our method makes no assumption on the form or smoothness of ϕ : it does not need to assume that it is of the form $\mathbb{I}_{\{\phi \geq 0\}}$, or to assume that $\nabla \phi$ is tractable. For completeness, whenever the algorithm of (Uribe et al. 2021) was applicable and computing the gradient of ϕ did not require any additional simulation budget, we have run it on the test cases considered here and found that it outperformed our algorithm. In more realistic settings, computing $\nabla \phi$ would likely increase the simulation budget, and it would be interesting to compare the two algorithms in more details to understand when this extra computation cost is worthwhile. We reserve such a question for future research and will not consider the algorithm of (Uribe et al. 2021) further, as our aim in this paper is to establish benchmark results for a general algorithm which works for any function ϕ .

3.2 Definition of the function ℓ

The statement of our result involves the following function ℓ , which is represented in **Fig-1**:

$$\ell : x \in (0, \infty) \mapsto -\log(x) + x - 1. \quad (6)$$

In the following, $(\lambda, \mathbf{d}) \in \mathbb{R} \times \mathbb{R}^n$ is an eigenpair of a matrix A if $A\mathbf{d} = \lambda\mathbf{d}$ and $\|\mathbf{d}\| = 1$. A diagonalizable matrix has n distinct eigenpairs, say $((\lambda_i^*, \mathbf{d}_i^*), i = 1, \dots, n)$, and we say that these eigenpairs are ranked in decreasing ℓ -order if $\ell(\lambda_1^*) \geq \dots \geq \ell(\lambda_n^*)$.

3.3 Main result of the paper

The main result of the present paper is to compute the exact value for \mathbf{m}_k^* in Equation 5, which therefore paves the way for efficient high-dimensional estimation schemes.

Theorem 3.1. *Let $(\lambda_i^*, \mathbf{d}_i^*)$ be the eigenpairs of A ranked in decreasing ℓ -order. Then for $1 \leq k \leq n$, the*

solution $(\mathbf{m}_k^*, \lambda_k^*)$ to Equation 5 is given by

$$\mathbf{m}_k^* = \mathbf{m}^* \text{ and } \lambda_k^* = I_n + \sum_{i=1}^k (\lambda_i^* - 1) \mathbf{d}_i^* (\mathbf{d}_i^*)^\top. \quad (7)$$

The proof of Theorem 3.1 is detailed in [?@sec-proof](#). For $k = 1$ for instance, the shape of the function ℓ depicted in [?@fig-1](#) implies that the matrix $\lambda_1^* = I_n + (\lambda^* - 1) \mathbf{d}^* (\mathbf{d}^*)^\top$ with $(\lambda^*, \mathbf{d}^*)$ the eigenpair of λ^* and λ^* either the largest or the smallest eigenvalue of λ^* , depending on which one maximizes ℓ .

This theoretical result therefore suggests to reduce dimension by estimating eigenpairs of λ^* , rank them in decreasing ℓ -order and then use the k first eigenpairs $((\hat{\lambda}_i^*, \hat{\mathbf{d}}_i^*), i = 1, \dots, k)$ to build the covariance matrix $\hat{\lambda}_k^* = \sum_{i=1}^k (\hat{\lambda}_i^* - 1) \hat{\mathbf{d}}_i^* (\hat{\mathbf{d}}_i^*)^\top + I_n$ and the corresponding auxiliary density. This scheme is summarized in Algorithm 1. The effective dimension k is obtained by Algorithm 2, see Section 3.4 below. The proof of the theorem is shown in [@#sec-proof](#).

Algorithm 1 Algorithm suggested by Theorem 1.

- 1: **Data:** Sample sizes N and M
 - 2: **Result:** Estimation $\widehat{\mathcal{E}}_N$ of integral \mathcal{E}
 - 3: - Generate a sample $\mathbf{X}_1, \dots, \mathbf{X}_M$ on \mathbb{R}^n independently according to g^*
 - 4: - Estimate $\hat{\mathbf{m}}^*$ and $\hat{\lambda}^*$ defined in Equation 10 and Equation 11 with this sample
 - 5: - Compute the eigenpairs $(\hat{\lambda}_i^*, \hat{\mathbf{d}}_i^*)$ of $\hat{\lambda}^*$ ranked in decreasing ℓ -order
 - 6: - Compute the matrix $\hat{\lambda}_k^* = \sum_{i=1}^k (\hat{\lambda}_i^* - 1) \hat{\mathbf{d}}_i^* (\hat{\mathbf{d}}_i^*)^\top + I_n$ with k obtained by applying Algorithm 2 with input $(\hat{\lambda}_1^*, \dots, \hat{\lambda}_n^*)$
 - 7: - Generate a new sample $\mathbf{X}_1, \dots, \mathbf{X}_N$ independently from $g' = g_{\hat{\mathbf{m}}^*, \hat{\lambda}_k^*}$
 - 8: - Return $\widehat{\mathcal{E}}_N = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{X}_i) \frac{f(\mathbf{X}_i)}{g'(\mathbf{X}_i)}$
-

Remark. The value 1 plays a particular role in Theorem 3.1, in that, as ℓ is minimized in 1, eigenvectors with eigenvalues 1 will only be selected once all other eigenvalues will have been picked: in other words, if $\lambda_i^* = 1$ then $\lambda_j^* = 1$ for all $j \geq i$. The reason why 1 plays this special role is due to the form of the covariance matrix that we impose. More precisely, looking for covariance matrices in the set $\mathcal{L}_{n,k}$ amounts to looking for covariance matrices which, once diagonalized, have one's on the diagonal except possibly for k values (the α_i 's). As k will be small, typically $k = 1$ or 2 , this amounts to looking for covariance matrices which are perturbation of the identity. This is particularly relevant as we assume that f is a standard Gaussian density. What Theorem 3.1 tells is that, when trying to approximate λ^* by such matrices, we should first consider eigenvectors with eigenvalues as different as possible from 1, the “distance” to 1 being measured by ℓ . If one was imposing a different form on λ_k^* (which can be interesting if the distribution f is not standardized), then a different result would arise. For instance, if one was looking for matrices where the “default” choice would be some $\lambda > 0$ for the diagonal entries that are not estimated, i.e., a matrix of the form $\sum_i (\alpha_i - \lambda) \mathbf{d}_i \mathbf{d}_i^\top + \lambda I_n$, then eigenpairs would be ranked according to the function $\ell(\cdot/\lambda)$, meaning that one would look for eigenvectors associated to eigenvalues as different as possible from λ .

In the first step of Algorithm 1, we assume g^* can be sampled independently. This is a reasonable assumption as classical techniques such as importance sampling with self-normalized weights or Markov Chain Monte Carlo (MCMC) can be applied in this case (see for instance (Chan and Kroese 2012), (Grace, Kroese, and Sandmann 2014)). In this paper, we choose to apply a basic rejection method that yields perfect independent samples from g^* , possibly at the price of a high computational cost. As the primary goal of this paper is to understand whether the \mathbf{d}_i^* 's are indeed good projection

directions, this cost will not be taken into account. Possible improvements to relax this assumption are discussed in the conclusion of the paper.

3.4 Choice of the number of dimensions k

The choice of the effective dimension k , i.e., the number of projection directions considered, is important. If it is close to n , then the matrix $\hat{\mathbf{x}}_k^*$ will be close to $\hat{\mathbf{x}}^*$ which is the situation we want to avoid in the first place. On the other hand, setting $k = 1$ in all cases may be too simple and lead to suboptimal results. In practice, however this is often a good choice. In order to adapt k dynamically, we consider a simple method based on the value of the KL divergence. Given the eigenvalues $\lambda_1, \dots, \lambda_n$ ranked in decreasing ℓ -order, we look for the maximal gap $\max_{i=1 \dots n-1} \{\delta_i = \ell(\lambda_{i+1}) - \ell(\lambda_i)\}$, in the sequence $(\ell(\lambda_1), \dots, \ell(\lambda_n))$. This allows to choose k such that $\sum_{i=1}^k \ell(\lambda_i)$ is close to $\sum_{i=1}^n \ell(\lambda_i)$ which is equal, up to an additive constant, to the minimal KL divergence (shown in Lemma 7.1). The precise method is described in Algorithm 2.

Algorithm 2 Choice of the number of dimensions

- 1: **Data:** Sequence of positive numbers $\lambda_1, \dots, \lambda_n$ in decreasing ℓ -order
 - 2: **Result:** Number of selected dimensions k
 - 3: - Compute the increments $\delta_i = \ell(\lambda_{i+1}) - \ell(\lambda_i)$ for $i = 1 \dots n - 1$
 - 4: - Return $k = \arg \max \delta_i$, the index of the maximum of the differences.
-

3.5 Theoretical result concerning the projection on \mathbf{m}^*

In (El Masri, Morio, and Simatos 2021), the authors propose to project on the mean \mathbf{m}^* of the optimal auxiliary density g^* . Numerically, this algorithm is shown to perform well, but only a very heuristic explanation based on the light tail of the Gaussian distribution is provided to motivate this choice. It turns out that the techniques used in the proof of Theorem 3.1 can shed light on why projecting on \mathbf{m}^* may indeed be a good idea. Let us first state our theoretical result, and then explain why it justifies the idea of projecting on \mathbf{m}^* .

