

# Learning from Distributions via Support Measure Machines

**Authored by:**

Kenji Fukumizu  
Bernhard Schölkopf  
Krikamol Muandet  
Francesco Dinuzzo

## **Abstract**

This paper presents a kernel-based discriminative learning framework on probability measures. Rather than relying on large collections of vectorial training examples, our framework learns using a collection of probability distributions that have been constructed to meaningfully represent training data. By representing these probability distributions as mean embeddings in the reproducing kernel Hilbert space (RKHS), we are able to apply many standard kernel-based learning techniques in straightforward fashion. To accomplish this, we construct a generalization of the support vector machine (SVM) called a support measure machine (SMM). Our analyses of SMMs provides several insights into their relationship to traditional SVMs. Based on such insights, we propose a flexible SVM (Flex-SVM) that places different kernel functions on each training example. Experimental results on both synthetic and real-world data demonstrate the effectiveness of our proposed framework.

## **1 Paper Body**

Discriminative learning algorithms are typically trained from large collections of vectorial training examples. In many classical learning problems, however, it is arguably more appropriate to represent training data not as individual data points, but as probability distributions. There are, in fact, multiple reasons why probability distributions may be preferable. Firstly, uncertain or missing data naturally arises in many applications. For example, gene expression data obtained from the microarray experiments are known to be very noisy due to various sources of variabilities [1]. In order to reduce uncertainty, and to allow for estimates of confidence levels, experiments are often replicated. Unfortunately, the feasibility of replicating the microarray experiments is often inhibited by cost constraints, as well as the amount of available mRNA. To cope with experimental uncertainty given a limited amount of data, it is natural to represent

each array as a probability distribution that has been designed to approximate the variability of gene expressions across slides. Probability distributions may be equally appropriate given an abundance of training data. In data-rich disciplines such as neuroinformatics, climate informatics, and astronomy, a high throughput experiment can easily generate a huge amount of data, leading to significant computational challenges in both time and space. Instead of scaling up one's learning algorithms, one can scale down one's dataset by constructing a smaller collection of distributions which represents groups of similar samples. Besides computational efficiency, aggregate statistics can potentially incorporate higher-level information that represents the collective behavior of multiple data points. 1

Previous attempts have been made to learn from distributions by creating positive definite (p.d.) kernels on probability measures. In [2], the probability product kernel (PPK) was proposed as a generalized inner product between two input objects, which is in fact closely related to well-known kernels such as the Bhattacharyya kernel [3] and the exponential symmetrized Kullback-Leibler (KL) divergence [4]. In [5], an extension of a two-parameter family of Hilbertian metrics of Topsøe was used to define Hilbertian kernels on probability measures. In [6], the semi-group kernels were designed for objects with additive semi-group structure such as positive measures. Recently, [7] introduced nonextensive information theoretic kernels on probability measures based on new Jensen-Shannon-type divergences. Although these kernels have proven successful in many applications, they are designed specifically for certain properties of distributions and application domains. Moreover, there has been no attempt in making a connection to the kernels on corresponding input spaces. The contributions of this paper can be summarized as follows. First, we prove the representer theorem for a regularization framework over the space of probability distributions, which is a generalization of regularization over the input space on which the distributions are defined (Section 2). Second, a family of positive definite kernels on distributions is introduced (Section 3). Based on such kernels, a learning algorithm on probability measures called support measure machine (SMM) is proposed. An SVM on the input space is provably a special case of the SMM. Third, the paper presents the relations between sample-based and distribution-based methods (Section 4). If the distributions depend only on the locations in the input space, the SMM particularly reduces to a more flexible SVM that places different kernels on each data point.

## 2

### Regularization on probability distributions

Given a non-empty set  $X$ , let  $\mathcal{P}$  denote the set of all probability measures  $P$  on a measurable space  $(X, \mathcal{A})$ , where  $\mathcal{A}$  is a  $\sigma$ -algebra of subsets of  $X$ . The goal of this work is to learn a function  $h : \mathcal{P} \rightarrow Y$  given a set of example pairs  $\{(P_i, y_i)\}_{i=1}^m$ , where  $P_i \in \mathcal{P}$  and  $y_i \in Y$ . In other words, we consider a supervised setting in which input training examples are probability distributions. In this paper, we focus on the binary classification problem, i.e.,  $Y = \{+1, -1\}$ . In order to learn from distributions, we employ a compact representation that not only preserves necessary information of individual distributions, but also

permits efficient computations. That is, we adopt a Hilbert space embedding to represent the distribution as a mean function in an RKHS [8, 9]. Formally, let  $H$  denote an RKHS of functions  $f : X \rightarrow \mathbb{R}$ , endowed with a reproducing kernel  $k : X \times X \rightarrow \mathbb{R}$ . The mean map from  $\mathcal{P}$  into  $H$  is defined as  $Z : \mathcal{P} \rightarrow H$ ,  $P \mapsto k(x, \cdot) dP(x)$ . (1)  $X$

We assume that  $k(x, \cdot)$  is bounded for any  $x \in X$ . It can be shown that, if  $k$  is characteristic, the map (1) is injective, i.e., all the information about the distribution is preserved [10]. For any  $P$ , letting  $\Phi_P = \Phi(P)$ , we have the reproducing property  $E_P[f] = \langle \Phi_P, \Phi f \rangle_H$ .

