

# Low Dimensional Explicit Feature Maps

Ondřej Chum

CMP, Faculty of Electrical Engineering, CTU in Prague

chum@cmp.felk.cvut.cz

## Abstract

*Approximating non-linear kernels<sup>1</sup> by finite-dimensional feature maps is a popular approach for speeding up training and evaluation of support vector machines or to encode information into efficient match kernels. We propose a novel method of data independent construction of low dimensional feature maps. The problem is cast as a linear program which jointly considers competing objectives: the quality of the approximation and the dimensionality of the feature map.*

*For both shift-invariant and homogeneous kernels the proposed method achieves a better approximations at the same dimensionality or comparable approximations at lower dimensionality of the feature map compared with state-of-the-art methods.*

## 1. Introduction

Kernel machines, such as support vector machines (SVMs), can approximate any function or decision boundary arbitrarily well when provided with enough training data. However, such methods scale poorly with the size of the training set. On the other hand, it was shown [6] that linear SVMs can be trained in linear time with the number of training examples, which allows its application to very large datasets. Approximate embeddings, or feature maps, can preserve the accuracy of kernel methods and enable scaling to large datasets at the same time.

The demand for linear approximations of non-linear kernel functions is not limited to SVM classification. The idea of efficient match kernels [2] has been used in various areas of computer vision. Examples where linear approximation of non-linear kernels plays an important role, are kernel descriptors of interest points proposed in [1]. In the domain of image retrieval, [15] proposes to encode the dominant orientation of regions of interest into aggregated image descriptors, such as VLAD [5] or Fisher vectors [12].

<sup>1</sup>I would like to thank Tomáš Werner for his valuable opinions and interesting discussions. This work was supported by MSMT LL1303 ERC-CZ grant.

Formally, for a positive definite kernel [14]  $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  there exists a Hilbert space  $\mathcal{H}$  and a mapping  $\Psi : \mathbb{R}^n \rightarrow \mathcal{H}$ , so that  $K(\mathbf{x}, \mathbf{y}) = \langle \Psi(\mathbf{x}), \Psi(\mathbf{y}) \rangle_{\mathcal{H}}$ , where  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  is a scalar product in  $\mathcal{H}$ . We address the problem of finding a low-dimensional mapping  $\hat{\Psi} : \mathbb{R}^n \rightarrow \mathbb{R}^D$  so that  $\hat{\Psi}(x)^\top \hat{\Psi}(y) \approx \langle \Psi(\mathbf{x}), \Psi(\mathbf{y}) \rangle_{\mathcal{H}}$ . A natural requirement from the accuracy point of view is to introduce as little error as possible by the approximation. On the other hand, from the practical point of view, the dimensionality of the approximate feature map should be as low as possible. These two criteria are clearly competing. We propose an optimization approach which is relaxed into a linear program, that trades off both criteria.

## 1.1. Related work

We briefly review the most relevant work on data independent (no training data needed) methods of kernel approximation. Random Fourier features were introduced by Rahimi and Recht in [13]. The feature map is a Monte Carlo approximation of the kernel where each dimension of the feature map represents a cosine function drawn from the distribution given by the spectrum of the kernel signature. The Monte Carlo approach requires relatively high number of samples to provide accurate approximation, however, unlike most other approaches, is directly applicable to very high dimensional input data. The idea has been extended from shift-invariant kernels to skewed multiplicative histogram kernels in [9]. Maji and Berg in [11] approximate the intersection kernel by a sparse feature map in closed-form. In [19] high dimensional sparse feature maps are derived and their relation to product quantisation is shown. In [18] Vedaldi and Zisserman introduced a generalization of explicit feature maps to the family of additive homogeneous kernels. In our paper, considerably better approximations with feature maps of the same dimensionality or equally good approximation with lower dimensionality of the feature maps is achieved compared with results of [18] implemented in [17]. We also show that the proposed method allows for optimization of meaningful errors measured on the homogeneous kernel output, rather than solely approximating the kernel signature.

## 2. Problem formulation

In this section, the problem of shift-invariant kernel approximation is outlined, and then the proposed approach is described. For now, we will focus only on one dimensional kernels  $K(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ . Kernels in more dimensions are discussed in section 6.

Consider a family of shift-invariant (or stationary) kernels

$$K(x, y) = K(x + c, y + c) \quad \forall x, y, c \in \mathbb{R}. \quad (1)$$

A signature  $k(\lambda) : \mathbb{R} \rightarrow \mathbb{R}$  of a shift invariant kernel  $K$  is defined as  $k(\lambda) = K(-\lambda/2, \lambda/2)$ . Such a one-dimensional signature function fully specifies the kernel, since  $K(x, y) = k(x - y)$ .

We will study approximations  $\hat{K} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  to shift-invariant kernels, in particular such approximations that can be written as an inner product of low-dimensional feature maps  $\hat{\Psi} : \mathbb{R} \rightarrow \mathbb{R}^D$

$$\hat{K}(x, y) = \hat{\Psi}(x)^\top \hat{\Psi}(y) \approx K(x, y).$$

The kernel  $K$  will be approximated via approximating the kernel signature  $k$  by  $\hat{k} : \mathbb{R} \rightarrow \mathbb{R}$  in the form

$$\hat{k}(\lambda) = \sum_{\omega \in \Omega} \alpha_\omega \cos(\omega\lambda), \quad (2)$$

where  $\omega$  is a frequency,  $\Omega$  is a finite set of frequencies  $\Omega \subset [0, \omega^{\max}]$ , and  $\alpha_\omega \in \mathbb{R}_0^+$  are non-negative weights. Kernels in the form of (2) can be directly converted into feature maps

$$\hat{\Psi}_\omega(x) = \begin{pmatrix} \sqrt{\alpha_\omega} \cos(\omega x) \\ \sqrt{\alpha_\omega} \sin(\omega x) \end{pmatrix}.$$

From the identity

$$\cos(x - y) = \cos(x) \cos(y) + \sin(x) \sin(y)$$

it follows that

$$\hat{\Psi}_\omega(x)^\top \hat{\Psi}_\omega(y) = \alpha_\omega \cos(\omega(x - y)).$$

