# **Learning Random Kernel Approximations for Object Recognition**

Eduard Gabriel Băzăvan <sup>1</sup>

Fuxin Li<sup>2</sup>

Cristian Sminchisescu<sup>2,1</sup>

eduard.bazavan@imar.ro

fli@cc.gatech.edu

cristian.sminchisescu@ins.uni-bonn.de

<sup>1</sup>Institute of Mathematics of the Romanian Academy
<sup>2</sup>Faculty of Mathematics and Natural Science, University of Bonn
<sup>3</sup>Georgia Institute of Technology

## **Abstract**

Approximations based on random Fourier features have recently emerged as an efficient and formally consistent methodology to design large-scale kernel machines [24]. By expressing the kernel as a Fourier expansion, features are generated based on a finite set of random basis projections, sampled from the Fourier transform of the kernel, with inner products that are Monte Carlo approximations of the original kernel. Based on the observation that different kernel-induced Fourier sampling distributions correspond to different kernel parameters, we show that an optimization process in the Fourier domain can be used to identify the different frequency bands that are useful for prediction on training data. Moreover, the application of group Lasso [37] to random feature vectors corresponding to a linear combination of multiple kernels, leads to efficient and scalable reformulations of the standard multiple kernel learning model [33]. In this paper we develop the linear Fourier approximation methodology for both single and multiple gradient-based kernel learning and show that it produces fast and accurate predictors on a complex dataset such as the Visual Object Challenge 2011 (VOC2011).

## 1. Introduction

The proper choice of kernel function and its hyperparameters are crucial to the success of applying kernel methods to practical applications. These selections span a number of different problems: from choosing a single width parameter in radial basis kernels, scaling different feature dimensions with different weights [6], to learning a linear or a nonlinear combination of multiple kernels (MKL) [17]. In complicated practical problems such as computer vision, sometimes the need of multiple kernels arises naturally [34, 14]. Images can be represented using descriptors based on shape, color and texture, and these descriptors have different roles

in classifying different categories. In such situations it is in principle easy to design a kernel classifier where each kernel represents one of the descriptors and the classifier would be based on a weighted combination of kernels with category dependent learnt weights.

A natural difficulty in kernel learning is scalability. Kernel methods scale with an already mediocre time complexity of at least  $O(N^{2.3})$  with respect to the size N of the training set, but combining multiple kernels or learning the hyperparameters of a single kernel significantly slows down training. Some speed-ups apply for specific kernels, but only in limited scenarios [22]. In consequence, most of the kernel learning approaches so far are only capable to handle a few thousand training examples at most. This is insufficient for the current age of massive datasets, such as a 11 million images ImageNet, or the 19 million articles within Wikipedia.

An emerging technique that can in principle speed up the costly kernel method while at the same time preserving its non-linear predictive power is the random Fourier feature methodology (RFF) [24, 35, 19]. By sampling components from the frequency space of the kernel using Monte Carlo methods, RFF obtains a bounded, approximate representation of a kernel embedding that may initially span an infinite-dimensional space. Many operations are simplified once such representation is available, the most notable being that any kernel learning algorithm would now scale as O(N), where N is the number of examples. This opens the path for applying kernel methods to the realm of massive datasets.

In the seminal work on random Fourier features [24], the methodology was developed primarily for radial basis kernels. Recent work [35, 19] focused on extending this technique to a number of other useful kernels defined on histogram features (empirical estimates of multinomial distributions), such as the  $\chi^2$  and histogram intersection measures. However, the potential of the linear random

Fourier methodology for kernel learning remains largely unexplored. In this paper we develop the methodology for learning both single kernel and for multiple kernel combinations in the Fourier domain and show that these produce accurate and efficient models. We conduct experiments in visual object recognition, using the difficult PASCAL Visual Object Challenges 2011 dataset, in order to demonstrate the performance of the proposed Fourier kernel learning methodology and compare against non-linear learning algorithms, designed to operate in the original kernel space.

#### 2. Related work

Approaches to kernel learning can be broadly classified into methods that estimate the hyper-parameters of a single kernel[25, 6, 15] and methods that learn the weights of a linear combinations of kernels and possibly their hyper-parameters – the so-called multiple kernel learning framework or MKL[31, 2, 16, 36].

A popular approach for single kernel learning is the gradient-based method pursued by Chapelle *et al.* [6, 15]. Keerthi *et al.* [15] give an efficient algorithm that alternates between learning an SVM and optimizing the feature weights. Cortes *et al.* [8] propose a two-stage method based on a modification of a classical kernel alignment [9] metric and prove a number of learning bounds. Approaches to single kernel learning under a kernel prior have been pursued both as semi-definite programs[17] and within unconstrained optimization formulations[11], although in both cases the optimization is involved and complexity is an issue. More recent methods attempt to track the entire regularization path in the non-linear case [27, 20].

Multiple kernel learning provides a powerful conceptual framework for both model combination and model selection and has attracted significant research recently. Initial approaches originating with work by Lanckriet et al. [17] estimated a linear combination of kernels using semi-definite programming. This was reformulated by Bach et al. [3] as block-norm regularization, reducing the optimization problem to a second order cone program applicable to medium scale problems. More recent methods [7] learn a polynomial combination of kernels. A number of approaches have pursued different  $l_p$ -norms for kernel selection [31, 16, 36]. Hierarchical kernel learning approaches like (HKL) [2] perform feature selection combinatorially, by choosing kernel combinations obtained by mapping to a directed acyclic graph. The search for such combinations can then be performed in polynomial time.