Theorem 3.2. Consider $\in \mathcal{L}_{n,1}$ of the form $= I_n + (\alpha - 1)\mathbf{d}\mathbf{d}^\top$ with $\alpha > 0$ and $\|\mathbf{d}\| = 1$. Then the minimizer in (α, \mathbf{d}) of the KL divergence between f and $g_{\mathbf{m}^*}$, is $(1 + \|\mathbf{m}^*\|^2, \mathbf{m}^*/\|\mathbf{m}^*\|)$:

$$(1 + \|\mathbf{m}^*\|^2, \mathbf{m}^*/\|\mathbf{m}^*\|) = \arg \min_{\alpha, \mathbf{d}} \{D(f, g_{\mathbf{m}^*, I_n + (\alpha - 1)\mathbf{d}\mathbf{d}^\top}) : \alpha > 0, \|\mathbf{d}\| = 1\}.$$

The proof of Theorem 3.2 is detailed in **@sec-proof**. In other words, \mathbf{m}^* appears as an optimal projection direction when one seeks to minimize the KL divergence between f and the Gaussian density with mean \mathbf{m}^* and covariance of the form $I_n + (\alpha - 1)\mathbf{d}\mathbf{d}^\top$. Let us now explain why this minimization problem is indeed relevant, and why choosing an auxiliary density which minimizes this KL divergence may indeed lead to an accurate estimation. The justification deeply relies on the recent results by (Chatterjee and Diaconis 2018).

As mentioned above, in a reliability context where one seeks to estimate a small probability $p = P(\mathbf{X} \in A)$, Theorem 1.3 in (Chatterjee and Diaconis 2018) shows that $D(g^*, g)$ governs the sample size required for an accurate estimation of p : more precisely, the estimation is accurate if the sample size is larger than $e^{D(g^*, g)}$, and inaccurate otherwise. This motivates the rationale for minimizing the KL divergence with g^* .

However, in high dimensions, importance sampling is known to fail because of the weight degeneracy problem whereby $\max_i L_i / \sum_i L_i \approx 1$, with the L_i 's the unnormalized importance weights, or likelihood ratios: $L_i = f(\mathbf{X}_i)/g(\mathbf{X}_i)$ with the \mathbf{X}_i 's i.i.d. drawn according to g . Theorem 2.3 in (Chatterjee and

Diaconis 2018) shows that the weight degeneracy problem is avoided if the empirical mean of the likelihood ratios is close to 1, and for this, Theorem 1.1 in (Chatterjee and Diaconis 2018) shows that the sample size should be larger than $e^{D(f,g)}$. In other words, these results suggest that the KL divergence with g^* governs the sample size for an accurate estimation of p , while the KL divergence with f governs the weight degeneracy problem.

In light of these results, it becomes natural to consider the KL divergence with f and not only g^* (Owen and Zhou 2000). Of course, minimizing $D(f, g_{\mathbf{m}})$ without constraints on \mathbf{m} and is trivial since $g_{\mathbf{m}} = f$ for $\mathbf{m} = 0$ and $= I_n$. However, these choices are the ones we want to avoid in the first place, and so it makes sense to impose some constraints on \mathbf{m} and . If one keeps in mind the other objective of getting close to g^* , then the choice $\mathbf{m} = \mathbf{m}^*$ becomes very natural, and we are led to considering the optimization problem of Theorem 3.2 (when $\in \mathcal{L}_{n,1}$ is a rank-1 perturbation of the identity).

4 Computational framework

4.1 General framework

The objective of the numerical simulations is to evaluate the impact of the choice of the covariance matrix on the estimation accuracy of a high dimensional integral \mathcal{E} . We compare in this section the estimation results for different choices of the auxiliary covariance matrix when the IS auxiliary density is Gaussian. To extend this comparison, we also compute the results when the IS auxiliary density is chosen with the von Mises–Fisher– Nakagami (vMFN) model recently proposed in (Papaioannou, Geyer, and Straub 2019) for high dimensional probability estimation.

In the following section we test these different models of auxiliary densities on five test cases, where f is a standard Gaussian density. This choice is not a theoretical limitation as we can in principle always come back to this case by transforming the vector \mathbf{X} with isoprobabilistic transformations (see for instance (Hohenbichler and Rackwitz 1981), (Liu and Der Kiureghian 1986)).

The precise numerical framework that we will consider to assess the efficiency of the different auxiliary models is as follows. We assume first that M i.i.d. random samples $\mathbf{X}_1, \dots, \mathbf{X}_M$ distributed from g^* are available from rejection sampling. From these samples, the parameters of the Gaussian and of the vMFN auxiliary density are computed to get an auxiliary density g' . Finally, N samples are generated from g' to provide an estimation of E with IS. This procedure is summarized by the following stages:

1. Generate a sample $\mathbf{X}_1, \dots, \mathbf{X}_M$ independently according to g^* ;
2. From $\mathbf{X}_1, \dots, \mathbf{X}_M$, compute the parameters of the auxiliary parametric density g' ;
3. Generate a new sample $\mathbf{X}_1, \dots, \mathbf{X}_N$ independently from g' ;
4. Estimate E with $\widehat{\mathcal{E}}_N = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{X}_i) \frac{f(\mathbf{X}_i)}{g'(\mathbf{X}_i)}$.

The number of samples M and N are respectively set to $M = 500$ and $N = 2000$. This procedure is then repeated 500 times to provide a mean estimation $\widehat{\mathcal{E}}$ of \mathcal{E} . In the result tables, for each auxiliary density g' we report the corresponding value for the relative error $\widehat{\mathcal{E}}/\mathcal{E} - 1$ and the coefficient of variation of the 500 iterations (the empirical standard deviation divided by $\widehat{\mathcal{E}}$). As was established in the proof of Theorem 3.1, the KL divergence is, up to an additive constant, equal to $D'() = \log\| + \text{tr}(*^{-1})$ which we will refer to as partial KL divergence. In the result tables, we also report thus the mean value of $D'()$ to analyse the relevance of the auxiliary density $g_{\widehat{\mathbf{m}}^*}$, for six choices of covariance matrix . The next sections specify the different parameters of g' for the Gaussian model and for the vMFN model we have considered in the simulations.

4.2 Choice of the auxiliary density g' for the Gaussian model

The goal is to get benchmark results to assess whether one can improve estimations of Gaussian IS auxiliary density by projecting the covariance matrix Σ^* in the proposed directions \mathbf{d}_i^* . The algorithm that we study here (Algorithms 1+2) aims more precisely at understanding whether:

- projecting can improve the situation with respect to the empirical covariance matrix;
- the \mathbf{d}_i^* 's are good candidates, in particular compared to the choice \mathbf{m}^* suggested in (El Masri, Morio, and Simatos 2021);
- what is the impact in making errors in estimating the eigenpairs $(\lambda_i^*, \mathbf{d}_i^*)$.

Let us define the estimate $\hat{\mathbf{m}}^*$ of \mathbf{m}^* from the M i.i.d. random samples $\mathbf{X}_1, \dots, \mathbf{X}_M$ distributed from g^* with

$$\hat{\mathbf{m}}^* = \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i. \quad (8)$$

In our numerical test cases, we will compare six different choices of Gaussian auxiliary distributions g' with mean $\hat{\mathbf{m}}^*$ and the following covariance matrices summarized in Table 1:

1. Σ^* : the optimal covariance matrix given by Equation 3;
2. $\hat{\Sigma}^*$: the empirical estimation of Σ^* given by

$$\hat{\Sigma}^* = \frac{1}{M} \sum_{i=1}^M (\mathbf{X}_i - \hat{\mathbf{m}}^*)(\mathbf{X}_i - \hat{\mathbf{m}}^*)^\top. \quad (9)$$

The four other covariance matrices considered in the numerical simulations are of the form $\sum_{i=1}^k (\nu_i - 1) \mathbf{d}_i \mathbf{d}_i^\top + I_n$ where ν_i is the variance of Σ^* in the direction \mathbf{d}_i , $\nu_i = \mathbf{d}_i^\top \Sigma^* \mathbf{d}_i$. The considered choice of k and \mathbf{d}_i gives the following covariance matrices:

3. $\hat{\Sigma}_{\text{opt}}^*$ is obtained by choosing $\mathbf{d}_i = \mathbf{d}_i^*$ of Theorem 3.1, which is supposed to be perfectly known from Σ^* and k is computed with Algorithm 2;
4. $\hat{\Sigma}_{\text{opt}}^{+d}$ is obtained by choosing $\mathbf{d}_i = \hat{\mathbf{d}}_i^*$ the i -th eigenvector of $\hat{\Sigma}^*$ (in ℓ -order), which is an estimation of \mathbf{d}_i^* , and k is computed with Algorithm 2;
5. $\hat{\Sigma}_{\text{mean}}^*$ is obtained by choosing $k = 1$ and $\mathbf{d}_1 = \mathbf{m}^* / \|\mathbf{m}^*\|$;
6. $\hat{\Sigma}_{\text{mean}}^{+d}$ is obtained by choosing $k = 1$ and $\mathbf{d}_1 = \hat{\mathbf{m}}^* / \|\hat{\mathbf{m}}^*\|$, where $\hat{\mathbf{m}}^*$ given by Equation 8.