(2)

That is, we can see the mean embedding  $\Phi_P$  as a feature map associated with the kernel  $K : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$ , defined as  $K(P, Q) = \langle \Phi_P, \Phi_Q \rangle_H$ . Since  $\langle \Phi_P, \Phi_Q \rangle_H = \int \int k(x, z) dP(x) dQ(z)$ , it also follows that  $K(P, Q) = \langle \int k(x, \cdot) dP(x), \int k(z, \cdot) dQ(z) \rangle_H$ , where the second equality follows from the reproducing property of  $H$ . It is immediate that  $K$  is a p.d. kernel on  $\mathcal{P}$ . The following theorem shows that optimal solutions of a suitable class of regularization problems involving distributions can be expressed as a finite linear combination of mean embeddings. Theorem 1. Given training examples  $(P_i, y_i) \in \mathcal{P} \times \mathbb{R}$ ,  $i = 1, \dots, m$ , a strictly monotonically increasing function  $\gamma : [0, +\infty) \rightarrow \mathbb{R}$ , and a loss function  $\ell : (\mathcal{P} \times \mathbb{R})^m \rightarrow \mathbb{R} \cup \{+\infty\}$ , any  $f \in H$  minimizing the regularized risk functional  $\gamma(P_1, y_1, E_{P_1}[f], \dots, P_m, y_m, E_{P_m}[f]) + \gamma(\langle f, \Phi \rangle_H)$  for some  $\gamma \in \mathbb{R}$ ,  $i = 1, \dots, m$ .

(3)

admits a representation of the form  $f =$

Theorem 1 clearly indicates how each distribution contributes to the minimizer of (3). Roughly speaking, the coefficients  $\gamma_i$  controls the contribution of the distributions through the mean embeddings  $\Phi_{P_i}$ . Furthermore, if we restrict  $\mathcal{P}$  to a class of Dirac measures  $\delta_x$  on  $X$  and consider

functional (3) reduces to the usual regularization functional [11] the training set  $\{(\delta_{x_i}, y_i)\}_{i=1}^m$ , the  $\mathcal{P}$  is  $\mathcal{P}_m$  and the solution reduces to  $f = \sum_{i=1}^m \gamma_i k(x_i, \cdot)$ . Therefore, the standard representer theorem is recovered as a particular case (see also [12] for more general results on representer theorem). Note that, on the one hand, the minimization problem (3) is different from minimizing the functional  $E_{P_1} \dots E_{P_m} \gamma(x_1, y_1, f(x_1), \dots, x_m, y_m, f(x_m)) + \gamma(\langle f, \Phi \rangle_H)$  for the special case of the additive loss  $\gamma$ . Therefore, the solution of our regularization problem is different from what one would get in the limit by training on an infinitely many points sampled from  $P_1, \dots, P_m$ . On the other hand, it is also different from minimizing the functional  $\gamma(M_1, y_1, f(M_1), \dots, M_m, y_m, f(M_m)) + \gamma(\langle f, \Phi \rangle_H)$  where  $M_i = E_{\delta_{x_i}} \Phi[x]$ . In a sense, our framework is something in between.

3

Kernels on probability distributions

As the map (1) is linear in  $P$ , optimizing the functional (3) amounts to finding  $R$  a function in  $H$  that approximate well functions from  $\mathcal{P}$  to  $\mathbb{R}$  in the function class  $F = \{P \in \mathcal{P} : \int g dP \in C(X) \text{ for all } g \in C(X)\}$  where  $C(X)$  is a class of bounded continuous functions on  $X$ . Since  $\delta_x \in \mathcal{P}$  for any  $x \in X$ , it follows that  $C(X) \subset$

$F \in C(P)$  where  $C(P)$  is a class of bounded continuous functions on  $P$  endowed with the topology of weak convergence and the associated Borel  $\sigma$ -algebra. The following lemma states the relation between the RKHS  $H$  induced by the kernel  $k$  and the function class  $F$ . Lemma 2. Assuming that  $X$  is compact, the RKHS  $H$  induced by a kernel  $k$  is dense in  $F$  if  $k$  is universal, i.e., for every function  $F \in F$  and every  $\epsilon > 0$  there exists a function  $g \in H$  with  $\sup_{P \times P} |F(P) - g(dP)| < \epsilon$ . Proof. Assume that  $k$  is universal. Then, for every function  $f \in C(X)$  and every  $\epsilon > 0$  there exists a function  $g \in H$  induced by  $k$  with  $\sup_{x \in X} |f(x) - g(x)| < \epsilon$  [13]. Hence, by linearity of  $F$ , for every  $F \in F$  and every  $\epsilon > 0$  there exists a function  $h \in H$  such that  $\sup_{P \times P} |F(P) - h(dP)| < \epsilon$ .