A difficulty with many single or multiple kernel learning formulations is their relatively unfavorable scaling properties. When kernels are used, such methods usually scale at least quadratically with the number of examples. Sometimes scaling with the number of kernel parameters and the number of kernels can also become an issue. A formulation

of kernel learning within a linear Fourier framework, as pursued in this paper, carries the promise of better scalability and wider applicability to large datasets, while at the same time preserving the non-linear predictive power that makes kernel methods attractive.

## 3. Learning Kernels Parameters

## 3.1. Random Fourier Approximations

Kernels offer substantial power and flexibility to predictive models. Central to their use is the 'kernel trick' which allows the design of algorithms that depend only on the inner product of their arguments. This is based on the property of positive definite kernel functions  $k(\cdot, \cdot)$  to define a dot product and a lifting  $\phi$  so that the dot product between lifted data points can be efficiently computed as  $\phi(\mathbf{x})^{\top}\phi(\mathbf{y}) = k(\mathbf{x}, \mathbf{y})$ . Algorithms become linear in a complex feature space induced by  $\phi$ , but access the data only through evaluations of the kernel k in input space. This makes possible to handle complex or infinite dimensional feature space mappings  $\phi$ , and indeed these usually give the best results in visual recognition datasets [13, 18]. However the very advantage that made kernel methods popular the kernel trick, requires the manipulation of matrices of pairwise examples. For large datasets, the scaling of kernel methods is at least quadratic in the number of examples. This makes their direct usage impractical beyond datasets of  $10^5$  elements. The methodology we pursue is to approximate the kernel function linearly using random feature maps.

Key to the random Fourier methodology is Bochner's theorem, that connects positive definite kernels and their Fourier transforms [28, 24, 19]. For positive definite translation-invariant kernels  $k(\mathbf{x},\mathbf{y}) = k(\mathbf{x} - \mathbf{y})$  on  $\mathbb{R}^m$ , Bochner's theorem guarantees that every kernel is the inverse Fourier transform of a proper probability distribution  $\mu$ . Defining  $\zeta(\mathbf{x}) = e^{j\mathbf{x}^\top\gamma}$  (with j the imaginary unit), the following equality holds:

$$k_{\sigma}(\mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^m} e^{j(\mathbf{x} - \mathbf{y})^{\top} \gamma} d\mu_k(\gamma)$$
$$= \mathbf{E}_{\mu} [\zeta_{\gamma}(\mathbf{x}) \zeta_{\gamma}(\mathbf{y})^*]$$
$$\approx \phi_{\Gamma}(\mathbf{x})^{\top} \phi_{\Gamma}(\mathbf{y})$$

where \* is the (complex) conjugate, and

$$\phi_{\Gamma}(\mathbf{x}) = \sqrt{\frac{2}{d}} \left[ \cos \left( \mathbf{x}^{\top} \boldsymbol{\gamma}_i + 2\pi b_i \right) \right]_{i=1,d}$$
 (1)

is the random feature map at frequencies  $\Gamma = \{\gamma_1, \ldots, \gamma_d\}$ , with  $b \sim \mathrm{U}[0, 2\pi]$ , the uniform distribution in the interval  $[0, 2\pi]$ . We approximate the expectation  $\mathbf{E}_{\mu}[\zeta_{\gamma}(\mathbf{x})\zeta_{\gamma}(\mathbf{y})^*]$  by means of a Monte-Carlo sample

drawn from the distribution  $\mu_k$ . This is an operation on linear functions with explicit features. The algorithm for the change of representation has two steps: i) Generate d random samples  $\Gamma$  from the distribution  $\mu_k$ ; ii) Compute the random projection  $\phi_{\Gamma}$  using (1), for all training examples. The approximation has the convergence rate of Monte Carlo methods,  $O(d^{-1/2})$ , dependent on the random sample size, but independent of the input dimension. One usually needs up to a few thousand dimensions to approximate the original kernel accurately. Datasets containing hundreds of thousands of examples can be trained in a few hours, for simpler models. This motivates our interest in advancing the fundamental theory and practice of kernel learning for an efficient class of approximations with randomized feature dimension, with linear scaling in the number of examples, and with predictable convergence rates.

## 3.2. Single Kernel Learning

First we consider learning strategies for the hyper-parameters  $\sigma$  of the kernel k based on the Fourier feature methodology. To achieve this goal we need a learning criterion and an algorithm for optimizing the hyper-parameters. Since model training is much faster due to the linear dependence on the training set size, we can perform more complex parameter learning than classic approaches that usually rely on a grid search procedure, like cross-validation. This efficiency of the linear formulation also allows us to scale parameter learning towards numbers of parameters which are unattainable with classical methods. Our approach to kernel learning will optimize the hyper-parameters with respect to the error on a held-out validation set, for models obtained on the training set. This has been shown to be a viable procedure to prevent overfitting[13].

