

# Combining Multiple Kernels for Efficient Image Classification

Behjat Siddiquie  
Dept. of Computer Science  
University of Maryland  
College Park, 20742, USA  
behjat@cs.umd.edu

## ABSTRACT

We investigate the problem of combining multiple feature channels for the purpose of efficient image classification. Discriminative kernel based methods, such as SVMs, have been shown to be quite effective for image classification. To use these methods with several feature channels, one needs to combine base kernels computed from them. Multiple kernel learning is an effective method for combining the base kernels. However, the cost of computing the kernel similarities of a test image with each of the support vectors for all feature channels is extremely high. We propose an alternate method, where training data instances are selected for each of the base kernels using boosting. A composite decision function is learnt, which can be evaluated by computing kernel similarities with respect to only these chosen instances. This method significantly reduces the number of kernel computations required during testing. Experimental results on the benchmark UCI datasets, as well as on two challenging painting and chart datasets, are included to demonstrate the effectiveness of our method.

## Categories and Subject Descriptors

I.5.2 [Pattern Recognition]: Design Methodology-Classifier design and evaluation; I.4.7 [Image Processing and Computer Vision]: Feature Measurement-Feature Representation

## General Terms

Algorithms, Performance, Experimentation

## Keywords

kernel learning, paintings, charts, boosting

## 1. INTRODUCTION

We address the problem of combining multiple heterogeneous features for image classification. Categorizing images based on stylistic variations such as scene content and painting

genre requires reliance on a rich feature repertoire. Classification is accomplished by comparing distributions of features, e.g., color, texture, gradient histograms [13, 25, 20]. For instance, Grauman and Darrell proposed the Pyramid Match Kernel (PMK) to compute Mercer kernels between feature distributions for Support Vector Machine (SVM) based classification. This has been shown to be effective for object categorization [12], scene analysis [20] and document analysis [29]. Approaches such as PMK would compute a kernel matrix for each feature distribution. We explore techniques for combining the kernels from multiple features for efficient and robust recognition.

A number of techniques have been proposed to learn the optimal combination of a set of kernels for SVM-based classification. Lanckriet et al. proposed an approach for Multiple Kernel Learning (MKL) through semi-definite programming [19]. Sonnenburg et al. generalized MKL to regression and one-class SVMs, and enhanced the ability to handle large scale problems. Rakotomamonjy et al. increased the efficiency of MKL and demonstrated its utility on several standard datasets including the UCI repository [24]. They compute multiple kernels by varying the parameters of polynomial and Gaussian kernels, and apply MKL to compute an optimal combination. Bosch et al. learn the optimal mixture between two kernels - shape and appearance - using a validation set [6]. Varma and Ray propose to minimize the number of kernels involved in the final classification by including  $L_1$  norms of the kernel weights in the SVM optimization function [27]. Bi et al. proposed a boosting-based classifier that combines multiple kernel matrices for regression and classification [5].

The efficiency of MKL-based SVM classifiers during the testing phase depends upon the number of support vectors and the number of features. In general, multi-class problems requiring subtle distinctions entail a large number of support vectors. The computational cost is substantial when the kernels are complex, e.g., matching similarity of feature distributions. Is it possible to reduce the number of complex kernel computations while maintaining performance? We propose an approach for combining multiple kernels through a feature selection process followed by SVM learning. Let  $K_m(\cdot, \cdot)$  be the kernel values for the  $m^{\text{th}}$  feature channel computed using approaches such as Pyramid Match Kernel. The columns of  $K_m$  are considered to be features embedding the images in a high-dimensional space based on similarity to training examples. During the training phase, a subset of

the columns are chosen using Gentle Boost [1] based on their discriminative power, and a new kernel  $K$  is constructed. This is provided as input to an SVM for final classification. Kernels of test images need to be computed for only the chosen set of columns - much smaller than the full set of kernel values. This results in substantial reductions in computational complexity during the testing phase. The consequent approach is simple and relies on well understood techniques of Boosting and SVMs. Boosting methods have previously been used for feature selection [30], to learn kernels directly from data [9, 16], and for selecting a subset of kernels for concept detection in [17].

We compare our Boosted Kernel SVM (BK-SVM) approach with the Efficient Multiple Kernel Learning (EMKL) approach proposed in [24]. EMKL has been shown to increase the efficiency of kernel learning while enabling the use of a large number of kernels within SVM. It uses all the kernel values for classification - a superset of the features obtained by the greedy Boosting-based selection. BK-SVM and EMKL are tested in three scenarios: standard datasets from the UCI repository [2], a novel Painting dataset and a previously reported database of computer generated charts [29]. Results indicate that BK-SVM's classification accuracy is comparable to that of EMKL, with the additional advantage of a much smaller number of complex kernel computations.

The Painting dataset consists of nearly 500 images downloaded from the Internet - the task being to classify images into 6 genres. This provides a good testbed as the classification is subtle, requiring a large variety of features. Recently, there have been studies on the classification of paintings based on their style, artist, period and brushwork [32, 18, 15, 22, 31]. A semi-supervised method employing a variety of feature channels to annotate painting brushwork was presented in [32]. In [18] paintings are classified according to artist. Li et al. [22] have used 2D multi-resolution HMMs with multi-level Debauchies wavelet coefficient features to identify the artists of ancient Chinese paintings. In [21], high level semantic concepts are combined with low level image features to annotate paintings based on period, style and artist. In some of these methods such as [32, 21] a high level of domain knowledge has been used to develop the hierarchy of classes and to select appropriate image features. We use a large repertoire of simple features and rely on machine learning to obtain the combination best suited for the classification. This provides the potential for application in other categorization tasks.

