

A Large Dimensional Analysis of Least Squares Support Vector Machines

Liao Zhenyu
Couillet Romain

Laboratoire de Signaux et Systèmes (L2S, UMR8506)
CNRS-CentraleSupélec-Université Paris-Sud
3 rue Joliot-Curie, 91192 Gif-sur-Yvette, France

ZHENYU.LIAO@L2S.CENTRALESUPELEC.FR
ROMAIN.COUILLET@CENTRALESUPELEC.FR

Editor:

Abstract

In this article, a large dimensional performance analysis of kernel least squares support vector machines (LS-SVMs) is provided under the assumption of a two-class Gaussian mixture model for the input data. Building upon recent random matrix advances, when both the dimension of data p and their number n grow large at the same rate, we show that the LS-SVM decision function converges to a normal-distributed variable, the mean and variance of which depend explicitly on a local behavior of the kernel function. This theoretical result is then applied to real data sets which, despite their non-Gaussianity, exhibit a surprisingly similar behavior. Our analysis provides a deeper understanding of the mechanism into play in SVM-type methods and in particular of the impact on the choice of the kernel function as well as some of their theoretical limits.

Keywords: kernel methods, random matrices, support vector machines

1. Introduction

In the past two decades, due to their surprising classification capability and simple implementation, kernel support vector machine (SVM) (Cortes and Vapnik, 1995) and its variants (Suykens and Vandewalle, 1999; Fung and Mangasarian, 2005; Lee and Mangasarian, 2001) have been used in a wide variety of classification applications, such as face detection (Osuna et al., 1997; Papageorgiou and Poggio, 2000), handwritten digit recognition (LeCun et al., 1995), and text categorization (Drucker et al., 1999; Joachims, 1998). In all aforementioned applications, the dimension of data p and their number n are large: in the hundreds and even thousands. The significance of working in this large n, p regime is even more convincing in the Big Data paradigm today where handling data which are both numerous and large dimensional becomes increasingly common.

Firmly grounded in the framework of statistical learning theory (Vapnik, 2013), support vector machine has two main features: (i) in SVM, the training data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ are mapped into some *feature space* through a non-linear function φ , which, thanks to the so-called “kernel trick” (Scholkopf and Smola, 2001), needs not be computed explicitly, so that some *kernel function* f is introduced in place of the inner product in the feature space as $f(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^\top \varphi(\mathbf{y})$, and (ii) a standard optimization method is used to find the classifier

that both minimizes the training error and yields a good generalization performance for unknown data.

As the training of SVMs involves a quadratic programming problem, the computation complexity of SVM training algorithms can be intensive when the number of training examples n becomes large (at least quadratic with respect to n). It is thus difficult to deal with large scale problems with traditional SVMs. To cope with this limitation, least squares SVM (LS-SVM) was proposed by [Suykens and Vandewalle \(1999\)](#), providing a more computationally efficient implementation of the traditional SVMs, by using equality optimization constraints instead of inequalities, which results in an explicit solution (from a set of linear equations) rather than an implicit one in SVMs. This article is mostly concerned with this particular type of SVMs.

Trained SVMs are strongly data-driven: the data with generally unknown statistics are passed through a non-linear kernel function f and standard optimization methods are used to find the best classifier. All these features make the performance of SVM hardly traceable (within the classical finite p and n regime). To understand the mechanism of SVMs, the notion of *VC dimension* was introduced to provide bounds on the generalization performance of SVM ([Vapnik, 2013](#)), while a probabilistic interpretation of LS-SVM was discussed by [Van Gestel et al. \(2002\)](#) through a Bayesian inference approach. In other related works, connections between LS-SVMs and SVMs were revealed by [Ye and Xiong \(2007\)](#), and more relationships were shown between SVM-type and other learning methods, e.g., ; LS-SVMs and extreme learning machines (ELMs) by [Huang et al. \(2012\)](#); SVMs and regularization networks (RNs) by [Poggio et al. \(2002\)](#), etc. But a theoretical analysis of LS-SVM for large dimensional data sets of central interest here is still missing. This is all the more frustrating that such an analysis may serve as the first essential step towards more complex machine learning methods as SVMs and neural networks.

Unlike the classical regime of $n \rightarrow \infty$ while p fixed, where the diversity of the number of data provides convergence through laws of large numbers, working in the large n, p regime by letting in addition $p \rightarrow \infty$ helps exploit the diversity offered by the size of each data point, providing us with another dimension to guarantee the convergence of some key objects in our analysis, and thus plays a crucial role in the analysis of the elusive *kernel matrix* $\mathbf{K} = \{f(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^n$. Recent breakthroughs in random matrix theory have allowed one to overtake the theoretical difficulty posed by the non-linearity of the aforementioned kernel function f ([El Karoui et al., 2010](#); [Couillet et al., 2016](#)) and thus make an in-depth analysis of LS-SVM possible in the large n, p regime. These tools were notably used to assess the performance of the popular Ng-Weiss-Jordan kernel spectral clustering methods for large data sets ([Couillet et al., 2016](#)), in the analysis of graphed-based semi-supervised learning ([Mai and Couillet, 2017](#)) or for the development of novel kernel subspace clustering methods ([Couillet and Kammoun, 2016](#)).

Similar to these works, in this article, we provide a performance analysis of LS-SVM, in the regime of $n, p \rightarrow \infty$ and $p/n \rightarrow c_0 \in (0, \infty)$, under the assumption of a two-class Gaussian mixture model of means $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$ and covariance matrices $\mathbf{C}_1, \mathbf{C}_2$ for the input data. The Gaussian assumption may seem artificial to the practitioners, but reveals first insights into how SVM-type methods deal with the information in means and covariances from a more quantitative point of view. Besides, the early investigations ([Couillet et al., 2016](#); [Mai and Couillet, 2017](#)) have revealed that the behavior of some machine learning methods

under Gaussian or deterministic practical input data sets are a close match, despite the obvious non-Gaussianity of the latter.

Our main finding is that, as by [Couillet et al. \(2016\)](#), in the large n, p regime and under suitable conditions on the Gaussian input statistics, a non-trivial asymptotic classification error rate (i.e., neither 0 nor 1) can be obtained and the decision function of LS-SVM converges to a Gaussian random variable whose mean and variance depend on the statistics of the two different classes and on the behavior of the kernel function f evaluated at $2(n_1 \text{tr } \mathbf{C}_1 + n_2 \text{tr } \mathbf{C}_2)/(np)$, with n_1 and n_2 the number of instances in each class. This brings novel insights into some key issues of SVM-type methods such as kernel function selection and parameter optimization (see [Van Gestel et al., 2002](#); [Cherkassky and Ma, 2004](#); [Chapelle et al., 2002](#); [Ayat et al., 2005](#); [Weston et al., 2000](#); [Huang and Wang, 2006](#)), but from a more theoretically grounded viewpoint. More importantly, we again confirm through simulations that our theoretical findings closely match the performance obtained on real data sets, which conveys a strong applicative motivation for this work.

In the remainder of the article, we provide a rigorous statement of our main results. The problem of LS-SVM is discussed in Section 2 and our model and main results presented in Section 3, while all proofs are deferred to the appendices. Then in Section 4, attention will be paid on some special cases that are more analytically tractable and thus convey a deeper understanding of the mechanisms of LS-SVM. Section 5 concludes the paper by summarizing the main results and outlining future research directions.

Notations: Boldface lowercase (uppercase) characters stand for vectors (matrices), and scalars non-boldface respectively. $\mathbf{1}_n$ is the column vector of ones of size n , $\mathbf{0}_n$ the column vector of zeroes, and \mathbf{I}_n the $n \times n$ identity matrix. The notation $(\cdot)^\top$ denotes the transpose operator. The norm $\|\cdot\|$ is the Euclidean norm for vectors and the operator norm for matrices. The notation $\mathbb{P}(\cdot)$ denotes the probability measure of a random variable. The notation \xrightarrow{d} denotes convergence in distribution and $\xrightarrow{\text{a.s.}}$ almost sure convergence, respectively. The operator $\mathcal{D}(\mathbf{v}) = \mathcal{D}\{v_a\}_{a=1}^k$ is the diagonal matrix having v_1, \dots, v_k as its ordered diagonal elements. We denote $\{v_a\}_{a=1}^k$ a column vector with a -th entry (or block entry) v_a (which may be a vector), while $\{V_{ab}\}_{a,b=1}^k$ denotes a square matrix with entry (or block-entry) (a, b) given by V_{ab} (which may be a matrix).

2. Problem Statement

Least squares support vector machines (LS-SVMs) are a modification of the standard SVM ([Vapnik, 2013](#)) introduced by [Suykens and Vandewalle \(1999\)](#) developed to overcome the drawbacks of SVM related to computational efficiency. The optimization problem has half the number of parameters and benefits from solving a linear system of equations instead of a quadratic programming problem as in standard SVM and is thus more practical for large dimensional learning tasks. In this article, we will focus on a binary classification problem using LS-SVM as described in the following paragraph.

Given a set of training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ of size n , where data $\mathbf{x}_i \in \mathbb{R}^p$ and labels $y_i \in \{-1, 1\}$, the objective of LS-SVM is to devise a decision function $g(\mathbf{x})$ that ideally maps all \mathbf{x}_i in the training set to y_i and subsequently all unknown data \mathbf{x} to their corresponding y value. Here we note $\mathbf{x}_i \in \mathcal{C}_1$ if $y_i = -1$ and $\mathbf{x}_i \in \mathcal{C}_2$ if $y_i = 1$ and shall say that \mathbf{x}_i belongs to class \mathcal{C}_1 or class \mathcal{C}_2 . Due to the often nonlinear separability

of these training data in the input space \mathbb{R}^p , in most cases, one associates the training data \mathbf{x}_i to some feature space \mathcal{H} through a nonlinear mapping $\varphi : \mathbf{x}_i \mapsto \varphi(\mathbf{x}_i) \in \mathcal{H}$. Constrained optimization methods are then used to define a separating hyperplane in space \mathcal{H} with direction vector \mathbf{w} and correspondingly find a function $g(\mathbf{x}) = \mathbf{w}^\top \varphi(\mathbf{x}) + b$ that minimizes the training errors $e_i = y_i - (\mathbf{w}^\top \varphi(\mathbf{x}_i) + b)$, and at the same time yields a good generalization performance by minimizing the norm of \mathbf{w} (Smola and Schölkopf, 2004). The LS-SVM approach consists in minimizing the squared errors e_i^2 , thus resulting in

$$\arg \min_{\mathbf{w}, b} J(\mathbf{w}, e) = \|\mathbf{w}\|^2 + \frac{\gamma}{n} \sum_{i=1}^n e_i^2 \quad (1)$$

$$\text{such that } y_i = \mathbf{w}^\top \varphi(\mathbf{x}_i) + b + e_i, \quad i = 1, \dots, n$$

where $\gamma > 0$ is a penalty factor that weights the structural risk $\|\mathbf{w}\|^2$ against the empirical risk $\frac{1}{n} \sum_{i=1}^n e_i^2$.

The problem can be solved by introducing Lagrange multipliers $\alpha_i, i = 1, \dots, n$ with solution $\mathbf{w} = \sum_{i=1}^n \alpha_i \varphi(\mathbf{x}_i)$, where, letting $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n]^\top$, we have

$$\begin{cases} \boldsymbol{\alpha} &= \mathbf{S}^{-1} \left(\mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \right) \mathbf{y} = \mathbf{S}^{-1} (\mathbf{y} - b \mathbf{1}_n) \\ b &= \frac{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{y}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \end{cases} \quad (2)$$

with $\mathbf{S} = \mathbf{K} + \frac{n}{\gamma} \mathbf{I}_n$ and $\mathbf{K} \triangleq \{\varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j)\}_{i,j=1}^n$ referred to as the kernel matrix (Suykens and Vandewalle, 1999).