The feature map  $\hat{\Psi}(x)$  defined by the signature  $\hat{k}$  is a concatenation of  $\hat{\Psi}_\omega$  for all  $\omega \in \Omega$ . The dimensionality  $D(\hat{k})$  of the feature map  $\hat{\Psi}(x)$  is

$$D(\hat{k}) = \sum_{\omega \in \Omega} \delta(\alpha_\omega) D_\omega, \quad (3)$$

where  $\delta(\alpha_\omega) = 0$  for  $\alpha_\omega = 0$ , and  $\delta(\alpha_\omega) = 1$  otherwise,

$$D_\omega = \begin{cases} 1 & \omega = 0 \\ 2 & \omega > 0. \end{cases}$$

Here  $D_\omega$  denotes the dimensionality of feature map for a particular frequency  $\omega$ . The value of  $D_\omega = 1$  for  $\omega = 0$  comes from the fact that  $\hat{\Psi}_0 = (\sqrt{\alpha_\omega}, 0)^\top$ , where the zero can be dropped from the embedding.

**Input domain.** The input  $x$  for the kernel function is typically some measurement, such as coordinates of a point in a canonical patch (of fixed size), angle of the dominant orientation, or an entry of a normalized histogram. We make the assumption that the measured features  $x$  come from a bounded interval  $x \in [a, b]$ . This assumption is natural for many practical problems. Given the properties of the shift-invariant kernels,  $x \in [a, b]$  implies that the kernel signature  $k$  needs to be approximated on interval  $[-M, M]$ ,  $M = b - a$ .

**Error function.** The similarity of the original signature function  $k$  and its approximation  $\hat{k}$ , is measured by an error function  $C(k, \hat{k}) \in \mathbb{R}_0^+$ . In order to use discrete optimization methods, the error function used in the paper will only depend on a finite number of points  $z$  from an evaluation set  $Z$ ,  $z \in Z \subset [0, M]$ . The points  $z$  are non-negative, as both  $k$  and  $\hat{k}$  are symmetric. The discretization of the input domain is optimal for quantities that are discrete, such as for pixel coordinates. In many domains, sufficiently fine discretization introduces negligible error compared to the error introduced by the measurement estimation, e.g., the angle of the dominant orientation of a feature patch. If a continuous input domain is essential, the number  $|Z|$  of the points  $z$  has to be adjusted with respect to the maximal frequency  $\omega^{\max}$  and the spectrum of the kernel signature  $k$ .

In the paper, the two following error functions will be used

$$C_1(k, \hat{k}) = \sum_{z \in Z} w(z) \cdot |k(z) - \hat{k}(z)|, \quad (4)$$

$$C_\infty(k, \hat{k}) = \max_{z \in Z} w(z) \cdot |k(z) - \hat{k}(z)|, \quad (5)$$

where  $w(z) \in \mathbb{R}_0^+$  are weights that adjust the relative importance of the approximation error at point  $z$ . For all  $w(z) = 1$  constant, (4) represents  $L_1$  norm and (5) represents  $L_\infty$  norm.

### 2.1. Optimization

Two antagonistic objectives have to be considered in the approximation task: keeping the dimensionality  $D(\hat{k})$  of the embedding low and obtaining as close an approximation, measured by  $C(k, \hat{k})$ , of the kernel as possible.

Since  $D(\hat{k})$  is not convex, not even continuous, we apply an LP relaxation [16] to make the optimization tractable. Instead of dealing with the dimensionality  $D(\hat{k})$ , which is a weighted  $L_0$  norm (3), a weighted  $L_1$  norm

$$\bar{D}(\hat{k}) = \sum_{\omega \in \Omega} D_\omega \alpha_\omega \quad (6)$$

is used, recall that  $\alpha_\omega \geq 0$ .

The task of finding approximation  $\hat{k}$  that minimizes  $\bar{D}(\hat{k})$  while preserving a defined quality of the approximation is formulated as a linear program

$$\min_{\hat{k}} \bar{D}(\hat{k}) \quad \text{subject to } C(k, \hat{k}) \leq C^{\max} \in \mathbb{R}^+.$$

Finding an approximation  $\hat{k}$  of fixed dimensionality  $D^{\max}$  of the feature map is sought while minimizing  $C(k, \hat{k})$  is approximated by a linear program

$$\min_{\hat{k}} \bar{D}(\hat{k}) + \gamma C(k, \hat{k}),$$

where  $\gamma \in \mathbb{R}^+$  is a constant controlling the trade-off between the quality of the approximation and the relaxed dimensionality  $\bar{D}$  of the feature map. A version of binary search for the appropriate weight  $\gamma$  is used: the LP is executed for the value of  $D$  (not  $\bar{D}$ ), if  $D \leq D^{\max}$  the value of  $\gamma$  is increased (higher importance to the fit cost), otherwise the value of  $\gamma$  is decreased (higher importance to the solution sparsity). The solution with the best fit is selected among LP outputs with  $D \geq D^{\max}$ , for outputs with  $D > D^{\max}$  considering only top largest values of  $\alpha_\omega$  so that the dimensionality is at most  $D^{\max}$ .

### 3. Periodic kernels

Let  $k$  be a kernel signature that is periodic with period  $2M$ . The task in this section is to approximate  $k$  on interval  $[-M, M]$ , which is equivalent to approximating on all of  $\mathbb{R}$  since  $k$  is periodic. The spectrum of  $k$  is restricted to harmonics of the base frequency  $\pi/M$ , and hence

$$\Omega_0 = \left\{ i \frac{\pi}{M} \mid i \in \mathbb{N}_0 \right\}. \quad (7)$$

A standard approach to this problem is to project the function  $k$  to an orthogonal basis  $\cos(i\pi/M\lambda)$ . The function  $k$  is then approximated using basis functions with the highest values of the coefficients. Such an approach is efficient for one dimensional kernels and the method proposed in this paper does not bring any contribution to this problem. Results for multiplicative kernels (Section 6) are applicable to periodic functions.

### 4. Aperiodic kernels

In this section, an approximation of kernels on interval  $[-M, M]$  with signature that is not periodic (or do not have period  $2M$ ) is derived. Many shift invariant kernels, including the RBF kernel, are not periodic.

