

---

# WASSERSTEIN EXPONENTIAL KERNELS

---

**Henri De Plaen**

ESAT-STADIUS

KU Leuven

Leuven, 3001 Belgium

henri.deplaen@esat.kuleuven.be

**Michaël Fanuel**

ESAT-STADIUS

KU Leuven

Leuven, 3001 Belgium

michael.fanuel@esat.kuleuven.be

**Johan A. K. Suykens**

ESAT-STADIUS

KU Leuven

Leuven, 3001 Belgium

johan.suykens@esat.kuleuven.be

February 6, 2020

## ABSTRACT

In the context of kernel methods, the similarity between data points is encoded by the kernel function which is often defined thanks to the Euclidean distance, a common example being the squared exponential kernel. Recently, other distances relying on optimal transport theory – such as the Wasserstein distance between probability distributions – have shown their practical relevance for different machine learning techniques. In this paper, we study the use of exponential kernels defined thanks to the regularized Wasserstein distance and discuss their positive definiteness. More specifically, we define Wasserstein feature maps and illustrate their interest for supervised learning problems involving shapes and images. Empirically, Wasserstein squared exponential kernels are shown to yield smaller classification errors on small training sets of shapes, compared to analogous classifiers using Euclidean distances.

## 1 Introduction

Contemporary machine learning methods frequently rely on neural networks, and shape recognition relies more specifically on convolutional neural networks. The big advantage of the latter is its ability to take the underlying structure of the data into account by treating neighboring pixels together. If these methods are very often impressive by their performance, they are also known for their drawbacks such as a weak robustness and a difficult explainability. On the other side, though not always being as accurate as neural networks, kernel methods are praised for their easy explainability and robustness. Another advantage of kernel methods is their versatility as they easily be used in supervised and unsupervised methods, as well as for generation [1]. We emphasize here the interest of choosing a particular kernel based on Wasserstein distance for classifying small datasets consisting of shapes.

In the context of kernel methods, squared exponential kernel functions are widely used, mainly because of their universal approximation properties and their empirical success. These Gaussians consist of the exponential of the negative Euclidean distance squared. However, the Euclidean distance might not always be appropriate to compare data points when data has some specific structure. Indeed, it measures the correspondence of each feature independently of the other features. For example, let's consider the case of two identical 2D-shapes. When the two shapes overlap, their Euclidean distance is zero. However, if they do not overlap, their relative Euclidean distance becomes large although the shapes are identical. In other words, the Euclidean distance only compares each pixel at the same place on the grid, not taking the neighbouring pixels into account. The general structure of the features is not taken into account, only their strict correspondence. However, another distance – the Wasserstein distance – gained popularity in recent years

since it can incorporate the structure of the data if the dataset can be processed so that the datapoints can be considered as probability distributions.

## Contributions

The contributions of this paper are the following. Empirically, we demonstrate that squared exponential kernels (1) based on a regularized Wasserstein distance are performant on small scale classification problems involving shape datasets, compared for instance to the popular Gaussian RBF kernel [2]. Also, an approximation technique is proposed, with the so-called Wasserstein feature map, so that a positive semi-definite (psd) kernel can be defined from the Wasserstein squared exponential kernel which is not necessarily psd.

## Notations and conventions

In the sequel, we denote vectors by bold lower case letters. Let  $\mathbf{1}$  be the all ones column vector. Also, we define  $\delta_y$  to be the Dirac measure at point  $y$ . A kernel  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is called positive semi-definite if all kernel matrices  $K = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^n$  are positive semi-definite.

## Wasserstein distances

The Wasserstein distance is a central notion in optimal transport theory. Also known as the *earth mover's distance*, it corresponds to the optimal transportation cost between two measures [3, 4]. Let  $p > 0$ . We then define two normalized empirical measures  $\alpha = \sum_{i=1}^m a_i \delta_{\mathbf{y}_i}$  and  $\beta = \sum_{j=1}^n b_j \delta_{\mathbf{z}_j}$  such that  $\alpha^\top \mathbf{1} = 1$  and  $\beta^\top \mathbf{1} = 1$ , and where  $\{\mathbf{y}_i \in \mathbb{R}^d\}_{i=1}^m, \{\mathbf{z}_j \in \mathbb{R}^d\}_{j=1}^n$  are support points. Also, we define an Euclidean distance matrix  $d_{ij} = \|\mathbf{y}_i - \mathbf{z}_j\|_2$ . Then, the  $p$ -Wasserstein distance is given by

$$\mathcal{W}_p(\alpha, \beta) = \left( \min_{\pi \in \Pi(\alpha, \beta)} \sum_{i,j} \pi_{ij} d_{ij}^p \right)^{1/p},$$