The matrices $\hat{\Sigma}_{\text{opt}}^*$ and $\hat{\Sigma}_{\text{mean}}^*$ use the estimation $\hat{\Sigma}^*$ but the actual directions \mathbf{d}_i^* or \mathbf{m}^* , while the matrices $\hat{\Sigma}_{\text{opt}}^{+d}$ and $\hat{\Sigma}_{\text{mean}}^{+d}$ involve an additional estimation of the directions. By definition, Σ^* will give optimal results, while results for $\hat{\Sigma}^*$ will deteriorate as the dimension increases, which is the well-known behavior which we try to improve. Moreover, for $\hat{\Sigma}_{\text{mean}}^*$ and $\hat{\Sigma}_{\text{opt}}^*$, the projection directions, if not known analytically, are obtained by a brute force Monte Carlo scheme with a very high simulation budget. Finally, we emphasize that Algorithm 1 corresponds to estimating and projecting on the \mathbf{d}_i^* 's, and so the matrix $\hat{\Sigma}_k^*$ of Algorithm 1 is equal to the matrix $\hat{\Sigma}_{\text{opt}}^{+d}$.

Except for Σ^* , the five other matrices involve one or two estimations: $\hat{\Sigma}^*$ is the empirical estimation of Σ^* given by Equation 9. The four others are obtained by projecting $\hat{\Sigma}^*$ on: (i) the optimal directions \mathbf{d}_i^* for $\hat{\Sigma}_{\text{opt}}^*$; (ii) estimations $\hat{\mathbf{d}}_i^*$ of the optimal directions \mathbf{d}_i^* for $\hat{\Sigma}_{\text{opt}}^{+d}$; (iii) \mathbf{m}^* for $\hat{\Sigma}_{\text{mean}}^*$; (iv) the estimation $\hat{\mathbf{m}}^*$ in Equation 3 of \mathbf{m}^* for $\hat{\Sigma}_{\text{mean}}^{+d}$. The subscript therefore indicates the choice for the projection direction, while the superscript +d indicates whether these directions are estimated or not.

Table 1: Presentation of the six covariance matrices considered in the numerical examples.

	*	^*	^ _{opt}	^ _{mean}	^+d _{opt}	^+d _{mean}
Initial covariance matrix	*	^*	^*	^*	^*	^*
Projection directions (exact or estimated)	-	-	Exact	Exact	Esti- mated	Esti- mated
Choice for the projection direction	None	None	Opt	Mean	Opt	Mean

4.3 Choice of the auxiliary density g' for the von Mises–Fisher–Nakagami model

Von Mises–Fisher–Nakagami (vMFN) distributions were proposed in (Papaioannou, Geyer, and Straub 2019) as an alternative to the Gaussian parametric family to perform IS for high dimensional probability estimation. A random vector \mathbf{X} drawn according to the vMFN distribution can be written as $\mathbf{X} = R\mathbf{A}$ where $\mathbf{A} = \frac{\mathbf{X}}{\|\mathbf{X}\|}$ is a unit random vector following the von Mises–Fisher distribution, and $R = \|\mathbf{X}\|$ is a positive random variable with a Nakagami distribution; further, R and \mathbf{A} are independent. The vMFN pdf can be written as

$$g_{\text{vMFN}}(\mathbf{x}) = g_{\text{N}}(\|\mathbf{x}\|, p, \omega) \times g_{\text{vMF}}\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}, \boldsymbol{\mu}, \kappa\right). \quad (10)$$

The density $g_{\text{N}}(\|\mathbf{x}\|, p, \omega)$ is the Nakagami distribution with shape parameter $p \geq 0.5$ and a spread parameter $\omega > 0$ defined by

$$g_{\text{N}}(\|\mathbf{x}\|, p, \omega) = \frac{2p^p}{\Gamma(p)\omega^p} \|\mathbf{x}\|^{2p-1} \exp\left(-\frac{p}{\omega} \|\mathbf{x}\|^2\right)$$

and the density $g_{\text{vMF}}\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}, \boldsymbol{\mu}, \kappa\right)$ is the von Mises–Fisher distribution, given by

$$g_{\text{vMF}}\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}, \boldsymbol{\mu}, \kappa\right) = C_n(\kappa) \exp\left(\kappa \boldsymbol{\mu}^T \frac{\mathbf{x}}{\|\mathbf{x}\|}\right),$$

where $C_n(\kappa)$ is a normalizing constant, $\boldsymbol{\mu}$ is a mean direction (with $\|\boldsymbol{\mu}\| = 1$) and $\kappa > 0$ is a concentration parameter.

Choosing a vMFN distribution therefore amounts to choosing the parameters p , ω , $\boldsymbol{\mu}$ and κ . There are therefore $n + 3$ parameters to estimate, which is a significant reduction compared to the $\frac{n(n+3)}{2}$ required parameters of the Gaussian model with full covariance matrix.

Following (Papaioannou, Geyer, and Straub 2019), given a sample $\mathbf{X}_1, \dots, \mathbf{X}_M$ distributed from g^* , the parameters ω , p , $\boldsymbol{\mu}$ and κ are set in the following way in order to define g' :

$$\hat{\omega} = \frac{1}{M} \sum_{i=1}^M \|\mathbf{X}_i\|^2 \quad \text{and} \quad \hat{p} = \frac{\hat{\omega}^2}{\hat{\tau} - \hat{\omega}^2} \quad \text{with} \quad \hat{\tau} = \frac{1}{M} \sum_{i=1}^M \|\mathbf{X}_i\|^4$$

and

$$\hat{\boldsymbol{\mu}} = \frac{\sum_{i=1}^M \frac{\mathbf{X}_i}{\|\mathbf{X}_i\|}}{\|\sum_{i=1}^M \frac{\mathbf{X}_i}{\|\mathbf{X}_i\|}\|} \quad \text{and} \quad \hat{\kappa} = \frac{n\hat{\chi} - \hat{\chi}^3}{1 - \hat{\chi}^2} \quad \text{with} \quad \hat{\chi} = \min\left(\left\|\frac{1}{M} \sum_{i=1}^M \frac{\mathbf{X}_i}{\|\mathbf{X}_i\|}\right\|, 0.95\right).$$

5 Numerical results on five test cases

The proposed numerical framework is applied on three examples that are often considered to assess the performance of importance sampling algorithms and also two test cases from the area of financial mathematics.

5.1 Test case 1: one-dimensional optimal projection

We consider a test case where all computations can be made exactly. This is a classical example of rare event probability estimation, often used to test the robustness of a method in high dimensions. It is given by $\phi(\mathbf{x}) = \mathbb{I}_{\{\varphi(\mathbf{x}) \geq 0\}}$ with φ the following affine function:

$$\varphi : \mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto \sum_{j=1}^n x_j - 3\sqrt{n}. \quad (11)$$

The quantity of interest \mathcal{E} is defined as $\mathcal{E} = \int_{\mathbb{R}^n} \phi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = \mathbb{P}_f(\varphi(\mathbf{X}) \geq 0) \simeq 1.35 \cdot 10^{-3}$ for all n where the density f is the standard n -dimensional Gaussian distribution. Here, the zero-variance density is $g^*(\mathbf{x}) = \frac{f(\mathbf{x}) \mathbb{I}_{\{\varphi(\mathbf{x}) \geq 0\}}}{\mathcal{E}}$, and the optimal parameters \mathbf{m}^* and Σ^* in Equation 3 can be computed exactly, namely $\mathbf{m}^* = \alpha \mathbf{1}$ with $\alpha = e^{-9/2} / (E(2\pi)^{1/2})$ and $\mathbf{1} = \frac{1}{\sqrt{n}}(1, \dots, 1) \in \mathbb{R}^n$ the normalized constant vector, and $\Sigma^* = (\nu - 1)\mathbf{1}\mathbf{1}^\top + I_n$ with $\nu = 3\alpha - \alpha^2 + 1$.

5.1.1 Evolution of the partial KL divergence and spectrum

Figure 1 represents the evolution as the dimension varies between 5 and 100 of the partial KL divergence D' for three different choices of covariance matrix: the optimal matrix Σ^* , its empirical estimation $\hat{\Sigma}^*$ and the estimation $\hat{\Sigma}_k^*$ of the optimal lower-dimensional covariance matrix. We can notice that the partial KL divergence for $\hat{\Sigma}^*$ grows much faster than the other two, and that the partial KL divergence for $\hat{\Sigma}_k^*$ remains very close to the optimal value $D'(\Sigma^*)$. As the KL divergence is a proxy for the efficiency of the auxiliary density (it is for instance closely related to the number of samples required for a given precision (Chatterjee and Diaconis 2018)), this suggests that using $\hat{\Sigma}_k^*$ will provide results close to optimal.

We now check this claim. As $\Sigma^* = (\nu - 1)\mathbf{1}\mathbf{1}^\top + I_n$, its eigenpairs are $(\nu, \mathbf{1})$ and $(1, \mathbf{d}_i)$ where the \mathbf{d}_i 's form an orthonormal basis of the space orthogonal to the space spanned by $\mathbf{1}$. In particular, $(\nu, \mathbf{1})$ is the largest (in ℓ -order) eigenpair of Σ^* and $\Sigma_k^* = \Sigma^*$ for any $k \geq 1$.