#### 4.1. Discrete frequencies

Following [18], for an aperiodic kernel signature  $k$ , there is a function  $g$  with period  $2M$  and  $g(\lambda) = k(\lambda)$  for

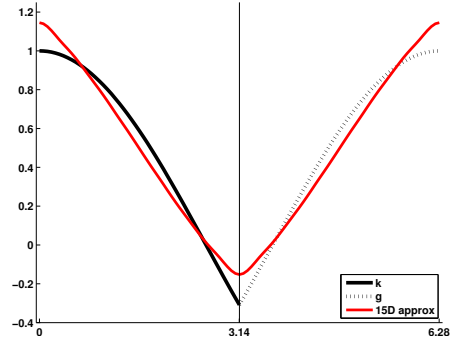


Figure 1. Kernel signature  $k = \cos(0.6\lambda)$  (solid black curve) that is not periodic on interval  $[-\pi, \pi]$  is approximated via approximating periodic function  $g$  (solid black and dashed black) using only harmonic angular frequencies  $\Omega = \mathbb{N}_0$ . Frequencies with negative coefficients are truncated which leads to poor approximation (red curve).

$\lambda \in [-M, M]$ . Then approximating periodic  $g$ , as in the previous section, using harmonic frequencies  $\Omega_0$  (7) only, approximates  $k$  on  $[-M, M]$ .

This approach has two drawbacks: First, even though  $k$  has a non-negative spectrum due to Bochner’s theorem, this does not hold for  $g$ . All frequencies with negative weights have to be left out [18]. As a consequence, the approximation of the signature function  $k$  cannot be arbitrarily precise, even for very high dimensional feature maps. Second, approximating  $g$  instead of  $k$  is not optimal with respect to the dimensionality of the feature map. To demonstrate these claims, consider the toy example in Figure 1. The kernel signature  $k(\lambda) = \cos(0.6\lambda)$  approximated on  $[-\pi, \pi]$  is not periodic with the period of  $2\pi$ . The approximation of periodic  $g$  by harmonic frequencies  $\omega \in \mathbb{N}_0$  with non-negative coefficients is not satisfactory. The exact feature map, originating from  $\hat{k} = \cos(0.6\lambda)$ , is two-dimensional, but the optimal frequency  $\omega^* = 0.6 \notin \Omega = \mathbb{N}_0$ <sup>2</sup>.

A simple generalization of the above approach increases the number of possible frequencies with increasing  $j \in \mathbb{N}$

$$\Omega_j = \left\{ i \frac{\pi}{2^j M} \mid i \in \mathbb{N}_0 \right\}. \quad (8)$$

Since  $\Omega_j \subset \Omega_{j+1}$ , approximation with frequencies  $\Omega_{j+1}$  will not be worse than with  $\Omega_j$ . With  $j$  approaching infinity, the set  $\Omega_j$  will contain frequencies arbitrarily close to any real-valued frequency. However, sets  $\Omega_j$  with large  $j$  are impractical in real problems. Sets  $\Omega_j$  of practical use lead to better approximations than  $\Omega_0$ , but still can only reach a discrete subset of possible frequencies.

Using the set of frequencies in (8) can be interpreted as approximating a periodic function  $g(\lambda)$  with period of

<sup>2</sup>The problem can be alleviated by approximating the signature  $k = \cos(0.6\lambda)$  on interval  $[-5\pi, 5\pi]$ . This toy example was selected as an extreme case to demonstrate the drawbacks of the standard approach.

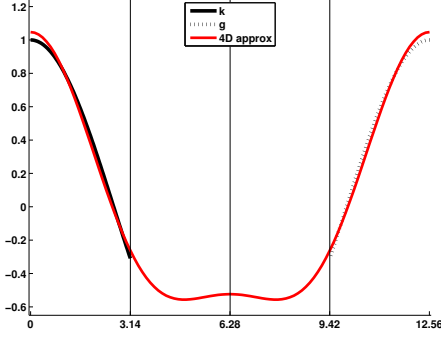


Figure 2. Approximating kernel signature  $k = \cos(0.6\lambda)$  (solid black curve) on  $[-\pi, \pi]$  via approximation of periodic function  $g$  (solid black and dashed black) with period  $4\pi$  using harmonic frequencies  $F = \{i/2 \mid i \in \mathbb{N}_0\}$ . There are no constraints imposed on  $g$  on  $(\pi, 3\pi)$ . Approximation by 4D feature map drawn in red.

$2^{j+1}M$ , where  $g(\lambda) = k(\lambda)$  for  $\lambda \in [-M, M]$ , and no constraints imposed on  $g(\lambda)$  in interval  $\lambda \in (M, 2^jM)$ . The situation is depicted in Figure 2 for  $j = 1$ .

## 4.2. Continuous frequencies

In the framework of discrete optimization used in this paper, the pool of frequencies  $\Omega$  is required to be finite and, for practical reasons, not extremely large. To access any real frequency while preserving finite  $\Omega$ , we will slightly modify the form of the approximation of the kernel signature  $\hat{k}$  to

$$\hat{k}(\lambda) = \sum_{\omega \in \Omega} \alpha_\omega \cos((\omega + d_\omega)\lambda), \quad (9)$$

where

$$|d_\omega| \leq d^{\max} \quad (10)$$

is a small difference in the frequency. The differences  $d_\omega$  are estimated jointly with the weights  $\alpha_\omega$  by the linear program. That is, instead of using exactly frequency  $\omega$  in the approximation, any frequency within interval  $[\omega - d^{\max}, \omega + d^{\max}]$  can be used. The first order Taylor expansion of the cosine function in the frequency variable  $\omega$  (not in  $\lambda$ ) reads

$$\cos((\omega + d_\omega)\lambda) \approx \cos(\omega\lambda) - d_\omega\lambda \sin(\omega\lambda). \quad (11)$$

Such an approximation is good only in a small neighbourhood of  $\omega$ , which is controlled by the size  $d^{\max}$  of the “trust region” (10). By substituting (11) into (9), we obtain

$$\hat{k}(\lambda) = \sum_{\omega \in \Omega} \alpha_\omega \cos(\omega\lambda) - \sum_{\omega \in \Omega} d_\omega \alpha_\omega \lambda \sin(\omega\lambda). \quad (12)$$

By introducing an auxiliary variable  $\beta_\omega = d_\omega \alpha_\omega$ , equation (10) transforms to

$$|\beta_\omega| \leq \alpha_\omega d^{\max}. \quad (13)$$

Both (12) and (13) in variables  $(\alpha_\omega, \beta_\omega)$  are in a form that can be written as a linear program. Compared to the original formulation,  $|\Omega|$  variables  $\beta_\omega$ , and  $2|\Omega|$  constraints (13) were introduced in the linear program.