with  $\Pi(\alpha, \beta) = \{\Pi \in \mathbb{R}^{m \times n} | \Pi \mathbf{1} = \alpha \text{ and } \Pi^\top \mathbf{1} = \beta\}$ , the set of joint distributions  $\pi$  with specified marginals given by  $\alpha$  and  $\beta$ . Intuitively, the optimal probability distribution  $\pi^*$  represents the optimal mass transportation scheme from  $\alpha$  to  $\beta$ . A particular result occurs in the one-dimensional ( $d = 1$ ) case assuming the support points are ordered, i.e.,  $y_1 \leq \dots \leq y_m$  and  $z_1 \leq \dots \leq z_n$ , where the Wasserstein distance reduces to an  $\ell^p$ -norm:  $\mathcal{W}_p^p \left( \frac{1}{n} \sum_{i=1}^n \delta_{y_i}, \frac{1}{n} \sum_{j=1}^n \delta_{z_j} \right) = \frac{1}{n} \|\mathbf{y} - \mathbf{z}\|_p^p$  [4]. This connection between  $\ell^p$ -norms and Wasserstein distances is only clear in one dimension, illustrating here again the fact that  $\ell^p$ -norms don't take the underlying structure into account. To take it into account, we need to consider the case  $d > 1$ . In this way we can define a new kernel function

$$k_W(\alpha, \beta) = \exp \left( -\frac{W_2^2(\alpha, \beta)}{2\sigma^2} \right). \quad (1)$$

However, this has some undesirable consequences concerning positive definiteness. A kernel  $k(\mathbf{x}, \mathbf{y}) = \exp(-t f(\mathbf{x}, \mathbf{y}))$  is positive semi-definite for all  $t > 0$  if and only if  $f(\mathbf{x}, \mathbf{y})$  is Hermitian and conditionally negative semi-definite [5]. Recall that a kernel is conditionally negative semi-definite if any Gram matrix  $F = [f(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^n$  (with  $n \geq 2$ ) built from a discrete sample satisfies  $\mathbf{c}^\top F \mathbf{c} \leq 0$  for all  $\mathbf{c}$  such that  $\mathbf{1}^\top \mathbf{c} = 0$ . However, the Wasserstein distance for  $d > 1$  is not necessarily conditionally negative definite [4]. The consequence is that we cannot guarantee that any resulting squared exponential kernel matrix built with the 2-Wasserstein distance is positive definite. This property is fundamental in kernel theory and more specifically for defining reproducing kernel Hilbert spaces (RKHS; see [6] for more details).

## 2 Dealing with indefinite exponential kernels

This restriction has lead authors to consider only some specific cases of Wasserstein distances which are known to be positive definite. The one-dimensional generic case is proven to be positive definite and has lead to the introduction of sliced Wasserstein distances [7, 8]. Another notable case is the Wasserstein distance between two Gaussians in more than one dimension, which can even be written in closed form [4].

Some kernel methods are still usable with non positive definite kernels, such as LS-SVMs [9, 10]. However, this leads to a slightly different interpretation of the global problem, using Kreĭn spaces for which a weaker version of the representer

theorem holds [11]. In this paper, we propose an alternative which allows us to still work with a positive definite kernel approximating the squared exponential kernel. If the Wasserstein exponential kernel can not be used, we can always find a parameter  $\sigma > 0$  and a finite dimensional feature map resulting in a positive definite kernel.

## 2.1 Positive definite squared exponential kernels and bandwidth choice

In this section, we show that for a given dataset, the corresponding Gram matrix of  $k_W$  is positive definite if the bandwidth parameter  $\sigma > 0$  is small enough.

**Definition 2.1.** Let  $d : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}_{\geq 0}$  be a symmetric function such that  $d(\mathbf{x}, \mathbf{x}) = 0$  and let  $\{\mathbf{x}_i \in \mathcal{D}\}_{i=1}^N$  be a dataset. A squared exponential kernel matrix is defined as

$$\mathbf{K}_{d,\sigma} = \left[ \exp \left( \frac{-d^2(\mathbf{x}_i, \mathbf{x}_j)}{2\sigma^2} \right) \right]_{i,j=1}^N.$$

By construction, this exponential kernel matrix will be symmetric and have a diagonal consisting only of ones. Its eigenvalues are real. To investigate its (semi)-definiteness, we have to investigate the sign of the minimum eigenvalue. The minimum eigenvalue  $\lambda_{\min}(\sigma)$  of  $\mathbf{K}_{d,\sigma}$  is the function  $\lambda_{\min} : \mathbb{R}_{>0} \rightarrow \mathbb{R}, \sigma \mapsto \min \{\lambda_1, \dots, \lambda_N\}$  where  $\lambda_1, \dots, \lambda_N$  are the eigenvalues of  $\mathbf{K}_{d,\sigma}$ . We can now prove the following result:

**Lemma 2.1.** The eigenvalues of the exponential kernel matrix  $\mathbf{K}_{d,\sigma}$  are continuous functions of  $\sigma$ . In particular,  $\lambda_{\min}(\sigma)$  is continuous.