The next section presents details of combining multiple kernels, followed by experiments on the UCI datasets. Section 4 presents the Painting dataset, the features used and the experimental results. Section 5 describes experiments on the Charts dataset, followed by concluding remarks.

## 2. LEARNING A MIXTURE OF KERNELS

Content-based image categorization typically represents images with histograms or distributions of features from channels such as texture, color and local gradients [10, 23]. Classification is performed by comparing such distributions. Grauman and Darrell [13] proposed the Pyramid Match Kernel (PMK) for efficiently computing Mercer kernels between

feature distributions and apply it to Support Vector Machine (SVM)-based object categorization [13]. A closely related approach used spatial distributions of features for scene recognition [20]. These techniques use SVM to learn the manifold of image categories and show good generalization. However, classifying images based on subtle style variations, e.g., painting genres, requires a large repertoire of feature channels. Techniques such as PMK would compute a kernel matrix for each feature channel. We are thus faced with the problem of determining the best mixture of the kernels for a given classification task.

A number of Multiple Kernel Learning (MKL) techniques have been proposed to compute linear combinations of kernels for classification by SVM [19, 24, 26]. Let  $\{K_1, K_2, \dots, K_M\}$  be the kernel matrices computed for various feature modalities. MKL computes an optimal classification kernel

$$K(q_i, q_j) = \sum_{m=1}^M \beta_m K_m(q_i, q_j) \quad (1)$$

where  $q_i$  are the training images. Recent MKL techniques have progressively improved training efficiency, e.g., [24, 7]. However, classifying a test image  $x$ , using a non-linear SVM, requires computing its kernel value with respect to the selected set of training support vectors  $S$  for all feature channels with  $\beta_m \neq 0$ , i.e.  $K_m(q, x) \forall q \in S$  and  $\forall m$  where  $\beta_m \neq 0$ . This has  $O(c\tilde{N}\tilde{M})$  computational complexity where

- $c$  is the complexity of computing the kernels. This is significant when matching the similarity of distributions.
- $\tilde{N}$  is the number of support vectors,  $|S|$ . Classification problems with difficult decision boundaries require a large set of support vectors. Some approaches propose to reduce this by approximating  $S$  with a reduced set of vectors, e.g., [8]. However, they are unsuitable for our case as each kernel is constructed from a different feature modality. Moreover, it is desirable to include as many training images as possible for good generalization (large  $N$ ).
- $\tilde{M} = |\{m | \beta_m \neq 0\}|$ . MKL methods reduce  $\tilde{M}$  by imposing sparsity constraints on the weights  $\beta_m$  [26]. However, this may not provide significant benefits when a large variety of features are required for classification.

Is it possible to reduce the number of kernel computations while maintaining performance?

Consider a vector constructed for a test image by concatenating its kernel values with all the training images. For an image  $x$ , this would be an  $NM$  dimensional vector

$$\mathbf{f}(x) = \langle K_1(q_1, x) \dots K_1(q_N, x) \dots K_M(q_1, x) \dots K_M(q_N, x) \rangle \quad (2)$$

We use Gentle Boost to determine the  $L$  most discriminative dimensions for a classification problem. The size of  $L$  is chosen such that  $|L| \ll \tilde{N}\tilde{M}$ . This results in a reduced dimensional vector for each image, denoted by  $\hat{\mathbf{f}}(x)$ . An SVM is trained to classify images based on the  $\hat{\mathbf{f}}$ 's. E.g., for

a linear SVM, the kernel between two images  $x$  and  $y$  would be

$$\Phi(x, y) = \sum_{\langle n, m \rangle \in L} K_m(x, q_n) K_m(q_n, y) \quad (3)$$

For each test image, this requires  $O(|L|)$  complex kernel  $K_m(\cdot, \cdot)$  computations, and  $O(N|L|)$  computations of a simpler kernel such as linear or RBF. This significantly reduces the computational complexity.

To better understand the nature of  $\Phi(\cdot, \cdot)$ , notice that the Pyramid Match Kernel between two images  $x$  and  $y$  can be abstracted as a dot-product between two bit-vectors,  $\psi_m(x)^T \psi_m(y)$ , where  $m$  is the feature channel [13]. Therefore, eq.(3) is equivalent to

$$\begin{aligned} \Phi(x, y) &= \sum_{\langle n, m \rangle \in L} \psi_m(x)^T \psi_m(q_n) \psi_m(q_n)^T \psi_m(y) \\ &= \sum_m \psi_m(x)^T \left[ \sum_{\langle n, m \rangle \in L} \psi_m(q_n) \psi_m(q_n)^T \right] \psi_m(y) \end{aligned} \quad (4)$$

The inner matrix,  $A_m = \sum \psi_m(q_n) \psi_m(q_n)^T$ , is a semi-definite matrix. It is easy to show that for a RBF SVM

$$\Phi(x, y) = \exp \frac{1}{\sigma^2} \sum_m \|\psi_m(x) - \psi_m(y)\|_{A_m}^2 \quad (5)$$

Intuitively,  $A$ 's are akin to covariance matrices of the exemplar images in  $L$ . When  $L$  is constructed to maximize discrimination between classes,  $A$  defines a discriminative projection.