Given $\boldsymbol{\alpha}$ and b , a new datum \mathbf{x} is then classified into class \mathcal{C}_1 or \mathcal{C}_2 depending on the value of the following decision function

$$g(\mathbf{x}) = \boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x}) + b \quad (3)$$

where $\mathbf{k}(\mathbf{x}) = \{\varphi(\mathbf{x})^\top \varphi(\mathbf{x}_j)\}_{j=1}^n \in \mathbb{R}^n$. More precisely, \mathbf{x} is associated to class \mathcal{C}_1 if $g(\mathbf{x})$ takes a small value (below a certain threshold ξ) and to class \mathcal{C}_2 otherwise.¹

Thanks to the so-called “kernel trick” (Schölkopf and Smola, 2001), as shown in (2) and (3) that, both in the “training” and “testing” steps, one only needs to evaluate the inner product $\varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j)$ or $\varphi(\mathbf{x})^\top \varphi(\mathbf{x}_j)$, and never needs to know explicitly the mapping $\varphi(\cdot)$. In the rest of this article, we assume that the kernel is *translation invariant* and focus on kernel functions $f : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ that satisfy $\varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j) = f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)$ and shall redefine \mathbf{K} and $\mathbf{k}(\mathbf{x})$ for datapoint \mathbf{x} as²

$$\mathbf{K} = \{f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)\}_{i,j=1}^n \quad (4)$$

$$\mathbf{k}(\mathbf{x}) = \{f(\|\mathbf{x} - \mathbf{x}_j\|^2/p)\}_{j=1}^n. \quad (5)$$

Some often used kernel functions are the Gaussian radial basis (RBF) kernel $f(x) = \exp(-\frac{x}{2\sigma^2})$ with $\sigma > 0$ and the polynomial kernel $f(x) = \sum_{i=0}^d a_i x^i$ with $d \geq 1$.

In the rest of this article, we will focus on the performance of LS-SVM, in the large n, p regime, by studying the asymptotic behavior of the decision function $g(\mathbf{x})$ given in (3), on a two-class classification problem with some statistical properties, whose model will be specified in the next section.

1. Since data from \mathcal{C}_1 are labeled -1 while data from \mathcal{C}_2 are labeled 1 .
 2. As shall be seen later, the division by p here is a convenient normalization in the large n, p regime.

3. Main Results

3.1 Model and Assumptions

Evaluating the performance of LS-SVM is made difficult by the heavily data-driven aspect of the method. As a first approximation, we shall assume in this article that all \mathbf{x}_i 's are extracted from a Gaussian mixture, thereby allowing for a thorough theoretical analysis.

Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ be independent vectors belonging to two distribution classes $\mathcal{C}_1, \mathcal{C}_2$, with $\mathbf{x}_1, \dots, \mathbf{x}_{n_1} \in \mathcal{C}_1$ and $\mathbf{x}_{n_1+1}, \dots, \mathbf{x}_n \in \mathcal{C}_2$ (so that class \mathcal{C}_1 has cardinality n_1 and class \mathcal{C}_2 has cardinality $n - n_1 = n_2$). We assume that $\mathbf{x}_i \in \mathcal{C}_a$ for $a \in \{1, 2\}$ if

$$\mathbf{x}_i = \boldsymbol{\mu}_a + \sqrt{p}\boldsymbol{\omega}_i$$

for some $\boldsymbol{\mu}_a \in \mathbb{R}^p$ and $\boldsymbol{\omega}_i \sim \mathcal{N}(0, p^{-1}\mathbf{C}_a)$, with $\mathbf{C}_a \in \mathbb{R}^{p \times p}$ some symmetric and nonnegative definite matrix.

To achieve an asymptotically non-trivial misclassification rate (i.e., neither 0 nor 1), similar to the work of [Couillet et al. \(2016\)](#), we shall consider the large dimensional regime where both n and p grow large simultaneously with the following growth rate assumption:

Assumption 1 (Growth Rate) *As $n, p \rightarrow \infty$, for $a \in \{1, 2\}$, the following conditions hold.*

- **Data scaling:** $p/n \triangleq c_0 \rightarrow \bar{c}_0 > 0$.
- **Class scaling:** $n_a/n \triangleq c_a \rightarrow \bar{c}_a > 0$.
- **Mean scaling:** $\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\| = O(1)$.
- **Covariance scaling:** $\|\mathbf{C}_a\| = O(1)$ and $\text{tr}(\mathbf{C}_2 - \mathbf{C}_1) = O(\sqrt{n})$.
- for $\mathbf{C}^\circ \triangleq \frac{n_1}{n}\mathbf{C}_1 + \frac{n_2}{n}\mathbf{C}_2$, $\frac{2}{p}\text{tr}\mathbf{C}^\circ \rightarrow \tau > 0$ as $n, p \rightarrow \infty$.

Aside from the last assumption, stated here mostly for technical convenience, if a single one of these assumptions would be relaxed, e.g., if $\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|$ scales faster than $O(1)$, then asymptotic perfect classification would be achieved. If instead $\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\| = o(1)$, $\text{tr}(\mathbf{C}_2 - \mathbf{C}_1) = o(\sqrt{n})$ and $\text{tr}(\mathbf{C}_a - \mathbf{C}^\circ)(\mathbf{C}_b - \mathbf{C}^\circ) = o(n)$ for all $a, b \in \{1, 2\}$, classification becomes asymptotically impossible. Refer to [Couillet et al. \(2016\)](#) for more details.

A key observation, also made by [Couillet et al. \(2016\)](#), is that, under Assumption 1, for all pairs $i \neq j$,

$$\|\mathbf{x}_i - \mathbf{x}_j\|^2/p \xrightarrow{\text{a.s.}} \tau \tag{6}$$

and the convergence is even uniform across all $i \neq j$. This remark is the crux of all subsequent results (note that, surprisingly at first, it states that all data are essentially equivalent, irrespective of classes and that the matrix \mathbf{K} defined in (4) has all its entries essentially equal “in the limit”).

The function f defining the kernel matrix \mathbf{K} in (4) shall be requested to satisfy the following additional technical assumption:

Assumption 2 (Kernel Function) *The function f is a three-times differentiable function in a neighborhood of τ .*

The objective of this article is to assess the performance of LS-SVM, under the setting of Assumptions 1 and 2, as $n, p \rightarrow \infty$, by studying the asymptotic behavior of the decision function $g(\mathbf{x})$ defined in (3). Following the work of El Karoui et al. (2010); Couillet et al. (2016), under our basic settings, the convergence in (6) makes it possible to linearize the kernel matrix \mathbf{K} around the matrix $f(\tau)\mathbf{1}_n\mathbf{1}_n^\top$, and thus the intractable non-linear kernel matrix \mathbf{K} can be asymptotically linearized in the large n, p regime. As such, since the decision function $g(\mathbf{x})$ is explicitly defined as a function of \mathbf{K} (through $\boldsymbol{\alpha}$ and b as defined in (2)), one can work out an asymptotic linearization of $g(\mathbf{x})$ as a function of the kernel function f and the statistics of the data \mathbf{x}_i 's. This analysis, presented in detail in Appendix A, allows one to reveal the relationship between the performance of LS-SVM and the kernel function f as well as the given learning task, in the case of Gaussian input data as $n, p \rightarrow \infty$.

3.2 Asymptotic Behavior of $g(\mathbf{x})$

Before going into our main results, a few notations need to be introduced. In the remainder of the article, we shall use the following deterministic and random elements notations

$$\begin{aligned}\mathbf{P} &\triangleq \mathbf{I}_n - \frac{\mathbf{1}_n\mathbf{1}_n^\top}{n} \in \mathbb{R}^{n \times n} \\ \boldsymbol{\Omega} &\triangleq [\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_n] \in \mathbb{R}^{p \times n} \\ \boldsymbol{\psi} &\triangleq \{\|\boldsymbol{\omega}_i\|^2 - \mathbb{E}[\|\boldsymbol{\omega}_i\|^2]\}_{i=1}^n \in \mathbb{R}^n.\end{aligned}$$

Under Assumptions 1 and 2, following up Couillet et al. (2016), one can approximate the kernel matrix \mathbf{K} by $\hat{\mathbf{K}}$ in such a way that

$$\|\mathbf{K} - \hat{\mathbf{K}}\| \xrightarrow{\text{a.s.}} 0$$

with $\hat{\mathbf{K}} = -2f'(\tau)(\mathbf{M} + \mathbf{V}\mathbf{V}^\top) + (f(0) - f(\tau) + \tau f'(\tau))\mathbf{I}_n$ for some matrices \mathbf{M} and \mathbf{V} , where \mathbf{M} is a standard random matrix model (of operator norm $O(1)$) and $\mathbf{V}\mathbf{V}^\top$ a small rank matrix (of operator norm $O(n)$), which depends both on $\mathbf{P}, \boldsymbol{\Omega}, \boldsymbol{\psi}$ and on the class features $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$ and $\mathbf{C}_1, \mathbf{C}_2$. The same analysis is applied to the vector $\mathbf{k}(\mathbf{x})$ by similarly defining the following random variables for a new datum $\mathbf{x} \in \mathcal{C}_a$, $a \in \{1, 2\}$

$$\begin{aligned}\boldsymbol{\omega}_{\mathbf{x}} &\triangleq \frac{\mathbf{x} - \boldsymbol{\mu}_a}{\sqrt{p}} \in \mathbb{R}^p \\ \boldsymbol{\psi}_{\mathbf{x}} &\triangleq \|\boldsymbol{\omega}_{\mathbf{x}}\|^2 - \mathbb{E}[\|\boldsymbol{\omega}_{\mathbf{x}}\|^2] \in \mathbb{R}.\end{aligned}$$

Based on the approximation $\mathbf{K} \approx \hat{\mathbf{K}}$, a Taylor expansion is then performed on the inverse $\mathbf{S}^{-1} = (\mathbf{K} + n\mathbf{I}_n/\gamma)^{-1}$ to obtain an approximation of \mathbf{S}^{-1} , and subsequently on $\boldsymbol{\alpha}$ and b which depend explicitly on the random matrix \mathbf{S}^{-1} . At last, plugging these results in $g(\mathbf{x}) = \boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x}) + b$, one finds the main technical result of this article as follows:

Theorem 1 (Random Equivalent) *Let Assumptions 1 and 2 hold, and $g(\mathbf{x})$ be defined by (3). Then, as $n, p \rightarrow \infty$, $n(g(\mathbf{x}) - \hat{g}(\mathbf{x})) \xrightarrow{\text{a.s.}} 0$, where*

$$\hat{g}(\mathbf{x}) = \begin{cases} c_2 - c_1 + \gamma (\mathfrak{P} - 2c_1c_2^2\mathfrak{D}), & \text{if } \mathbf{x} \in \mathcal{C}_1 \\ c_2 - c_1 + \gamma (\mathfrak{P} + 2c_1^2c_2\mathfrak{D}), & \text{if } \mathbf{x} \in \mathcal{C}_2 \end{cases} \quad (7)$$

with

$$\mathfrak{P} = -\frac{2f'(\tau)}{n} \mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \boldsymbol{\omega}_x - \frac{4c_1 c_2 f'(\tau)}{\sqrt{p}} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \boldsymbol{\omega}_x + 2c_1 c_2 f''(\tau) \psi_x \frac{\text{tr}(\mathbf{C}_2 - \mathbf{C}_1)}{p} \quad (8)$$

$$\mathfrak{D} = -\frac{2f'(\tau)}{p} \|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2 + \frac{f''(\tau)}{p^2} (\text{tr}(\mathbf{C}_2 - \mathbf{C}_1))^2 + \frac{2f''(\tau)}{p^2} \text{tr}((\mathbf{C}_2 - \mathbf{C}_1)^2). \quad (9)$$

Leaving the proof to Appendix A, Theorem 1 tells us that the decision function $g(\mathbf{x})$ has a random equivalent $\hat{g}(\mathbf{x})$ that consists of three parts:

1. the deterministic term $c_2 - c_1$ of order $O(1)$ that depends on the number of instances in each class, which essentially comes from the term $\mathbf{1}_n^\top \mathbf{y}/n$ in b ;
2. the “noisy” term \mathfrak{P} of order $O(n^{-1})$ which is a function of the zero mean random variables $\boldsymbol{\omega}_x$ and ψ_x , thus in particular $\mathbb{E}[\mathfrak{P}] = 0$;
3. the “informative” term containing \mathfrak{D} , also of order $O(n^{-1})$, which features the deterministic differences between the two classes.

From Theorem 1, under the basic settings of Assumption 1, as $n \rightarrow \infty$, for Gaussian data $\mathbf{x} \in \mathcal{C}_a$, $a \in \{1, 2\}$, we can show that $\hat{g}(\mathbf{x})$ (and thus $g(\mathbf{x})$) converges to a random Gaussian variable the mean and variance of which are given by the following theorem. The proof is deferred to Appendix B.