In practice, we do not use this theoretical knowledge and $\Sigma^*, \hat{\Sigma}_k^*$ and the eigenpairs are estimated. The six covariance matrices introduced in Section 4.2 and in which we are interested are as follows:

- $\Sigma^* = (\nu - 1)\mathbf{1}\mathbf{1}^\top + I_n$;
- $\hat{\Sigma}^*$ given by Equation 9;
- $\hat{\Sigma}_{\text{opt}}^*$ and $\hat{\Sigma}_{\text{mean}}^*$ are equal and given by $(\hat{\lambda} - 1)\mathbf{1}\mathbf{1}^\top + I_n$ with $\hat{\lambda} = \mathbf{1}^\top \hat{\Sigma}^* \mathbf{1}$. This amounts to assuming that the projection direction $\mathbf{1}$ is perfectly known, whereas the variance in this direction is estimated;
- $\hat{\Sigma}_{\text{opt}}^{+d} = (\hat{\lambda} - 1)\hat{\mathbf{d}}\hat{\mathbf{d}}^\top + I_n$ with $(\hat{\lambda}, \hat{\mathbf{d}})$ the smallest eigenpair of $\hat{\Sigma}^*$. The difference with the previous case is that we do not assume anymore that the optimal projection direction $\mathbf{1}$ is known, and so it needs to be estimated;
- $\hat{\Sigma}_{\text{mean}}^{+d} = (\hat{\lambda} - 1)\frac{\hat{\mathbf{m}}^*(\hat{\mathbf{m}}^*)^\top}{\|\hat{\mathbf{m}}^*\|^2} + I_n$ with $\hat{\mathbf{m}}^*$ given by Equation 8 and $\hat{\lambda} = \frac{(\hat{\mathbf{m}}^*)^\top \hat{\Sigma}^* \hat{\mathbf{m}}^*}{\|\hat{\mathbf{m}}^*\|^2}$. Here we assume that \mathbf{m}^* is a good projection direction, but is unknown and therefore needs to be estimated.

Note that in the particularly simple case considered here, both $\hat{\mathbf{m}}^*/\|\hat{\mathbf{m}}^*\|$ and $\hat{\mathbf{d}}$ are estimators of $\mathbf{1}$ but they are obtained by different methods. In the next example we will consider a case where \mathbf{m}^* is not an optimal projection direction as given by Theorem 3.1.

?@fig-eigsum-2 represents the images by ℓ of the eigenvalues of * and ** . This picture carries a very important insight. We notice that the estimation of most eigenvalues is poor: indeed, all the blue crosses except the leftmost one are meant to be estimator of 1, whereas we see that they are more or less uniformly spread between 0.2 and 2.5. This means that the variance terms in the corresponding directions are poorly estimated, which could be the explanation on why the use of ** gives an inaccurate estimation. But what we remark also is that the function ℓ is quite flat around one: as a consequence, although the eigenvalues offer significant variability, this variability is smoothed by the action of ℓ . Indeed, the images of the eigenvalues by ℓ take values between 0 and 0.8 and have smaller variability. Moreover, $\ell(x)$ increases sharply as x approaches 0 and thus efficiently distinguishes between the two leftmost estimated eigenvalues and is able to separate them.

5.1.2 Numerical results

We report in ?@tbl-sum the numerical results for the six different matrices and the vMFN model for the dimension $n = 100$. The column * gives the optimal results, while the column ** corresponds to the results that we are trying to improve. Comparing these two columns, we notice as expected that the estimation of \mathcal{E} with ** is significantly degraded. Compared to the first column * , the third and fourth columns with $\hat{\mathbf{opt}} = \hat{\mathbf{mean}}$ correspond to the best projection direction $\mathbf{1}$ (as for *) but estimating the variance in this direction (instead of the true variance) with $\mathbf{1}^{\top} \mathbf{1}$. This choice performs very well, with numerical results similar to the optimal ones. This can be understood since in this case, both $\hat{\mathbf{opt}}$ and * are of the form $\alpha \mathbf{1}\mathbf{1}^{\top} + I_n$ and so estimating $\hat{\mathbf{opt}}$ requires only a one-dimensional estimation (namely, the estimation of α). Next, the last two columns $\hat{\mathbf{opt}}^{+d}$ and $\hat{\mathbf{mean}}^{+d}$ highlight the impact of having to estimate the projection directions in addition to the variance since these two matrices are of the form $\hat{\alpha} \hat{\mathbf{1}}\hat{\mathbf{1}}^{\top} + I_n$ with both $\hat{\alpha}$ (the variance term) and $\hat{\mathbf{1}}$ (the direction) being estimated. We observe that these matrices yield results which are close to optimal and greatly improve the estimation obtained using ** . In dimension 100, the coefficient of variation is around 5.2% for $\hat{\mathbf{mean}}^{+d}$, and around 11.1% for $\hat{\mathbf{opt}}^{+d}$, compared to 2.6% for ** .

Moreover, we observe that $\hat{\mathbf{mean}}^{+d}$ gives better results than $\hat{\mathbf{opt}}^{+d}$. We suggest that this is because $\hat{\mathbf{m}}^*/\|\hat{\mathbf{m}}^*\|$ is a better estimator of $\mathbf{1}$ than the eigenvector of ** . Indeed, evaluating $\hat{\mathbf{m}}^*$ requires the estimation of n parameters, whereas ** needs around $n^2/2$ parameters to estimate, so the eigenvector is finally more noisy than the mean vector. In the last column, we present the vMFN estimation that is slightly more efficient than the estimation obtained with $\hat{\mathbf{mean}}^{+d}$.

Thus, the proposed idea improves significantly the probability estimation in high dimensions. But we see that the method taken in (El Masri, Morio, and Simatos 2021) with the projection \mathbf{m}^* is at least as much efficient in this example where we need only a one-dimensional projection. The next case shows that the projection on more than one direction can outperform the one-dimensional projection on \mathbf{m}^* .

5.2 Test case 2: projection in 2 directions

The second test case is again a probability estimation, i.e., it is of the form $\phi = \mathbb{I}_{\{\phi \geq 0\}}$ with now the function ϕ having some quadratic terms:

$$\phi : \mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto x_1 - 25x_2^2 - 30x_3^2 - 1. \quad (12)$$

The quantity of interest \mathcal{E} is defined as $\mathcal{E} = \int_{\mathbb{R}^n} \phi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = \mathbb{P}_f(\phi(\mathbf{X}) \geq 0)$ for all n where the density f is the standard n -dimensional Gaussian distribution. This function is motivated in part

because \mathbf{m}^* and \mathbf{d}_1^* are different and also because Algorithm 2 chooses two projection directions. Thus, this is an example where $\hat{\mathbf{m}}_{\text{mean}}$ and $\hat{\mathbf{m}}_{\text{opt}}$ are significantly different.

5.2.1 Evolution of the partial KL divergence and spectrum

We check on **fig-inefficiency-parab-1** that the partial KL divergence obeys the same behavior as for the previous example, namely the one associated with $\hat{\mathbf{m}}^*$ increases much faster than the ones associated with $\hat{\mathbf{m}}^*$ and $\hat{\mathbf{d}}_k^*$, which again suggests that projecting can improve the situation. Since the function φ only depends on the first three variables and is even in x_2 and x_3 , one gets that $\mathbf{m}^* = \alpha \mathbf{e}_1$ with $\alpha = \mathbb{E}(X_1 \mid X_1 \geq 25X_2^2 + 30X_3^2 + 1) \approx 1.9$ (here and in the sequel, \mathbf{e}_i denotes the i th canonical vector of \mathbb{R}^n , i.e., all its coordinates are 0 except the i -th one which is equal to one), and that Σ^* is diagonal with

$$\Sigma^* = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & 0 & \cdots & 0 \\ 0 & 0 & \lambda_3 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{pmatrix}.$$

Note that the off-diagonal elements of the submatrix $(\Sigma^*)_{1 \leq i, j \leq 3}$ are indeed 0 since they result from integrating an odd function of an odd random variable with an even conditioning. For instance, if $F(x) = \mathbb{P}(30X_3^2 + 1 \leq x)$, then by conditioning on (X_1, X_3) we obtain

$$\Sigma_{12}^* = \mathbb{E}((X_1 - \alpha)X_2 \mid X_1 - 25X_2^2 \geq 30X_3^2 + 1) = \frac{1}{\mathcal{E}} \mathbb{E}[(X_1 - \alpha) \mathbb{E}(X_2 F(X_1 - 25X_2^2) \mid X_1)]$$

which is 0 as $x_2 F(x_1 - x_2^2)$ is an odd function of x_2 for fixed x_1 , and X_2 has an even density.

We can numerically compute $\lambda_1 \approx 0.28$, $\lambda_2 \approx 0.009$ and $\lambda_3 \approx 0.008$. These values correspond to the red squares in **fig-inefficiency-parab-2** which shows that the smallest eigenvalues are properly estimated. Moreover, Algorithm 2 selects the two largest eigenvalues, which have the highest ℓ -values. These two eigenvalues thus correspond to the eigenvectors \mathbf{e}_2 and \mathbf{e}_3 , and so we see that on this example, the optimal directions predicted by Theorem 3.1 are significantly different (actually, orthogonal) from \mathbf{m}^* which is proportional to \mathbf{e}_1 .