*Proof.* This is a direct consequence of the continuity of the roots of a polynomial under continuously varying coefficients. Therefore, we have to prove that the coefficients of the characteristic polynomial of the exponential kernel matrix  $\mathbf{K}_{d,\sigma}$  is continuous in function of  $\sigma$ . The characteristic polynomial is given by  $\det(\mathbf{K}_{d,\sigma} - \lambda \mathbf{I})$  and by the formula of Leibniz, we ultimately have that the characteristic polynomial is a sum of products of elements of  $\mathbf{K}_{d,\sigma} - \lambda \mathbf{I}$ , which are continuous in function of  $\sigma$ . Hence, the coefficients are continuous and so are the eigenvalues.  $\square$

**Lemma 2.2.**  $\lim_{\sigma \rightarrow 0} \mathbf{K}_{d,\sigma} = \text{id}$  and thus  $\lambda_{\min}(0) = 1$ .

*Proof.* From Definition 2.1, we know that  $[\mathbf{K}_{d,\sigma}]_{i,j} = \exp \left( \frac{-d^2(\mathbf{x}_i, \mathbf{x}_j)}{2\sigma^2} \right)$  with  $d^2(\mathbf{x}_i, \mathbf{x}_i) = 0$  and  $d^2(\mathbf{x}_i, \mathbf{x}_j) > 0$  for  $i \neq j$ . Denote  $C_{i,j} = d^2(\mathbf{x}_i, \mathbf{x}_j)$  for simplicity. We have thus  $\lim_{\sigma \rightarrow 0} \exp \left( \frac{0}{2\sigma^2} \right) = 1$  and  $\lim_{\sigma \rightarrow 0} \exp \left( -\frac{C_{i,j}}{2\sigma^2} \right) = 0$  with  $C_{i,j} > 0$  for  $i \neq j$ , hence the identity matrix. By consequence, all the eigenvalues are equal to 1.  $\square$

**Lemma 2.3.** We have  $\lim_{\sigma \rightarrow \infty} \mathbf{K}_{d,\sigma} = \mathbf{1}\mathbf{1}^T$  and thus  $\lim_{\sigma \rightarrow \infty} \lambda_{\min}(\sigma) = 0$ .

*Proof.* Similarly as before, we have  $\lim_{\sigma \rightarrow +\infty} [\mathbf{K}_{d,\sigma}]_{i,j} = 1$  everywhere. By consequence, we have  $\lambda_{\max} = N$  and all others equal to zero, hence  $\lambda_{\min} = 0$ .  $\square$

**Proposition 2.4.** There exists a  $\sigma_{\text{PSD}} \in \mathbb{R}_+$  such that  $\mathbf{K}_{d,\sigma}$  is positive semi-definite for all  $\sigma \leq \sigma_{\text{PSD}}$ .

*Proof.* Let's proceed *ad absurdum* and suppose this is not the case. We consider the sequence  $(\sigma_n)_n$  converging to 0 with  $\sigma_0 = \sigma_{\text{PSD}}$ . There must exist some subsequence  $(\sigma_{n_j})_j$  such that  $(\lambda_{\min}(\sigma_{n_j}))_j < 0$ . If this sequence is finite, then it suffices to consider a new sequence with  $\sigma_{\text{PSD}} = \sigma_{n_{j_{\max}}+1}$ . If this subsequence is infinite, then  $(\lambda_{\min}(\sigma_n))_n$  cannot converge to 1. This is impossible because of the continuity of  $\lambda_{\min}(\sigma)$  (lemma 2.1) and its convergence to 1 (lemma 2.2). Hence, there exist some  $\sigma_{\text{PSD}} > 0$  such that  $\lambda_{\min}(\sigma) \geq 0$  for all  $\sigma \leq \sigma_{\text{PSD}}$ . This proves our proposition.  $\square$

We can empirically see the result of Proposition 2.4 in Fig. 1, where all eigenvalues are positive. To give some intuition, decreasing the  $\sigma$  tends to make the smallest distances more predominant, pushing the smallest eigenvalue progressively to the positive side. In this sense, an indefinite kernel matrix with  $\sigma$  close to  $\sigma_{\text{PSD}}$  will lead to very proportionally very small negative eigenvalues in magnitude. In this case, a finite positive definite approximation can be justified.