In the sequel we denote X and y the matrix of inputs or covariates (row-wise organized) and y the matrix of targets on the training set, respectively, whereas U and v are the validation inputs and targets, respectively. We use  $\phi_{\Gamma}$ , introduced earlier, as our random Fourier feature map but we will drop  $\Gamma$  to simplify notation. We define  $\phi(X)$  to be the feature map applied to all the rows of the matrix X.

We will learn the hyper-parameters of a kernel ridge regression model  $\beta \equiv \beta(\sigma)$  (other margin-based training costs, such as hinge loss, logistic loss, etc., can be similarly adopted into the framework)

$$\min_{\beta} \frac{1}{2} \| \boldsymbol{\phi}(\mathbf{X}) \boldsymbol{\beta} - \mathbf{y} \|_{2}^{2} + \lambda \| \boldsymbol{\beta} \|_{2}^{2}. \tag{2}$$

For this problem the optimum can be obtained in closed form as

$$\boldsymbol{\beta} = \left(\boldsymbol{\phi}(\mathbf{X})^{\top} \boldsymbol{\phi}(\mathbf{X}) + \lambda \mathbf{I}_{d}\right)^{-1} \boldsymbol{\phi}(\mathbf{X})^{\top} \mathbf{y}. \tag{3}$$

To learn the hyper-parameters  $\sigma$  we optimize the squared  $l_2$  validation error. Given  $f = \phi(\mathbf{U})^T \boldsymbol{\beta} - \mathbf{v}$  and  $r(\boldsymbol{\sigma})$  a

regularizer for the hyper-parameters, e.g.  $r(\sigma) = \|\sigma\|_2^2$  we optimize

$$\min_{\boldsymbol{\sigma}} \|f\|_2^2 + r(\boldsymbol{\sigma}) \tag{4}$$

Note that we can easily obtain a gradient with regard to the kernel parameters for this optimization. The random feature map  $\phi$  establishes a connection between the original input representation  $\mathbf{X}$  and the approximation of its lifting  $\phi(\mathbf{X})$ . This can be used to derive analytical expressions not only for kernel expansions in the input space but also for the gradient with regard to the kernel parameters. We can minimize the loss using a local-descent based optimization.

Given  $\sigma_i$  the *i*'th dimension of the kernel parameter vector  $\sigma$ , then

$$\frac{\partial \|f\|_{2}^{2}}{\partial \sigma_{i}} = \operatorname{Tr}\left[\left(\frac{\partial \|f\|_{2}^{2}}{\partial f}\right)^{\top} \frac{\partial f}{\partial \sigma_{i}}\right]$$

$$= f^{\top} \frac{\partial f}{\partial \sigma_{i}} = f^{\top} \left(\frac{\partial \phi(\mathbf{U})}{\partial \sigma_{i}} \boldsymbol{\beta} + \phi(\mathbf{U}) \frac{\partial \boldsymbol{\beta}}{\partial \sigma_{i}}\right)$$
(5)

For translation invariant kernels, we can easily compute  $\phi_{\Gamma}$ . Manipulating the sampling distribution  $\mu_k$  associated with kernel k in an optimization process identifies the different frequency bands that are useful for prediction on training data. This effectively leads to a posterior frequency distribution, given  $\mu_k$  as a prior. Generating samples from this posterior and optimizing with respect to their expectation on the training set gives a direct way of learning the parameters of the kernel. However as the kernel hyper-parameters  $\sigma$  change,  $\Gamma$  needs to be re-sampled. Although this is feasible, changing the sampling distribution introduces conceptual difficulties during optimization as the objective function change, and sampling becomes a source of noise and non-smoothness for the optimization. Importance sampling can be used to avoid resampling at each step, but if the parametrized frequency distribution drifts far away from the starting distribution, then the original samples will have little importance weights, and resampling would potentially be needed nevertheless for convergence. Such difficulties can be overcome for several interesting kernels that belong to a class of functions where the samples from  $\mu_k$  can be written as

$$\gamma = \boldsymbol{\sigma} \cdot h(\boldsymbol{\omega}) \tag{6}$$

where h is the quantile function,  $\omega$  are uniformly sampled and fixed, and  $\cdot$  is the Hadamard product of two vectors. In this case, throughout the entire optimization process, sampling needs to be done only once from the uniform distribution for  $\omega$ . When  $\sigma$  is changed, samples from the new distribution are generated automatically from  $\omega$ . Examples of such kernels are the *Gaussian*, the *generalized skewed* $\chi_2$  and the *generalized skewed intersection*[19] (see table 1).

kernel	parametric form (k)	probability distribution $(\mu)$	quantile $(\gamma)$	
gaussian	$\frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{\ x-y\ ^2}{2\sigma^2}}$	$\frac{\sigma}{\sqrt{2\pi}}e^{-\frac{\sigma^2\omega}{2}}$	$\sigma\sqrt{2}\mathrm{erf}^{-1}(2u-1)$	
skewed $-\chi_2$	$\frac{2(x+c)^{\sigma}(y+c)^{\sigma}}{(x+c)^{2\sigma}+(y+c)^{2\sigma}}$	$\frac{1}{2\sigma}\operatorname{sech}\left(\frac{\pi\omega}{2\sigma}\right)$	$\sigma \frac{2}{\pi} \log \left( \tan \left( \frac{\pi}{2} u \right) \right)$	
skewed intersection	$\min \left\{ \frac{(x+c)^{\sigma}}{(y+c)^{\sigma}}, \frac{(y+c)^{\sigma}}{(x+c)^{\sigma}} \right\}$	$\frac{1}{\pi} \frac{\sigma}{\sigma^2 + \omega^2}$	$\sigma \tan \left(\pi (u - \frac{1}{2})\right)$	