5.2.2 Numerical results

The numerical results of our simulations are presented in **tbl-parabol**. We remark as before that, when using $\hat{\mathbf{m}}^*$, the accuracy quickly deteriorates as the dimension increases as shows the coefficient of variation of 396% in dimension $n = 100$. In contrast, $\hat{\mathbf{m}}_{\text{opt}}$ leads to very accurate results, which remain close to optimal up to the same dimension $n = 100$. This behavior is to compare with the evolution of the relative KL divergence: contrary to $\hat{\mathbf{m}}^*$, $\hat{\mathbf{m}}_{\text{opt}}$ gives a partial KL divergence close to optimal in dimension $n = 100$. This confirms that the KL divergence is indeed a good proxy to assess the relevance of an auxiliary density.

It is also interesting to note that the direction \mathbf{m}^* improves the situation compared to not projecting (column $\hat{\mathbf{m}}_{\text{mean}}$ compared to $\hat{\mathbf{m}}^*$), but using $\hat{\mathbf{m}}_{\text{opt}}$ gives significantly better results, with for instance a coefficient of variation around 3.7% for $\hat{\mathbf{m}}_{\text{opt}}$ and around 28.8% for $\hat{\mathbf{m}}_{\text{mean}}$ in dimension $n = 100$. Thus, this confirms our theoretical result that the \mathbf{d}_i^* 's are good directions on which to project.

Finally, we notice that performing estimations of the projection directions instead of taking the true ones (columns $\hat{\mathbf{m}}_{\text{opt}}^{\text{d}}$ vs $\hat{\mathbf{m}}_{\text{opt}}$) slightly degrades the situation, making the coefficient of variation increase from 3.7 to 7.2% even if the accuracy remains satisfactory. The vMFN model is also not really adapted to this example as it gives results similar to $\hat{\mathbf{m}}_{\text{mean}}$. Gaussian density family are more able to fit g^* than vMFN parametric model in this test case.

Remark. For the two test cases studied so far, projecting $\hat{\mathbf{m}}^*$ in the Failure-Informed Subspace (FIS) of (Uribe et al. 2021) (see the introduction) would outperform our method with $\hat{\mathbf{m}}_k^*$, leading to results close to those obtained with \mathbf{m}^* . However, computing the FIS relies on the knowledge of the gradient of the function ϕ , which is straightforward to compute in these two test cases, and the method of (Uribe et al. 2021) can be applied because they are rare-event problems (i.e., ϕ is of the form $\phi = \mathbb{I}_{\{\phi \geq 0\}}$). In the next section we present other applications where the evaluation of the FIS is not feasible since either the function is not differentiable (test case of Section 5.4) or the example is not a rare event simulation problem (test cases of Section 5.3 and Section 5.5).

5.3 Test case 3: banana shape distribution

The third test case we consider is the integration of the banana shape distribution h , which is a classical test case in importance sampling (Cornuet et al. 2012), (Elvira et al. 2019). The banana shape distribution is the following pdf

$$h(\mathbf{x}) = g_{\mathbf{0},C}(x_1, x_2 + b(x_1^2 - \sigma^2), x_3, \dots, x_n). \quad (13)$$

The term $g_{\mathbf{0},C}$ represents the pdf of a Gaussian distribution of mean $\mathbf{0}$ and diagonal covariance matrix $C = \text{diag}(\sigma^2, 1, \dots, 1)$. The value of b and σ^2 are respectively set to $b = 800$ and $\sigma^2 = 0.0025$. We choose ϕ such that the optimal IS density g^* is equal to h , i.e., we choose $\phi = h/f$ so that the integral \mathcal{E} that we are trying to estimate is equal to $\mathcal{E} = \int \phi f = 1$. This choice is made in order to have an optimal covariance matrix \mathbf{m}^* whose two largest eigenvalues (in ℓ -order) correspond to the smallest and largest eigenvalues, as can be seen in ?@fig-inefficiency-banana-2. More formally, the optimal value of the Gaussian parameters are given by $\mathbf{m}^* = \mathbf{0}$ and \mathbf{m}^* is diagonal with

$$\mathbf{m}^* = \begin{pmatrix} 0.0025 & 0 & 0 & 0 & \dots & 0 \\ 0 & 9 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{pmatrix}.$$

The evolution of the KL partial divergence is given in ?@fig-inefficiency-banana-1. As the optimal mean \mathbf{m}^* is equal to $\mathbf{0}$, we cannot project on \mathbf{m}^* and so the matrix $\hat{\mathbf{m}}_{\text{mean}}^*$ is not defined. However, the numerical estimation $\hat{\mathbf{m}}^*$ will not be equal to 0 and so the approach proposed in (El Masri, Morio, and Simatos 2021) with $\hat{\mathbf{m}}_{\text{mean}}^{+d}$ is still applicable numerically.

The simulation results for the different covariance matrices and the vMFN density are given in ?@tbl-banana. The matrices $\hat{\mathbf{m}}_{\text{opt}}^*$ and $\hat{\mathbf{m}}_{\text{opt}}^{+d}$ perform very well for the estimation of E with an accuracy of the same order as the optimal covariance matrix \mathbf{m}^* . The effect of estimating the $k = 2$ main projection directions does not affect much the estimation performance as $\hat{\mathbf{m}}_{\text{opt}}^{+d}$ is still efficient compared to $\hat{\mathbf{m}}_{\text{opt}}^*$. The estimation results with $\hat{\mathbf{m}}_{\text{mean}}^{+d}$ are not really accurate and this choice is in fact roughly equivalent to choosing a random projection direction. The vMFN parametric model is not adapted to this test case as the vMFN estimate is not close to 1.

5.4 Application 1: large portfolio losses

The next example is a rare event application in finance, taken from (Bassamboo, Juneja, and Zeevi 2008), (Chan and Kroese 2012). The unknown integral is $\mathcal{E} = \int_{\mathbb{R}^{n+2}} \phi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = \mathbb{P}_f(\phi(\mathbf{X}) \geq 0)$, with $\phi = \mathbb{I}_{\{\phi \geq 0\}}$ and f is the standard $n + 2$ -dimensional Gaussian distribution. The function ϕ is the portfolio loss function defined as:

$$\phi(\mathbf{x}) = \sum_{j=3}^{n+2} \mathbb{I}_{\{\Psi(x_1, x_2, x_j) \geq 0.5\sqrt{n}\}} - bn, \quad (14)$$

with

$$\Psi(x_1, x_2, x_j) = \left(qx_1 + 3(1 - q^2)^{1/2} x_j \right) \left[F_{\Gamma}^{-1} (F_{\mathcal{N}}(x_2)) \right]^{-1/2},$$

where F_{Γ} and $F_{\mathcal{N}}$ are the cumulative distribution functions of Gamma(6, 6) and $\mathcal{N}(0, 1)$ random variables respectively. The constant b is chosen such that the probability is of the order of 10^{-3} in all dimension, then we have $b = 0.45$ when $n \leq 30$, $b = 0.3$ when $30 < n \leq 70$, and $b = 0.25$ when $n > 70$.

The reference value of this probability \mathcal{E} is reported in `?@tbl-portfolio` for dimension $n = 100$. The optimal parameters \mathbf{m}^* and \mathbf{d}_1^* cannot be computed analytically, but they are accurately estimated by Monte Carlo with a large sample. It turns out that \mathbf{m}^* and the first eigenvector \mathbf{d}_1^* of \mathbf{C}^* are numerically indistinguishable and that Algorithm 2 selects $k = 1$ projection direction, so that numerically, the choices $\hat{\mathbf{g}}_{\text{opt}}$ and $\hat{\mathbf{g}}_{\text{mean}}$ are indistinguishable and gives the same estimation results. Actually, the fact that these two estimators behave similarly does not seem to come from the fact that \mathbf{m}^* and \mathbf{d}_1^* are close: this relation can be broken for instance by a simple translation argument (see remark after `?@tbl-payoff`), but even then they behave similarly. The KL partial divergence and the spectrum with the associated ℓ -order are presented respectively in `?@fig-inefficiency-portfolio-1` and in `?@fig-inefficiency-portfolio-2`.

The results of `?@tbl-portfolio` show similar trends as for the first test case of Section 5.1. First, projecting seems indeed a relevant idea, as using $\hat{\mathbf{g}}_{\text{opt}}$ or $\hat{\mathbf{g}}_{\text{mean}}$ greatly improves the situation compared to $\hat{\mathbf{g}}^*$. This is particularly salient as $\hat{\mathbf{g}}^*$ yields an important bias and a coefficient of variation of 370%, whereas projecting on \mathbf{d}_1^* or \mathbf{m}^* yields a coefficient of variation of 14% and of 6.9% respectively. This improvement is still true even when the projection directions are estimated: compared to $\hat{\mathbf{g}}_{\text{opt}}$ and $\hat{\mathbf{g}}_{\text{mean}}$, $\hat{\mathbf{g}}_{\text{opt}}^{+d}$ and $\hat{\mathbf{g}}_{\text{mean}}^{+d}$ give coefficients of variation between 11.8 and 14.6%. Finally, $\hat{\mathbf{g}}_{\text{opt}}^{+d}$ seems to behave better than $\hat{\mathbf{g}}_{\text{mean}}^{+d}$.