Theorem 2 (Gaussian Approximation) *Under the setting of Theorem 1, $n(g(\mathbf{x}) - G_a) \xrightarrow{d} 0$, where*

$$G_a \sim \mathcal{N}(\mathbf{E}_a, \text{Var}_a)$$

with

$$\mathbf{E}_a = \begin{cases} c_2 - c_1 - 2c_2 \cdot c_1 c_2 \gamma \mathfrak{D}, & a = 1 \\ c_2 - c_1 + 2c_1 \cdot c_1 c_2 \gamma \mathfrak{D}, & a = 2 \end{cases}$$

$$\text{Var}_a = 8\gamma^2 c_1^2 c_2^2 (\mathcal{V}_1^a + \mathcal{V}_2^a + \mathcal{V}_3^a)$$

and

$$\mathcal{V}_1^a = \frac{(f''(\tau))^2}{p^4} (\text{tr}(\mathbf{C}_2 - \mathbf{C}_1))^2 \text{tr} \mathbf{C}_a^2$$

$$\mathcal{V}_2^a = \frac{2(f'(\tau))^2}{p^2} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \mathbf{C}_a (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$$

$$\mathcal{V}_3^a = \frac{2(f'(\tau))^2}{np^2} \left(\frac{\text{tr} \mathbf{C}_1 \mathbf{C}_a}{c_1} + \frac{\text{tr} \mathbf{C}_2 \mathbf{C}_a}{c_2} \right).$$

Theorem 2 is our main practical result as it allows one to evaluate the large n, p performance of LS-SVM for Gaussian data. While dwelling on the implications of Theorem 1 and 2, several remarks and discussions are in order.

Remark 3 (Dominant Bias) From Theorem 1, under the key Assumption 1, both the random “noisy” \mathfrak{P} and the deterministic “informative” \mathfrak{D} terms are of order $O(n^{-1})$, which means for the decision function $g(\mathbf{x}) = c_2 - c_1 + O(n^{-1})$. This result somehow contradicts the classical decision criterion proposed by [Suykens and Vandewalle \(1999\)](#), based on the sign of $g(\mathbf{x})$, i.e., \mathbf{x} is associated to class \mathcal{C}_1 if $g(\mathbf{x}) < 0$ and to class \mathcal{C}_2 otherwise. When $c_1 \neq c_2$, this would lead to an asymptotic classification of all new data \mathbf{x} ’s in the same class as $n \rightarrow \infty$. Instead, a first result of Theorem 1 is that the decision threshold ξ should be taken as $\xi = \xi_n = c_2 - c_1 + O(n^{-1})$ when $c_1 \neq c_2$.

This finding was mentioned by [Van Gestel et al. \(2002\)](#) through a Bayesian inference analysis: the term $c_2 - c_1$ appears in the “bias term” b under the form of prior class probabilities $P(y = -1)$, $P(y = 1)$ and allows for adjusting classification problems with different prior class probabilities in the training and test sets. This idea of a (static) bias term correction has also been applied in the work of [Evgeniou et al. \(2000\)](#) in order to improve the validation set performance. Here we confirm this finding by Figure 1 with $c_1 = 1/4$ and $c_2 = 3/4$, where the histograms of $g(\mathbf{x})$ for $\mathbf{x} \in \mathcal{C}_1$ and \mathcal{C}_2 center somewhere close to $c_2 - c_1 = 0.5$, thus resulting in a trivial classification if one takes $\xi = 0$ because $P(g(\mathbf{x}) < \xi \mid \mathbf{x} \in \mathcal{C}_1) \rightarrow 0$ and $P(g(\mathbf{x}) > \xi \mid \mathbf{x} \in \mathcal{C}_2) \rightarrow 1$ as $n, p \rightarrow \infty$ (the convergence being in fact an equality for finite n, p in this particular figure).

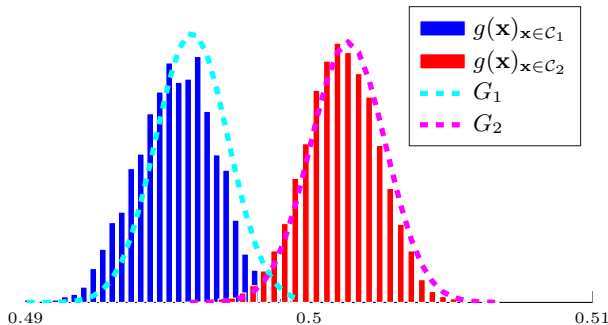


Figure 1: Gaussian approximation of $g(\mathbf{x})$, $n = 256$, $p = 512$, $c_1 = 1/4$, $c_2 = 3/4$, $\gamma = 1$, Gaussian kernel with $\sigma^2 = 1$, $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$ with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 3; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{5}{\sqrt{p}})$.

An alternative to alleviate the bias issue is to normalize the label vector \mathbf{y} . From the proof of Theorem 1 in Appendix A we see the bias of $c_2 - c_1$ is due to the fact that in b one has $\mathbf{1}_n^T \mathbf{y} = n_2 - n_1 \neq 0$. Thus, one may normalize the labels y_i as $y_i^* = -1/c_1$ if $\mathbf{x}_i \in \mathcal{C}_1$ and $y_i^* = 1/c_2$ if $\mathbf{x}_i \in \mathcal{C}_2$, so that the relation $\mathbf{1}_n^T \mathbf{y}^* = 0$ is satisfied. This formulation is also referred to as the *Fishers targets*: $\{-n/n_1, n/n_2\}$ in the context of Kernel Fisher Discriminant Analysis ([Baudat and Anouar, 2000](#); [Scholkopf and Mullert, 1999](#)). With the aforementioned normalized labels \mathbf{y}^* , we have the following lemma that reveals the connection between the corresponding decision function $g^*(\mathbf{x})$ and $g(\mathbf{x})$.

Lemma 4 Let $g(\mathbf{x})$ be defined by (3) and $g^*(\mathbf{x})$ be defined as $g^*(\mathbf{x}) = (\boldsymbol{\alpha}^*)^\top \mathbf{k}(\mathbf{x}) + b^*$, with $(\boldsymbol{\alpha}^*, b^*)$ given by (2) with \mathbf{y}^* in the place of \mathbf{y} . Then, the following relation holds true

$$g(\mathbf{x}) - (c_2 - c_1) = 2c_1c_2g^*(\mathbf{x}).$$

Proof From (2) and (3) we get

$$g(x) = \mathbf{y}^\top \left(\mathbf{S}^{-1} - \frac{\mathbf{S}^{-1} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \right) \mathbf{k}(\mathbf{x}) + \frac{\mathbf{y}^\top \mathbf{S}^{-1} \mathbf{1}_n}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} = \mathbf{y}^\top \boldsymbol{\varpi}$$

with $\boldsymbol{\varpi} = \left(\mathbf{S}^{-1} - \frac{\mathbf{S}^{-1} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \right) \mathbf{k}(\mathbf{x}) + \frac{\mathbf{S}^{-1} \mathbf{1}_n}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n}$. Besides, note that $\mathbf{1}_n^\top \boldsymbol{\varpi} = 1$. We thus have

$$\begin{aligned} g(\mathbf{x}) - (c_2 - c_1) &= \mathbf{y}^\top \boldsymbol{\varpi} - (c_2 - c_1) \mathbf{1}_n^\top \boldsymbol{\varpi} = 2c_1c_2 \left(\frac{\mathbf{y} - (c_2 - c_1) \mathbf{1}_n}{2c_1c_2} \right)^\top \boldsymbol{\varpi} \\ &= 2c_1c_2 (\mathbf{y}^*)^\top \boldsymbol{\varpi} = 2c_1c_2 g^*(\mathbf{x}) \end{aligned}$$

which concludes the proof. ■

As a consequence of Lemma 4, instead of Theorem 2 for standard labels \mathbf{y} , one would have the following corollary for the corresponding Gaussian approximation of $g^*(\mathbf{x})$ when normalized labels \mathbf{y}^* are used.

Corollary 5 (Gaussian Approximation of $g^*(\mathbf{x})$) Under the setting of Theorem 1, and with $g^*(\mathbf{x})$ defined in Lemma 4, $n(g^*(\mathbf{x}) - G_a^*) \xrightarrow{d} 0$, where

$$G_a^* \sim \mathcal{N}(\mathbf{E}_a^*, \text{Var}_a^*)$$

with

$$\begin{aligned} \mathbf{E}_a^* &= \begin{cases} -c_2\gamma\mathcal{D}, & a = 1 \\ +c_1\gamma\mathcal{D}, & a = 2 \end{cases} \\ \text{Var}_a^* &= 2\gamma^2 (\mathcal{V}_1^a + \mathcal{V}_2^a + \mathcal{V}_3^a) \end{aligned}$$

and \mathcal{D} is defined by (9), $\mathcal{V}_1^a, \mathcal{V}_2^a$ and \mathcal{V}_3^a as in Theorem 2.

Figure 2 illustrates this result in the same settings as Figure 1. Compared to Figure 1, one can observe that in Figure 2 both histograms are now centered close to 0 (at distance $O(n^{-1})$ from zero) instead of $c_2 - c_1 = 1/2$. Still, even in the case where normalized labels \mathbf{y}^* are used as observed in Figure 2 (where the histograms cross at about $-0.004 \approx 1/n$), taking $\xi = 0$ as a decision threshold may not be an appropriate choice, as $\mathbf{E}_1^* \neq -\mathbf{E}_2^*$.

Remark 6 (Insignificance of γ) As a direct result of Theorem 1 and Remark 3, note in (7) that $\hat{g}(\mathbf{x}) - (c_2 - c_1)$ is proportional to the hyperparameter γ , which indicates that, rather surprisingly, the tuning of γ is (asymptotically) of no importance when $n, p \rightarrow \infty$ since it does not alter the classification statistics when one uses the sign of $g(\mathbf{x}) - (c_2 - c_1)$ for the decision.³

3. This remark is only valid only under Assumption 1 and $\gamma = O(1)$, i.e., γ is considered to remain a constant as $n, p \rightarrow \infty$.

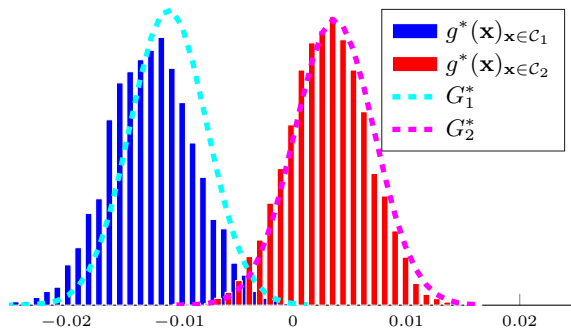


Figure 2: Gaussian approximation of $g^*(\mathbf{x})$, $n = 256$, $p = 512$, $c_1 = 1/4$, $c_2 = 3/4$, $\gamma = 1$, Gaussian kernel with $\sigma^2 = 1$, $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 3; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{5}{\sqrt{p}})$.

Letting $Q(x) = \frac{1}{2\pi} \int_x^\infty \exp(-t^2/2) dt$, from Theorem 2 and Corollary 5, we now have the following immediate corollary for the (asymptotic) classification error rate:

Corollary 7 (Asymptotic Error Rate) *Under the setting of Theorem 1, for a threshold ξ_n possibly depending on n , as $n \rightarrow \infty$,*

$$\mathbb{P}(g(\mathbf{x}) > \xi_n \mid \mathbf{x} \in \mathcal{C}_1) - Q\left(\frac{\xi_n - E_1}{\sqrt{\text{Var}_1}}\right) \rightarrow 0 \quad (10)$$

$$\mathbb{P}(g(\mathbf{x}) < \xi_n \mid \mathbf{x} \in \mathcal{C}_2) - Q\left(\frac{E_2 - \xi_n}{\sqrt{\text{Var}_2}}\right) \rightarrow 0 \quad (11)$$

with E_a and Var_a given in Theorem 2.

Obviously, Corollary 7 is only meaningful when $\xi_n = c_2 - c_1 + O(n^{-1})$ as recalled earlier. Besides, it is clear from Lemma 4 and Corollary 5 that $\mathbb{P}(g(\mathbf{x}) > \xi_n \mid \mathbf{x} \in \mathcal{C}_a) = \mathbb{P}(g^*(\mathbf{x}) > \xi_n - (c_2 - c_1) \mid \mathbf{x} \in \mathcal{C}_a)$, so that Corollary 7 extends naturally to $g^*(\mathbf{x})$ when normalized labels \mathbf{y}^* are applied.

Corollary 7 allows one to compute the asymptotic error rate of classification as a function of E_a , Var_a and the threshold ξ_n . Combined with Theorem 2, one may note the significance of a proper choice of the kernel function f . For instance, if $f'(\tau) = 0$, the term $\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1$ vanishes from the mean and variance of G_a , meaning that the classification of LS-SVM will not rely (at least asymptotically and under Assumption 1) on the differences in means of the two classes. Figure 3 corroborates this finding with the same theoretical Gaussian approximations G_1 and G_2 in subfigures (a) and (b). When $\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2$ varies from 0 in (a) to 18 in (b), the distribution of $g(\mathbf{x})$ remains almost the same for $\mathbf{x} \in \mathcal{C}_1$ and \mathcal{C}_2 respectively.

More traceable special cases and discussions around the choice of kernel function f will be given in the next section.