Nonlinear kernels on  $P$  can be defined in an analogous way to nonlinear kernels on  $X$ , by treating mean embeddings  $\mu_P$  of  $P \in P$  as its feature representation. First, assume that the map (1) is injective and let  $\langle \cdot, \cdot \rangle_P$  be an inner product on  $P$ . By linearity, we have  $\langle \mu_P, \mu_Q \rangle_P = \langle \mu_P, \mu_Q \rangle_H$  (cf. [8] for more details). Then, the nonlinear kernels on  $P$  can be defined as  $K(P, Q) = \langle \mu_P, \mu_Q \rangle_P = \langle \mu_P, \mu_Q \rangle_H$  where  $\langle \cdot, \cdot \rangle_H$  is a p.d. kernel. As a result, many standard nonlinear kernels on  $X$  can be used to define nonlinear kernels on  $P$  as long as the kernel evaluation depends entirely on the inner product  $\langle \mu_P, \mu_Q \rangle_H$ , e.g.,  $K(P, Q) = \langle \mu_P, \mu_Q \rangle_H + c$ . Although requiring more computational effort, their practical use is simple and flexible. Specifically, the notion of p.d. kernels on distributions proposed in this work is so generic that standard kernel functions can be reused to derive kernels on distributions that are different from many other kernel functions proposed specifically for certain distributions. It has been recently proved that the Gaussian RBF kernel given by  $K(P, Q) = \exp(-\frac{1}{2} \| \mu_P - \mu_Q \|_H^2)$ ,  $\mu_P, \mu_Q \in P$  is universal w.r.t  $C(P)$  given that  $X$  is compact and the map (1) is injective [14]. Despite its success in real-world applications, the theory of kernel-based classifiers beyond the input space  $X \subset \mathbb{R}^d$ , as also mentioned by [14], is still incomplete. It is therefore of theoretical interest to consider more general classes of universal kernels on probability distributions. 3.1

#### Support measure machines

This subsection extends SVMs to deal with probability distributions, leading to support measure machines (SMMs). In its general form, an SMM amounts to solving an SVM problem with the expected kernel  $K(P, Q) = \mathbb{E}_{x \sim P, z \sim Q} [k(x, z)]$ . This kernel can be computed in closed-form for certain classes of distributions and kernels  $k$ . Examples are given in Table 1. Alternatively, one can approximate the kernel  $K(P, Q)$  by the empirical estimate:  $n$

$$\hat{K}_n(P, Q) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, z_j) \quad (4)$$

$\hat{P}_n$  and  $\hat{Q}_n$  are empirical distributions of  $P$  and  $Q$  given random samples  $\{x_i\}_{i=1}^n$  and  $\{z_j\}_{j=1}^m$ , respectively. A finite sample of size  $n$  from a distribution  $P$  suffices (with high probability) [3]

Table 1: the analytic forms of expected kernels for different choices of kernels and distributions.

Embedding kernel  $k(x, y)$

$K(P_i, P_j) = h^T P_i, P_j i H$   
 Arbitrary  $P(m; ?)$  Gaussian  $N(m; ?)$   
 Linear  $h_x, y_i$  Gaussian RBF  $\exp(-\frac{1}{2} \|x - y\|^2)$   
 Gaussian  $N(m; ?)$  Gaussian  $N(m; ?)$   
 Polynomial degree 2  $(h_x, y_i + 1)^2$  Polynomial degree 3  $(h_x, y_i + 1)^3$   
 $m^T i m_j + ?_{ij} \text{tr} ?_i \exp(-\frac{1}{2} (m_i - m_j)^T (?_i + ?_j + ?^{-1} I)^{-1} (m_i - m_j))$   
 $1 / \sqrt{?_i + ?_j + I} - \frac{1}{2} m^T T(h_{mi}, m_{ji} + 1) + \text{tr} ?_i ?_j + m^T i ?_j m_i + m_j ?_i m_j$   
 $(h_{mi}, m_{ji} + 1)^3 + 6m^T ? ? m_i i j j T + 3(h_{mi}, m_{ji} + 1)(\text{tr} ?_i ?_j + m^T i ?_j$   
 $m_i + m_j ?_i m_j)$

1  
 to compute an approximation within an error of  $O(m^{-2})$ . Instead, if the sample set is sufficiently large, one may choose to approximate the true distribution by simpler probabilistic models, e.g., a mixture of Gaussians model, and choose a kernel  $k$  whose expected value admits an analytic form. Storing only the parameters of probabilistic models may save some space compared to storing all data points. Note that the standard SVM feature map  $\phi(x)$  is usually nonlinear in  $x$ , whereas  $\phi P$  is linear in  $P$ . Thus, for an SMM, the first level kernel  $k$  is used to obtain a vectorial representation of the measures, and the second level kernel  $K$  allows for a nonlinear algorithm on distributions. For clarity, we will refer to  $k$  and  $K$  as the embedding kernel and the level-2 kernel, respectively

4

#### Theoretical analyses

This section presents key theoretical aspects of the proposed framework, which reveal important connection between kernel-based learning algorithms on the space of distributions and on the input space on which they are defined.