Table 1. Examples of kernels that can be efficiently optimized within our framework, presented for the unidimensional case. For the multidimensional case, the Fourier methodology decomposes, hence it requires a simple multiplication of parameters corresponding to each dimension. In the above the random variable u is drawn from the uniform distribution.

Based on this property, by differentiating the feature map (1) we obtain

$$\frac{\partial \phi(\mathbf{u}_k)}{\partial \sigma_i} = \sqrt{\frac{2}{d}} \left[ \frac{\partial \cos(\mathbf{u}_k^{\top} (\boldsymbol{\sigma} \cdot h(\boldsymbol{\omega}_j)) + 2\pi b_j)}{\partial \sigma_i} \right]_{j=1,d}$$

$$= \sqrt{\frac{2}{d}} \left[ -\mathbf{u}_{k,i}^{\top} h(\boldsymbol{\omega}_{j,i}) \sin(\mathbf{u}_k^{\top} (\boldsymbol{\sigma} \cdot h(\boldsymbol{\omega}_j)) + 2\pi b_j) \right]_{i=1,d}$$
(7)

To compute the second term  $\frac{\partial \boldsymbol{\beta}}{\partial \sigma_i}$ , we first define the matrix  $\mathbf{Q} = \phi(\mathbf{X})^{\top} \phi(\mathbf{X}) + \lambda \mathbf{I}_d$  and using the standard result  $\frac{\partial \mathbf{Q}^{-1}}{\partial \sigma_i} = -\mathbf{Q}^{-1} \frac{\partial \mathbf{Q}}{\partial \sigma_i} \mathbf{Q}^{-1}$ , we obtain

$$\frac{\partial \boldsymbol{\beta}}{\partial \sigma_i} = \frac{\partial \mathbf{Q}^{-1}}{\partial \sigma_i} \boldsymbol{\phi}(\mathbf{X})^{\top} \mathbf{y} + \mathbf{Q}^{-1} \left( \frac{\partial \boldsymbol{\phi}(\mathbf{X})}{\partial \sigma_i} \right)^{\top} \mathbf{y} \tag{8}$$

$$= -\mathbf{Q}^{-1} \frac{\partial \mathbf{Q}}{\partial \sigma_i} \mathbf{Q}^{-1} \boldsymbol{\phi}(\mathbf{X})^{\top} \mathbf{y} + \mathbf{Q}^{-1} \left( \frac{\partial \boldsymbol{\phi}(\mathbf{X})}{\partial \sigma_i} \right)^{\top} \mathbf{y}$$

It is easy to see that  $\frac{\partial \phi(\mathbf{X})}{\partial \sigma_i}$  can be computed in the same way as in equation (7), as the gradient of  $\mathbf{Q}$  with respect to  $\sigma_i$  can be obtained as

$$\frac{\partial \mathbf{Q}}{\partial \sigma_i} = \left(\frac{\partial \phi(\mathbf{X})}{\partial \sigma_i}\right)^{\top} \phi(\mathbf{X}) + \phi(\mathbf{X})^{\top} \frac{\partial \phi(\mathbf{X})}{\partial \sigma_i}.$$
 (9)

Computing the gradient of  $r(\sigma)$  will depend on the type of regularization chosen. In general this is a smooth function of  $\sigma$  (e.g.  $l_2$  norm) so it will be straightforward to compute  $\frac{\partial r(\sigma)}{\partial \sigma_i}$ .

Now that we have a closed form formula for the gradient with respect to all kernel parameters in the Fourier domain we can plug it into a non-linear optimizer and estimate the parameters.

This overall philosophy bears a certain similarity to the objective introduced in [6] where the authors use a gradient descent learning technique for kernel ridge regression based on the exact kernel matrix. It is also similiar to the multiple kernel learning technique for products of kernels introduced in [33]. However, besides the technical differences that are introduced by the usage of a Fourier embedding map, an important advantage for our methodology is that we do not have to store the kernel matrices in memory

which has  $O(N^2)$  memory cost. Our memory requirement is just  $O(Nd+d^2)$  where d is the size of the Fourier embedding. The computational complexity of our method is dominated by  $O(i_{skl}(Nd^2+d^3+rNd))$  where d is the size of the random Fourier features, N is the number of training samples, r is the number of parameters and  $i_{skl}$  is the number of iterations to convergence.  $O(Nd^2+d^3)$  is the cost of computing matrix  $\mathbf{Q}$  and inverting it.