**Implementation details.** In our experiments, we first apply an approximation with a discrete set  $\Omega$  of frequencies equally spaced in  $[0, \omega^{\max}]$ , with spacing at most  $d^{\max} = 0.1$ , as described in section 2.1. Then, an iterative process is executed. Each iteration is composed of execution of the LP approximation using the first order Taylor expansion formulation (9) followed by a frequency update

$$\begin{aligned} d_\omega &= \beta_\omega / \alpha_\omega \quad \dots \text{ compute } d\text{'s} \\ \omega &\leftarrow \omega + d_\omega \quad \dots \text{ update frequencies} \\ d^{\max} &\leftarrow d^{\max} / 2 \quad \dots \text{ reduce the max step.} \end{aligned}$$

The iteration is used to eliminate the approximation error introduced by the Taylor expansion. In each step, the allowed difference in frequency  $d^{\max}$  is halved, which guarantees the convergence.

## 5. Homogeneous kernels

A homogeneous kernel is a positive definite kernel  $K : \mathbb{R}_0^+ \times \mathbb{R}_0^+ \rightarrow \mathbb{R}_0^+$  satisfying

$$K_h(cx, cy) = cK_h(x, y) \quad \forall x, y, c \geq 0.$$

Following [18], by setting  $c = \sqrt{xy}$ , any homogeneous kernel can be decomposed as

$$\begin{aligned} K_h(x, y) &= \sqrt{xy} \cdot K_h\left(\sqrt{x/y}, \sqrt{y/x}\right) = \\ &= \sqrt{xy} \cdot k_h(\log y - \log x), \end{aligned} \quad (14)$$

where  $k_h(\lambda)$  is a signature of  $K_h$

$$k_h(\lambda) = K_h(e^{-\lambda/2}, e^{\lambda/2}).$$

The signature of the homogeneous kernel (14) resembles the signature of the shift-invariant kernel after transforming the input domain into log-space. The homogeneous kernel can be approximated [18] via approximating the signature  $k_h(\lambda)$  by  $\hat{k}(\lambda)$  in a similar manner as in (2). The resulting feature map is then

$$\hat{\Psi}_\omega(x) = \begin{pmatrix} \sqrt{\alpha_\omega x} \cdot \cos(\omega \log x) \\ \sqrt{\alpha_\omega x} \cdot \sin(\omega \log x) \end{pmatrix}.$$

While for shift-invariant kernels optimizing the signature approximation was equivalent to optimizing the kernel approximation, we show that for homogeneous kernels the situation is different. We will demonstrate it on the  $L_\infty$  error measure, as it is independent of the data distribution. Derivation for other error measures is straightforward.

We first derive minimization  $\min_{\hat{k}} \varepsilon_A$  of the absolute  $L_\infty$  error

$$\begin{aligned} \varepsilon_A &= \max_{x, y \in (0, b]} |K_h(x, y) - \hat{K}_h(x, y)| = \\ &= \max_{x, y \in (0, b]} \sqrt{xy} \cdot \left| k_h\left(\log \frac{y}{x}\right) - \hat{k}\left(\log \frac{y}{x}\right) \right|. \end{aligned} \quad (15)$$

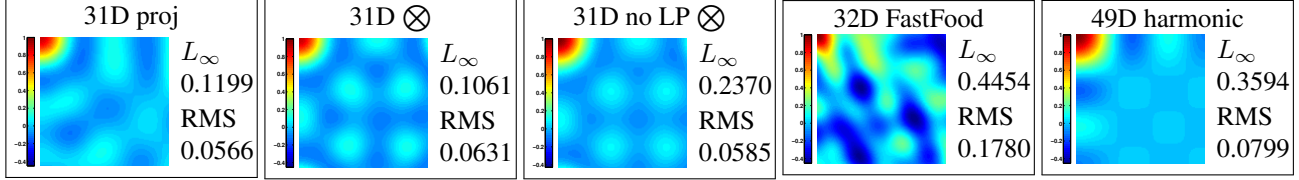


Figure 3. Comparison of kernel signature approximations of a symmetric 2D RBG kernel by projections, modulation, FastFood [7], and modulation of harmonic frequencies . Only one quadrant of the kernel signature is shown.

Let  $\lambda = \log y - \log x \geq 0$  (which is equivalent to  $y \geq x$ ) without a loss of generality, as  $k_h$  is symmetric. The error  $\varepsilon_A$  can be written as

$$\begin{aligned} \varepsilon_A &= \max_{y \in (0, b], \lambda \geq 0} y e^{-\lambda/2} |k_h(\lambda) - \hat{k}(\lambda)| \\ &= b \cdot \max_{\lambda \geq 0} e^{-\lambda/2} |k_h(\lambda) - \hat{k}(\lambda)|. \end{aligned}$$

This is achieved by optimizing the approximation of signature  $k_h$  with weighted  $C_\infty$  error function (5) with weight

$$w_A(\lambda) = e^{-\lambda/2}. \quad (16)$$

Similarly, we derive the minimization  $\min_{\hat{k}} \varepsilon_R$  of the relative  $L_\infty$  error

$$\begin{aligned} \varepsilon_R &= \max_{x, y \in (0, b]} \frac{|K_h(x, y) - \hat{K}_h(x, y)|}{K_h(x, y)} = \quad (17) \\ &= \max_{\lambda} \frac{1}{k_h(\lambda)} |k_h(\lambda) - \hat{k}(\lambda)|. \end{aligned}$$

Optimizing for  $L_\infty$  of the kernel relative error (17) is equivalent to optimizing weighted  $C_\infty$  error of the kernel signature approximation with weight

$$w_R(\lambda) = \frac{1}{k_h(\lambda)}. \quad (18)$$

While  $w_A$  is decreasing, that is, the fit should be tighter for small  $\lambda$ ,  $w_R$  is increasing and a better fit should be at the tail of the kernel signature.