#### Risk deviation bound

Given a training sample  $\{(P_i, y_i)\}_{i=1}^m$  drawn i.i.d. from some unknown probability distribution  $P$  on  $P \times Y$ , a loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ , and a function class  $\mathcal{H}$ , the goal of statistical learning is to find the function  $f \in \mathcal{H}$  that minimizes the expected risk functional  $R(f) = \mathbb{E}_{P \times Y} \ell(y, f(x)) dP(x)$ . Since  $P$  is unknown, the empirical risk  $R_{\text{emp}}(f) = \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i))$  is considered instead. Furthermore,  $\frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i)) dP_i(x)$  based on the training sample is  $\frac{1}{m} \sum_{i=1}^m \int \ell(y_i, f(x)) dP_i(x)$  the risk functional can be simplified further by considering  $m \times n$   $i=1 \sum_{j=1}^n x_{ij} ? P_i ?(y_i, f(x_{ij}))$  based on  $n$  samples  $x_{ij}$  drawn from each  $P_i$ . Our framework, on the other hand, alleviates the problem by minimizing the risk functional  $R_{\mathcal{H}}(f) = \int \ell(y, \mathbb{E}_P[f(x)]) dP(P, y)$  for  $f \in \mathcal{H}$  with corresponding empirical risk functional  $\frac{1}{m} \sum_{i=1}^m R_{\text{emp}}^i(f) = \frac{1}{m} \sum_{i=1}^m \int \ell(y_i, \mathbb{E}_{P_i}[f(x)]) dP_i(P, y)$  (cf. the discussion at the end of Section 2). It is often easier to optimize  $R_{\text{emp}}(f)$  as the expectation can be computed exactly for certain choices of  $P_i$  and  $\mathcal{H}$ . Moreover, for universal  $\mathcal{H}$ , this simplification preserves all information of the distributions. Nevertheless, there is still a loss of information due to the loss function  $\ell$ . Due to the i.i.d. assumption, the analysis of the difference between  $R$  and  $R_{\mathcal{H}}$  can be simplified w.l.o.g. to the analysis of the difference between  $\mathbb{E}_P[\ell(y, f(x))]$  and  $\int \ell(y, \mathbb{E}_P[f(x)]) dP(P, y)$  for a particular distribution  $P \in \mathcal{P}$ . The theorem below provides a bound on the difference between  $\mathbb{E}_P[\ell(y, f(x))]$  and  $\int \ell(y, \mathbb{E}_P[f(x)]) dP(P, y)$ . Theorem 3. Given an arbitrary probability distribution  $P$  with variance

$\frac{1}{2}$ , a Lipschitz continuous function  $f: \mathcal{R} \rightarrow \mathcal{R}$  with constant  $C_f$ , an arbitrary loss function  $\ell: \mathcal{R} \rightarrow \mathcal{R} \rightarrow \mathcal{R}$  that is Lipschitz continuous in the second argument with constant  $C_\ell$ , it follows that  $|\mathbb{E}_x \ell(y, f(x)) - \ell(y, \mathbb{E}_x f(x))| \leq 2C_\ell C_f$  for any  $y \in \mathcal{R}$ . Theorem 3 indicates that if the random variable  $x$  is concentrated around its mean and the function  $f$  and  $\ell$  are well-behaved, i.e., Lipschitz continuous, then the loss deviation  $|\mathbb{E}_P \ell(y, f(x)) - \ell(y, \mathbb{E}_P f(x))|$  will be small. As a result, if this holds for any distribution  $P_i$  in the training set  $\{(P_i, y_i)\}_{i=1}^m$ , the true risk deviation  $|\mathbb{R} - \mathbb{R}|$  is also expected to be small. <sup>4</sup>

#### 4.2

##### Flexible support vector machines

It turns out that, for certain choices of distributions  $P$ , the linear SMM trained using  $\{(P_i, y_i)\}_{i=1}^m$  is equivalent to an SVM trained using some samples  $\{(x_i, y_i)\}_{i=1}^m$  with an appropriate choice of kernel function. **Lemma 4.** Let  $k(x, z)$  be a bounded p.d. kernel on a measure space  $\mathcal{R}$  such that  $\int \int k(x, z)^2 dx dz < \infty$ , and  $g(x, x')$  be a square integrable function such that  $\int g(x, x')^2 dx < \infty$  for all  $x$ . Given a sample  $\{(P_i, y_i)\}_{i=1}^m$  where each  $P_i$  is assumed to have a density given by  $g(x_i, x)$ , the linear SMM is equivalent to the SVM on the training sample  $\{(x_i, y_i)\}_{i=1}^m$  with kernel  $K_g(x, z) = \int k(x, z')g(x, z')g(z, z')dz'$ .

Note that the important assumption for this equivalence is that the distributions  $P_i$  differ only in their location in the parameter space. This need not be the case in all possible applications of SMMs.