We note that the approach does not restrict the number of support vectors chosen by the SVM. It only restricts the SVM's kernel to be based on a limited number of base kernel columns.

## 2.1 Boosting for Feature Selection

Discriminative feature selection is a well studied problem in machine learning, e.g., Xiao et al. propose a variant of boosting called Joint Boost for feature selection [30]. We use Gentle Boost for its simplicity and robustness [1, 11]. Let  $\mathbf{f}$ 's be  $D$  dimensional vectors.  $D$  is typically large; in our case  $D = NM$ . The basic version of Gentle Boost defines a set of weak learners  $h(\mathbf{f})$  where each  $h(\cdot)$  is a linear classifier along a single dimension. The algorithm iteratively chooses a set of weak learners to maximize classification accuracy. The weak learner chosen at the  $t^{\text{th}}$  iteration, namely  $h_t(\cdot)$ , is the one providing maximal increase in classification accuracy with respect to the set of previously chosen classifiers  $h_1, \dots, h_{t-1}$ . Thus, the choice of dimensions depends upon the incremental benefit relative to previous choices. In spite of the greedy nature of the selection process, Boosting has been shown to perform well in many classification tasks [11]. Outline of Gentle Boost:

- Given:  $(x_1, y_1), \dots, (x_n, y_n)$  where  $x_i \in X$  and  $y_i \in \{-1, 1\}$
- Initialize the weights  $D(i) = \frac{1}{m}$
- For  $t = 1, \dots, T$

- Choose confidence value  $\alpha_t \in R$
- Find the classifier  $h_t$  which minimizes the classification error with respect to the distribution  $D_t$
- Update the weights  $D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$  where  $Z_t$  is a normalization factor.

- $\{h_t\}$  are the selected features.

## 3. EXPERIMENTS WITH UCI DATASETS

The boosting-based feature selection is an efficient but greedy approach. To observe its performance penalties, BK-SVM was applied to four datasets from the UCI repository, specifically the Liver, Ionosphere, Pima and Sonar datasets. The kernels were simple polynomial and Gaussian functions. Here, the motivation was solely to empirically observe the performance for standard datasets. The efficiency gains become evident for more complicated kernel functions used later in the Painting and Chart datasets.

The classification results were compared with those of the Efficient Multiple Kernel Learning (EMKL) algorithm described in [24]. For each dataset, a large number of Gaussian and polynomial kernels are computed as described in [24]. The base kernels include Gaussian kernels with 10 different bandwidths  $\sigma$  on all variables and on each single variable, and also polynomial kernels of degree 1 to 3 on all variables and each single variable. EMKL and BK-SVM are used to learn a mixture of the kernels appropriate for classification. During the testing phase, the number of kernels computation required in EMKL is a product of the number of kernels selected and the number of support vectors. In case of BK-SVM, the complexity depends upon the number of kernel columns chosen by boosting.

The classification results are summarized in Table 1. They indicate that BK-SVM performs close to the baseline EMKL approach even though the number of kernel computations is more than an order of magnitude lower. The loss of performance of approximately 2% may be ascribed to the greedy selection of kernel columns. The results also demonstrate the scalability of our method, which performs comparably to EMKL even in the case of the Sonar dataset where a large number of kernels(793) are used for learning with only a small number of training examples(104). These trends are reflected in the experiments with painting and computer-generated chart datasets - described in the next sections. The modest loss in performance is outweighed by the large decrease in computational complexity, especially for complex kernel functions.

## 4. PAINTING DATASET

BK-SVM was applied to painting genre classification. A dataset of 81 Abstract Expressionist, 84 Baroque, 84 Cubist, 82 Graffiti, 89 Impressionist and 78 Renaissance (total of 498) paintings was collected from the Internet. The painting styles along with the painters of each style are listed in Table 2. Some of the public domain images are shown in Figure 1. The distinguishing features for painting styles are not clearly defined due to its abstract nature. There is high intra-class variation due to differences between the painters of a particular style and also between the different paintings of individual painters [3]. The content in the paintings varies

Table 1: Experiments on UCI Dataset

Dataset			BK-SVM		EMKL	
name	size	kernels	accuracy	kernel computations	accuracy	kernel computations
Liver	345	91	$66.0 \pm 5.0$	40	$65.0 \pm 2.3$	$1607 \pm 324$
Ionosphere	351	442	$92.0 \pm 4.0$	40	$92.3 \pm 1.4$	$1496 \pm 266$
Pima	768	117	$73.0 \pm 7.0$	60	$75.8 \pm 1.6$	$3123 \pm 526$
Sonar	208	793	$75.0 \pm 5.0$	20	$78.6 \pm 4.2$	$2538 \pm 351$

Table 2: Painting Classes

Painting Style	Artist
Abstract Expressionist	Arshile Gorky, Helen Frankenthaler, James Brooks, Jane Frank, Jean Paul Riopelle, Kenzo Okada, Paul Jenkins
Baroque	Anthony Van Dyck, Artemisia Gentileschi, Carravagio, Diego Velazquez, Jan Vermeer, Nicholas Poussin, Peter Paul Reubens, Rembrandt
Cubist	Fernand Leger, Georges Barque, Gino Severini, Jacques Villon, Juan Gris, Lyonel Feininger, Pablo Picasso
Graffiti	Anonymous
Impressionist	Alfred Sisley, Camille Pissarro, Claude Monet, Frederic Bazille, Mary Cassatt, Pierre Auguste Renoir, Edouard Manet
Renaissance	Correggio, Raphael, Leonardo Da Vinci, Sandro Botticelli, Titan, Giorgione, Pieter Brueghel, Michelangelo

significantly and occasionally paintings of different styles depict the same scene or object, further complicating the problem. Having been compiled from a variety of sources, the images have variations in scale and illumination as well. The classification task is complex, requiring a rich set of features. This makes the dataset a good testbed for BK-SVM.