To apply the proposed approximation method, we need to select the size of the interval  $[-M, M]$  where the kernel signature should be approximated. The optimal choice is  $M = \log(b/m)$ , where  $b$  is the largest expected input value and  $m$  is the smallest non-zero input value. For instance, for histograms with 8bit entries, such as SIFT [10],  $M = \log(255/1)$ .

## 6. Kernels in more dimensions

Measurements in many applications take the form of high dimensional vectors  $\mathbf{x}, \mathbf{y} \in \mathcal{X}^n$ , where  $\mathcal{X}$  is  $\mathbb{R}$  or  $\mathbb{R}_0^+$  depending on the type of the input data. We will use  $x^i$  to denote  $i$ -th component of vector  $\mathbf{x}$ . So far only one dimensional kernels have been considered. These kernels can be

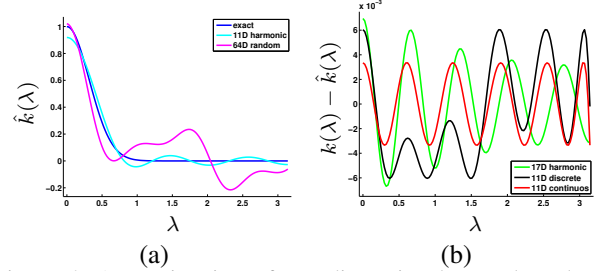


Figure 4. Approximation of one-dimensional RBF kernel: (a) the shape of the approximate kernel signature for 11D feature map via orthogonal projection onto angular harmonics and  $\hat{\Psi}_{\text{RF}}(0)^\top \hat{\Psi}_{\text{RF}}(\lambda)$  for random feature maps [13], (b) approximation error for 17 dimensional feature map via orthogonal projection, the proposed 11 dimensional feature map for the discrete and continuous methods respectively.

extended to higher dimensional input by either additive or multiplicative combination.

Additive kernels are defined as

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n K(x^i, y^i).$$

In computer vision, the following homogeneous additive kernels are commonly used:  $\chi^2$ , intersection, Hellinger, and Jensen-Shannon kernels [18]. The feature map for the additive kernel is a concatenation of feature maps for each dimension. The additive construction of the feature map increases the dimensionality  $D$  of the feature map linearly with the input dimension  $n$ , which is acceptable and we will not study the multi-dimensional additive feature maps further.

Multiplicative kernels, such as multi-dimensional RBF with diagonal  $\Sigma$ , can be written as

$$K(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^n K(x^i, y^i).$$

The feature map is a tensor (Kronecker) product of the feature maps for each input dimension  $\hat{\Psi}(\mathbf{x}) = \otimes \hat{\Psi}(x^i)$ . The construction of the multiplicative kernel, often called modulation, is increasing the dimensionality of the final feature map  $\hat{\Psi}(\mathbf{x})$  exponentially with the number of input dimensions. Therefore, multiplicative kernels constructed

by modulation are suitable for only low-dimensional input data. We will discuss multiplicative kernels in detail in section 6.2.

The proposed method of discrete optimization is not suitable for approximation of kernels with high-dimensional input data. In practice, direct application is tractable for kernels with up to 3 dimensions. Nevertheless, even with this limitation, there are practical applications that would benefit from our approach. Consider problems where low dimensional geometric data, such as the position of a point in a patch and the orientation of the gradient at that point, are to be encoded, *e.g.*, in interest point descriptors such as in [1] or [3]. Two different approaches using different forms of functions  $\hat{k}$  approximating the kernel signature will be considered. The optimization formulation is essentially identical to the one-dimensional case, including the extension exploiting entire continuous spectrum from section 4.2. The difference is in the construction of the frequency pool  $\Omega$  and the discrete evaluation domain  $Z$ . The size of these sets is the bottle neck of the discrete optimization approach, as both sets grow fast with the increasing dimensionality  $n$  of the input data. Two different forms of functions  $\hat{k}$  approximating the kernel signature will be considered. The approaches are compared in section 7.3.

### 6.1. Approximation by projections

A general method directly approximating the  $n$ -dimensional kernel signature uses form of the approximation  $\hat{k}_P$  similar to [13]

$$\hat{k}_P(\boldsymbol{\lambda}) = \sum_{\boldsymbol{\omega} \in \Omega} \alpha_{\boldsymbol{\omega}} \cos(\boldsymbol{\omega}^\top \boldsymbol{\lambda}), \quad (19)$$

where  $\boldsymbol{\lambda} = \mathbf{x} - \mathbf{y} \in \mathbb{R}^n$ ,  $\Omega \subset \mathbb{R}^n$ . Geometrically,  $\hat{k}_P$  can be seen as projecting  $\boldsymbol{\lambda}$  onto  $\boldsymbol{\omega}$  and then encoding the projection by a cosine with frequency  $\|\boldsymbol{\omega}\|$ . Since (19) is only symmetric on each line passing through the origin, that is  $\hat{k}_P(\boldsymbol{\lambda}) = \hat{k}_P(-\boldsymbol{\lambda})$ , the evaluation set has to be constructed as  $Z \subset \prod_{i=1}^{n-1} [-M_i, M_i] \times [0, M_n]$ . The finite pool of frequencies is  $\Omega \subset \mathbb{R}^n$ . The projection method is capable of approximating multi-dimensional kernels, even those that are not multiplicative.