4. Special Cases and Further Discussions

In this section, some special cases will be investigated in detail, during which we will reveal the significance of the kernel function f as well as the training set.

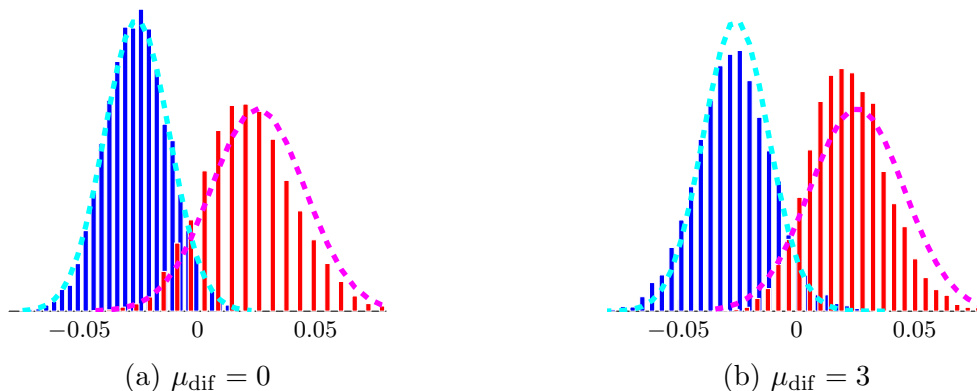


Figure 3: Gaussian approximation of $g(\mathbf{x})$, $n = 256$, $p = 512$, $c_1 = c_2 = 1/2$, $\gamma = 1$, polynomial kernel with $f(\tau) = 4$, $f'(\tau) = 0$, and $f''(\tau) = 2$. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; \mu_{\text{dif}}; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{5}{\sqrt{p}})$.

4.1 More Discussions Over the Kernel Function f

Following the discussion at the end of Section 3, if $f'(\tau) = 0$, the information about the statistical means of the two different classes is somehow lost and will not help perform the classification. Nonetheless, we find that, rather surprisingly, if one further assumes $\text{tr } \mathbf{C}_1 = \text{tr } \mathbf{C}_2 + o(\sqrt{n})$, using a kernel f that satisfies $f'(\tau) = 0$ results in $\text{Var}_a = 0$ while \mathbf{E}_a may remain non-zero, thereby ensuring a vanishing error rate (as long as $f''(\tau) \neq 0$). Intuitively speaking, the kernels that satisfy $f'(\tau) = 0$ play an important role in extracting the information of “shape” of both classes, making the classification extremely accurate even in cases that are deemed impossible to classify according to our comments following Assumption 1. This phenomenon was also remarked by Couillet et al. (2016) and deeply investigated by Couillet and Kammoun (2016). Figure 4 substantiates this finding for $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}$, for which $\text{tr } \mathbf{C}_1 = \text{tr } \mathbf{C}_2 = p$.

Remark 8 (Condition on Kernel Function f) *From Theorem 2 and Corollary 5, one observes that $|\mathbf{E}_1 - \mathbf{E}_2|$ is always proportional to the “informative” term \mathfrak{D} and should, for fixed Var_a , be made as large as possible to avoid the overlap of $g(\mathbf{x})$ for \mathbf{x} from different classes. Since Var_a does not depend on the signs of $f'(\tau)$ and $f''(\tau)$, it is easily deduced that, to achieve optimal classification performance, one needs to choose the kernel function f such that $f(\tau) > 0$, $f'(\tau) < 0$ and $f''(\tau) > 0$.*

Incidentally, the condition in Remark 8 is naturally satisfied for Gaussian kernel $f(x) = \exp(-x/(2\sigma^2))$ for any σ , meaning that, even without specific tuning of the kernel parameter σ through cross validation or other techniques, LS-SVM is expected to perform rather well with a Gaussian kernel (as shown in Figure 5), which is not always the case for polynomial kernels. This especially entails, for a second-order polynomial kernel given by

4. Unless particularly stated, the classification error will be understood as $c_1\text{P}(g(\mathbf{x}) > \xi_n | \mathbf{x} \in \mathcal{C}_1) + c_2\text{P}(g(\mathbf{x}) < \xi_n | \mathbf{x} \in \mathcal{C}_2)$.

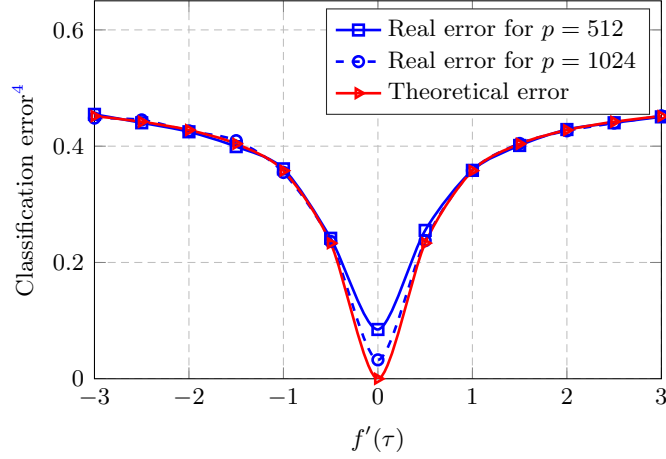


Figure 4: Performance of LS-SVM, $c_0 = 1/4$, $c_1 = c_2 = 1/2$, $\gamma = 1$, polynomial kernel with $f(\tau) = 4$, $f''(\tau) = 2$. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{0}_p$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}$.

$f(x) = a_2x^2 + a_1x + a_0$, that special attention should be paid to meeting the aforementioned condition when tuning the kernel parameters a_2 , a_1 and a_0 . Figure 6 attests of this finding with Gaussian input data. A rapid increase in classification error rate can be observed both in theory and in practice as soon as the condition $f'(\tau) < 0$, $f''(\tau) > 0$ is no longer satisfied.

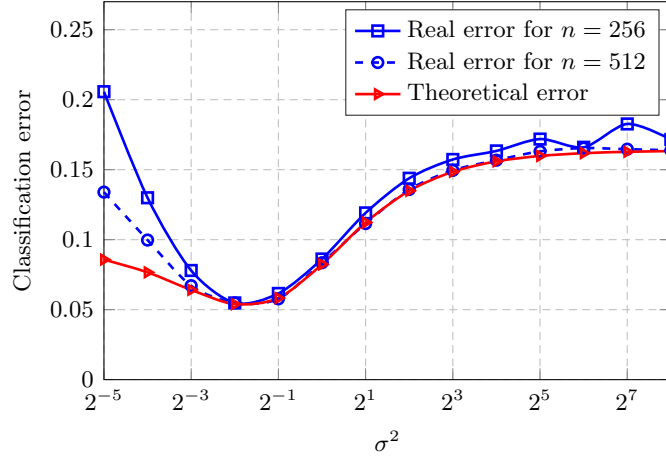


Figure 5: Performance of LS-SVM, $c_0 = 2$, $c_1 = c_2 = 1/2$, $\gamma = 1$, Gaussian kernel. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 2; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{4}{\sqrt{p}})$.

Note also from both Figure 4 and Figure 5 that, when n, p are doubled (from 2048, 512 to 4096, 1024 in Figure 4 and from 256, 512 to 512, 1024 in Figure 5), the real error becomes closer to the theoretical one, which confirms the asymptotic result as $n, p \rightarrow \infty$.

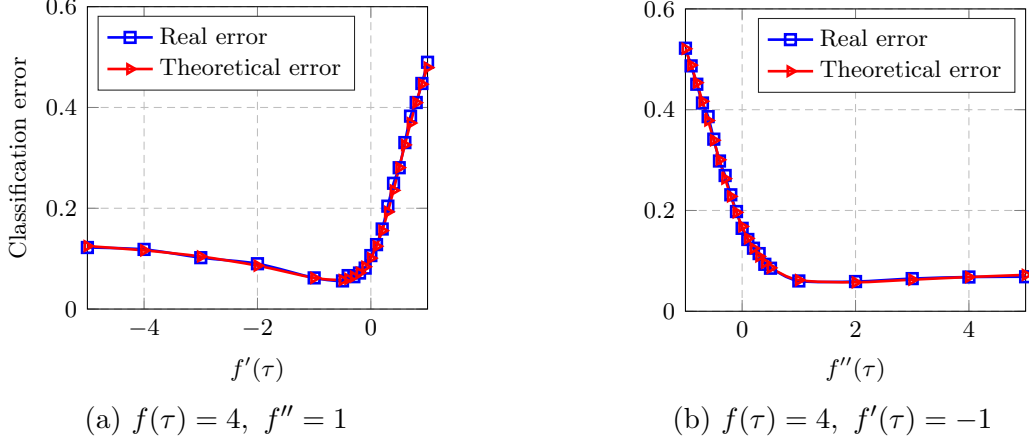


Figure 6: Performance of LS-SVM, $n = 256$, $p = 512$, $c_1 = c_2 = 1/2, \gamma = 1$, polynomial kernel. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 2; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{4}{\sqrt{p}})$.

Clearly, for practical use, one needs to know in advance the value of τ before training so that the kernel f can be properly chosen during the training step. The estimation of τ is possible, in the large n, p regime, with the following lemma:

Lemma 9 *Under Assumptions 1 and 2, as $n \rightarrow \infty$,*

$$\frac{2}{n} \sum_{i=1}^n \frac{\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{p} \xrightarrow{\text{a.s.}} \tau \quad (12)$$

with $\bar{\mathbf{x}} \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

Proof

$$\frac{2}{n} \sum_{i=1}^n \frac{\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{p} = \frac{2c_1c_2\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2}{p} + \frac{2}{n} \sum_{i=1}^n \|\boldsymbol{\omega}_i - \bar{\boldsymbol{\omega}}\|^2 + \kappa$$

with $\kappa = \frac{4}{n\sqrt{p}}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \left(-c_2 \sum_{\mathbf{x}_i \in \mathcal{C}_1} \boldsymbol{\omega}_i + c_1 \sum_{\mathbf{x}_j \in \mathcal{C}_2} \boldsymbol{\omega}_j \right)$ and $\bar{\boldsymbol{\omega}} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\omega}_i$.

According to Assumption 1 we have $\frac{2c_1c_2}{p}\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2 = O(n^{-1})$. The term κ is a linear combination of independent zero-mean Gaussian variables and thus $\kappa \sim \mathcal{N}(0, \text{Var}[\kappa])$ with $\text{Var}[\kappa] = \frac{16c_1c_2}{np^2}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top (c_2\mathbf{C}_1 + c_1\mathbf{C}_2)(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) = O(n^{-3})$. We thus deduce from Chebyshev's inequality and Borel-Cantelli lemma that $\kappa \xrightarrow{\text{a.s.}} 0$.

We then work on the last term $\frac{2}{n} \sum_{i=1}^n \|\boldsymbol{\omega}_i - \bar{\boldsymbol{\omega}}\|^2$ as

$$\frac{2}{n} \sum_{i=1}^n \|\boldsymbol{\omega}_i - \bar{\boldsymbol{\omega}}\|^2 = \frac{2}{n} \sum_{i=1}^n \|\boldsymbol{\omega}_i\|^2 - 2\|\bar{\boldsymbol{\omega}}\|^2.$$

Since $\bar{\boldsymbol{\omega}} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}^\circ/np)$, we deduce that $\|\bar{\boldsymbol{\omega}}\|^2 \xrightarrow{\text{a.s.}} 0$. Ultimately by the strong law of large numbers, the term $\frac{2}{n} \sum_{i=1}^n \|\boldsymbol{\omega}_i\|^2 \xrightarrow{\text{a.s.}} \tau$, which concludes the proof. \blacksquare

4.2 Some Limiting Cases

4.2.1 DOMINANT DEVIATION IN MEANS

When $\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2$ is largely dominant over $(\text{tr}(\mathbf{C}_2 - \mathbf{C}_1))^2/p$ and $\text{tr}((\mathbf{C}_2 - \mathbf{C}_1)^2)/p$, from Theorem 2, both $E_a - (c_2 - c_1)$ and $\sqrt{\text{Var}_a}$ are (approximately) proportional to $f'(\tau)$, which eventually makes the choice of the kernel irrelevant (as long as $f'(\tau) \neq 0$). This result also holds true for E_a^* and $\sqrt{\text{Var}_a^*}$ when normalized labels \mathbf{y}^* are applied, as a result of Lemma 4.