**Proof.**

Furthermore, we have  $K_g(x, z) = \int k(x, z')g(x, z')g(z, z')dz' = \int k(z', z)g(z', x)g(z', z)dz' = \int k(z, z')g(z, x)g(z, z')dz'$ . Thus, it is clear that the feature map of  $x$  depends not only on the kernel  $k$ , but also on the density  $g(x, x')$ . Consequently, by virtue of Lemma 4, the kernel  $K_g$  allows the SVM to place different kernels at each data point. We call this algorithm a flexible SVM (Flex-SVM). <sup>2</sup> Consider for example the linear SMM with Gaussian distributions  $N(x_1; \mu_1, \Sigma_1), \dots, N(x_m; \mu_m, \Sigma_m)$  and Gaussian RBF kernel  $k_2$  with bandwidth parameter  $\gamma$ . The convolution theorem of Gaussian distributions implies that this SMM is equivalent to a flexible SVM that places a data-dependent kernel  $k_2 + 2\gamma^2(x_i, x_j)$  on training example  $x_i$ , i.e., a Gaussian RBF kernel with larger bandwidth. <sup>5</sup>

##### Related works

The kernel  $K(P, Q) = \mathbb{E}_{x \sim P, y \sim Q} k(x, y)$  is in fact a special case of the Hilbertian metric [5], with the associated kernel  $K(P, Q) = \mathbb{E}_{P, Q} [k(x, x')]$ , and a generative mean map kernel (GMMK) proposed by [15]. In the GMMK, the kernel between two objects  $x$  and  $y$  is defined via  $p^x$  and  $p^y$ , which are estimated probabilistic models of  $x$  and  $y$ , respectively. That is, a probabilistic model  $p^x$  is learned for each example and used as a surrogate to construct the kernel between those examples. The idea of surrogate kernels has also been adopted by the Probability Product Kernel (PPK) [2]. In this case, we have  $K(p, p') = \int p(x)p'(x)dx$ , which has been shown to be a special case of GMMK when  $\gamma = 1$  [15]. Consequently, GMMK, PPK with  $\gamma = 1$ , and our linear

kernels are equivalent when the embedding kernel is  $k(x, x') = \phi(x) \cdot \phi(x')$ . More recently, the empirical kernel (4) was employed in an unsupervised way for multi-task learning to generalize to a previously unseen task [16]. In contrast, we treat the probability distributions in a supervised way (cf. the regularized functional (3)) and the kernel is not restricted to only the empirical kernel. The use of expected kernels in dealing with the uncertainty in the input data has a connection to robust SVMs. For instance, a generalized form of the SVM in [17] incorporates the probabilistic uncertainty into the maximization of the margin. This results in a second-order cone programming (SOCP) that generalizes the standard SVM. In SOCP, one needs to specify the parameter  $\gamma_i$  that reflects the probability of correctly classifying the  $i$ th training example. The parameter  $\gamma_i$  is therefore closely related to the parameter  $\sigma_i$ , which specifies the variance of the distribution centered at the  $i$ th example. [18] showed the equivalence between SVMs using expected kernels and SOCP when  $\gamma_i = 0$ . When  $\gamma_i \neq 0$ , the mean and covariance of missing kernel entries have to be estimated explicitly, making the SOCP more involved for nonlinear kernels. Although achieving comparable performance to the standard SVM with expected kernels, the SOCP requires a more computationally extensive SOCP solver, as opposed to simple quadratic programming (QP).

6

#### Experimental results

In the experiments, we primarily consider three different learning algorithms: i) SVM is considered as a baseline algorithm. ii) Augmented SVM (ASVM) is an SVM trained on augmented samples drawn according to the distributions  $\{P_i\}_{i=1}^m$ . The same number of examples are drawn from each distribution. iii) SMM is distribution-based method that can be applied directly on the distributions  $P_1, \dots, P_m$ .

We used the LIBSVM implementation.

5

100

80 60 40 Embedding RBF 1 Level-2 RBF Embedding RBF 2 Level-2 Poly

20 0 0

(a) decision boundaries.

1

2

3

4

5

6

Parameters

7

8

Embedding RBF 2

Embedding RBF 1

Accuracy(%)

100

90 80 70 60 50 40

Level-2 POLY

Level-2 RBF

(b) sensitivity of kernel parameters

Figure 1: (a) the decision boundaries of SVM, ASVM, and SMM. (b) the heatmap plots of average accuracies of SMM over 30 experiments using POLY-RBF (center) and RBF-RBF (right) kernel combinations with the plots of average accuracies at different parameter values (left).

Level-2 kernels

Table 2: accuracies (%) of SMM on synthetic data with different combinations of embedding and level-2 kernels.