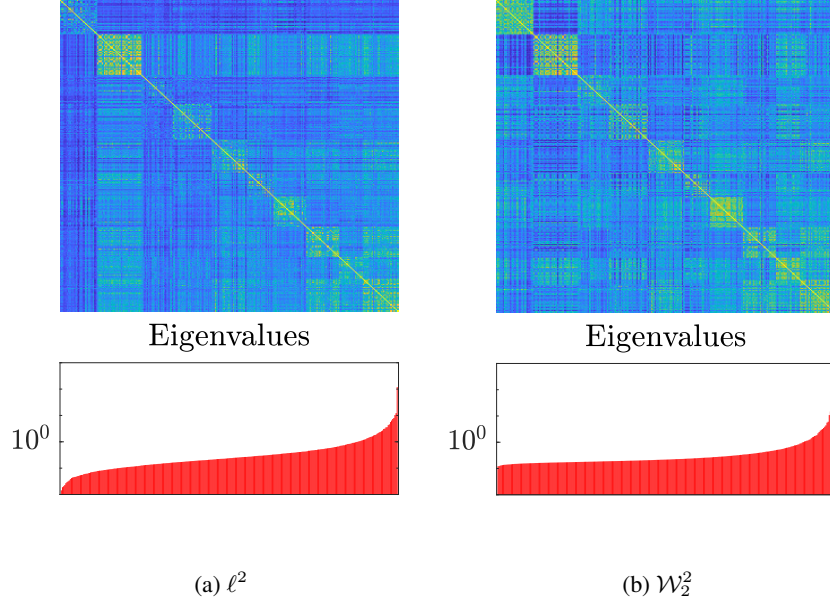


Figure 1: Comparison of the classical squared exponential kernel matrix (based on a  $\ell^2$ -distance) and the introduced Wasserstein exponential kernel matrix on 250 normalized digits of the MNIST dataset [12]. The digits are ordered by class in ascending order.

## 2.2 Wasserstein features

We can consider a finite dimensional feature map  $\phi(\mathbf{x})$  such that the positive semi-definite kernel  $\phi(\mathbf{x})^\top \phi(\mathbf{y})$  approximates  $k_W(\mathbf{x}, \mathbf{y})$  given in (1). This finite approximation is based on a training dataset  $\{\mathbf{x}_i\}_{i=1}^N$  for constructing an original kernel matrix  $\mathbf{K} = [k_W(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^N \in \mathbb{R}^{N \times N}$ . It suffices to truncate the spectral decomposition of the kernel matrix  $\mathbf{K} = \sum_{l=1}^N \lambda_l \mathbf{v}_l \mathbf{v}_l^\top$  to the  $\ell$  largest strictly positive eigenvalues. This will result in a new positive definite kernel matrix  $\mathbf{K}^{(\ell)} \stackrel{\text{def}}{=} \sum_{l=1}^{\ell} \lambda_l \mathbf{v}_l \mathbf{v}_l^\top \succ 0$  with  $\lambda_1 \geq \dots \geq \lambda_N$ . We can now reconstruct the different components of an approximate feature map

$$\phi_l(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\lambda_l}} \mathbf{k}_x^\top \mathbf{v}_l, \quad \text{for } i = 1, \dots, \ell, \quad (2)$$

with  $\mathbf{k}_x \stackrel{\text{def}}{=} [k_W(\mathbf{x}, \mathbf{x}_1) \dots k_W(\mathbf{x}, \mathbf{x}_N)]^\top$ . We refer to these different components as the *Wasserstein features* as they compose the approximate feature map  $\phi(\mathbf{x}) \stackrel{\text{def}}{=} [\phi_1(\mathbf{x}) \dots \phi_\ell(\mathbf{x})]^\top$  of the Wasserstein exponential kernel. This approximate feature map is constructed by using a training dataset, but can afterwards be evaluated at any out-of-sample point. By construction, we can verify that the Wasserstein features evaluated on the training dataset result in the truncated kernel matrix:

**Proposition 2.5.** *We have  $[\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)]_{i,j=1}^N = \mathbf{K}^{(\ell)}$ .*

*Proof.* It suffices to observe that  $\mathbf{k}_{\mathbf{x}_i} = \sum_{l=1}^N \lambda_l \mathbf{v}_l [\mathbf{v}_l]_i$ . By consequence, we have  $\phi_l(\mathbf{x}_i) = \sqrt{\lambda_l} [\mathbf{v}_l]_i$ .  $\square$

Proposition 2.4 suggests that even if no suitable  $\sigma$  can be found such that the kernel matrix is positive, the negative eigenvalues will remain very small in magnitude. By consequence, we can suppress them without losing much information. A truncated kernel is thus very close to the original one in spectral norm. This justifies the Wasserstein features in this sense that they are very close to the Wasserstein exponential kernel as well as being positive definite by construction. This fact can be visualized on Fig. 2.

Clearly, the *Wasserstein features* yield a positive semi-definite kernel. Moreover, it is also advantageous to work with finite dimensional feature maps to reduce the training time. Indeed, the computation of the Wasserstein distance (or an approximation with e.g. Sinkhorn's algorithm [13]) is still relatively expensive compared to  $\ell^2$  distance.