### 6.2. Approximation by modulation

For multiplicative kernels, the following form of  $\hat{k}_M$  is useful

$$\hat{k}_M(\boldsymbol{\lambda}) = \sum_{\boldsymbol{\omega} \in \Omega} \alpha_{\boldsymbol{\omega}} \prod_i^n \cos(\omega^i \lambda^i). \quad (20)$$

Since (20) is symmetric along all axes, that is  $\hat{k}_M((\lambda^1, \dots, \lambda^n)^\top) = \hat{k}_M((\pm\lambda^1, \dots, \pm\lambda^n)^\top)$ , it is sufficient to optimize only in  $Z \subset \prod_{i=1}^n [0, M_i]$ . Let  $\hat{\Psi}^i$  be a feature map optimized separately over the  $i$ -th dimension

from frequencies  $\Omega^i$  and corresponding weights  $\alpha_{\omega^i}$ , and let

$$\Omega_{\otimes} = \{\boldsymbol{\omega} = (\omega^1, \dots, \omega^n) \mid \omega^i \in \Omega^i\}, \quad (21)$$

$$\alpha_{\boldsymbol{\omega}} = \prod_{i=1}^n \alpha_{\omega^i}, \quad D_{\boldsymbol{\omega}} = \prod_{i=1}^n D_{\omega^i} \quad \text{for } \boldsymbol{\omega} = (\omega^1, \dots, \omega^n).$$

The feature map constructed from frequencies  $\Omega_{\otimes}$  and weights  $\alpha_{\boldsymbol{\omega}}$  is equivalent to  $\hat{\Psi}(\mathbf{x}) = \bigotimes \hat{\Psi}(x^i)$ .

The dimensionality of feature map  $\hat{\Psi}$  can be reduced by dropping frequencies  $\boldsymbol{\omega}$  with small coefficient  $\alpha_{\boldsymbol{\omega}}$ . Even though frequencies with small  $\alpha_{\omega^i}$  may be still important for approximation in dimension  $i$ , the product of such weights can exponentially lower the impact. We refer to this greedy method as ' $\bigotimes$  no LP' in the experiments. Better approximation results are achieved by executing the proposed LP optimization. Using the frequency pool  $\Omega_{\otimes}$  (21) significantly increases the speed of the algorithm.

Note that frequencies from the modulation approach can be transformed into frequencies of the projection approach using identity  $\cos(x) \cos(y) = \cos(x+y)/2 + \cos(x-y)/2$ . However, while the left-hand side of the equation represents one entry in the frequency pool  $\Omega$  for modulation, it generates two entries for the projection case. Overall, the projection approach is more general at the cost of larger LP problem (larger in both, the size of  $\Omega$  and in the size of the evaluation set  $Z$ , as discussed in section 6.1). Visual difference in the approximation is shown in Figure 3 (three leftmost plots).

## 7. Experimental comparison

In this section, we evaluate the quality of the proposed feature map construction for an RBF kernel as a representative of aperiodic functions, homogeneous kernels represented by  $\chi^2$  kernel, and two dimensional RBF kernel.

### 7.1. RBF kernel

A number of different feature maps approximating a one-dimensional RBF kernel with  $\sigma = 0.2$  on interval  $x, y \in [0, \pi]$  were compared. Figure 4(a) shows two examples of rather poor approximations of the underlying kernel signature: an orthogonal projection onto a cosine basis with angular frequencies  $\Omega = \{0, \dots, 5\}$  resulting in 11D feature map, and random explicit feature map of 64 dimensions [8]. Since random explicit feature maps are only approximately shift-invariant, the plot shows  $\hat{\Psi}_{\text{RF}}(0)^\top \hat{\Psi}_{\text{RF}}(\lambda)$ . Figure 4(b) shows the values of the absolute error  $k(\lambda) - \hat{k}(\lambda)$  for three comparable feature maps: an orthogonal projection onto a cosine basis with angular frequencies  $\Omega = \{0, \dots, 8\}$  resulting in 17D feature map (labelled 17D harmonic), the 11D feature map by the proposed discrete method (section 4.1), and the 11D feature map by the proposed continuous method (section 4.2). All three approximations would be indistinguishable from the exact  $k$

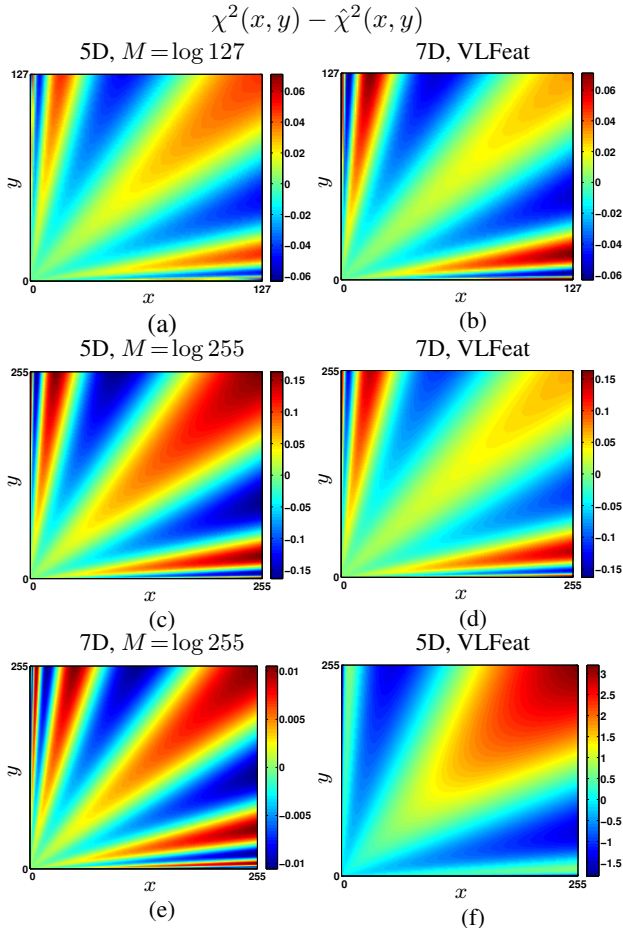


Figure 5. Comparison of the absolute error of the  $\chi^2$  approximation. Left column is the proposed method optimizing (15), right column shows results for Vedaldi [18], VLFeat implementation [17]. The first two rows compare the proposed 5D mapping to 7D mapping of VLFeat on  $x, y \in \{0, \dots, 127\}$  and  $x, y \in \{0, \dots, 255\}$  respectively. The third row shows the error of the 7D proposed mapping and 5D mapping of VLFeat for comparison.

on Figure 4(a). It can be seen, that the proposed continuous method is superior to the orthogonal projection, despite the fact that it uses only 11 dimensions compared to 17 of the orthogonal projection. The continuous method reduces the optimized  $C_\infty$  error function over the discrete method to approximately one half in this case (from  $6.7 \cdot 10^{-3}$  to  $3.3 \cdot 10^{-3}$ ).