4.2.2 c_0 LARGE OR SMALL

Note that, differently from both \mathcal{V}_1 and \mathcal{V}_2 , \mathcal{V}_3 is a function of c_0 as it can be rewritten as

$$\mathcal{V}_3^a = \frac{2c_0 (f'(\tau))^2}{p^3} \left(\frac{\text{tr} \mathbf{C}_1 \mathbf{C}_a}{c_1} + \frac{\text{tr} \mathbf{C}_2 \mathbf{C}_a}{c_2} \right)$$

which indicates that the variance of $g(\mathbf{x})$ grows as c_0 becomes large. This result is easily understood since, with p fixed, a small c_0 means a larger n , and with more training samples, one may “gain” more information of the two different classes, which reduces the “uncertainty” of the classifier. When $n \rightarrow \infty$ with a fixed p , the LS-SVM is considered “well-trained” and its performance can be described with Theorem 2 by taking $\mathcal{V}_3 = 0$. On the contrary, when $c_0 \rightarrow \infty$, with few training data, LS-SVM does not sample sufficiently the high dimensional space of the \mathbf{x} ’s, thus resulting in a classifier with arbitrarily large variance (for fixed means). Figure 7 confirms this result with p fixed to 256 while n varies from 8 to 8192.

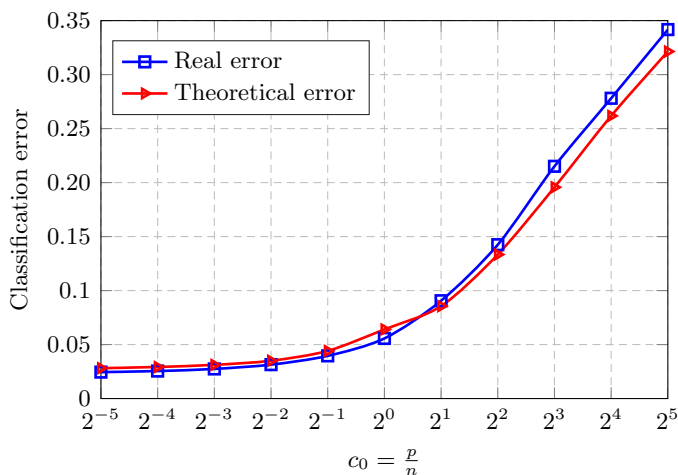


Figure 7: Performance of LS-SVM, $p = 256$, $c_1 = c_2 = 1/2$, $\gamma = 1$, Gaussian kernel with $\sigma^2 = 1$. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 2; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{4}{\sqrt{p}})$.

4.2.3 $c_1 \rightarrow 0$

As revealed in Remark 3, the number of instances in each class n_a in the training set plays a significant role in the performance of classification. A natural question arises: what happens when one class is strongly dominant over the other? Take the case of $c_1 \rightarrow 0, c_2 \rightarrow 1$. From Corollary 5, one has $E_1^* \rightarrow -\gamma \mathcal{D}$ and $E_2^* \rightarrow 0$ while $\mathcal{V}_3^a \rightarrow \infty$ because of $c_1 \rightarrow 0$ in the denominator, which then makes the ratio $\frac{E_a^*}{\sqrt{\text{Var}_a^*}}$ (and thus $\frac{E_a - (c_2 - c_1)}{\sqrt{\text{Var}_a}}$) go to zero, resulting in a poorly-performing LS-SVM. The same occurs when $c_1 \rightarrow 1$ and $c_2 \rightarrow 0$. Figure 8 collaborates this finding with $c_1 = 1/32$ in subfigure (a) and $1/2$ in (b). Note that in subfigure (a), even with a smartly chosen threshold ξ , LS-SVM shall not perform as well as in the case $c_1 = c_2$, as a result of the significant overlap of the two histograms.

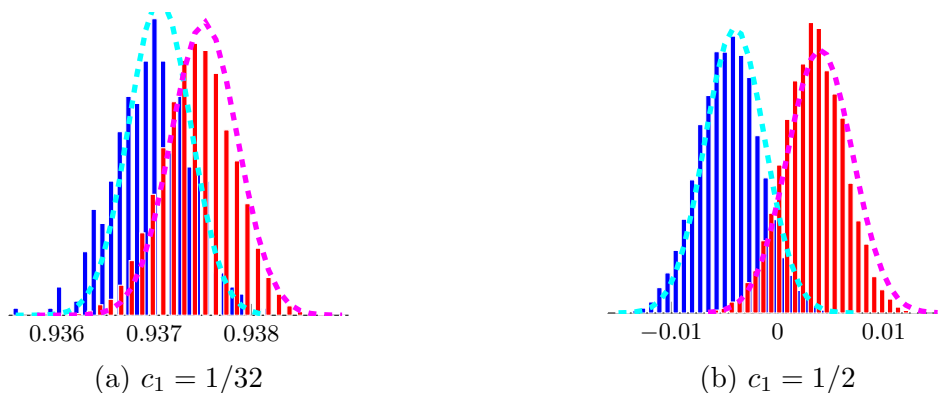


Figure 8: Gaussian approximation of $g(\mathbf{x})$, $n = 256$, $p = 512$, $c_2 = 1 - c_1, \gamma = 1$, Gaussian kernel with $\sigma^2 = 1$. $\mathbf{x} \in \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$, with $\boldsymbol{\mu}_a = [\mathbf{0}_{a-1}; 2; \mathbf{0}_{p-a}]$, $\mathbf{C}_1 = \mathbf{I}_p$ and $\{\mathbf{C}_2\}_{i,j} = .4^{|i-j|}(1 + \frac{4}{\sqrt{p}})$.

4.3 Applying to Real-world Data sets

When the classification performance of real-world data sets is concerned, our theory is limited by: i) the fact that it is an asymptotic result and allows for an estimation error of order $O(1/\sqrt{p})$ between theory and practice and ii) the strong Gaussian assumption for the input data.

However, when applied to real-world data sets, here to the popular MNIST database (LeCun et al., 1998), our asymptotic results, which are theoretically only applicable for Gaussian data, show an unexpectedly similar behavior. Here we consider a two-class classification problem with a training set of $n = 256$ vectorized images of size $p = 784$ randomly selected from the MNIST database (numbers 8 and 9). Then a test set of $n_{\text{test}} = 256$ is used to evaluate the classification performance. Means and covariances are empirically obtained from the full set of 11 800 MNIST images (5 851 images of number 8 and 5 949 of number 9). Despite the obvious non-Gaussianity of the input data, the distribution of $g(\mathbf{x})$ is still surprisingly close to its Gaussian approximation computed from Theorem 2, as shown in

Figure 9 (a) as well as in (b) when Gaussian white noise is artificially added to the image vectors.

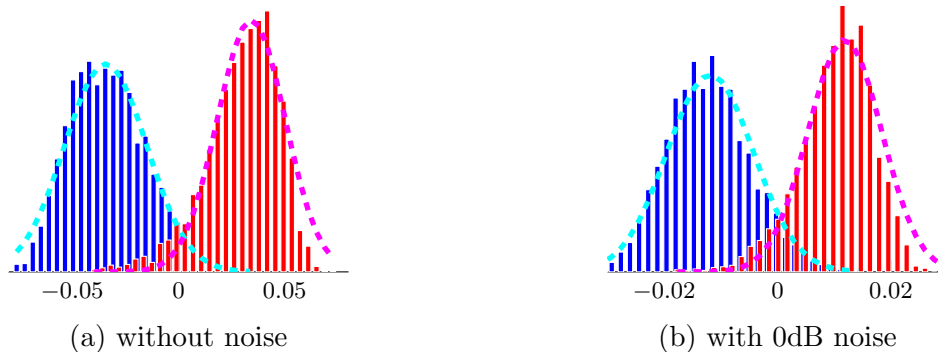


Figure 9: Gaussian approximation of $g(\mathbf{x})$, $n = 256$, $p = 784$, $c_1 = c_2 = \frac{1}{2}$, $\gamma = 1$, Gaussian kernel with $\sigma = 1$, MNIST data (numbers 8 and 9) without and with 0dB noise.

In Figure 10 we evaluated the performance of LS-SVM on the MNIST data set (with and without noise) as a function of the kernel parameter σ of Gaussian kernel $f(x) = \exp(-\frac{x}{2\sigma^2})$. Surprisingly, compared to Figure 5, we face the situation where there is little difference in the performance of LS-SVM as soon as σ^2 is away from 0, which likely comes from the fact that the difference in means $\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1$ is so large that it becomes predominant over the influence of covariances as mentioned in the special case of Section 4.2.1. This argument is numerically sustained by Table 1. The gap between theory and practice observed as $\sigma^2 \rightarrow 0$ is likely a result of the finite n, p (as in Figure 5) rather than of the Gaussian assumption of the input data, since we observe a similar behavior even when Gaussian white noise is added.

Table 1: Empirical estimation of differences in means and covariances of MNIST data (numbers 8 and 9)

	without noise	with 0dB noise
$\ \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\ ^2$	251	96
$(\text{tr}(\mathbf{C}_2 - \mathbf{C}_1))^2 / p$	19	3
$\text{tr}((\mathbf{C}_2 - \mathbf{C}_1)^2) / p$	30	5

5. Concluding remarks

In this work, through a performance analysis of LS-SVM for large dimensional data, we reveal the significant choice of training data with $n_1 = n_2$, as well as the interplay between the pivotal kernel function f and the statistical structure of the data. The normalized labels $y_i^* \in \{-1/c_1, 1/c_2\}$ are proposed to mitigate the damage of $c_2 - c_1$ in the decision

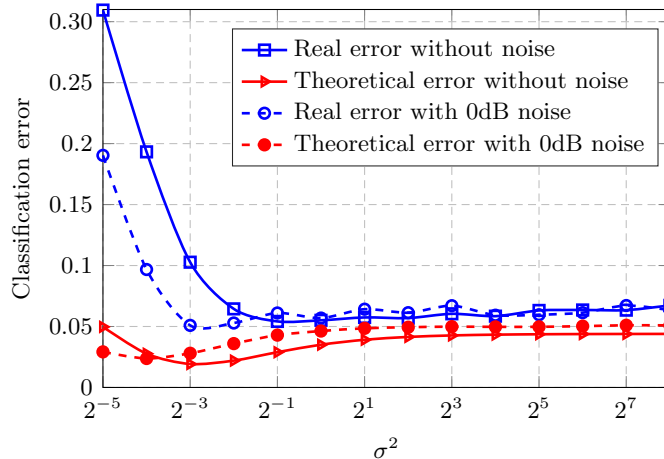


Figure 10: $n = 256$, $p = 784$, $c_1 = c_2 = 1/2$, $\gamma = 1$, Gaussian kernel, MNIST data (numbers 8 and 9) with and without noise.

function, although the best way is always to use balanced training sets. We have proved the irrelevance of γ when it is considered to remain constant in the large n, p regime; however, this argument is not guaranteed to hold true when γ scales with n, p . Also, even though our theoretical results are built upon the assumption of Gaussian data, similar behavior is observed with real-world large dimensional data sets, which offers a possible application despite the strong Gaussian assumption in the general context of large scale supervised learning.

The extension of the present work to the asymptotic performance analysis of the classical SVM requires more efforts since, there, the decision function $g(\mathbf{x})$ depends implicitly (through the solution to a quadratic programming problem) rather than explicitly on the underlying kernel matrix \mathbf{K} . Additional technical tools are thus required to cope with this dependence structure.

The link between LS-SVM and extreme learning machine (ELM) was brought to light by [Huang et al. \(2012\)](#) and the performance analysis of ELM in large dimension has been investigated by [Cosme and Couillet \(2017\)](#). Together with these works, we have the possibility to identify the tight but subtle relation between the kernel function and the activation function in the context of some simple structured neural networks. This is notably of interest when the data sets are so large that computing \mathbf{K} and the decision function g becomes prohibitive, a problem largely alleviated by neural networks with controllable number of neurons. This link also generally opens up a possible direction of research into the complex neural networks realm.

Acknowledgments

We would like to acknowledge this work is supported by the ANR Project RMT4GRAPH (ANR-14-CE28-0006). This paper was presented in part at the 42nd IEEE International

Conference on Acoustics, Speech and Signal Processing (ICASSP'17), New Orleans, USA, March 2017.

Appendix A. Proof of Theorem 1

Our key interest here is on the decision function of LS-SVM: $g(\mathbf{x}) = \boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x}) + b$ with $(\boldsymbol{\alpha}, b)$ given by

$$\begin{cases} \boldsymbol{\alpha} &= \mathbf{S}^{-1} \left(\mathbf{I}_n - \frac{\mathbf{1}_n^\top \mathbf{1}_n \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \right) \mathbf{y} \\ b &= \frac{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{y}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \end{cases}$$

and $\mathbf{S}^{-1} = \left(\mathbf{K} + \frac{n}{\gamma} \mathbf{I}_n \right)^{-1}$.