## 3.3. Multiple Kernel Learning

Previous work has proven an underlying connection between multiple kernel learning and group Lasso[1]. Although this is an interesting property, it has found relatively limited practical applications so far. In this section we show that by using the random Fourier methodology, we can do just the opposite and lift group Lasso to the kernel domain. We propose a multiple kernel learning formulation where the features are initially transformed using the random Fourier framework, concatenated, and group Lasso is applied to them (RFF-GL). We prove that this new formulation is equivalent with the multiple kernel learning formulation introduced in [33] and then compare the two methodologies. Experiments show that both approaches have similar performance. In contrast, group Lasso based on random Fourier features runs faster and scales better since we do not need to compute or store the Gramm matrices associated with the features.

Let X be the input matrix for the training set organized row-wise, y the associated targets and  $\{k_1,\ldots,k_r\}$  a set of kernels. For example X could be a concatenation of multiple image features such as SIFT or HOG and therefore should be represented as  $\{X_1,\ldots,X_r\}$ . But to simplify the notation we would not refer to the individial components and will just use X when we refer to the input. We denote  $\mathcal{F}_i$  the matrix of random Fourier features obtained from approximating  $k_i$  on inputs X. We concatenate  $\mathcal{F}_i$  columnwise to obtain a matrix  $\mathcal{F}$  on which we can apply group Lasso with r non-overlapping groups, each corresponding to the Fourier embedding of a different kernel. We use the  $l_1$ - $l_2$  formulation [37] which can be written as:

$$\min_{\mathbf{w}} \lambda \sum_{i=1}^{r} \|\mathbf{w}_{\mathcal{F}_i}\|_2 + l(\mathbf{y}, \mathcal{F}\mathbf{w}), \tag{10}$$

where  $\mathbf{w}_{\mathcal{F}_i}$  are the weights applied to the features  $\mathcal{F}_i$  and l(y, f(x)) is a loss function, which can be a quadratic loss or an approximation for the  $\epsilon$ -insensitive regression loss. In our experiments the optimization was performed with a group Lasso solver [21] which was adapted for our specific loss functions.

We compare the above formulation with the standard multiple kernel learning method (GMKL) presented in [33] where we consider the regression problem under  $l_1$  regularization. For GMKL we have to build the Gram matrices  $\{\mathbf{K}_i\}_{i=1..r}$  for each kernel  $k_i$  on  $\mathbf{X}$ . We focus on the regression problem and define  $\mathbf{K_d} = \sum_{i=1..r} \mathbf{d}_i \mathbf{K}_i$ . The al-

gorithm will output the optimized values for d and a set of support vectors and both these elements will be leveraged in the model used for testing.

#### 3.3.1 Proof of equivalence

We now show that GMKL is equivalent as an optimization procedure with RFF-GL. For GMKL we have the following primal problem

$$\min_{\mathbf{w}, \mathbf{d}} \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} + C \sum_{i=1}^{n} l(y_i, \psi(\mathbf{x}_i)^{\top} \mathbf{w}) + \sum_{t=1}^{r} d_t$$
subject to  $\mathbf{d} \ge 0$  (11)

where

$$\psi(\mathbf{x}_i) = [\sqrt{d_1}\psi_1(\mathbf{x}_i)^\top, \dots, \sqrt{d_r}\psi_r(\mathbf{x}_i)^\top]^\top$$
 (12)

and we defined  $\psi_t(\mathbf{x})$  as the feature embedding associated to kernel  $k_t$ .

From the Representer Theorem [29] we now show that the solution for w will be a linear combination of the basis functions

$$\mathbf{w} = \sum_{i=1}^{n} \alpha_i \psi(\mathbf{x}_i)$$

$$= \left[ \sqrt{d_1} \mathbf{w}_1^{\top}, \dots, \sqrt{d_k} \mathbf{w}_k^{\top} \right]^{\top}$$
(13)

where we have dropped the bias term for simplicity. We can rewrite the above optimization problem

$$\min_{\mathbf{w}, \mathbf{d}} \sum_{t=1}^{k} \frac{d_t}{2} \|\mathbf{w}_t\|_2^2 + C \sum_{i=1}^{n} l(y_i, \sum_{t=1}^{k} d_t \psi_t(\mathbf{x}_i)^\top \mathbf{w}_t) + \sum_{t=1}^{k} d_t$$
subject to  $\mathbf{d} \ge 0$  (14)

Following a standard trick from the multiple kernel learning literature[1] we make the substitutions  $\mathbf{w}_t \to d_t \mathbf{w}_t$  to obtain

$$\min_{\mathbf{w}, \mathbf{d}} \sum_{t=1}^{k} \frac{\|\mathbf{w}_t\|_2^2}{2d_l} + C \sum_{i=1}^{n} l(y_i, \sum_{t=1}^{k} \psi_t(\mathbf{x}_i)^\top \mathbf{w}_t) + \sum_{t=1}^{k} d_t$$
subject to  $\mathbf{d} \ge 0$  (15)

It can easily be shown that

$$\frac{1}{2} \frac{\|\mathbf{w}_t\|_2^2}{d_t} + d_t \ge \sqrt{2} \|\mathbf{w}_t\|_2 \tag{16}$$

and we observe that the loss function no longer depends on d. Therefore the primal multiple kernel learning formulation is equivalent to the following group Lasso formulation

$$\min_{\mathbf{w}} \lambda \sum_{t=1}^{k} \|\mathbf{w}_t\|_2 + \sum_{i=1}^{n} l(y_i, \sum_{t=1}^{k} \psi_t(\mathbf{x}_i)^{\top} \mathbf{w}_l)$$
 (17)

where we set  $\lambda = \frac{\sqrt{2}}{C}$ . The equality happens when  $d_t = \frac{1}{\sqrt{2}} \|\mathbf{w}_t\|_2$ . This approach is suitable for other types of regularization for the parameters  $\mathbf{d}$  as well, not only the  $l_1$  norm which we have used. The only restriction is that we should be able to find a closed form solution for (16).