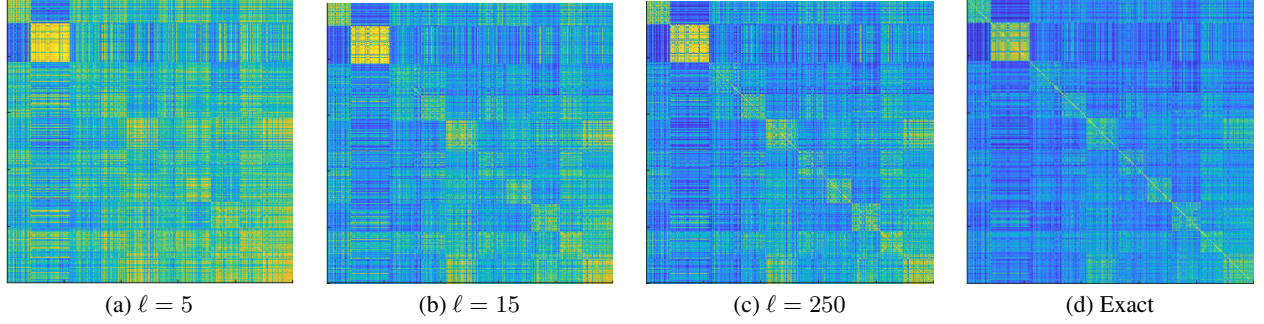


Figure 2: Kernel matrices constructed as the inner products of a different number Wasserstein features of a test set. These matrices are compared with the exact Wasserstein squared exponential kernel matrix of the test set. Both the training set and the test are of size  $N = 500$ .

Table 1: Percentage of classification error on the test set of three datasets. The standard deviation is given in parenthesis. The number of repeated simulations is 7 for MNIST, 8 for Quickdraw and 6 for USPS.

Dataset	MNIST		Quickdraw		USPS	
Method	Avg.	Best	Avg.	Best	Avg.	Best
Wass. LS-SVM (Core+OOS)	3.95 ( $\pm 0.18$ )	3.74	11.45 ( $\pm 0.39$ )	10.97	6.77 ( $\pm 0.52$ )	6.20
Wass. LS-SVM (Core)	3.81 ( $\pm 0.34$ )	3.28	10.80 ( $\pm 0.19$ )	10.52	7.93 ( $\pm 1.45$ )	6.35
Wass. LS-SVM (Indef.)	<b>3.40</b> ( $\pm 0.11$ )	<b>3.23</b>	<b>10.75</b> ( $\pm 0.27$ )	10.35	6.15 ( $\pm 0.67$ )	5.45
R. Wass. LS-SVM (Core+OOS)	3.91 ( $\pm 0.27$ )	3.45	11.79 ( $\pm 0.48$ )	10.95	6.68 ( $\pm 0.80$ )	5.70
R. Wass. LS-SVM (Core)	3.71 ( $\pm 0.15$ )	3.46	10.99 ( $\pm 0.44$ )	<b>10.07</b>	6.35 ( $\pm 0.11$ )	6.20
R. Wass. LS-SVM (Indef.)	3.48 ( $\pm 0.13$ )	3.29	12.43 ( $\pm 0.43$ )	11.95	<b>5.70</b> ( $\pm 0.29$ )	<b>5.40</b>
Wass. kNN	6.31 ( $\pm 0.33$ )	5.81	12.26 ( $\pm 0.33$ )	11.91	6.60 ( $\pm 0.44$ )	6.00
RBF LS-SVM	4.26 ( $\pm 0.10$ )	4.07	11.46 ( $\pm 0.20$ )	11.23	6.75 ( $\pm 0.04$ )	6.70
$\ell^2$ kNN	7.20 ( $\pm 0.15$ )	6.95	15.32 ( $\pm 0.40$ )	14.68	7.52 ( $\pm 0.38$ )	7.20
Set size	Core + OOS	Others	Core + OOS	Others	Core + OOS	Others
Training	1500 + 2500	4000	500 + 750	1250	1000 + 1500	2500
Validation	5000	5000	5000	5000	2000	2000
Test	10 000	10 000	10 000	10 000	2000	2000