Before going into the detailed proof, as we will frequently deal with random variables evolving as n, p grow large, we shall use the extension of the $O(\cdot)$ notation introduced by Couillet et al. (2016): for a random variable $x \equiv x_n$ and $u_n \geq 0$, we write $x = O(u_n)$ if for any $\eta > 0$ and $D > 0$, we have $n^D \mathbb{P}(x \geq n^\eta u_n) \rightarrow 0$.

When multidimensional objects are concerned, $\mathbf{v} = O(u_n)$ means the maximum entry of a vector (or a diagonal matrix) \mathbf{v} in absolute value is $O(u_n)$ and $\mathbf{M} = O(u_n)$ means that the operator norm of \mathbf{M} is $O(u_n)$. We refer the reader to the work of Couillet et al. (2016) for more discussions on these practical definitions.

Under the growth rate settings of Assumption 1, from Couillet et al. (2016), the approximation of the kernel matrix \mathbf{K} is given by

$$\mathbf{K} = -2f'(\tau) \left(\mathbf{P} \boldsymbol{\Omega}^\top \boldsymbol{\Omega} \mathbf{P} + \mathbf{A} \right) + \beta \mathbf{I}_n + O(n^{-\frac{1}{2}}) \quad (13)$$

with $\beta = f(0) - f(\tau) + \tau f'(\tau)$ and $\mathbf{A} = \mathbf{A}_n + \mathbf{A}_{\sqrt{n}} + \mathbf{A}_1$, $\mathbf{A}_n = -\frac{f(\tau)}{2f'(\tau)} \mathbf{1}_n \mathbf{1}_n^\top$ and $\mathbf{A}_{\sqrt{n}}, \mathbf{A}_1$ given by (14) and (15) at the top of next page, where we denote

$$\begin{aligned} t_a &\triangleq \frac{\text{tr}(\mathbf{C}_a - \mathbf{C}^\circ)}{\sqrt{p}} = O(1) \\ (\boldsymbol{\psi})^2 &\triangleq [(\boldsymbol{\psi}_1)^2, \dots, (\boldsymbol{\psi}_n)^2]^\top. \end{aligned}$$

$$\mathbf{A}_{\sqrt{n}} = -\frac{1}{2} \left[\boldsymbol{\psi} \mathbf{1}_n^\top + \mathbf{1}_n \boldsymbol{\psi}^\top + \left\{ t_a \frac{\mathbf{1}_{n_a}}{\sqrt{p}} \right\}_{a=1}^2 \mathbf{1}_n^\top + \mathbf{1}_n \left\{ t_b \frac{\mathbf{1}_{n_b}}{\sqrt{p}} \right\}_{b=1}^2 \right] \quad (14)$$

$$\begin{aligned}
 \mathbf{A}_1 = & -\frac{1}{2} \left[\left\{ \left\| \boldsymbol{\mu}_a - \boldsymbol{\mu}_b \right\|^2 \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p} \right\}_{a,b=1}^2 + 2 \left\{ \frac{(\boldsymbol{\Omega} \mathbf{P})_a^\top (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a) \mathbf{1}_{n_b}^\top}{\sqrt{p}} \right\}_{a,b=1}^2 \right. \\
 & \left. - 2 \left\{ \frac{\mathbf{1}_{n_a} (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)^\top (\boldsymbol{\Omega} \mathbf{P})_b}{\sqrt{p}} \right\}_{a,b=1}^2 \right] - \frac{f''(\tau)}{4f'(\tau)} \left[(\boldsymbol{\psi})^2 \mathbf{1}_n^\top + \mathbf{1}_n [(\boldsymbol{\psi})^2]^\top + \left\{ t_a^2 \frac{\mathbf{1}_{n_a}}{p} \right\}_{a=1}^2 \mathbf{1}_n^\top \right. \\
 & \left. + \mathbf{1}_n \left\{ t_b^2 \frac{\mathbf{1}_{n_b}^\top}{p} \right\}_{b=1}^2 + 2 \left\{ t_a t_b \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p} \right\}_{a,b=1}^2 + 2\mathcal{D} \{ t_a \mathbf{I}_{n_a} \}_{a=1}^2 \boldsymbol{\psi} \frac{\mathbf{1}_n^\top}{\sqrt{p}} + 2\boldsymbol{\psi} \left\{ t_b \frac{\mathbf{1}_{n_b}^\top}{\sqrt{p}} \right\}_{b=1}^2 \right. \\
 & \left. + 2 \frac{\mathbf{1}_n}{\sqrt{p}} (\boldsymbol{\psi})^\top \mathcal{D} \{ t_a \mathbf{1}_{n_a} \}_{a=1}^2 + 2 \left\{ t_a \frac{\mathbf{1}_{n_a}}{\sqrt{p}} \right\}_{a=1}^2 (\boldsymbol{\psi})^\top + 4 \left\{ \text{tr}(\mathbf{C}_a \mathbf{C}_b) \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p^2} \right\}_{a,b=1}^2 + 2\boldsymbol{\psi} (\boldsymbol{\psi})^\top \right] \quad (15)
 \end{aligned}$$

We start with studying the inverse \mathbf{S}^{-1} . The terms of leading order of \mathbf{K} , i.e., $-2f'(\tau)\mathbf{A}_n = f(\tau)\mathbf{1}_n \mathbf{1}_n^\top$ and $\frac{n}{\gamma}\mathbf{I}_n$ are both of operator norm $O(n)$. Therefore a Taylor expansion can be performed as

$$\begin{aligned}
 \mathbf{S}^{-1} &= \left(\mathbf{K} + \frac{n}{\gamma} \mathbf{I}_n \right)^{-1} = \frac{1}{n} \left[\mathbf{L}^{-1} - \frac{2f'(\tau)}{n} \left(\mathbf{A}_{\sqrt{n}} + \mathbf{A}_1 + \mathbf{P} \boldsymbol{\Omega}^\top \boldsymbol{\Omega} \mathbf{P} \right) + \frac{\beta \mathbf{I}_n}{n} + O(n^{-\frac{3}{2}}) \right]^{-1} \\
 &= \frac{\mathbf{L}}{n} + \frac{2f'(\tau)}{n^2} \mathbf{L} \mathbf{A}_{\sqrt{n}} \mathbf{L} + \mathbf{L} \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{L} + O(n^{-\frac{5}{2}})
 \end{aligned}$$

with $\mathbf{L} = \left(f(\tau) \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} + \frac{\mathbf{I}_n}{\gamma} \right)^{-1}$ of order $O(1)$ and $\mathbf{Q} = \frac{2f'(\tau)}{n^2} \left(\mathbf{A}_1 + \mathbf{P} \boldsymbol{\Omega}^\top \boldsymbol{\Omega} \mathbf{P} + \frac{2f'(\tau)}{n} \mathbf{A}_{\sqrt{n}} \mathbf{L} \mathbf{A}_{\sqrt{n}} \right)$.

With the Sherman-Morrison formula we are able to compute explicitly \mathbf{L} as

$$\mathbf{L} = \left(f(\tau) \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} + \frac{\mathbf{I}_n}{\gamma} \right)^{-1} = \gamma \left(\mathbf{I}_n - \frac{\gamma f(\tau)}{1 + \gamma f(\tau)} \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} \right) = \frac{\gamma}{1 + \gamma f(\tau)} \mathbf{I}_n + \frac{\gamma^2 f(\tau)}{1 + \gamma f(\tau)} \mathbf{P} = O(1).$$

Writing \mathbf{L} as a linear combination of \mathbf{I}_n and \mathbf{P} is useful when computing $\mathbf{L} \mathbf{1}_n$ or $\mathbf{1}_n^\top \mathbf{L}$, because by the definition of $\mathbf{P} = \mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n}$, we have $\mathbf{1}_n^\top \mathbf{P} = \mathbf{P} \mathbf{1}_n = \mathbf{0}$.

We shall start with the term $\mathbf{1}_n^\top \mathbf{S}^{-1}$, since it is the basis of several other terms appearing in $\boldsymbol{\alpha}$ and b :

$$\mathbf{1}_n^\top \mathbf{S}^{-1} = \frac{\gamma \mathbf{1}_n^\top}{1 + \gamma f(\tau)} \left[\frac{\mathbf{I}_n}{n} + \frac{2f'(\tau)}{n^2} \mathbf{A}_{\sqrt{n}} \mathbf{L} + \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{L} \right] + O(n^{-\frac{3}{2}})$$

since $\mathbf{1}_n^\top \mathbf{L} = \frac{\gamma}{1 + \gamma f(\tau)} \mathbf{1}_n^\top$.

With $\mathbf{1}_n^\top \mathbf{S}^{-1}$ at hand, we next obtain

$$\begin{aligned}
 \mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1} &= \frac{\gamma}{1 + \gamma f(\tau)} \left[\underbrace{\frac{\mathbf{1}_n \mathbf{1}_n^\top}{n}}_{O(1)} + \underbrace{\frac{2f'(\tau)}{n^2} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{A}_{\sqrt{n}} \mathbf{L}}_{O(n^{-1/2})} + \underbrace{\mathbf{1}_n \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{L}}_{O(n^{-1})} \right] + O(n^{-\frac{3}{2}}) \quad (16) \\
 \mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{y} &= \frac{\gamma}{1 + \gamma f(\tau)} \left[\underbrace{c_2 - c_1}_{O(1)} + \underbrace{\frac{2f'(\tau)}{n^2} \mathbf{1}_n^\top \mathbf{A}_{\sqrt{n}} \mathbf{L} \mathbf{y}}_{O(n^{-1/2})} + \underbrace{\mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{L} \mathbf{y}}_{O(n^{-1})} \right] + O(n^{-\frac{3}{2}})
 \end{aligned}$$

$$\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n = \frac{\gamma}{1 + \gamma f(\tau)} \left[\underbrace{1}_{O(1)} + \underbrace{\frac{2f'(\tau) \gamma \mathbf{1}_n^\top \mathbf{A} \sqrt{n} \mathbf{1}_n}{n^2 (1 + \gamma f(\tau))}}_{O(n^{-1/2})} + \underbrace{\frac{\gamma}{1 + \gamma f(\tau)} \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{1}_n}_{O(n^{-1})} \right] + O(n^{-\frac{3}{2}}).$$

The inverse of $\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n$ can consequently be computed using a Taylor expansion around its leading order, allowing an error term of $O(n^{-\frac{3}{2}})$ as

$$\frac{1}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} = \frac{1 + \gamma f(\tau)}{\gamma} \left[\underbrace{1}_{O(1)} - \underbrace{\frac{2f'(\tau) \gamma \mathbf{1}_n^\top \mathbf{A} \sqrt{n} \mathbf{1}_n}{n^2 (1 + \gamma f(\tau))}}_{O(n^{-1/2})} - \underbrace{\frac{\gamma}{1 + \gamma f(\tau)} \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{1}_n}_{O(n^{-1})} \right] + O(n^{-\frac{3}{2}}). \quad (17)$$

Combing (16) with (17) we deduce

$$\frac{\mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} = \underbrace{\frac{\mathbf{1}_n \mathbf{1}_n^\top}{n}}_{O(1)} + \underbrace{\frac{2f'(\tau)}{n^2} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{A} \sqrt{n} \left[\mathbf{L} - \frac{\gamma \mathbf{1}_n \mathbf{1}_n^\top}{1 + \gamma f(\tau)} \right]}_{O(n^{-1/2})} + \underbrace{\mathbf{1}_n \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \left[\mathbf{L} - \frac{\gamma \mathbf{1}_n \mathbf{1}_n^\top}{1 + \gamma f(\tau)} \right]}_{O(n^{-1})} + O(n^{-\frac{3}{2}}) \quad (18)$$

and similarly the expression of b with (14)–(15) as

$$\begin{aligned} b &= \underbrace{c_2 - c_1}_{O(1)} - \underbrace{\frac{2\gamma}{\sqrt{p}} c_1 c_2 f'(\tau) (t_2 - t_1)}_{O(n^{-1/2})} - \underbrace{\frac{\gamma f''(\tau)}{2n} \mathbf{y}^\top \mathbf{P} (\boldsymbol{\psi})^2 + \frac{4\gamma c_1 c_2}{p} [c_1 T_1 + (c_2 - c_1) D - c_2 T_2]}_{O(n^{-1})} \\ &\quad - \underbrace{\frac{\gamma f'(\tau)}{n} \mathbf{y}^\top \mathbf{P} \boldsymbol{\psi}}_{O(n^{-1})} + O(n^{-\frac{3}{2}}) \end{aligned} \quad (19)$$

with

$$\begin{aligned} D &= \frac{f'(\tau)}{2} \|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2 + \frac{f''(\tau)}{4} (t_1 + t_2)^2 + f''(\tau) \frac{\text{tr} \mathbf{C}_1 \mathbf{C}_2}{p} \\ T_a &= f''(\tau) t_a^2 + f''(\tau) \frac{\text{tr} \mathbf{C}_1 \mathbf{C}_2}{p}. \end{aligned}$$