## 4.1 Features

Inspired by previous studies on painting classification, a large variety of features are computed. Each feature channel produces a distribution of filter responses for a given image. The similarity of images is defined as the match between the distributions. To reduce the computational cost of extracting the features the images were resized such that their smaller side was a maximum of 1000 pixels while maintaining the aspect ratio.

### 4.1.1 Texture

Texture features capture brushwork and characteristics of the depicted scene. They have been shown to be effective in classification of paintings [22, 32, 15]. We employ the MR8 filter bank [28] as it responds to both isotropic and anisotropic textures and was observed to perform better than Gabor filter banks. The MR8 filter bank consists of a Gaussian and a Laplacian of Gaussian with  $\sigma = 10$  and oriented edge and bar filters at 3 scales  $(\sigma_x, \sigma_y) = \{(1, 3), (2, 6), (4, 12)\}$  and 6 orientations. Only the maximum response is recorded at each scale for each of the edge and bar filters across all orientations. This provides 8 responses at each pixel. The responses at all the pixels are combined to form a set of vectors, denoted by  $F_{\text{texture}}$ .

### 4.1.2 Histograms of Oriented Gradients (HOG)

HOG based descriptors have been extensively used for representing local shape in object recognition [10, 4, 23]. They

have some degree of invariance to illumination and geometric transformations. We compute two types of features using HOG:

1.  $F_{\text{HOGdense}}$ : set of HOG features on overlapping  $8 \times 8$  sized patches placed on a dense regular grid with a spacing of 4 pixels - similar to [10].
2.  $F_{\text{HOGsparse}}$ : sparse set of HOG features computed on  $8 \times 8$  patches centered on all edge points. This was inspired by [4].

### 4.1.3 Color

Color is probably the most important aspect of paintings. Color features have been previously employed for classifying paintings [15, 31]. We use local histograms to represent color features consisting of 10 bins of the pixel intensities of each color channel. The histograms are computed in  $8 \times 8$  sized patches centered on a dense grid over the image. This generates a set of vectors denoted by  $F_{\text{color}}$ . The histograms of different color channels were concatenated because the joint histograms are quite sparse and have a high dimensionality. Experiments indicated that RGB, HSV and LUV had similar performance. Only results for RGB color-space are presented here.

### 4.1.4 Edge Continuity

Edge Continuity is used to enhance the saliency of long continuous curves relative to scattered and cluttered edges. We use the technique described in [14] for computing the saliency maps of the images. HOG features are extracted from these saliency maps from patches centered on edges having high saliency. The obtained set of HOG vectors is denoted by  $F_{\text{HOGsal}}$ .