### 3 Experiments

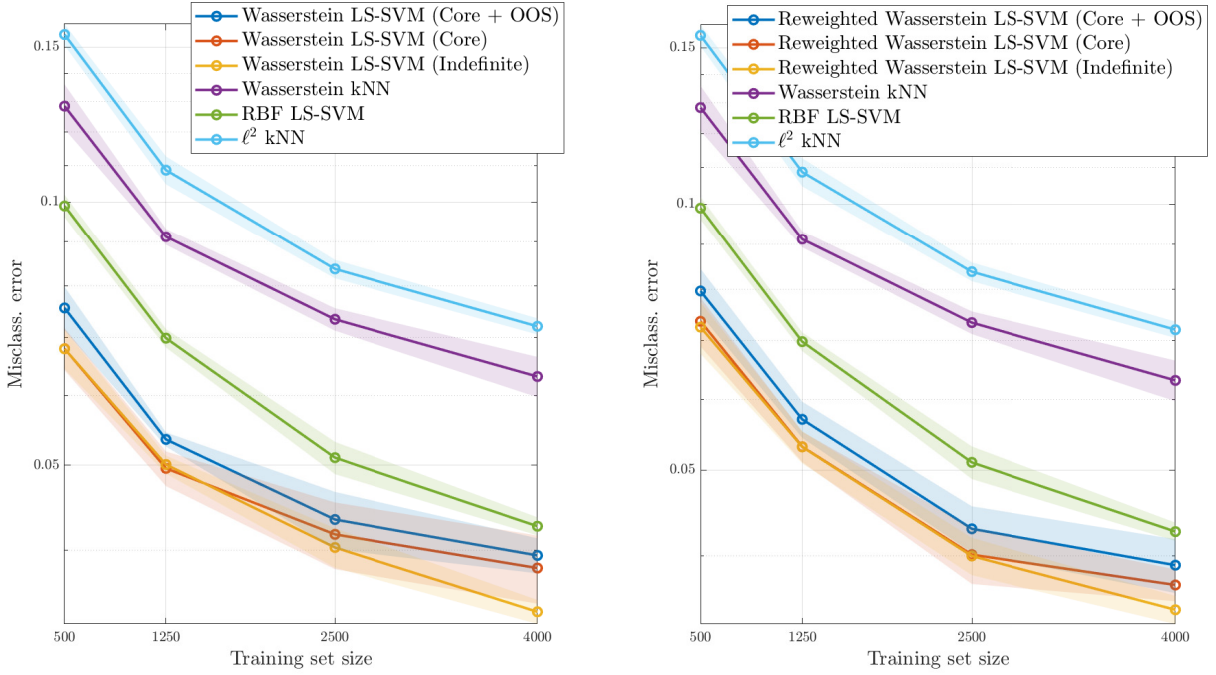
#### 3.1 Setup for 2D shape classification

Let  $u$  be a greyscale image that we unfold as a vector of length  $m$  and so that  $u_i > 0$  is the “grey” value at the pixel  $y_i$  of a pixel grid. It is mapped to a probability  $\alpha = \sum_{i=1}^m a_i \delta_{y_i}$  by defining  $a_i = u_i / \|u\|_1$ , so that the mass of  $\alpha$  is one. In practice, the  $p = 2$  Wasserstein distance is calculated in this paper with the help of the well-known entropic regularization, namely

$$\mathcal{W}_2^2(\alpha, \beta, \epsilon) = \min_{\pi \in \Pi(\alpha, \beta)} \sum_{i,j} \pi_{ij} d_{ij}^2 + \epsilon \pi_{ij} \log \pi_{ij},$$

where  $\epsilon > 0$  is a small regularization term and  $d_{ij}$  is the Euclidean distance between pixels located at  $y_i$  and  $y_j$  in a pixel grid. The advantage of this regularized problem is that its solution can be efficiently obtained thanks to the Sinkhorn algorithm, which can be parallelized. For more details, we refer to [4]. All the simulations used  $\epsilon = 0.4$  and the diagonal of the distance matrix set to zero.





(a) Comparison of Wasserstein exponential kernels with other similar methods.

(b) Comparison of *reweighted* Wasserstein exponential kernels with other similar methods.

Figure 3: Mean misclassification rates for various subset sizes of the MNIST dataset, computed on 7 simulations. The standard deviation is given by the errors bars. For the specific case of “Core + OOS”, the out-of-sample subset represents 300 datapoints on 500, 750 on 1250, 1500 on 2500 and 2500 on 4000. The size of validation set is always 5000 and of the test set always 10 000.

### 3.2 Shape recognition

We illustrate the use of the Wasserstein based kernels in the context of shape classifications. Namely, we train a Least Squares Support Vector Machine [14] classifier on subsets of the MNIST [12], Quickdraw<sup>1</sup> and USPS [15] datasets, which are sampled uniformly at random. These three datasets contain handwritten digits and shapes. The multiclass problem is solved by a one-versus-one encoding. One instance of these binary classifiers  $f(\mathbf{x}) = \text{sign}(\mathbf{w}^{\top} \phi(\mathbf{x}) + b^*)$  is obtained by solving

$$\min_{\substack{\mathbf{w} \in \mathbb{R}^{\ell}; b \in \mathbb{R} \\ e_i \in \mathbb{R}}} \mathbf{w}^{\top} \mathbf{w} + \frac{\gamma}{N} \sum_{i=1}^N e_i^2 \text{ s.t. } e_i = y_i - \mathbf{w}^{\top} \phi(\mathbf{x}_i) - b, \quad (3)$$