5.5 Application 2: discretized Asian payoff

Our last numerical experiment is a mathematical finance example coming from (Kawai 2018), representing a discrete approximation of a standard Asian payoff under the Black–Scholes model. The goal is to estimate the integral $\mathcal{E} = \int_{\mathbb{R}^n} \phi(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$ with f the standard n -dimensional Gaussian distribution and the following function ϕ :

$$\phi : \mathbf{x} = (x_1, \dots, x_n) \mapsto e^{-rT} \left[\frac{S_0}{n} \sum_{i=1}^n \exp \left(i \left(r - \frac{\sigma^2}{2} \right) \frac{T}{n} + \sigma \sqrt{\frac{T}{n}} \sum_{k=1}^i x_k \right) - K \right]_+ \quad (15)$$

where $[y]_+ = \max(y, 0)$, for a real number y . The constants are taken from (Kawai 2018): $S_0 = 50$, $r = 0.05$, $T = 0.5$, $\sigma = 0.1$, $K = 55$, where they test the function for dimension $n = 16$. In our contribution, we test this example in dimension 100. Concerning \mathbf{m}^* and the \mathbf{d}_i^* 's, the situation is the same as in the previous example: they are not available analytically but can be estimated numerically by Monte Carlo with a large simulation budget. And again, it turns out that \mathbf{m}^* and the first eigenvector \mathbf{d}_1^* of \mathbf{C}^* are numerically indistinguishable and that Algorithm 2 selects $k = 1$ projection direction, so that $\hat{\mathbf{g}}_{\text{opt}}$ and $\hat{\mathbf{g}}_{\text{mean}}$ yield results that are numerically indistinguishable. The KL partial divergence and the spectrum with the associated ℓ -order are respectively presented in `?@fig-inefficiency-kawai-1` and `?@fig-inefficiency-kawai-2`.

The results of this example are given in `?@tbl-payoff`. The insight gained in the previous examples is confirmed. Projecting on \mathbf{m}^* or \mathbf{d}_1^* in dimension $n = 100$ enables to reach convergence and reduces (compared to $\hat{\mathbf{g}}^*$) the coefficient of variation from 559% to nearly 2%. Moreover, this improvement goes through even when projection directions are estimated, with again $\hat{\mathbf{g}}_{\text{mean}}^{+d}$ behaving better than $\hat{\mathbf{g}}_{\text{opt}}^{+d}$.

Remark. As already mentioned, the two directions \mathbf{m}^* and \mathbf{d}_1^* are numerically indistinguishable in the two examples of Section 5.4 and Section 5.5. However, we do not believe this relation to be highly relevant. For instance, this symmetry can be broken by changing ϕ into $\phi' = \phi(\cdot - \mu)$ and f into $f' = f(\cdot - \mu)$ for some $\mu \in \mathbb{R}^n$. Since $g^* \propto \phi f$, this amounts to translating g^* which thus changes \mathbf{m}^* into $\mathbf{m}^{*'} = \mathbf{m}^* + \mu$, but which does not change the covariance matrix (and therefore its leading eigenvector \mathbf{d}_1^*) which is translation-invariant. Note that this translation does not affect the integral $\mathcal{E} = \int \phi' f' = \int \phi f$, and so this modification leads to a new estimator $\widehat{\mathcal{E}}_\mu$ of the same quantity \mathcal{E} . However, it can be shown that $\widehat{\mathcal{E}}_\mu$ and $\widehat{\mathcal{E}}$ (the estimators considered throughout the paper) have the same law so that this translation, although it does break the relation $\mathbf{m}^* \approx \mathbf{d}_1^*$, does not change the law of the estimators. This suggests that, if the estimators based on $\hat{\Sigma}_{\text{opt}}$ and $\hat{\Sigma}_{\text{mean}}$ do behave similarly on these examples, this is not due to the fact that \mathbf{m}^* and \mathbf{d}_1^* are close but rather to Theorem 3.1 and Theorem 3.2. However, the fact that \mathbf{m}^* and \mathbf{d}_1^* are close bears some insight into the importance of the quality of the estimation of the projection direction as we now elaborate in the conclusion.

6 Conclusion

The goal of this paper is to assess the efficiency of projection methods in order to overcome the curse of dimensionality for importance sampling. Based on a new theoretical result (Theorem 3.1), we propose to project on the subspace spanned by the eigenvectors \mathbf{d}_i^* 's corresponding to the largest eigenvalues of the optimal covariance matrix Σ^* , where eigenvalues are ranked based on their image by some explicit function ℓ . Our numerical results show that if the \mathbf{d}_i^* 's were perfectly known, then projecting on them (column $\hat{\Sigma}_{\text{opt}}$ in the result tables of section 6) would greatly improve the final estimation compared to using the empirical estimation of the covariance matrix (column $\hat{\Sigma}^*$) and actually lead to results which are comparable to those obtained with the optimal covariance matrix (column Σ^*). Moreover, we show that this improvement goes through when one estimates the \mathbf{d}_i^* 's (column $\hat{\Sigma}_{\text{opt}}^{\text{+d}}$) by computing the eigenpairs of Σ^* : indeed, using $\hat{\Sigma}_{\text{opt}}^{\text{+d}}$ as covariance matrix instead of $\hat{\Sigma}^*$ gives results which remain accurate up to the dimension $n = 100$ considered in the present paper.

Moreover, we compare the directions \mathbf{d}_i^* 's, which are justified by Theorem 3.1, with the algorithm proposed in (El Masri, Morio, and Simatos 2021) which amounts to projecting on \mathbf{m}^* (column $\hat{\Sigma}_{\text{mean}}$ when \mathbf{m}^* is assumed to be known, or column $\hat{\Sigma}_{\text{mean}}^{\text{+d}}$ when one uses the estimation $\hat{\mathbf{m}}^*$ of \mathbf{m}^*). On three out of the five numerical examples considered, it turns out that \mathbf{m}^* and \mathbf{d}_1^* can be proved to be equal, or are numerically indistinguishable, and Algorithm 2 selects one projection direction. In these cases, the choices $\hat{\Sigma}_{\text{opt}}$ and $\hat{\Sigma}_{\text{mean}}$ lead to similar numerical results. The second and third test cases of Section 5.2 break the relation between \mathbf{m}^* and \mathbf{d}_1^* , and in this case, $\hat{\Sigma}_{\text{opt}}^{\text{+d}}$ clearly outperforms $\hat{\Sigma}_{\text{mean}}^{\text{+d}}$. It would be interesting to delve deeper into the relation between \mathbf{m}^* and the eigenvectors of Σ^* , and try and understand when estimating the \mathbf{d}_i^* 's instead of the simpler \mathbf{m}^* is indeed worthwhile.

These theoretical and numerical results show that the \mathbf{d}_i^* 's of Theorem 3.1 are good directions in which to estimate variance terms. With the insight gained, we see several ways to extend our results. Two in particular stand out:

- study different ways of estimating the eigenpairs $(\lambda_i^*, \mathbf{d}_i^*)$;
- incorporate this method in adaptive importance sampling schemes, in particular the cross-entropy method (Rubinstein and Kroese 2017).

For the first point, remember that we made the choice to estimate the eigenpairs of Σ^* by computing the eigenpairs of Σ^* . Moreover, in the numerical examples of Section 5.1, Section 5.4 and Section 5.5 where \mathbf{m}^* and \mathbf{d}_1^* are equal or indistinguishable, we saw that $\hat{\Sigma}_{\text{mean}}^{\text{+d}}$ performed better than $\hat{\Sigma}_{\text{opt}}^{\text{+d}}$ and we conjecture that this is because $\hat{\mathbf{m}}^*$ is a better estimator than $\hat{\mathbf{d}}_1^*$ (recall that $\mathbf{m}^* = \mathbf{d}_1^*$ for the example of Section 5.1, while in Section 5.4 and Section 5.5 they are numerically indistinguishable and so, for all

practical purposes, $\hat{\mathbf{m}}^*$ and $\hat{\mathbf{d}}_1^*$ can be considered estimators of the same direction). This suggests that improving the estimation of the \mathbf{d}_i^* 's can indeed improve the final estimation of E . Possible ways to do so consist in adapting existing results on the estimation of covariance matrices (for instance (Ledoit and Wolf 2004)) or even directly results on the estimation of eigenvalues of covariance matrices such as (Benaych-Georges and Nadakuditi 2011), (Mestre 2008a), (Mestre 2008b), (Nadakuditi and Edelman 2008), which we plan to do in future work. Moreover, it would be interesting to relax the assumption that one can sample from g^* in order to estimate g^* . For the second point, we plan to investigate how the idea of the present paper can improve the efficiency of adaptive importance sampling schemes in high dimensions. In this case, there is an additional difficulty, namely the introduction of likelihood ratios can lead to the problem of weight degeneracy which is another reason why performance of such schemes degrades in high dimensions (Bengtsson, Bickel, and Li 2008)).