## 4.2 Pyramid Match Kernel

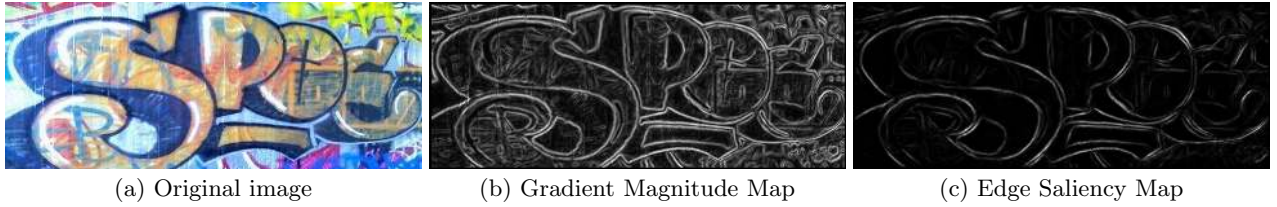


Figure 2: Salient Edges

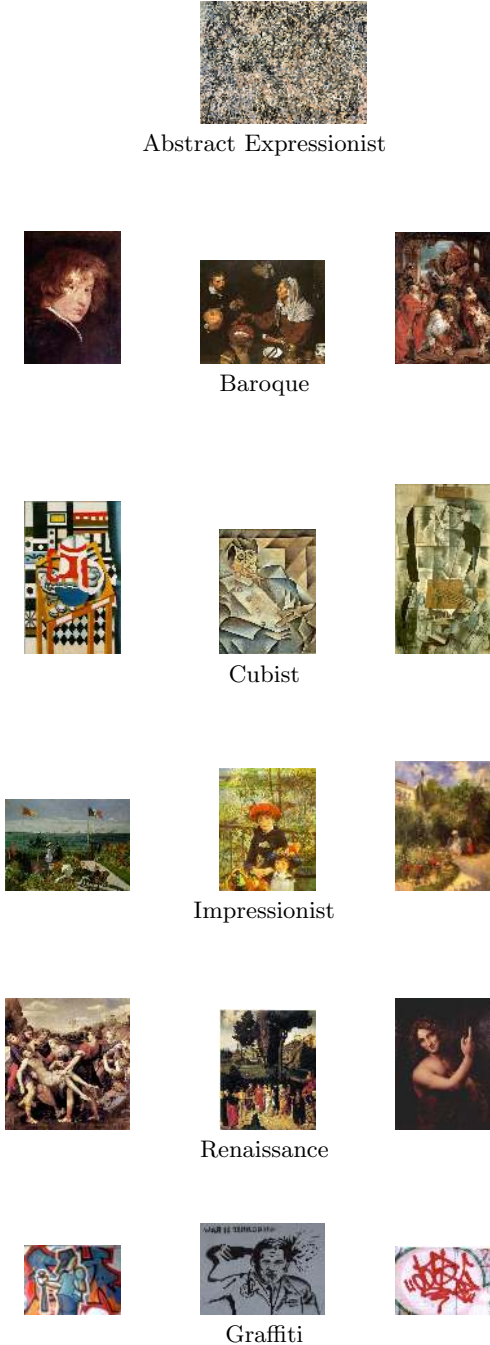


Figure 1: Examples of Paintings

Each of the described features produces a set of vectors for a given image. For each feature channel, similarity between images is computed based on the similarity between the two sets of vectors, computed using Pyramid Match Kernel (PMK) [12]. The sets can have different cardinalities. The approach has been shown to be efficient and effective for image classification. In this section we briefly describe the kernel. Let  $X$  and  $Y$  be two sets of feature vectors in a  $d$ -dimensional feature space. Now consider  $L + 1$  levels of histograms  $H^0, H^1, \dots, H^L$ . The level 0 of the histogram consists of just 1 bin which is the entire space, the level 1 of the histogram consists of  $2^d$  bins equally dividing the feature space into two parts along all dimensions. Similarly the level  $l$  of the histogram consists of  $D = 2^{dl}$  bins. Let  $H_X^l$  and  $H_Y^l$  denote the histograms of  $X$  and  $Y$  at level  $l$  with  $H_X^l(i)$  and  $H_Y^l(i)$  being the number of feature vectors of  $X$  and  $Y$  respectively falling into the  $i$ th bin at level  $l$ . A histogram intersection gives the number of matches at this level.

$$I(H_X^l, H_Y^l) = \sum_{i=1}^D \min(H_X^l(i), H_Y^l(i))$$

But note that all the matches at level  $l + 1$  are also matches at this level and hence the number of new matches at level  $l$  is  $I(H_X^l, H_Y^l) - I(H_X^{l+1}, H_Y^{l+1})$ . The matches at level  $l$  are weighted by  $\frac{1}{2^{L-l}}$  in order to give higher weights to matches which happen at smaller bin sizes and hence have a higher similarity. The total match between  $X$  and  $Y$  is defined as the similarity between  $X$  and  $Y$

$$K(X, Y) = I(H_X^L, H_Y^L) + \sum_{l=0}^{L-1} \frac{1}{2^{L-l}} I(H_X^l, H_Y^l) - I(H_X^{l+1}, H_Y^{l+1})$$

To avoid biasing the kernel toward larger input sets it is normalized

$$K(X, Y) = \frac{K(X, Y)}{\sqrt{K(X, X)K(Y, Y)}}$$

This normalization also ensures that  $\forall X, Y K(X, Y) \in [0, 1]$ . It has been shown that PMK may be abstracted as the dot product between two bit vectors. Therefore, it is a Mercer kernel and can be directly used in an SVM.

### 4.3 Classification Results

Training the BK-SVM consists of the following steps:

- Each of the  $M$  described features is extracted for all training images  $q_i$ .
- PMK was used to compute kernel values  $K_m(q_i, q_j)$ ,  $\forall q_i, q_j, m$ , producing  $M$  kernel matrices,  $K_1, \dots, K_M$ .

- A vector  $\mathbf{f}_i$  is constructed for each  $q_i$  by concatenating the kernel values as defined in eq.(2).
- Boosting is used to select a set of  $L$  dimensions that best classify  $\mathbf{f}_i$ 's into the painting genres. The number of exemplar images selected is equal to the number of iterations of boosting and thus can be easily controlled.
- A new RBF kernel matrix  $\Phi$  is constructed from the selected dimensions (i.e. columns of  $K_m$ 's) through the relation in eq.(5). A one-vs-all multi-class SVM is trained on  $\Phi$ .

During the testing phase, PMK is computed between a given test image and the  $L$  selected training images. Classification is performed through the trained RBF SVM.

For comparison, EMKL was employed for the same classification task. For EMKL, we learn separate kernels for each individual classifier, using the same parameters that were used for the UCI datasets ( $C = 100$ , maximal number of iterations=500, duality gap=0.01).

The experiments were repeated 10 times with a 5-fold cross-validation (80% training data and 20% test data) and each time the training and test sets were chosen randomly. The results are listed in Table 3. They indicate that both the EMKL and our method perform much better than each of the individual feature channels. Our method used on average 490 kernel similarity computations whereas the EMKL method requires about 2400 kernel similarity computations. Despite this reduction in complexity, the performance of our method is quite close to that of the EMKL method. Figure 5 plots the performance of our method as the number of selected features is varied. It can be seen that the performance is better than the individual feature channels even when there are very few exemplar images and it increases as the number of selected images increases.