Now we consider the loss function which has to be differentiable in order to be suitable for a group Lasso formulation. The  $\epsilon$ -insensitive loss ( $\epsilon$ -IL) used in standard support vector regression is not differentiable but we can use a smooth approximation instead. In our case we have selected an  $\epsilon$ -insensitive  $\gamma$  logistic loss ( $\epsilon$ -IGLL) function defined as [10, 26]

$$l_{\gamma,\epsilon}(y_i, f(\mathbf{x})) = \frac{1}{\gamma} \log \left( 1 + e^{\gamma(f(\mathbf{x}) - y_i - \epsilon)} \right)$$

$$+ \frac{1}{\gamma} \log \left( 1 + e^{\gamma(-f(\mathbf{x}) + y_i - \epsilon)} \right)$$

$$- \frac{2}{\gamma} \log (1 + e^{-\gamma \epsilon})$$
(18)

We could also consider quadratic or logistic losses for solving the group Lasso.

## 3.3.2 Computational complexity

We assume an average time complexity for a support vector machine algorithm [5] to be  $O(N_{sv}^3)$  where  $N_{sv}$  is the number of support vectors. Since GMKL is based on several evaluations of the standard SVM solver we can show the time complexity of the multiple kernel learning framework is  $O\left(rN^2m^2+i_{GMKL}(N^2r+N_{sv}^3)+rN_{sv}Nm^2\right)$  where r is the number of kernels, N is the number of training samples, m is the size of the training input (we assume it is the same for all data for simplicity) and  $i_{GMKL}$  is the

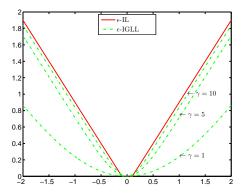


Figure 1. The  $\epsilon$ -insensitive  $\gamma$  loss function we use for optimization provides a good approximation to the exact  $\epsilon$ -insensitive loss. Here we show values for  $\epsilon = 0.1$  and  $\gamma = 1, 5, 10$ 

maximum number of calls to and SVM solver. The complexity is dominated by the term  $O(rN^2m^2)$  which is the cost of computing the kernel matrices.

On the other hand for the group Lasso formulation we have complexity  $O(rNm^2d+i_{gl}Ndr)$  where d is the size of the random Fourier features and  $i_{gl}$  is the number of iterations required for group Lasso. If we use an approximated kernel as in [30] another  $(2p+1)^2$  cost is added to the preprocessing step. The complexity becomes  $O(rNm^2(2p+1)^2d)$  where 2p+1 is the dimension of the approximated kernel as discussed in [35]. The memory complexity is O(Nrd) which is definitely smaller than  $O(rN^2)$  for GMKL.

We see that our algorithm scales linearly, whilst the non-linear kernel method is quadratic. The constants matter because for a small number of samples our method could run slower given that d is in the range of thousands,  $i_{GMKL}$  is usually  $10^2$  and  $i_{al}$  is in the order of d.

## 4. Experiments

We present experiments for single kernel learning and multiple kernel learning, comparing the random linear Fourier methodology with its non-linear kernel counterparts both in terms of running times and in terms of accuracy.

We have experimented with the PASCAL VOC2011 segmentation challenge. This consists of 2223 images for training with ground truth split into halves for training and validation. Another 1111 images are given for testing without ground truth. Following the standard procedure we have trained the methods on the training set (further split, internally into training and validation for kernel hyper-parameter learning) and tested on the PASCAL VOC2011 validation set. We have used the methodology presented in Li *et al.* [18] and relied on a segmentation algorithm from Carreira and Sminchisescu [4] to filter the initial pool of segments to around 100 segments per image. For each of these segments

we extracted 8 types of features among which two bag of visual words for color SIFT [32] and two dense gray scale SIFT descriptors one on each segment and one on the background of each segment, three types of phog descriptors two on the pb edges given by [12] computed at different scales, one on the contour and one on the foreground. The third phog descriptor uses no pb. The last descriptor is a pyramid of locally binary pattern features [23] for texture classification. Because the number of segments (around  $10^5$ ) was still too large for any kernel support vector based algorithm to cope with, we have chosen up to  $10^4$  segments for each class. The segments were chosen based on their individual scores and we balanced the examples to have a fair split between positive and negative segments. More details on how the scores are defined and computed could be found in [18].

## 4.1. Single kernel learning

We ran a set of experiments for our random Fourier features single kernel learning technique (RFF-SKL) and compared in terms of accuracy with the single kernel learning technique introduced by Chapelle  $et\ al.\ [6]$  (KRR-GD). We want to predict the class for more than  $10^5$  segments. We expect from RFF-SKL to give comparable results in terms of accuracy to KRR-GD. Results are shown in Table 2.