6.1

LIN POLY RBF

Embedding kernels POLY3 RBF

LIN

POLY2

85.20 $\pm$ 2.20 83.95 $\pm$ 2.11 87.80 $\pm$ 1.96

81.04 $\pm$ 3.11 81.34 $\pm$ 1.21 73.12 $\pm$ 3.29

81.10 $\pm$ 2.76 82.66 $\pm$ 1.75 78.28 $\pm$ 2.19

87.74 $\pm$ 2.19 88.06 $\pm$ 1.73 89.65 $\pm$ 1.37

URBF 85.39 $\pm$ 2.56 86.84 $\pm$ 1.51 86.86 $\pm$ 1.88

Synthetic data

Firstly, we conducted a basic experiment that illustrates a fundamental difference between SVM, ASVM, and SMM. A binary classification problem of 7 Gaussian distributions with different means and covariances was considered. We trained the SVM using only the means of the distributions, ASVM with 30 virtual examples generated from each distribution, and SMM using distributions as training examples. A Gaussian RBF kernel with  $\gamma = 0.25$  was used for all algorithms. Figure 1a shows the resulting decision boundaries. Having been trained only on means of the distributions, the SVM classifier tends to overemphasize the regions with high densities and underrepresent the lower density regions. In contrast, the ASVM is more expensive and sensitive to outliers, especially when learning on heavy-tailed distributions. The SMM treats each distribution as a training example and implicitly incorporates properties of the distributions, i.e., means and covariances, into the classifier. Note that the SVM can be trained to achieve a similar result to the SMM by choosing an appropriate value for  $\gamma$  (cf. Lemma 4). Nevertheless, this becomes more difficult if the training distributions are, for example, nonisotropic and have different covariance matrices. Secondly, we evaluate the performance of the SMM for different combinations of embedding and level-2 kernels. Two classes of synthetic Gaussian distributions on R10 were generated. The mean parameters of the positive and negative distributions are normally distributed with means  $m_+ = (1, \dots, 1)$  and  $m_- = (2, \dots, 2)$  and identical covariance matrix  $\Sigma = 0.5 \cdot I_{10}$ , respectively. The covariance matrix for each distribution is generated according to two Wishart distributions with covariance matrices given by  $\Sigma_+ = 0.6 \cdot I_{10}$  and  $\Sigma_- = 1.2 \cdot I_{10}$  with 10 degrees of freedom. The training set consists of



500 distributions from the positive class and 500 distributions from the negative class. The test set consists of 200 distributions with the same class proportion. The kernels used in the experiment include linear kernel (LIN), polynomial kernel of degree 2 (POLY2), polynomial kernel of degree 3 (POLY3), unnormalized Gaussian RBF kernel (RBF), and normalized Gaussian RBF kernel (NRBF). To fix parameter values of both kernel functions and SMM, 10-fold cross-validation (10-CV) is performed on a parameter grid,  $C \in \{2^{-3}, 2^{-2}, \dots, 2^7\}$  for SMM, bandwidth parameter  $\gamma \in \{10^{-3}, 10^{-2}, \dots, 10^2\}$  for Gaussian RBF kernels, and degree parameter  $d \in \{2, 3, 4, 5, 6\}$  for polynomial kernels. The average accuracy and  $\pm 1$  standard deviation for all kernel combinations over 30 repetitions are reported in Table 2. Moreover, we also investigate the sensitivity of kernel parameters for two kernel combinations: RBF-RBF and POLYRBF. In this case, we consider the bandwidth parameter  $\gamma = \{10^{-3}, 10^{-2}, \dots, 10^3\}$  for Gaussian 6

90 100  
90  
100 95 90 85  
100  
100  
100 95  
95 100  
95 100 95 90 100 90 80 70  
95 95 95  
90 85  
10  
20  
30  
10  
20 30 10 20 Number of virtual examples  
Scaling  
95  
100  
30  
10  
20  
30  
102 101 100 2000  
4000  
6000  
Number of virtual examples

Figure 3: relative computational cost of ASVM and SMM (baseline: SMM with 2000 virtual examples). 70 65 60 55 50

Figure 2: the performance of SVM, ASVM, and SMM algorithms on handwritten digits constructed using three basic transformations.

ASVM  
103

10-1  
 Accuracy (%)  
 Accuracy (%)  
 100 95 90 85  
 100  
 95  
 6 vs 9  
 Translation  
 100  
 3 vs 8  
 Rotation  
 3 vs 4  
 Relative comp. cost  
 SMM  
 1 vs 8  
 pLSA  
 SVM  
 LSMM NLSMM

Figure 4: accuracies of four different techniques for natural scene categorization.

RBF kernels and degree parameter  $d = \{2, 3, \dots, 8\}$  for polynomial kernels. Figure 1b depicts the accuracy values and average accuracies for considered kernel functions. Table 2 indicates that both embedding and level-2 kernels are important for the performance of the classifier. The embedding kernels tend to have more impact on the predictive performance compared to the level-2 kernels. This conclusion also coincides with the results depicted in Figure 1b.