## 7.2. Homogeneous kernels

We thoroughly evaluate the proposed method on a  $\chi^2$  kernel approximation  $\chi^2(x, y) = 2xy/(x + y)$ . We compared approximation by the proposed method with the state-of-the-art method of [18] available in VLFeat [17]. Our approach allows the error function on the kernel to be optimized, while the competing method approximates the ker-

	$\chi^2$		intersect		J-S	
	$L_\infty$	RMS	$L_\infty$	RMS	$L_\infty$	RMS
Ours 5D	<b>0.163</b>	<b>0.081</b>	<b>10.922</b>	<b>5.376</b>	<b>0.019</b>	<b>0.009</b>
VLFeat 5D	3.205	1.251	30.119	6.679	2.911	1.203
Ours 7D	<b>0.011</b>	<b>0.005</b>	<b>8.238</b>	<b>4.053</b>	<b><math>9e^{-4}</math></b>	<b><math>3e^{-4}</math></b>
VLFeat 7D	0.143	0.053	22.287	4.436	0.127	0.070

Table 1. Comparison of approximation precision of different homogeneous feature maps. Maximal error  $L_\infty$  and root mean square RMS error are compared on  $x, y \in \{0, \dots, 255\}$ .

nel signature. The comparison on the commonly used homogeneous kernels is summarized in Table 1.

First, we compare the absolute error of the approximation. For this experiment, the  $\varepsilon_A$  error (equation (15)) was minimized for the proposed method. The approximation errors are plotted in Figure 5. The first row compares our 5D feature map to the 7D feature map of VLFeat [17] on input data  $x, y \in \{0, \dots, 127\}$ . Note that the input values can be arbitrarily scaled, only the smallest non-zero ratio of the values is relevant. The kernel signature for the proposed method was optimized on interval  $[-M, M]$ ,  $M = \log 127$ . Even though the proposed method provides a lower dimensional feature map and  $L_\infty$  error was optimized, it outperforms (on this interval) the method VLFeat [17] in  $L_\infty$  error ( $\max_{x,y} |\chi^2(x, y) - \hat{\chi}^2(x, y)|$  : 0.048 vs. 0.071) as well as in  $L_2$  error ( $\sum_{x,y} (\chi^2(x, y) - \hat{\chi}^2(x, y))^2$  : 9.121 vs. 11.272).

The middle row of Figure 5 compares the proposed 5D feature map ( $M = \log 255$ ) and the 7D feature map of VLFeat [17] on input data  $x, y \in \{0, \dots, 255\}$ . The 7D feature map provides a slightly better approximation than the 5D map; however, the error range of the two feature maps is approximately the same. Replacing a 7D feature map by 5D feature map reduces the memory requirements and kernel evaluation time by 28%.

For a full comparison, we have included (bottom row of Figure 5) the proposed 7D feature map ( $M = \log 255$ ) with error order of magnitude lower than the two previously compared feature maps, and also the 5D feature map of VLFeat [17] with an order of magnitude higher error than the two previous feature maps. From this experiment, we see that (1) the proposed approximation outperforms the state-of-the-art feature maps, and (2) the feature map should be optimized for the input domain of particular application.

In the next experiment, we study how the approximation behaves outside the region for which it was optimized. The approximation error for different methods is plotted in Figure 6 for values of the ratio  $x/y$  up to  $1/10^7$ . It can be observed that outside the optimal region, the error of the kernel signature approximation  $|\hat{k}(\lambda) - k(\lambda)|$  increases and thus the error of the kernel  $|\sqrt{xy}[\hat{k}(\log y/x) - k(\log y/x)]$  also increases. Since  $k(\lambda)$  decays and  $\hat{k}(\lambda)$  is bounded,  $|\hat{k}(\lambda) - k(\lambda)|$  is also bounded. As a result, for sufficiently

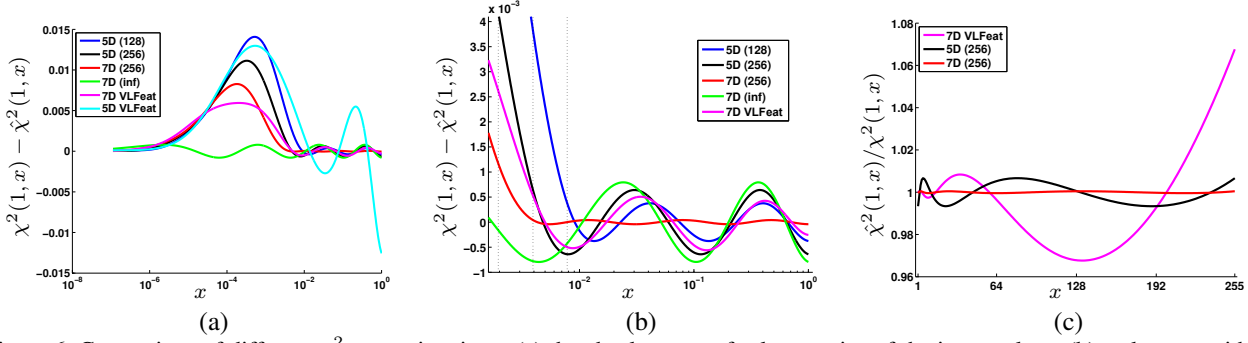


Figure 6. Comparison of different  $\chi^2$  approximations: (a) the absolute error for large ratios of the input values; (b) a close up with three dotted vertical lines at  $1/512$ ,  $1/256$  and  $1/128$  respectively; (c) relative error of the approximation. The error range colour mapping is fixed for first and second row. Logarithm of the number in brackets states the size  $M$  of the interval on which the kernel signatures were approximated.

large  $|\lambda|$  the error of the kernel is dominated by  $\sqrt{xy}$  and approaches zero. In Figure 6, the green curve corresponds to a kernel signature that has been optimized for a large enough interval so that the upper bound on the error is less than the optimized  $L_\infty$  inside the interval, thus having optimal error bound everywhere. This is only possible for error measures, such as  $\varepsilon_A$  (15) with decreasing weight  $w(\lambda)$  (16).