where  $y_i \in \{-1, 1\}$  and  $\phi(\mathbf{x}) \in \mathbb{R}^{\ell}$  is a feature map obtained for instance thanks to (2). The solution is obtained by solving

$$\begin{bmatrix} \sum_i \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^{\top} + \frac{N}{\gamma} \mathbb{I} & \sum_i \phi(\mathbf{x}_i) \\ \sum_i \phi(\mathbf{x}_i)^{\top} & N \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} \sum_i y_i \phi(\mathbf{x}_i) \\ \sum_i y_i \end{bmatrix},$$

which is a  $(\ell + 1) \times (\ell + 1)$  linear system. A classifier can also be obtained by solving the dual problem of (3). The optimality conditions of this dual problem yield the following  $(N + 1) \times (N + 1)$  linear system

$$\begin{bmatrix} K + \frac{N}{\gamma} \mathbb{I} & \mathbf{1} \\ \mathbf{1}^{\top} & 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}. \quad (4)$$

The resulting classifier has then the expression  $f(\mathbf{x}) = \text{sign}(\sum_{i=1}^N \alpha_i^* k(\mathbf{x}, \mathbf{x}_i) + b^*)$ . The hyperparameters  $\sigma > 0$  and  $\gamma > 0$  are chosen by validation. The final classification is done by minimizing the hamming distance on the

<sup>1</sup><https://quickdraw.withgoogle.com/data>

one-versus-one outputs [16]. In order to account for the amount of ink in the grey images  $\mathbf{u}$  and  $\mathbf{v}$ , we also introduce a reweighted kernel that is defined as

$$k_{RW}(\mathbf{u}, \mathbf{v}) = \|\mathbf{u}\|_1 \|\mathbf{v}\|_1 k_W \left( \frac{\mathbf{u}}{\|\mathbf{u}\|_1}, \frac{\mathbf{v}}{\|\mathbf{v}\|_1} \right). \quad (5)$$

Notice that a similar kernel has been defined with the Euclidean distance in [17, 18].

In our experiments, we compare several methods based on  $k_W$  and  $k_{RW}$ , Wasserstein and Euclidean distances.

### 3.2.1 Core Wasserstein kernel

The ‘‘Core’’ method consists in solving (3) thanks to the feature map (2) associated to  $K^{(\ell)}$ . The parameter  $\ell$  is chosen such that all the selected eigenvalues are larger than  $10^{-6}$  to avoid numerical instabilities. The optimal  $\mathbf{w}^*$  and  $\mathbf{b}^*$  are then obtained by solving a linear system.

### 3.2.2 Core Wasserstein kernel with out-of-sample

Our second method named ‘‘Core + OOS’’ uses almost the same methodology as ‘‘Core’’. However, a subset of the training set is used to construct the truncated Wasserstein kernel of Proposition 2.5. Then the out-of-sample (OOS) formula (2) is used to construct an approximation of the kernel matrix on the full training dataset. The advantage of this approximation is that it can avoid the full eigendecomposition of the kernel matrix which is necessary for the ‘‘Core’’ method.

### 3.2.3 Indefinite Wasserstein kernel

For this second method, we simply use for the kernel matrix the indefinite Gram matrix associated to (5) and solve the system (4) associated to the dual formulation of LS-SVM. While the associated optimization problem is not necessarily bounded in that case, the linear system (4) still has a solution in practice (almost surely if  $\gamma$  selected uniformly at random.). We name this method ‘‘Indefinite Wasserstein’’ in Figure 3.

### 3.2.4 Gaussian RBF

The previous methods are compared with a classical LS-SVM classifier with kernel

$$k(\mathbf{u}, \mathbf{v}) = \exp \left( \frac{-\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2} \right).$$

The parameters  $\sigma$  and  $\gamma$  are obtained by validation in the same spirit as above.

### 3.2.5 KNN

The same task is also performed for a kNN classifiers defined both with Euclidean and Wasserstein distances [19]. Those two methods are considered as benchmarks to assess the accuracy of the kernel methods hereabove. Notice that the number of nearest neighbours  $k$  is selected by validation.

## 3.3 Description of the simulations

The simulations are repeated several times and the mean classification error rate is given as well as the standard deviation. We emphasize that the classes are balanced in each of the datasets. The coded is provided on GitHub<sup>2</sup>.

## 3.4 Discussion

The results obtained by classifiers defined with Wasserstein exponential kernel  $k_W$  outperform the Euclidean and Wasserstein kNN classifiers, as well as LS-SVM with a Gaussian RBF kernel (see Fig. 3 and Table 1). The latter is especially outperformed when the number of training data points is limited to a few thousands. We observed empirically that the advantage of  $k_W$  is indeed reduced as the size of the training set further increases. Surprisingly, the classifier obtained for the indefinite  $k_W$  kernel yields the best performance when the training set is larger. For moderate size training sets, LS-SVM classifiers can be competitive with respect to other methods that do not rely on convolutional neural networks. The latter are known to be performant for relatively large training datasets. While an advantage of Wasserstein based methods is an increased accuracy in the classification tasks of this paper, a main disadvantage is the increased training time.