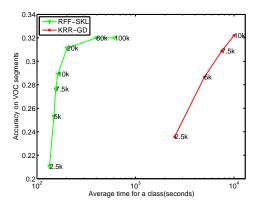


Figure 2. Time versus accuracy for classifying the VOC segments for RFF-SKL and KRR-GD. We varied the number of training points from  $10^3$  up to  $10^5$  and measured the average time for a class, for 20 classes. For KRR-GD it was not feasible to go beyond  $10^4$  training samples. x axis is log scale.

We vary the size of the training set from  $10^3$  to  $10^5$  and measured both the running times of RFF-SKL and the accuracy. In Figure 2 we present how the accuracy depends on the size of the training data and the average time required for RFF-SKL and KRR-GD to run on a class. We observe that RFF-SKL running time scales linearly in the number of examples. The number of random Fourier samples we have chosen was d=3000. This is consistent with our computational complexity discussed in Section 3.2. We are able to tune the hyper-parameters of a model trained with more

than  $10^5$  samples, each with 3000 attributes in less than 15 minutes.

Samples	$2.5 \cdot 10^{3}$	$5 \cdot 10^3$	$7.5 \cdot 10^4$	$10^{4}$	$2 \cdot 10^{4}$	$6 \cdot 10^{4}$
RFF-SKL	0.21	0.25	0.28	0.29	0.31	0.32
KRR-GD	0.24	0.29	0.31	0.32	*	*

Table 2. Accuracy of RFF-SKL versus KRR-GD. For a large number of training samples the nonlinear method could not be run due to memory limits.

## 4.2. Multiple Kernel Learning

We also report experiments for multiple kernel learning. For GMKL we have evaluated an exponentiated  $\chi^2$  kernel for each image feature. We set the scaling parameter to be the mean of the chi-square distance matrix following the procedure from [13]. The Gramm matrices were created for each class since for different classes we had to select different representative samples. For our random Fourier features within the group Lasso framework (RFF-GL) we have approximated the kernel as in [30] using the recommended settings. In Figure 3 we compare the accuracy of predicting the right class for the segments. We see that for a small number of training samples GMKL is slightly superior, but RFF-GL catches up due to its scalability. Working with kernel matrices larger than  $10^4$  was not feasible for GMKL.

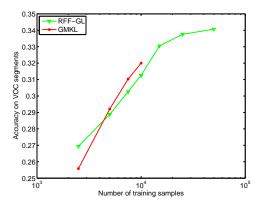


Figure 3. Accuracy for classifying the VOC segments for RFF-GL and GMKL, as a function of the number of training examples. RFF-GL scales significantly better than GMKL.

For group Lasso we have adapted the implementation presented in [21] whereas for comparisons with standard multiple kernel learning, we have used the GMKL implementation presented by Varma and Babu [33].

In Figure 4 we show the average running times for a class. We see that RFF-GL scales linearly and GMKL scales quadratically.

Following the relation given by eq. (16) in Table 3 we compare the weights given by the two methods on one of the classes on which we have performed regression – in this

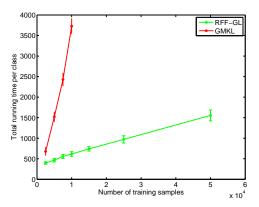


Figure 4. Running times expressed in seconds for RFF-GL and GMKL. Notice that RFF-GL scales linearly whereas GMKL scales quadratically. Error-bars are obtained by averaging over the number of classes (20 in this case).

case the aeroplane class. We see that both RFF-GL and GMKL favour the SIFT over pHOG.

## 5. Conclusions

Fourier methodology is a powerful and formally consistent class of linear approximation techniques for nonlinear kernel machines that carries the promise of combining good model scalability and non-linear prediction power. This has motivated research in extending the class of useful kernels that can be approximated, e.g. Chi-square[35, 19], but leaves ample space for reformulating standard problems like single or multiple kernel learning in the linear Fourier domain. In this paper we have developed gradient-based methods for single and multiple-kernel learning in the Fourier domain and showed that these are efficient and produce accurate results on a complex computer vision dataset like VOC2011. In future work we plan to explore alternative kernel basis expansions, feature selection and nonconvex optimization techniques for learning.

## References

- F. Bach. Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research*, (9):1179– 1225, 2008. 4, 5
- [2] F. Bach. Exploring large feature spaces with hierarchical multiple kernel learning. In *NIPS 21*, 2009. 2
- [3] F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the smo algorithm. In ICML, 2004. 2
- [4] J. Carreira and C. Sminchisescu. Constrained Parametric Min-Cuts for Automatic Object Segmentation. In CVPR, June 2010. 6
- [5] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. ACM Transactions on Intel-

	pHOG 1	pHOG 2	pHOG no pb	SIFT fg	SIFT bgnd	color SIFT fg	color SIFT bgnd	LBP
RFF-GL	0.2398	0.2260	0.5652	1.7903	1.0495	1.1287	0.8505	0.5847
GMKL	0.3862	0.1901	0.6428	1.5679	0.5333	0.4243	0.3319	0.1081

Table 3. For the RFF-GL we show the  $l_2$  norm of the coefficients corresponding to each group, whereas for GMKL we show the coefficients d mentioned in section 3.3. These are for the aeroplane class from VOC. We see both methods tend to associate similar weights to the different image features, and also the quantitative relationship given by eq. (16) is preserved.