The features, in general, perform quite well individually and also complement each other resulting in a significant improvement in performance when combined. In Baroque and Renaissance paintings, darker colors are used and this makes color histograms particularly useful for discriminating them from the other classes (Fig. 3). Color features are also useful for identifying Impressionist paintings to some extent as they tend to depict outdoor scenes with sunlight, landscapes and greenery.

The texture feature proved useful for distinguishing Impressionist images as they have distinctive brush strokes. Baroque paintings being darker, generate low responses with the filter banks and are again easily identified. Texture also distinguishes Renaissance from Graffiti paintings to some extent.

The cubist paintings are composed of dense geometrical structures such as straight lines, cubes, cylinders. Consequently, local shape features such as the dense HOG are useful in distinguishing them. The sparse HOG features encode the local shape around the edge points and prove useful for identifying Impressionist paintings.

Perceptually salient curves can be enhanced using edge continuity techniques[14]. Graffiti paintings tend to have smooth

continuous contours, which get enhanced in the saliency maps(Fig. 2) and the local shape features around these contours help discriminate them from other paintings. Saliency based features also help in identifying Renaissance paintings to some extent. The corresponding kernel similarity matrix is shown in Fig. 4.

Combining the feature channels helps reduce confusion caused by individual features. For instance, on the sole basis of color, a dark colored graffiti painting may be confused as a baroque painting. However, local shape information provided by saliency maps helps reduce this confusion. The confusion matrix obtained after combining features using BK-SVM is shown in Fig. 6. There is some degree of confusion between abstract expressionist and cubist paintings and upon examining the misclassifications, it was found out that most of the abstract expressionist paintings wrongly identified as cubist had the geometrical structures characteristic to cubist paintings. Some of the cubist paintings misclassified as abstract expressionist lacked these geometrical shapes. There are also some errors between impressionist and renaissance paintings.

To gain further insight into the construction of the individual one-vs-all classifiers, we looked at the average weights allocated by EMKL to the kernels for each individual classifier(Fig. 7). Color being an important feature was assigned a high weight in each of the individual classifiers and as expected, it turned out to be the most dominant feature for distinguishing Baroque paintings. Similarly the saliency kernel is weighted relatively high in the Cubist and Graffiti classifiers. Texture is also important to some extent in case of Baroque, Impressionist and Renaissance classes. We also observed that sparse HOG features are assigned extremely high weights in all the classifiers, indicating the significance of the local shape information represented by them. Dense HOG features are allocated high weights in the Cubist classifier as expected. On the whole, the weights seemed quite intuitive with features that distinguished a particular class well, being assigned a higher weight in the respective classifier. However, texture was weighted relatively low which is surprising, given the fact that it performs quite well individually. We conjecture that since both texture and HOG are based on local edges, they contain redundant information resulting in texture being ignored.

We did a similar study for BK-SVM, where we examined the proportion of exemplar images selected from each kernel for the individual classifiers(Fig. 8). Though some of the above mentioned trends were observed, like color and saliency being important for the baroque and graffiti paintings respectively, no single feature dominated the individual classifiers. We hypothesize that this is a result of the lack of any external constraints imposed by our method unlike the sparsity constraints imposed by EMKL. We are currently exploring the possibility of constraining the underlying boosting algorithm for improving the BK-SVM method.

## 5. CHART DATASET

We also compare the performance of our method to EMKL on the chart dataset [29] which consists of more than 650 images of computer generated charts belonging to 5 classes. Please refer to [29] for further details on the dataset as well

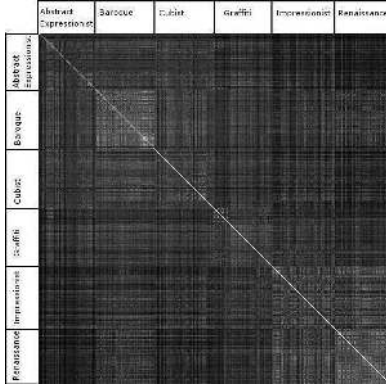


Figure 3: Color Histogram Kernel

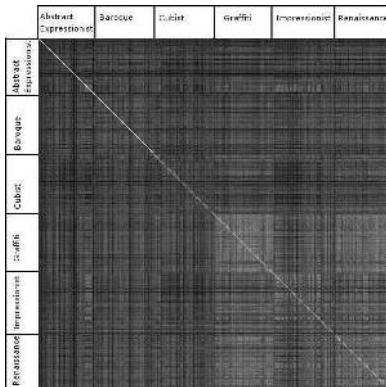


Figure 4: Edge Saliency Kernel

Table 3: Painting Classification Results

Feature	Accuracy
texture	$73.5 \pm 1.1$
color	$70.6 \pm 1.1$
dense HOG	$69.3 \pm 1.2$
sparse HOG	$69.0 \pm 1.0$
Saliency	$62.2 \pm 0.7$
Combined EMKL	$82.4 \pm 0.9$
Combined Our Method	$80.9 \pm 0.6$

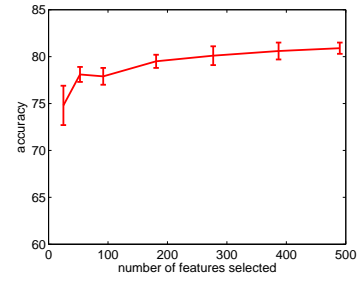


Figure 5: Variation in performance with the number of features for the painting dataset