<sup>2</sup>[https://github.com/hdeplaen/Exponential\\_Wasserstein\\_Kernels](https://github.com/hdeplaen/Exponential_Wasserstein_Kernels)

## 4 Conclusion

In this paper, we proposed the use of Wasserstein squared exponential kernels for classifying shapes given relatively small training datasets. Although the computation of Wasserstein distances is expensive, it can be made possible thanks to the entropic regularization and the Sinkhorn algorithm, as it is well known. The so-called Wasserstein features are also proposed to serve as an approximation of the Wasserstein squared exponential kernel which is not necessarily positive semidefinite. In particular, this construction is possible if the bandwidth parameter is small enough as it is explained by elementary theoretical results. These theoretical results also open a door to more general exponential kernels based on any measure of similarity.

## Acknowledgment

EU: The research leading to these results has received funding from the European Research Council under the European Union’s Horizon 2020 research and innovation program / ERC Advanced Grant E-DUALITY (787960). This paper reflects only the authors’ views and the Union is not liable for any use that may be made of the contained information. Research Council KUL: Optimization frameworks for deep kernel machines C14/18/068. Flemish Government: FWO: projects: GOA4917N (Deep Restricted Kernel Machines: Methods and Foundations), PhD/Postdoc grant. This research received funding from the Flemish Government under the “Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen” programme. Ford KU Leuven Research Alliance Project KUL0076 (Stability analysis and performance improvement of deep reinforcement learning algorithms). The computational resources and services used in this work were provided by the VSC (Flemish Supercomputer Center), funded by the Research Foundation - Flanders (FWO) and the Flemish Government – department EWI.

## References

- [1] Arun Pandey, Joachim Schreurs, and Johan A. K. Suykens. Generative restricted kernel machines, arxiv:1906.08144.
- [2] Carl E. Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [3] Cédric Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008.
- [4] Gabriel Peyré and Marco Cuturi. *Computational Optimal Transport: With Applications to Data Science*. Now Foundations and Trends, 2019.
- [5] Christian Berg, Jens Peter Reus Christensen, and Paul Ressel. *Harmonic Analysis on Semigroups*, volume 100 of *Graduate Texts in Mathematics*. Springer New York, New York, NY, 1984.
- [6] Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.
- [7] Mathieu Carrière, Marco Cuturi, and Steve Oudot. Sliced wasserstein kernel for persistence diagrams. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, page 664–673. JMLR.org, 2017.
- [8] Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. Generalized sliced wasserstein distances. In *Advances in Neural Information Processing Systems 32*, pages 261–272. Curran Associates, Inc., 2019.
- [9] Johan A. K. Suykens, Tony Van Gestel, Jos De Brabanter, Bart De Moor, and Joos Vandewalle. *Least Squares Support Vector Machines*. World Scientific, Singapore, 2002.
- [10] Xiaolin Huang, Andreas Maier, Joachim Hornegger, and Johan A. K. Suykens. Indefinite kernels in least squares support vector machines and principal component analysis. *Applied and Computational Harmonic Analysis*, 43(1):162–172, 2017.
- [11] Cheng Soon Ong, Xavier Mary, Stéphane Canu, and Alexander J Smola. Learning with non-positive kernels. In *Twenty-first international conference on Machine learning - ICML ’04*, page 81, New York, New York, USA, 2004. ACM Press.
- [12] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [13] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems 26*, pages 2292–2300. Curran Associates, Inc., 2013.
- [14] Johan A. K. Suykens and Joos Vandewalle. Least squares support vector machine classifiers. *Neural Process. Lett.*, 9(3):293–300, June 1999.
- [15] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, May 1994.
- [16] Johan A. K. Suykens and Joos Vandewalle. Multiclass least squares support vector machines. In *IJCNN’99. International Joint Conference on Neural Networks. Proceedings (Cat. No.99CH36339)*, volume 2, pages 900–903 vol.2, July 1999.
- [17] Julien Mairal. End-to-end kernel learning with supervised convolutional kernel networks. In *Advances in Neural Information Processing Systems 29*, pages 1399–1407. Curran Associates, Inc., 2016.



- [18] Dexiong Chen, Laurent Jacob, and Julien Mairal. Biological sequence modeling with convolutional kernel networks. *Bioinformatics*, 35(18):3294–3302, 02 2019.
- [19] Michael Snow and Jan Van Lent. Monge’s optimal transport distance for image classification, arxiv:1612.00181.