We thus have the random equivalent of b as $n \rightarrow \infty$. Moving to $\boldsymbol{\alpha}$, note from (16) that $\mathbf{L} - \frac{\gamma}{1 + \gamma f(\tau)} \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} = \gamma \mathbf{P}$, and we can thus rewrite

$$\frac{\mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} = \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} + \frac{2\gamma f'(\tau)}{n^2} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{A} \sqrt{n} \mathbf{P} + \gamma \mathbf{1}_n \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{P} + O(n^{-\frac{3}{2}}).$$

At this point, for $\boldsymbol{\alpha} = \mathbf{S}^{-1} \left(\mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^\top \mathbf{S}^{-1}}{\mathbf{1}_n^\top \mathbf{S}^{-1} \mathbf{1}_n} \right) \mathbf{y}$, we thus have

$$\boldsymbol{\alpha} = \mathbf{S}^{-1} \left[\mathbf{I}_n - \frac{2\gamma f'(\tau)}{n^2} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{A} \sqrt{n} - \gamma \mathbf{1}_n \mathbf{1}_n^\top \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \right] \mathbf{P} \mathbf{y} + O(n^{-\frac{5}{2}}).$$

Here again, we use $\mathbf{1}_n^\top \mathbf{L} = \frac{\gamma}{1+\gamma f(\tau)} \mathbf{1}_n^\top$ and $\mathbf{L} - \frac{\gamma}{1+\gamma f(\tau)} \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} = \gamma \mathbf{P}$, to eventually get

$$\boldsymbol{\alpha} = \underbrace{\frac{\gamma}{n} \mathbf{P} \mathbf{y}}_{O(n^{-1})} + \underbrace{\gamma^2 \mathbf{P} \left(\mathbf{Q} - \frac{\beta}{n^2} \mathbf{I}_n \right) \mathbf{P} \mathbf{y}}_{O(n^{-2})} - \underbrace{\frac{\gamma^2}{1+\gamma f(\tau)} \left(\frac{2f'(\tau)}{n^2} \right)^2 \mathbf{L} \mathbf{A} \sqrt{n} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{A} \sqrt{n} \mathbf{P} \mathbf{y}}_{O(n^{-2})} + O(n^{-\frac{5}{2}}). \quad (23)$$

Note here the absence of a term of order $O(n^{-3/2})$ in the expression of $\boldsymbol{\alpha}$ since $\mathbf{P} \mathbf{A} \sqrt{n} \mathbf{P} = 0$ from (14).

We shall now work on the ‘‘information vector’’ $\mathbf{k}(\mathbf{x})$, following the same analysis as in the work of Couillet et al. (2016) for the kernel matrix \mathbf{K} , assuming that $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_a, \mathbf{C}_a)$. Recalling the random variables definitions

$$\boldsymbol{\omega}_{\mathbf{x}} \triangleq \frac{\mathbf{x} - \boldsymbol{\mu}_a}{\sqrt{p}}$$

$$\psi_{\mathbf{x}} \triangleq \|\boldsymbol{\omega}_{\mathbf{x}}\|^2 - \mathbb{E} [\|\boldsymbol{\omega}_{\mathbf{x}}\|^2]$$

we can show that every entry of $\mathbf{k}(\mathbf{x})$ can be written as

$$\begin{aligned} \{\mathbf{k}(\mathbf{x})\}_j &= \underbrace{f(\tau)}_{O(1)} + f'(\tau) \left[\underbrace{\frac{t_a + t_b}{\sqrt{p}} + \psi_x + \psi_j - 2(\boldsymbol{\omega}_{\mathbf{x}})^\top \boldsymbol{\omega}_j}_{O(n^{-1/2})} + \underbrace{\frac{\|\boldsymbol{\mu}_b - \boldsymbol{\mu}_a\|^2}{p} + \frac{2}{\sqrt{p}} (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)^\top (\boldsymbol{\omega}_j - \boldsymbol{\omega}_{\mathbf{x}})}_{O(n^{-1})} \right] \\ &+ \frac{f''(\tau)}{2} \left[\underbrace{\left(\frac{t_a + t_b}{\sqrt{p}} + \psi_j + \psi_{\mathbf{x}} \right)^2 + \frac{4}{p^2} \text{tr} \mathbf{C}_a \mathbf{C}_b}_{O(n^{-1})} \right] + O(n^{-\frac{3}{2}}). \end{aligned} \quad (24)$$

Combining (23) and (24), we deduce

$$\boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x}) = \underbrace{\frac{2\gamma}{\sqrt{p}} c_1 c_2 f'(\tau) (t_2 - t_1)}_{O(n^{-1/2})} + \underbrace{\frac{\gamma}{n} \mathbf{y}^\top \mathbf{P} \tilde{\mathbf{k}}(\mathbf{x})}_{O(n^{-1})} + \underbrace{\frac{\gamma f'(\tau)}{n} \mathbf{y}^\top \mathbf{P} (\boldsymbol{\psi} - 2\mathbf{P} \boldsymbol{\Omega}^\top \boldsymbol{\omega}_{\mathbf{x}})}_{O(n^{-1})} + O(n^{-\frac{3}{2}}) \quad (25)$$

with $\tilde{\mathbf{k}}(\mathbf{x})$ given by

$$\begin{aligned} \tilde{\mathbf{k}}(\mathbf{x}) &= f'(\tau) \left[\left\{ \frac{\|\boldsymbol{\mu}_b - \boldsymbol{\mu}_a\|^2}{p} \mathbf{1}_{n_b} \right\}_{b=1}^2 - \frac{2}{\sqrt{p}} \left\{ \mathbf{1}_{n_b} (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)^\top \right\}_{b=1}^2 \boldsymbol{\omega}_{\mathbf{x}} + \frac{2}{\sqrt{p}} \mathcal{D} \left(\left\{ \mathbf{1}_{n_b} (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)^\top \right\}_{b=1}^2 \boldsymbol{\Omega} \right) \right] \\ &+ \frac{f''(\tau)}{2} \left[\left\{ \frac{(t_a + t_b)^2}{p} \mathbf{1}_{n_b} \right\}_{b=1}^2 + 2\mathcal{D} \left(\left\{ \frac{t_a + t_b}{\sqrt{p}} \mathbf{1}_{n_b} \right\}_{b=1}^2 \right) \boldsymbol{\psi} + 2 \left\{ \frac{t_a + t_b}{\sqrt{p}} \mathbf{1}_{n_b} \right\}_{b=1}^2 \psi_{\mathbf{x}} \right. \\ &\left. + (\boldsymbol{\psi})^2 + 2\psi_{\mathbf{x}} \boldsymbol{\psi} + \psi_{\mathbf{x}}^2 \mathbf{1}_n + \left\{ \frac{4}{p^2} \text{tr}(\mathbf{C}_a \mathbf{C}_b) \mathbf{1}_{n_b} \right\}_{b=1}^2 \right] \end{aligned} \quad (26)$$

At this point, note that the term of order $O(n^{-\frac{1}{2}})$ in the final object $g(\mathbf{x}) = \boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x}) + b$ disappears because in both (19) and (25) the term of order $O(n^{-1/2})$ is $\frac{2\gamma}{\sqrt{p}} c_1 c_2 f'(\tau) (t_2 - t_1)$ but of opposite signs. Also, we can say that the leading term $c_2 - c_1$ in b will remain in $g(\mathbf{x})$ as stated in Remark 3.

The development of $\mathbf{y}^\top \mathbf{P} \tilde{\mathbf{k}}(\mathbf{x})$ induces many simplifications, since i) $\mathbf{P} \mathbf{1}_n = \mathbf{0}$ and ii) random variables as $\boldsymbol{\omega}_\mathbf{x}$ and $\boldsymbol{\psi}$ in $\tilde{\mathbf{k}}(\mathbf{x})$, once multiplied by $\mathbf{y}^\top \mathbf{P}$, thanks to probabilistic averaging of independent zero-mean terms, are of smaller order and thus become negligible. We thus get

$$\begin{aligned} \frac{\gamma}{n} \mathbf{y}^\top \mathbf{P} \tilde{\mathbf{k}}(\mathbf{x}) &= 2\gamma c_1 c_2 f'(\tau) \left[\frac{\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_a\|^2 - \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_a\|^2}{p} - 2(\boldsymbol{\omega}_\mathbf{x})^\top \frac{\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1}{\sqrt{p}} \right] \\ &+ \frac{\gamma f''(\tau)}{2n} \mathbf{y}^\top \mathbf{P} (\boldsymbol{\psi})^2 + \gamma c_1 c_2 f''(\tau) \left[2 \left(\frac{t_a}{\sqrt{p}} + \boldsymbol{\psi}_\mathbf{x} \right) \frac{t_2 - t_1}{\sqrt{p}} + \frac{t_2^2 - t_1^2}{p} \right. \\ &\left. + \frac{4}{p^2} \text{tr}(\mathbf{C}_a \mathbf{C}_2 - \mathbf{C}_a \mathbf{C}_1) \right] + O(n^{-\frac{3}{2}}). \end{aligned} \quad (27)$$

This result, together with (25), completes the analysis of the term $\boldsymbol{\alpha}^\top \mathbf{k}(\mathbf{x})$. Combining (25)-(27) and (19) concludes the proof.

Appendix B. Proof of Theorem 2

This section is dedicated to the proof of the central limit theorem for

$$\hat{g}(\mathbf{x}) = c_2 - c_1 + \gamma (\mathfrak{P} + c_\mathbf{x} \mathfrak{D})$$

with the notation $c_\mathbf{x} = -2c_1 c_2^2$ for $\mathbf{x} \in \mathcal{C}_1$ and $c_\mathbf{x} = 2c_1^2 c_2$ for $\mathbf{x} \in \mathcal{C}_2$, and $\mathfrak{P}, \mathfrak{D}$ as defined in (8) and (9).

Our objective is to show that for $a \in \{1, 2\}$, $n(\hat{g}(\mathbf{x}) - G_a) \xrightarrow{d} 0$ for

$$G_a \sim \mathcal{N}(E_a, \text{Var}_a)$$

with E_a and Var_a given in Theorem 2. We recall here that $\mathbf{x} = \boldsymbol{\mu}_a + \sqrt{p} \boldsymbol{\omega}_\mathbf{x}$ with $\boldsymbol{\omega}_\mathbf{x} \sim \mathcal{N}(0, \frac{\mathbf{C}_a}{p})$.

Letting $\mathbf{z}_\mathbf{x}$ such that $\boldsymbol{\omega}_\mathbf{x} = \frac{(\mathbf{C}_a)^{1/2}}{\sqrt{p}} \mathbf{z}_\mathbf{x}$, we have $\mathbf{z}_\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_n)$ and we can rewrite $\hat{g}(\mathbf{x})$ in the following quadratic form (of $\mathbf{z}_\mathbf{x}$) as

$$\hat{g}(\mathbf{x}) = \mathbf{z}_\mathbf{x}^\top \mathbf{A} \mathbf{z}_\mathbf{x} + \mathbf{z}_\mathbf{x}^\top \mathbf{b} + c$$

with

$$\begin{aligned} \mathbf{A} &= 2\gamma c_1 c_2 f''(\tau) \frac{\text{tr}(\mathbf{C}_2 - \mathbf{C}_1)}{p} \frac{\mathbf{C}_a}{p} \\ \mathbf{b} &= -\frac{2\gamma f'(\tau)}{n} \frac{(\mathbf{C}_a)^{\frac{1}{2}}}{\sqrt{p}} \boldsymbol{\Omega} \mathbf{P} \mathbf{y} - \frac{4c_1 c_2 \gamma f'(\tau)}{\sqrt{p}} \frac{(\mathbf{C}_a)^{\frac{1}{2}}}{\sqrt{p}} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \\ c &= c_2 - c_1 + \gamma c_\mathbf{x} \mathfrak{D} - 2\gamma c_1 c_2 f''(\tau) \frac{\text{tr}(\mathbf{C}_2 - \mathbf{C}_1)}{p} \frac{\text{tr} \mathbf{C}_a}{p}. \end{aligned}$$

Since $\mathbf{z}_\mathbf{x}$ is Gaussian and has the same distribution as $\mathbf{U} \mathbf{z}_\mathbf{x}$ for any orthogonal matrix \mathbf{U} (i.e., such that $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}_n$), we choose \mathbf{U} that satisfies $\mathbf{A} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$, with $\boldsymbol{\Lambda}$ diagonal so that $\hat{g}(\mathbf{x})$ and $\tilde{g}(\mathbf{x})$ have the same distribution where

$$\tilde{g}(\mathbf{x}) = \mathbf{z}_\mathbf{x}^\top \boldsymbol{\Lambda} \mathbf{z}_\mathbf{x} + \mathbf{z}_\mathbf{x}^\top \tilde{\mathbf{b}} + c = \sum_{i=1}^n \left(z_i^2 \lambda_i + z_i \tilde{b}_i + \frac{c}{n} \right)$$

and $\tilde{\mathbf{b}} = \mathbf{U}^\top \mathbf{b}$, λ_i the diagonal elements of $\mathbf{\Lambda}$ and z_i the elements of $\mathbf{z}_\mathbf{x}$.