The last experiment with  $\chi^2$  kernel considers the relative error  $\varepsilon_R$  of the kernel fit (17). Three feature maps are compared: proposed 5D and 7D constructed to minimize the relative error, and 7D by VLFeat [17]. Plot in Figure 6 (c) show that the proposed method significantly outperforms its competitor.

### 7.3. Symmetric RBF kernel in 2D

Four methods approximating the symmetric 2D RBF kernel with  $\sigma = 0.2$  with kernel input variables  $x, y \in [0, \pi]^2$  were compared: two using the projection method described in section 6.1, and two using the modulation method (section 6.2). For the projection method, two different initialization of the frequency pool  $\Omega$  were used: a general initialization by discretization of angle and frequency  $\|\omega\|$  of  $\omega$ , referred to as 'Proj'; and by frequencies equivalent to  $\Omega_\otimes$  (21), obtained as a Cartesian product  $\Omega_\otimes = \Omega_{1D} \times \Omega_{1D}$ , where  $\Omega_{1D}$  correspond the the 11D feature map form section 7.1 (referred to as 'Proj  $\otimes$ '). For both methods, full LP optimization on 2D input, including the continuous extension was executed.

Both modulation methods (section 6.2) were initialized by  $\Omega_\otimes$ . One method (' $\otimes$ ') exploits the full LP optimization on 2D input, including the continuous extension. For the last method (' $\otimes$  no LP') the feature map is selected greedily based on the estimate  $\alpha_\omega = \alpha_{\omega^1} \alpha_{\omega^2}$ . The quantitative result are summarized in Table 2, and the qualitative results of approximations by 31 dimensional feature maps is shown in Figure 3 leftmost column.

The fastest approach ' $\otimes$  no LP' is the least precise. The

$D(\hat{k})$	Proj	Proj $\otimes$	$\otimes$	$\otimes$ no LP
31	0.1199	0.0984	0.1061	0.2370
73	0.0292	0.0148	0.0137	0.0439
101	0.0092	0.0038	0.0038	0.0122

Table 2. Comparison of the  $L_\infty$  error of 2D RBF kernel approximation.

most general approach 'Proj' is the slowest and performs slightly worse than the two approaches initialized by results of 1D optimization.

We made the following observations for the modulation methods: (1) the linear program selects different components than the greedy approach, (2) after the continuous extension, the frequencies in  $\Omega$  are no longer on a grid, originally defined by the Cartesian product  $\Omega^1 \times \Omega^2$ .

Finally, the comparison with other methods show in Figure 3 demonstrates that any of the proposed methods is superior to existing methods.

## 8. Conclusions

A novel method of data independent construction of low dimensional feature maps was proposed. The problem is cast as a linear program that jointly considers competing objectives: quality of the approximation and the dimensionality of the feature map. The proposed discrete optimization exploits the entire continuous spectrum of frequencies and achieves considerably better approximations with feature maps of the same dimensionality or equally good approximations with lower dimensional feature maps compared with the state-of-the-art methods. It was also demonstrated that the proposed method allows for optimization of meaningful errors measured on the homogeneous kernel output, rather than solely approximating the kernel signature.

Any application that uses explicit features maps would benefit from the results of this paper. The code is available [4].



## References

- [1] L. Bo, X. Ren, and D. Fox. Kernel descriptors for visual recognition. In *NIPS*, pages 244–252, 2010. 1, 6
- [2] L. Bo and C. Sminchisescu. Efficient match kernel between sets of features for visual recognition. In *NIPS*, pages 135–143, 2009. 1
- [3] A. Bursuc, G. Tolias, and H. Jégou. Kernel local descriptors with implicit rotation matching. In *ICMR*, 2015. 6
- [4] O. Chum. Implementation of low dimensional explicit feature maps. <http://cmp.felk.cvut.cz/~chum/code/ld-efm.html>. 8
- [5] H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *CVPR*, 2010. 1
- [6] T. Joachims. Training linear SVMs in linear time. In *KDD*, pages 217–226. ACM, 2006. 1
- [7] Q. Le, T. Sarlós, and A. Smola. Fastfood-approximating kernel expansions in loglinear time. In *ICML*, 2013. 5
- [8] F. Li, C. Ionescu, and C. Sminchisescu. Randfeat: Random fourier approximations for skewed multiplicative histogram kernels, release 1. <http://sminchisescu.ins.uni-bonn.de/code/randfeat/randfeat-release1.tar.gz>. 6
- [9] F. Li, C. Ionescu, and C. Sminchisescu. Random fourier approximations for skewed multiplicative histogram kernels. In *DAGM*, 2010. 1
- [10] D. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004. 5
- [11] S. Maji and A. C. Berg. Max-margin additive classifiers for detection. In *ICCV*, pages 40–47, 2009. 1
- [12] F. Perronnin, Y. Liu, J. Sanchez, and H. Poirier. Large-scale image retrieval with compressed fisher vectors. In *CVPR*, 2010. 1
- [13] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *NIPS*, pages 1177–1184, 2007. 1, 5, 6
- [14] B. Scholkopf and A. Smola. *Learning with Kernels*. MIT Press, 2002. 1
- [15] G. Tolias, T. Furon, and H. Jégou. Orientation covariant aggregation of local descriptors with embeddings. In *ECCV*, 2014. 1
- [16] R. J. Vanderbei. *Linear Programming: Foundations and Extensions*. Kluwer Academic Publishers, Boston, 1996. 2
- [17] A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/>, 2008. 1, 7, 8
- [18] A. Vedaldi and A. Zisserman. Efficient additive kernels via explicit feature maps. *TPAMI*, 34(3), 2012. 1, 3, 4, 5, 7
- [19] A. Vedaldi and A. Zisserman. Sparse kernel approximations for efficient classification and detection. In *CVPR*, 2012. 1