## 6.2 Handwritten digit recognition

In this section, the proposed framework is applied to distributions over equivalence classes of images that are invariant to basic transformations, namely, scaling, translation, and rotation. We consider the handwritten digits obtained from the USPS dataset. For each  $16 \times 16$  image, the distribution over the equivalence class of the transformations is determined by a prior on parameters associated with such transformations. Scaling and translation are parametrized by the scale factors  $(s_x, s_y)$  and displacements  $(t_x, t_y)$  along the  $x$  and  $y$  axes, respectively. The rotation is parametrized by an angle  $\theta$ . We adopt Gaussian distributions as prior distributions, including  $N([1, 1], 0.1I_2)$ ,  $N([0, 0], 5I_2)$ , and  $N(0; ?)$ . For each image, the virtual examples are obtained by sampling parameter values from the distribution and applying the transformation accordingly. Experiments are categorized into simple and difficult binary classification tasks. The former consists of classifying digit 1 against digit 8 and digit 3 against digit 4. The latter considers classifying digit 3 against digit 8 and digit 6 against digit 9. The initial dataset for each task is constructed by randomly selecting 100 examples from each class. Then, for each example in the initial dataset, we generate 10, 20, and 30 virtual examples using the aforementioned transformations to construct virtual data sets consisting of 2,000, 4,000,

and 6,000 examples, respectively. One third of examples in the initial dataset are used as a test set. The original examples are excluded from the virtual datasets. The virtual examples are normalized such that their feature values are in  $[0, 1]$ . Then, to reduce computational cost, principle component analysis (PCA) is performed to reduce the dimensionality to 16. We compare the SVM on the initial dataset, the ASVM on the virtual datasets, and the SMM. For SVM and ASVM, the Gaussian RBF kernel is used. For SMM, we employ the empirical kernel (4) with Gaussian RBF kernel as a base kernel. The parameters of the algorithms are fixed by 10-CV over parameters  $C \in \{2^{-3}, 2^{-2}, \dots, 2^7\}$  and  $\gamma \in \{0.01, 0.1, 1\}$ . The results depicted in Figure 2 clearly demonstrate the benefits of learning directly from the equivalence classes of digits under basic transformations<sup>2</sup>. In most cases, the SMM outperforms both the SVM and the ASVM as the number of virtual examples increases. Moreover, Figure 3 shows the benefit of the SMM over the ASVM in term of computational cost<sup>3</sup>.<sup>2</sup> While the reported results were obtained using virtual examples with Gaussian parameter distributions (Sec. 6.2), we got similar results using uniform distributions. TM R 3 Core 2 Duo CPU E8400 at The evaluation was made on a 64-bit desktop computer with Intel 3.00GHz<sup>2</sup> and 4GB of memory.

7

### 6.3

#### Natural scene categorization

This section illustrates benefits of the nonlinear kernels between distributions for learning natural scene categories in which the bag-of-word (BoW) representation is used to represent images in the dataset. Each image is represented as a collection of local patches, each being a codeword from a large vocabulary of codewords called codebook. Standard BoW representations encode each image as a histogram that enumerates the occurrence probability of local patches detected in the image w.r.t. those in the codebook. On the other hand, our setting represents each image as a distribution over these codewords. Thus, images of different scenes tends to generate distinct set of patches. Based on this representation, both the histogram and the local patches can be used in our framework. We use the dataset presented in [19]. According to their results, most errors occurs among the four indoor categories (830 images), namely, bedroom (174 images), living room (289 images), kitchen (151 images), and office (216 images). Therefore, we will focus on these four categories. For each category, we split the dataset randomly into two separate sets of images, 100 for training and the rest for testing. A codebook is formed from the training images of all categories. Firstly, interesting keypoints in the image are randomly detected. Local patches are then generated accordingly. After patch detection, each patch is transformed into a 128-dim SIFT vector [20]. Given the collection of detected patches, K-means clustering is performed over all local patches. Codewords are then defined as the centers of the learned clusters. Then, each patch in an image is mapped to a codeword and the image can be represented by the histogram of the codewords. In addition, we also have an  $M \times 128$  matrix of SIFT vectors where  $M$  is the number of codewords. We compare the performance of a Probabilistic Latent Semantic Analysis (pLSA) with the standard BoW rep-

resentation, SVM, linear SMM (LSMM), and nonlinear SMM (NLSMM). For SMM, we use the empirical embedding kernel with Gaussian RBF base kernel  $k: K(h_i, h_j) = \frac{1}{P} \sum_{r=1}^P \sum_{s=1}^P h_i(c_r) h_j(c_s) k(c_r, c_s)$  where  $h_i$  is the histogram of the  $i$ th image and  $c_r$  is the  $r$ th SIFT vector. A Gaussian RBF kernel is also used as the level-2 kernel for nonlinear SMM. For the SVM, we adopt a Gaussian RBF kernel with  $\ell_2$ -distance between the histograms [21], i.e.,

$\frac{1}{P} \sum_{r=1}^P (h_i(c_r) - h_j(c_r))^2$   $K(h_i, h_j) = \exp(-\gamma \sum_{r=1}^P (h_i(c_r) - h_j(c_r))^2)$  where  $\gamma(h_i, h_j) = \frac{1}{\sum_{r=1}^P (h_i(c_r) - h_j(c_r))^2}$ . The parameters of the algorithms are fixed by 10-CV over parameters  $C \in \{2^{-3}, 2^{-2}, \dots, 2^7\}$  and  $\gamma \in \{0.01, 0.1, 1\}$ . For NLSMM, we use the best  $\gamma$  of LSMM in the base kernel and perform 10-CV to choose  $\gamma$  parameter only for the level-2 kernel. To deal with multiple categories, we adopt the pairwise approach and voting scheme to categorize test images. The results in Figure 4 illustrate the benefit of the distribution-based framework. Understanding the context of a complex scene is challenging. Employing distribution-based methods provides an elegant way of utilizing higher-order statistics in natural images that could not be captured by traditional sample-based methods.