We note finally that it would be interesting to consider multimodal failure functions ϕ . Indeed, with unimodal functions, the light tail of the Gaussian random variable implies that the conditional variance decreases which explains why, in all our numerical examples with an indicator function, the highest eigenvalues ranked in ℓ -order are simply the smallest eigenvalues. However, for multimodal failure functions, we may expect the conditional variance to increase and that the highest eigenvalues ranked in ℓ -order are actually the largest ones. For multimodal problems, one may want to consider different parametric families of auxiliary densities, and so it would be interesting to see whether Theorem 3.1 can be extended to more general cases.

Acknowledgement

The first author was enrolled in a PhD program co-funded by ISAE-SUPAERO and ONERA—The French Aerospace Lab. Their financial support is gratefully acknowledged. This work is part of the activities of ONERA - ISAE - ENAC joint research group.

7 Proof of Theorem 3.1 and Theorem 3.2 {#sec-proof}

We begin with a preliminary lemma.

Lemma 7.1. *Let f be the density of the standard Gaussian vector in dimension n , $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_+$ and $g_* = f\phi/\mathcal{E}$ with $\mathcal{E} = \int f\phi$. Then for any \mathbf{m} and any of the form $\Sigma = I_n + \sum_i (\alpha_i - 1)\mathbf{d}_i\mathbf{d}_i^\top$ with $\alpha_i > 0$ and the \mathbf{d}_i 's orthonormal, we have*

$$\begin{aligned} D(g^*, g_{\mathbf{m},}) &= \frac{1}{2} \sum_i \left(\log \alpha_i - \left(1 - \frac{1}{\alpha_i}\right) \mathbf{d}_i^\top \mathbf{d}_i \right) + \frac{1}{2} (\mathbf{m} - \mathbf{m}^*)^\top \Sigma^{-1} (\mathbf{m} - \mathbf{m}^*) \\ &\quad - \frac{1}{2} \|\mathbf{m}^*\|^2 - \log \mathcal{E} + \mathbb{E}_{g^*}(\log \phi(\mathbf{X})). \end{aligned} \tag{16}$$

Proof. (of Lemma 7.1) For any $\mathbf{m} \in \mathbb{R}^n$ and $\Sigma \in \mathcal{S}_n^+$, we have by definition

$$D(g^*, g_{\mathbf{m},}) = \mathbb{E}_{g^*} \left(\log \left(\frac{g^*(\mathbf{X})}{g_{\mathbf{m},}(\mathbf{X})} \right) \right) = \mathbb{E}_{g^*} \left(\log \left(\frac{\frac{\phi(\mathbf{X})e^{-\frac{1}{2}\|\mathbf{X}\|^2}}{\mathcal{E}(2\pi)^{n/2}}}{\frac{e^{-\frac{1}{2}(\mathbf{X}-\mathbf{m})^\top \Sigma^{-1}(\mathbf{X}-\mathbf{m})}}{(2\pi)^{n/2}\|\Sigma\|^{1/2}}} \right) \right)$$

and so

$$\begin{aligned} D(g^*, g_{\mathbf{m},}) &= -\frac{1}{2} \mathbb{E}_{g^*}(\|\mathbf{X}\|^2) + \frac{1}{2} \mathbb{E}_{g^*}((\mathbf{X} - \mathbf{m})^\top \Sigma^{-1} (\mathbf{X} - \mathbf{m})) \\ &\quad + \frac{1}{2} \log \|\Sigma\| - \log \mathcal{E} + \mathbb{E}_{g^*}(\log \phi(\mathbf{X})). \end{aligned}$$

Because $\mathbb{E}_{g^*}(\mathbf{X}) = \mathbf{m}^*$, we have $\mathbb{E}_{g^*}(\|\mathbf{X}\|^2) = \mathbb{E}_{g^*}(\|\mathbf{X} - \mathbf{m}^*\|^2) + \|\mathbf{m}^*\|^2$ and

$$\mathbb{E}_{g^*}((\mathbf{X} - \mathbf{m})^{\top-1}(\mathbf{X} - \mathbf{m})) = \mathbb{E}_{g^*}((\mathbf{X} - \mathbf{m}^*)^{\top-1}(\mathbf{X} - \mathbf{m}^*)) + (\mathbf{m} - \mathbf{m}^*)^{\top-1}(\mathbf{m} - \mathbf{m}^*).$$

In the following derivations, we use the linearity of the trace and of the expectation, which make these two operators commute, as well as the identity $a^{\top}b = \text{tr}(ab^{\top})$ for any two vectors a and b . With this caveat, we obtain

$$\mathbb{E}_{g^*}[\|\mathbf{X} - \mathbf{m}^*\|^2] = \mathbb{E}_{g^*}[\text{tr}((\mathbf{X} - \mathbf{m}^*)(\mathbf{X} - \mathbf{m}^*)^{\top})] = \text{tr}^*$$

and we obtain with similar arguments $\mathbb{E}_{g^*}((\mathbf{X} - \mathbf{m}^*)^{\top-1}(\mathbf{X} - \mathbf{m}^*)) = \text{tr}^{(-1*)}$. Consider now $= I_n + \sum_i (\alpha_i - 1) \mathbf{d}_i \mathbf{d}_i^{\top}$ with $\alpha_i > 0$ and the \mathbf{d}_i 's orthonormal. Then the eigenvalues of potentially different from 1 are the α_i 's (α_i is the eigenvalue associated with \mathbf{d}_i), so that

$$\log\| = \sum_i \log \alpha_i.$$

Moreover, we have $^{-1} = I_n - \sum_i \beta_i \mathbf{d}_i \mathbf{d}_i^{\top}$ with $\beta_i = 1 - 1/\alpha_i$ and so

$$\text{tr}^{(-1*)} = \text{tr}^* - \sum_i \beta_i \mathbf{d}_i^{\top*} \mathbf{d}_i.$$

Gathering the previous relation, we finally obtain the desired result. \square

Proof. (of Theorem 3.1) From Equation 16 we see that the only dependency of $D(g^*, g_{\mathbf{m}})$ in \mathbf{m} is in the quadratic term $(\mathbf{m} - \mathbf{m}^*)^{\top-1}(\mathbf{m} - \mathbf{m}^*)$. As is definite positive, this term is ≥ 0 , and so it is minimized for $\mathbf{m} = \mathbf{m}^*$. Next, we see that the derivative in α_i is given by (here and in the sequel, we see $D(g^*, g_{\mathbf{m}})$ as a function of $\mathbf{v} = (\alpha_i)_i$ and $\mathbf{d} = (\mathbf{d}_i)_i$)

$$\frac{\partial D}{\partial \alpha_i}(\mathbf{v}, \mathbf{d}) = \frac{1}{\alpha_i} - \frac{1}{\alpha_i^2} \mathbf{d}_i^{\top*} \mathbf{d}_i = \frac{1}{\alpha_i^2} (\alpha_i - \mathbf{d}_i^{\top*} \mathbf{d}_i).$$

Thus, for fixed \mathbf{d} , D is decreasing in α_i for $\alpha_i < \mathbf{d}_i^{\top*} \mathbf{d}_i$ and then increasing for $\alpha_i > \mathbf{d}_i^{\top*} \mathbf{d}_i$, which shows that, for fixed \mathbf{d} , it is minimized for $\alpha_i = \mathbf{d}_i^{\top*} \mathbf{d}_i$. For this value (and $\mathbf{m} = \mathbf{m}^*$) we have

$$D(g^*, g_{\mathbf{m}^*}) = \sum_{i=1}^k [\log(\mathbf{d}_i^{\top*} \mathbf{d}_i) + 1 - \mathbf{d}_i^{\top*} \mathbf{d}_i] + C = - \sum_{i=1}^k \ell(\mathbf{d}_i^{\top*} \mathbf{d}_i) + C \quad (17)$$

with $C = -\frac{1}{2}\|\mathbf{m}^*\|^2 - \log \mathcal{E} + \mathbb{E}_{g^*}(\log \phi(\mathbf{X}))$ independent from the \mathbf{d}_i 's. Since ℓ is decreasing and then increasing, it is clear from this expression that in order to minimize D , one must choose the \mathbf{d}_i 's in order to either maximize or minimize $\mathbf{d}_i^{\top*} \mathbf{d}_i$, whichever maximizes ℓ . Since the variational characterization of eigenvalues shows that eigenvectors precisely solve this problem, we get the desired result. \square

Proof. (of Theorem 3.2) In Equation 16, the \mathbf{m}^* and the * that appear in the right-hand side are the mean and variance of the density g^* considered in the first argument of the Kullback–Leibler divergence. In particular, if we apply Equation 16 with $\phi \equiv 1$, we have $g^* = f$, and the \mathbf{m}^* and * of the right-hand side become 0 and I_n , respectively, so that

$$D(f, g_{\mathbf{m}}) = \frac{1}{2} \sum_i \left(\log \alpha_i - \left(1 - \frac{1}{\alpha_i} \right) \right) + \frac{1}{2} \mathbf{m}^{\top-1} \mathbf{m}.$$