Conditioning on $\mathbf{\Omega}$, we thus result in the sum of independent but not identically distributed random variables $r_i = z_i^2 \lambda_i + z_i \tilde{b}_i + \frac{c}{n}$. We then resort to the Lyapunov CLT (Billingsley, 2008, Theorem 27.3).

We begin by estimating the expectation and the variance

$$\begin{aligned}\mathbb{E}[r_i|\mathbf{\Omega}] &= \lambda_i + \frac{c}{n} \\ \text{Var}[r_i|\mathbf{\Omega}] &= \sigma_i^2 = 2\lambda_i^2 + \tilde{b}_i^2\end{aligned}$$

of r_i , so that

$$\begin{aligned}\sum_{i=1}^n \mathbb{E}[r_i|\mathbf{\Omega}] &= c_2 - c_1 + \gamma c_\mathbf{x} \mathfrak{D} = E_a \\ s_n^2 &= \sum_{i=1}^n \sigma_i^2 = 2 \text{tr}(\mathbf{A}^2) + \mathbf{b}^\top \mathbf{b} \\ &= 8\gamma^2 c_1^2 c_2^2 (f''(\tau))^2 \frac{(\text{tr}(\mathbf{C}_2 - \mathbf{C}_1))^2 \text{tr} \mathbf{C}_a^2}{p^2} + 4\gamma^2 \left(\frac{f'(\tau)}{n}\right)^2 \mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} \mathbf{P} \mathbf{y} \\ &\quad + \frac{16\gamma^2 c_1^2 c_2^2 (f'(\tau))^2}{p} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \frac{\mathbf{C}_a}{p} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + O(n^{-\frac{5}{2}}).\end{aligned}$$

We shall rewrite $\mathbf{\Omega}$ into two blocks as

$$\mathbf{\Omega} = \left[\frac{(\mathbf{C}_1)^{\frac{1}{2}}}{\sqrt{p}} \mathbf{Z}_1, \quad \frac{(\mathbf{C}_2)^{\frac{1}{2}}}{\sqrt{p}} \mathbf{Z}_2 \right]$$

where $\mathbf{Z}_1 \in \mathbb{R}^{p \times n_1}$ and $\mathbf{Z}_2 \in \mathbb{R}^{p \times n_2}$ with i.i.d. Gaussian entries with zero mean and unit variance. Then

$$\mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} = \frac{1}{p^2} \begin{bmatrix} \mathbf{Z}_1^\top (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{Z}_1 & \mathbf{Z}_1^\top (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{Z}_2 \\ \mathbf{Z}_2^\top (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{Z}_1 & \mathbf{Z}_2^\top (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{Z}_2 \end{bmatrix}$$

and with $\mathbf{P} \mathbf{y} = \mathbf{y} - (c_2 - c_1) \mathbf{1}_n$, we deduce

$$\begin{aligned}\mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} \mathbf{P} \mathbf{y} &= \frac{4}{p^2} \left(c_2^2 \mathbf{1}_{n_1}^\top \mathbf{Z}_1^\top (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_1)^{\frac{1}{2}} b \mathbf{Z}_1 \mathbf{1}_{n_1} - 2c_1 c_2 \right. \\ &\quad \left. \mathbf{1}_{n_1}^\top \mathbf{Z}_1^\top (\mathbf{C}_1)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{Z}_2 \mathbf{1}_{n_2} + c_2^2 \mathbf{1}_{n_1}^\top \mathbf{Z}_2^\top (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{C}_a (\mathbf{C}_2)^{\frac{1}{2}} \mathbf{Z}_2 \mathbf{1}_{n_2} \right).\end{aligned}$$

Since $\mathbf{Z}_i \mathbf{1}_{n_i} \sim \mathcal{N}(\mathbf{0}, n_i \mathbf{I}_{n_i})$, by applying the trace lemma (Bai and Silverstein, 2010, Lemma B.26) we get

$$\mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} \mathbf{P} \mathbf{y} - \frac{4nc_1^2 c_2^2}{p^2} \left(\frac{\text{tr} \mathbf{C}_1 \mathbf{C}_a}{c_1} + \frac{\text{tr} \mathbf{C}_2 \mathbf{C}_a}{c_2} \right) \xrightarrow{\text{a.s.}} 0. \quad (28)$$

Consider now the events

$$\begin{aligned}E &= \left\{ \left| \mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} \mathbf{P} \mathbf{y} - \rho \right| < \epsilon \right\} \\ \bar{E} &= \left\{ \left| \mathbf{y}^\top \mathbf{P} \mathbf{\Omega}^\top \frac{\mathbf{C}_a}{p} \mathbf{\Omega} \mathbf{P} \mathbf{y} - \rho \right| > \epsilon \right\}\end{aligned}$$

for any fixed ϵ with $\rho = \frac{4nc_1^2c_2^2}{p^2} \left(\frac{\text{tr } \mathbf{C}_1 \mathbf{C}_a}{c_1} + \frac{\text{tr } \mathbf{C}_2 \mathbf{C}_a}{c_2} \right)$ and write

$$\begin{aligned} \mathbb{E} \left[\exp \left(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n} \right) \right] &= \mathbb{E} \left[\exp \left(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n} \right) \middle| E \right] \\ &P(E) + \mathbb{E} \left[\exp \left(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n} \right) \middle| \bar{E} \right] P(\bar{E}) \end{aligned} \quad (29)$$

We start with the variable $\tilde{g}(\mathbf{x})|E$ and check that Lyapunov's condition for $\bar{r}_i = r_i - \mathbb{E}[r_i]$, conditioning on E ,

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n \mathbb{E}[|\bar{r}_i|^4] = 0$$

holds by rewriting

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n \mathbb{E}[|\bar{r}_i|^4] = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{60\lambda_i^4 + 12\lambda_i^2 \tilde{b}_i^2 + 3\tilde{b}_i^4}{s_n^4} = 0$$

since both λ_i and \tilde{b}_i are of order $O(n^{-3/2})$.

As a consequence of the above, we have the CLT for the random variable $\tilde{g}(\mathbf{x})|E$, thus

$$\mathbb{E} \left[\exp \left(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n} \right) \middle| E \right] \rightarrow \exp\left(-\frac{u^2}{2}\right).$$

Next, we see that the second term goes to zero because $|\mathbb{E}[\exp(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n}) | \bar{E}]| \leq 1$ and $P(\bar{E}) \rightarrow 0$ from (28) and we deduce

$$\mathbb{E} \left[\exp \left(iun \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n} \right) \right] \rightarrow \exp\left(-\frac{u^2}{2}\right).$$

With the help of Lévy's continuity theorem, we thus prove the CLT of the variable $n \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{s_n}$. Since $s_n^2 \rightarrow \text{Var}_a$, with Slutsky's theorem, we have the CLT for $n \frac{\tilde{g}(\mathbf{x}) - \mathbf{E}_a}{\sqrt{\text{Var}_a}}$ (thus for $n \frac{\hat{g}(\mathbf{x}) - \mathbf{E}_a}{\sqrt{\text{Var}_a}}$), and eventually for $n \frac{g(\mathbf{x}) - \mathbf{E}_a}{\sqrt{\text{Var}_a}}$ by Theorem 1 which completes the proof.

References

Nedjem-Eddine Ayat, Mohamed Cheriet, and Ching Y Suen. Automatic model selection for the optimization of SVM kernels. *Pattern Recognition*, 38(10):1733–1745, 2005.

Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.

Gaston Baudat and Fatiha Anouar. Generalized discriminant analysis using a kernel approach. *Neural computation*, 12(10):2385–2404, 2000.

Patrick Billingsley. *Probability and measure*. John Wiley & Sons, 2008.

- Olivier Chapelle, Vladimir Vapnik, Olivier Bousquet, and Sayan Mukherjee. Choosing multiple parameters for support vector machines. *Machine learning*, 46(1-3):131–159, 2002.
- Vladimir Cherkassky and Yunqian Ma. Practical selection of SVM parameters and noise estimation for SVM regression. *Neural networks*, 17(1):113–126, 2004.
- Corinna Cortes and Vladimir Vapnik. Support-Vector Networks. *Machine learning*, 20(3):273–297, 1995.
- Louart Cosme and Romain Couillet. Harnessing Neural Networks: A Random Matrix Approach. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- Romain Couillet and Abba Kammoun. Random Matrix Improved Subspace Clustering. In *2016 Asilomar Conference on Signals, Systems, and Computers*, 2016.
- Romain Couillet, Florent Benaych-Georges, et al. Kernel spectral clustering of large dimensional data. *Electronic Journal of Statistics*, 10(1):1393–1454, 2016.
- Harris Drucker, Donghui Wu, and Vladimir N Vapnik. Support vector machines for spam categorization. *IEEE Transactions on Neural networks*, 10(5):1048–1054, 1999.
- Noureddine El Karoui et al. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1):1–50, 2010.
- Theodoros Evgeniou, Massimiliano Pontil, Constantine Papageorgiou, and Tomaso Poggio. Image representations for object detection using kernel classifiers. In *Asian Conference on Computer Vision*, pages 687–692. Citeseer, 2000.
- Glenn M Fung and Olvi L Mangasarian. Multicategory proximal support vector machine classifiers. *Machine learning*, 59(1-2):77–97, 2005.
- Cheng-Lung Huang and Chieh-Jen Wang. A GA-based feature selection and parameters optimization for support vector machines. *Expert Systems with applications*, 31(2):231–240, 2006.
- Guang-Bin Huang, Hongming Zhou, Xiaojian Ding, and Rui Zhang. Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 42(2):513–529, 2012.
- Thorsten Joachims. Text categorization with support vector machines: Learning with many relevant features. In *European conference on machine learning*, pages 137–142. Springer, 1998.
- Yann LeCun, LD Jackel, Léon Bottou, Corinna Cortes, John S Denker, Harris Drucker, Isabelle Guyon, UA Muller, E Sackinger, Patrice Simard, et al. Learning algorithms for classification: A comparison on handwritten digit recognition. *Neural networks: the statistical mechanics perspective*, 261:276, 1995.

- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yuh-Jye Lee and Olvi L Mangasarian. RSVM: Reduced Support Vector Machines. In *SDM*, volume 1, pages 325–361, 2001.
- Xiaoyi Mai and Romain Couillet. The counterintuitive mechanism of graph-based semi-supervised learning in the big data regime. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- Edgar Osuna, Robert Freund, and Federico Girosit. Training support vector machines: an application to face detection. In *Computer vision and pattern recognition, 1997. Proceedings., 1997 IEEE computer society conference on*, pages 130–136. IEEE, 1997.
- Constantine Papageorgiou and Tomaso Poggio. A trainable system for object detection. *International Journal of Computer Vision*, 38(1):15–33, 2000.
- T Poggio, S Mukherjee, R Rifkin, A Raklin, and A Verri. B. uncertainty in geometric computations. *Kluwer Academic Publishers*, 22:131–141, 2002.
- Bernhard Scholkopf and Alexander J Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2001.
- Bernhard Scholkopf and Klaus-Robert Mullert. Fisher discriminant analysis with kernels. *Neural networks for signal processing IX*, 1(1):1, 1999.
- Alex J Smola and Bernhard Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222, 2004.
- Johan AK Suykens and Joos Vandewalle. Least squares support vector machine classifiers. *Neural processing letters*, 9(3):293–300, 1999.
- Tony Van Gestel, Johan AK Suykens, Gert Lanckriet, Annemie Lambrechts, Bart De Moor, and Joos Vandewalle. Bayesian framework for least-squares support vector machine classifiers, gaussian processes, and kernel fisher discriminant analysis. *Neural computation*, 14(5):1115–1147, 2002.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer Science & Business Media, 2013.
- Jason Weston, Sayan Mukherjee, Olivier Chapelle, Massimiliano Pontil, Tomaso Poggio, and Vladimir Vapnik. Feature selection for SVMs. 2000.
- Jieping Ye and Tao Xiong. SVM versus Least Squares SVM. In *AISTATS*, pages 644–651, 2007.