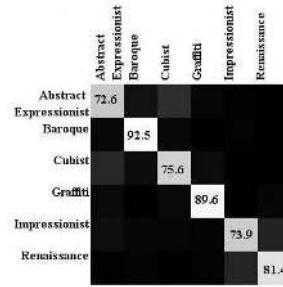


Figure 6: Confusion Matrix for the painting dataset

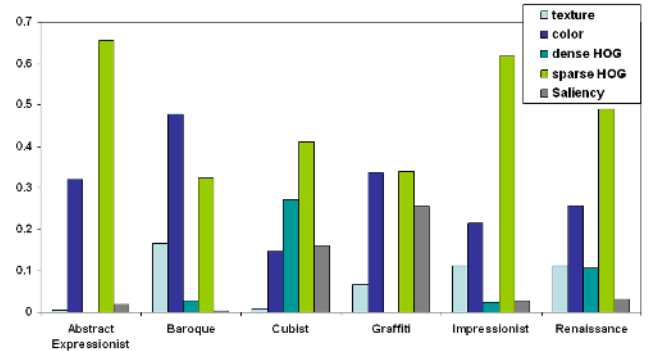


Figure 7: Average kernel weights learnt by EMKL for each classifier

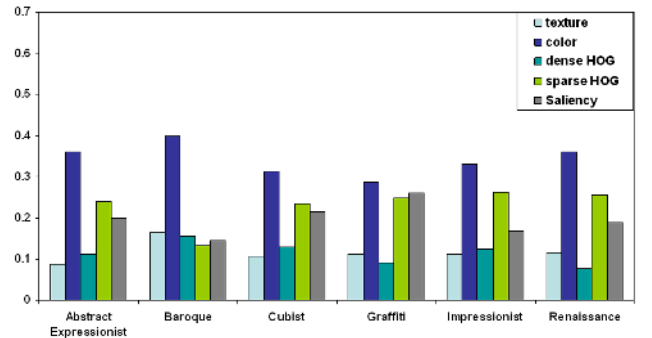


Figure 8: Average proportion of exemplar images selected from the feature channels for each classifier



**Table 4: Chart Classification**

Feature	Accuracy
HOG	$83.0 \pm 0.8$
Edge Continuity HOG	$67.2 \pm 0.9$
Edge Orient Histograms	$57.8 \pm 0.9$
Distance map	$82.2 \pm 0.6$
Edge Continuity locn.	$72.5 \pm 0.8$
Region Segmentation	$29.2 \pm 0.7$
Combined EMKL	$87.6 \pm 0.5$
Combined Our Method	$86.3 \pm 0.7$

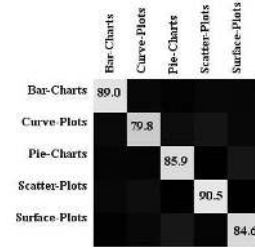
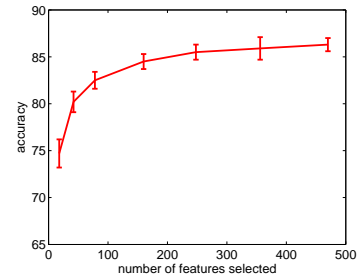
as the features used for the classification. As in the painting classification, we compute kernels for the individual features using PMK from which the composite kernels are learnt. The experiments are performed using the same protocol as the painting classification. The results are summarized in Table 4. Again, the performance of our method is quite close to that of EMKL and both of them are significantly better than the individual feature channels. The average number of kernel computations required by our method is 470 compared to 3200 kernel computations required by the EMKL method. Figure 10 plots the variation of performance with the number of features selected and again the performance increases with increasing number of features. The confusion matrix is shown in Fig. 9.

In both the Painting and the Chart datasets, BK-SVM requires nearly 5 times fewer kernel computations than EMKL for achieving comparable classification accuracy. This reduction, though substantial, is less compared to the 40-120 time reduction achieved on the UCI datasets. There are two plausible explanations:

- The painting and chart datasets have multiple classes, which makes the decision boundaries more complex than in case of the UCI datasets, which have only two classes.
- The UCI dataset experiments use base kernels produced by varying the parameters of Gaussian and polynomial kernels, many of which are likely to be redundant. Hence, a relatively sparse set of features selected by Boosting is sufficient to accurately approximate the optimal kernel. In case of the paintings and the chart datasets, each of the base kernels are computed from different feature channels and contain complementary information. Consequently, a number of exemplar instances are selected for each of the base kernels.

## 6. SUMMARY

This paper has presented a simple and efficient approach for learning a mixture of kernels. Our method, which works by greedily selecting exemplar data instances corresponding to each kernel using adaboost, has been shown to compare well to multiple kernel learning methods, while simultaneously reducing the number of kernel similarity computations required. The effectiveness of our method has been demonstrated on some of the benchmark UCI datasets. We have also tested our method on two extremely diverse and challenging image datasets, where a single feature channel is

**Figure 9: Confusion Matrix for the chart dataset****Figure 10: Variation in performance with the number of features for chart dataset**

not adequate for classification. We combine multiple kernels computed from different feature channels using the pyramid match method[13], obtaining results comparable to the MKL method. The results provide evidence that our approach is almost as accurate as the multiple kernel learning method, while being computationally much more efficient. We are also looking into combining these kernels in an unsupervised or semi-supervised manner to perform clustering.