Now, if we consider $\mathbf{m} = \mathbf{m}^*$ and $\mathbf{d} = I + (\alpha - 1)\mathbf{d}\mathbf{d}^\top$, we obtain (using $\alpha^{-1} = I - (1 - 1/\alpha)\mathbf{d}\mathbf{d}^\top$ as mentioned in the proof of Lemma 7.1)

$$D(f, g_{\mathbf{m}^*}) = \frac{1}{2} \left(\log \alpha - \left(1 - \frac{1}{\alpha}\right) (1 + (\mathbf{d}^\top \mathbf{m}^*)^2) \right) + \frac{1}{2} \|\mathbf{m}^*\|^2.$$

Then the function $x \mapsto \log x + (1/x - 1)\gamma$ is minimized for $x = \gamma$ where it takes the value $-\ell(\gamma)$: $D(f, g_{\mathbf{m}^*})$ is therefore minimized for $\alpha = 1 + (\mathbf{d}^\top \mathbf{m}^*)^2$ and for this value, we have

$$D(f, g_{\mathbf{m}^*}) = -\frac{1}{2} \ell(1 + (\mathbf{d}^\top \mathbf{m}^*)^2) + \frac{1}{2} \|\mathbf{m}^*\|^2.$$

As ℓ is increasing in $[1, \infty)$, this last quantity is minimized by maximizing $(\mathbf{d}^\top \mathbf{m}^*)^2$, which is obtained for $\mathbf{d} = \mathbf{m}^* / \|\mathbf{m}^*\|$. The result is proved. □

References

- Agapiou, Sergios, Omiros Papaspiliopoulos, Daniel Sanz-Alonso, and Andrew M Stuart. 2017. “Importance Sampling : Intrinsic Dimension and Computational Cost.” *Statistical Science* 32 (3): 405–31. <https://doi.org/10.1214/17-STS611>.
- Ashurbekova, Karina, Antoine Usseglio-Carleve, Florence Forbes, and Sophie Achard. 2020. “Optimal Shrinkage for Robust Covariance Matrix Estimators in a Small Sample Size Setting.”
- Au, S. K., and J. L. Beck. 2003. “Important Sampling in High Dimensions.” *Structural Safety* 25 (2): 139–63. [https://doi.org/10.1016/S0167-4730\(02\)00047-4](https://doi.org/10.1016/S0167-4730(02)00047-4).
- Bassamboo, Achal, Sandeep Juneja, and Assaf Zeevi. 2008. “Portfolio Credit Risk with Extremal Dependence: Asymptotic Analysis and Efficient Simulation.” *Operations Research* 56 (3): 593–606. <https://doi.org/10.1287/opre.1080.0513>.
- Benaych-Georges, Florent, and Raj Rao Nadakuditi. 2011. “The Eigenvalues and Eigenvectors of Finite, Low Rank Perturbations of Large Random Matrices.” *Advances in Mathematics* 227 (1): 494–521. <https://doi.org/10.1016/j.aim.2011.02.007>.
- Bengtsson, Thomas, Peter Bickel, and Bo Li. 2008. “Curse-of-Dimensionality Revisited: Collapse of the Particle Filter in Very Large Scale Systems.” In *Institute of Mathematical Statistics Collections*, 316–34. Beachwood, Ohio, USA: Institute of Mathematical Statistics. <https://doi.org/10.1214/193940307000000518>.
- Bucklew, James. 2013. *Introduction to Rare Event Simulation*. Springer Science & Business Media. <https://doi.org/10.1007/978-1-4757-4078-3>.
- Bugallo, Monica F., Victor Elvira, Luca Martino, David Luengo, Joaquin Miguez, and Petar M. Djuric. 2017. “Adaptive Importance Sampling: The Past, the Present, and the Future.” *IEEE Signal Processing Magazine* 34 (4): 60–79. <https://doi.org/10.1109/MSP.2017.2699226>.
- Cappé, Olivier, Randal Douc, Arnaud Guillin, Jean-Michel Marin, and Christian P. Robert. 2008. “Adaptive Importance Sampling in General Mixture Classes.” *Statistics and Computing* 18 (4): 447–59. <https://doi.org/10.1007/s11222-008-9059-x>.
- Chan, Joshua C. C., and Dirk P. Kroese. 2012. “Improved Cross-Entropy Method for Estimation.” *Statistics and Computing* 22 (5): 1031–40. <https://doi.org/10.1007/s11222-011-9275-7>.
- Chatterjee, Sourav, and Persi Diaconis. 2018. “The Sample Size Required in Importance Sampling.” *The Annals of Applied Probability* 28 (2): 1099–1135. <https://doi.org/10.1214/17-AAP1326>.
- Cornuet, Jean-Marie, Jean-Michel Marin, Antonietta Mira, and Christian P. Robert. 2012. “Adaptive Multiple Importance Sampling.” *Scandinavian Journal of Statistics* 39 (4): 798–812. <https://doi.org/10.1111/j.1467-9469.2011.00756.x>.
- El Masri, Maxime, Jérôme Morio, and Florian Simatos. 2021. “Improvement of the Cross-Entropy Method in High Dimension for Failure Probability Estimation Through a One-Dimensional

- Projection Without Gradient Estimation.” *Reliability Engineering & System Safety* 216: 107991. <https://doi.org/10.1016/j.ress.2021.107991>.
- El-Laham, Yousef, Victor Elvira, and Mónica Bugallo. 2019. “Recursive Shrinkage Covariance Learning in Adaptive Importance Sampling.” In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, 624–28. IEEE. <https://doi.org/10.1109/CAMSAP45676.2019.9022450>.
- Elvira, Victor, Luca Martino, David Luengo, and Mónica F. Bugallo. 2019. “Generalized Multiple Importance Sampling.” *Statistical Science* 34 (1): 129–55. <https://doi.org/10.1214/18-STS668>.
- Grace, Adam W., Dirk P. Kroese, and Werner Sandmann. 2014. “Automated State-Dependent Importance Sampling for Markov Jump Processes via Sampling from the Zero-Variance Distribution.” *Journal of Applied Probability* 51 (3): 741–55. <https://doi.org/10.1239/jap/1409932671>.
- Hohenbichler, Michael, and Rüdiger Rackwitz. 1981. “Non-Normal Dependent Vectors in Structural Safety.” *Journal of the Engineering Mechanics Division* 107 (6): 1227–38. <https://doi.org/10.1061/JMCEA3.0002777>.
- Kawai, Reiichiro. 2018. “Optimizing Adaptive Importance Sampling by Stochastic Approximation.” *SIAM Journal on Scientific Computing* 40 (4): A2774–2800. <https://doi.org/10.1137/18M1173472>.
- Ledoit, Olivier, and Michael Wolf. 2004. “A Well-Conditioned Estimator for Large-Dimensional Covariance Matrices.” *Journal of Multivariate Analysis* 88 (2): 365–411. [https://doi.org/10.1016/S0047-259X\(03\)00096-4](https://doi.org/10.1016/S0047-259X(03)00096-4).
- Liu, Pei-Ling, and Armen Der Kiureghian. 1986. “Multivariate Distribution Models with Prescribed Marginals and Covariances.” *Probabilistic Engineering Mechanics* 1 (2): 105–12. [https://doi.org/10.1016/0266-8920\(86\)90033-0](https://doi.org/10.1016/0266-8920(86)90033-0).
- Mestre, Xavier. 2008a. “Improved Estimation of Eigenvalues and Eigenvectors of Covariance Matrices Using Their Sample Estimates.” *IEEE Transactions on Information Theory* 54 (11): 5113–29. <https://doi.org/10.1109/TIT.2008.929938>.
- . 2008b. “On the Asymptotic Behavior of the Sample Estimates of Eigenvalues and Eigenvectors of Covariance Matrices.” *IEEE Transactions on Signal Processing* 56 (11): 5353–68. <https://doi.org/10.1109/TSP.2008.929662>.
- Nadakuditi, Raj Rao, and Alan Edelman. 2008. “Sample Eigenvalue Based Detection of High-Dimensional Signals in White Noise Using Relatively Few Samples.” *IEEE Transactions on Signal Processing* 56 (7): 2625–38. <https://doi.org/10.1109/TSP.2008.917356>.
- Owen, Art, and Yi Zhou. 2000. “Safe and Effective Importance Sampling.” *Journal of the American Statistical Association* 95 (449): 135–43. <https://doi.org/10.1080/01621459.2000.10473909>.
- Papaoannou, Iason, Sebastian Geyer, and Daniel Straub. 2019. “Improved Cross Entropy-Based Importance Sampling with a Flexible Mixture Model.” *Reliability Engineering & System Safety* 191 (November): 106564. <https://doi.org/10.1016/j.ress.2019.106564>.
- Rubinstein, Reuven Y., and Dirk P Kroese. 2011. *The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation and Machine Learning*. New York; London: Springer. <https://doi.org/10.1007/978-1-4757-4321-0>.
- Rubinstein, Reuven Y., and Dirk P. Kroese. 2017. *Simulation and the Monte Carlo Method*. Third edition. Wiley Series in Probability and Statistics. Hoboken, New Jersey: Wiley. <https://doi.org/10.1002/9781118631980>.
- Uribe, Felipe, Iason Papaoannou, Youssef M. Marzouk, and Daniel Straub. 2021. “Cross-Entropy-Based Importance Sampling with Failure-Informed Dimension Reduction for Rare Event Simulation.” *SIAM/ASA Journal on Uncertainty Quantification* 9 (2): 818–47. <https://doi.org/10.1137/20M1344585>.