- ligent Systems and Technology, 2:27:1-27:27, 2011. Software available at http://www.csie.ntu.edu.tw/
  ~cjlin/libsvm. 5
- [6] O. Chapelle, V. Vapnik, O. Bousquet, and S. Mukherjee. Choosing multiple parameters for support vector machines. *Machine Learning*, 46:131–159, 2002. 1, 2, 4, 6
- [7] C. Cortes, M. Mohri, and A. Rostamizadeh. Learning non-linear combinations of kernels. In *NIPS 21*, 2009. 2
- [8] C. Cortes, M. Mohri, and A. Rostamizadeh. Two-stage learning kernel algorithms. In *International Conference on Machine Learning*, 2010. 2
- [9] N. Cristianini, J. Kandola, A. Elisseeff, and J. Shawe-Taylor. On kernel-target alignment. In NIPS 14, 2002.
- [10] Ofer Dekel, Shai Shalev-Shwartz, and Yoram Singer. Smooth epsilon-insensitive regression by loss symmetrization, 2003. 5
- [11] Li F., Fu Y., Dai Y.H., Sminchisescu C., and Wang J. Kernel learning by unconstrained optimization. In AISTATS, 2009.
- [12] M. Maire P. Arbelaez C. Fowlkes and J. Malik. Using contours to detect and localize junctions in natural images. In CVPR 2008. 6
- [13] P. Gehler and S. Nowozin. On feature combination for multiclass object classification. In *ICCV*, 2009. 2, 3, 7
- [14] A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell. Gaussian processes for object categorization. *International Journal of Computer Vision*, 88:169–188, 2010.
- [15] S. Keerthi, V. Sindhwani, and O. Chapelle. An efficient method for gradient-based adaptation of hyperparameters in svm models. In NIPS 19, 2007. 2
- [16] M. Kloft, U. Brefeld, S., P. Laskov, K.-R. Müller, and A. Zien. Efficient and accurate lp-norm multiple kernel learning. In NIPS 21, 2009. 2
- [17] G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. El Ghaoui, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, 5:27–72, 2004. 1, 2
- [18] F. Li, J. Carreira, and C. Sminchisescu. Object Recognition as Ranking Holistic Figure-Ground Hypotheses. In CVPR, June 2010. 2, 6
- [19] F. Li, C. Ionescu, and C. Sminchisescu. Random Fourier approximations for skewed multiplicative histogram kernels. In *DAGM*, September 2010. 1, 2, 3, 7
- [20] F. Li and C. Sminchisescu. The Feature Selection Path in Kernel Methods. In Artificial Intelligence and Statistics. 2

- [21] J. Liu, S. Ji, and J. Ye. SLEP: Sparse Learning with Efficient Projections. Arizona State University, 2009. 5, 7
- [22] S. Maji, A. C. Berg, and J. Malik. Classification using intersection kernel support vector machines is efficient. In CVPR, 2008. 1
- [23] Timo Ojala, Matti Pietikinen, and Topi Menp. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE TPAMI*, 24(7):971–987, 2002. 6
- [24] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In NIPS, 2007. 1, 2
- [25] C. E. Rasmussen and C. K. I. Williams. Gaussian Processes for Machine Learning. MIT Press, 2006. 2
- [26] Jason D. M. Rennie. Maximum-margin logistic regression. http://people.csail.mit.edu/jrennie/writing, February 2004. 5
- [27] S Rosset. Tracking curved regularized optimization solution paths. In NIPS 17, 2005.
- [28] Walter Rudin. Fourier Analysis on Groups. Wiley-Interscience, 1990. 2
- [29] Bernhard Schölkopf and Alexander J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond.* Adaptive computation and machine learning. MIT Press, 2002. 5
- [30] V. Sreekanth, A. Vedaldi, C. V. Jawahar, and A. Zisserman. Generalized RBF feature maps for efficient detection. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2010. 6, 7
- [31] M. Szafranski, Y. Grandvalet, and A. Rakotomamonjy. Composite kernel learning. In *International Conference on Machine Learning*, 2008. 2
- [32] K. E. A. van de Sande, T. Gevers, and C. G. M. Snoek. Evaluating color descriptors for object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1582–1596, 2010. 6
- [33] M. Varma and B. R. Babu. More generality in efficient multiple kernel learning. In *International Conference on Machine Learning*, pages 1065–1072, June 2009. 1, 4, 5, 7
- [34] A. Vedaldi, V. Gulshan, M. Varma, and A. Zisserman. Multiple kernels for object detection. In *ICCV*, 2009. 1
- [35] A. Vedaldi and A. Zisserman. Efficient additive kernels via explicit feature maps. In *CVPR*, 2010. 1, 6, 7
- [36] S.V.N. Visnwanathan, Z. Sun, N. Theera-Ampornpunt, and Manik Varma. Multiple kernel learning and the smo algorithm. In NIPS 22, 2010. 2

[37] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68:49–67, 2006. 1, 4