7

## Conclusions

This paper proposes a method for kernel-based discriminative learning on probability distributions. The trick is to embed distributions into an RKHS, resulting in a simple and efficient learning algorithm on distributions. A family of linear and nonlinear kernels on distributions allows one to flexibly choose the kernel function that is suitable for the problems at hand. Our analyses provide insights into the relations between distribution-based methods and traditional sample-based methods, particularly the flexible SVM that allows the SVM to place different kernels on each training example. The experimental results illustrate the benefits of learning from a pool of distributions, compared to a pool of examples, both on synthetic and real-world data.

Acknowledgments KM would like to thank Zoubin Ghahramani, Arthur Gretton, Christian Walder, and Philipp Hennig for a fruitful discussion. We also thank all three insightful reviewers for their invaluable comments. 8

## 2 References

- [1] Y. H. Yang and T. Speed. Design issues for cDNA microarray experiments. *Nat. Rev. Genet.*, 3(8):579–588, 2002.
- [2] T. Jebara, R. Kondor, A. Howard, K. Bennett, and N. Cesa-bianchi. Probability product kernels. *Journal of Machine Learning Research*, 5:819–844, 2004.
- [3] A. Bhattacharyya. On a measure of divergence between two statistical populations defined by their probability distributions. *Bull. Calcutta Math Soc.*, 1943.
- [4] P. J. Moreno, P. P. Ho, and N. Vasconcelos. A Kullback-Leibler divergence based kernel for SVM classification in multimedia applications. In *Proceedings of Advances in Neural Information Processing Systems*. MIT Press, 2004.
- [5] M. Hein and O. Bousquet. Hilber-

tian metrics and positive definite kernels on probability. In Proceedings of The 12th International Conference on Artificial Intelligence and Statistics, pages 136?143, 2005. [6] M. Cuturi, K. Fukumizu, and J-P. Vert. Semigroup kernels on measures. *Journal of Machine Learning Research*, 6:1169?1198, 2005. [7] Andr e F. T. Martins, Noah A. Smith, Eric P. Xing, Pedro M. Q. Aguiar, and M ario A. T. Figueiredo. Nonextensive information theoretic kernels on measures. *Journal of Machine Learning Research*, 10:935?975, 2009. [8] A. Berline and Thomas C. Agnan. Reproducing kernel Hilbert spaces in probability and statistics. Kluwer Academic Publishers, 2004. [9] A. Smola, A. Gretton, L. Song, and B. Sch olkopf. A hilbert space embedding for distributions. In Proceedings of the 18th International Conference on Algorithmic Learning Theory, pages 13?31. Springer-Verlag, 2007. [10] B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Sch olkopf, and Gert R. G. Lanckriet. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 9:1517?1561, 2010. [11] B. Sch olkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In COLT ’01/EuroCOLT ’01, pages 416?426. Springer-Verlag, 2001. [12] F. Dinuzzo and B. Sch olkopf. The representer theorem for Hilbert spaces: a necessary and sufficient condition. In Advances in Neural Information Processing Systems 25, pages 189? 196. 2012. [13] I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2:67?93, 2001. [14] A. Christmann and I. Steinwart. Universal kernels on non-standard input spaces. In Proceedings of Advances in Neural Information Processing Systems, pages 406?414. 2010. [15] N. A. Mehta and A. G. Gray. Generative and latent mean map kernels. CoRR, abs/1005.0188, 2010. [16] G. Blanchard, G. Lee, and C. Scott. Generalizing from several related classification tasks to a new unlabeled sample. In Advances in Neural Information Processing Systems 24, pages 2178?2186. 2011. [17] P. K. Shivaswamy, C. Bhattacharyya, and A. J. Smola. Second order cone programming approaches for handling missing and uncertain data. *Journal of Machine Learning Research*, 7:1283?1314, 2006. [18] H.S. Anderson and M.R. Gupta. Expected kernel for missing features in support vector machines. In Statistical Signal Processing Workshop, pages 285?288, 2011. [19] L. Fei-fei. A bayesian hierarchical model for learning natural scene categories. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 524?531, 2005. [20] D. G. Lowe. Object recognition from local scale-invariant features. In Proceedings of the International Conference on Computer Vision, pages 1150?1157, Washington, DC, USA, 1999. [21] A. Vedaldi, V. Gulshan, M. Varma, and A. Zisserman. Multiple kernels for object detection. In Proceedings of the International Conference on Computer Vision, pages 606?613, 2009.