## 7. REFERENCES

- [1] <http://research.graphicon.ru/machine-learning/gml-adaboost-matlab-toolbox.html>.
- [2] <http://archive.ics.uci.edu/ml/>.
- [3] R. Arnheim. Art and visual perception. a psychology of the creative eye. 1955.
- [4] S. Belongie, J. Malik, and J. Puzicha. Shape matching and object recognition using shape contexts. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(4):509–522, Apr 2002.
- [5] J. Bi, T. Zhang, and K. P. Bennett. Column-generation boosting methods for mixture of kernels. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 521–526, New York, NY, USA, 2004. ACM.
- [6] A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. In *CIVR '07: Proceedings of the 6th ACM international conference on Image and video retrieval*, pages 401–408, New York, NY, USA, 2007. ACM.



- [7] A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. In *CIVR '07: Proceedings of the 6th ACM international conference on Image and video retrieval*, pages 401–408, New York, NY, USA, 2007. ACM.
- [8] C. J. C. Burges. Simplified support vector decision rules. In *International Conference on Machine Learning*, pages 71–77, 1996.
- [9] K. Crammer, J. Keshet, and Y. Singer. Kernel design using boosting. In *NIPS*, pages 537–544, 2002.
- [10] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, 1:886–893, 2005.
- [11] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: A statistical view of boosting. *The Annals of Statistics*, 28(2):337–374, 2000.
- [12] K. Grauman and T. Darrell. The pyramid match kernel: discriminative classification with sets of image features. *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, 2:1458–1465 Vol. 2, 17-21 Oct. 2005.
- [13] K. Grauman and T. Darrell. The pyramid match kernel: Efficient learning with sets of features. *J. Mach. Learn. Res.*, 8:725–760, 2007.
- [14] G. Guy and G. Medioni. Inferring global perceptual contours from local features. *IJCV*, 20(1/2):113–133, 1996.
- [15] H. v. d. Herik and E. O. Postma. Discovering the visual signature of painters. *Future Directions for Intelligent Systems and Information Sciences*, pages 129–147, 2000.
- [16] T. Hertz, A. B. Hillel, and D. Weinshall. Learning a kernel function for classification with small training samples. In *ICML '06: Proceedings of the 23rd international conference on Machine learning*, pages 401–408, New York, NY, USA, 2006. ACM.
- [17] W. Jiang, S.-F. Chang, and A. Loui. Kernel sharing with joint boosting for multi-class concept detection. *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pages 1–8, 17-22 June 2007.
- [18] D. Keren. Recognizing image "style" and activities in video using local features and naive bayes. *Pattern Recogn. Lett.*, 24(16):2913–2922, 2003.
- [19] G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *J. Mach. Learn. Res.*, 5:27–72, 2004.
- [20] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR '06: Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2169–2178, Washington, DC, USA, 2006. IEEE Computer Society.
- [21] L. Leslie, T.-S. Chua, and J. Ramesh. Annotation of paintings with high-level semantic concepts using transductive inference and ontology-based concept disambiguation. In *MULTIMEDIA '07: Proceedings of the 15th international conference on Multimedia*, pages 443–452, New York, NY, USA, 2007. ACM.
- [22] J. Li and J. Z. Wang. Studying digital imagery of ancient paintings by mixtures of stochastic models. *IEEE Transactions on Image Processing*, 13(3):340–353, 2004.
- [23] D. Lowe. Distinctive image features from scale-invariant keypoints. In *International Journal of Computer Vision*, volume 20, pages 91–110, 2003.
- [24] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. More efficiency in multiple kernel learning. In *ICML '07: Proceedings of the 24th international conference on Machine learning*, pages 775–782, New York, NY, USA, 2007. ACM.
- [25] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover's distance as a metric for image retrieval. *Int. J. Comput. Vision*, 40(2):99–121, 2000.
- [26] S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf. Large scale multiple kernel learning. *J. Mach. Learn. Res.*, 7:1531–1565, 2006.
- [27] M. Varma and D. Ray. Learning the discriminative power-invariance trade-off. *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pages 1–8, 14-21 Oct. 2007.
- [28] M. Varma and A. Zisserman. Classifying images of materials: Achieving viewpoint and illumination independence. *European Conference on Computer Vision*, pages 255–271, 2002.
- [29] S. Vitaladevuni, B. Siddiquie, J. Golbeck, and L. Davis. Classifying computer generated charts. *Content-Based Multimedia Indexing, 2007. CBMI '07. International Workshop on*, pages 85–92, 25-27 June 2007.
- [30] R. Xiao, W. Li, Y. Tian, and X. Tang. Joint boosting feature selection for robust face recognition. *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, 2:1415–1422, 2006.
- [31] M. Yelizaveta, C. Tat-Seng, and A. Irina. Analysis and retrieval of paintings using artistic color concepts. *Multimedia and Expo, 2005. ICME 2005. IEEE International Conference on*, pages 1246–1249, 6-6 July 2005.
- [32] M. Yelizaveta, C. Tat-Seng, and J. Ramesh. Semi-supervised annotation of brushwork in paintings domain using serial combinations of multiple experts. In *MULTIMEDIA '06: Proceedings of the 14th annual ACM international conference on Multimedia*, pages 529–538, 2006.