

Two-stage Augmented Kernel Matrix for Object Recognition

Muhammad Awais, Fei Yan, Krystian Mikolajczyk, and Josef Kittler

Centre for Vision, Speech and Signal Processing (CVSSP),
University of Surrey, UK

{m.rana, f.yan, k.mikolajczyk, j.kittler}@surrey.ac.uk

Abstract. Multiple Kernel Learning (MKL) has become a preferred choice for information fusion in image recognition problem. The aim of MKL is to learn the optimal combination of kernels formed from different features, thus, to learn the importance of different feature spaces for classification. Augmented Kernel Matrix (AKM) has recently been proposed to accommodate for the fact that a single training example may have different importance in different feature spaces, in contrast to the MKL that assigns the same weight to all examples in one feature space. However, the AKM approach is limited to small datasets due to its memory requirements.

We propose a novel two stage technique to make AKM applicable to large data problems. In the first stage various kernels are combined into different groups automatically using kernel alignment. Next, the most influential training examples are identified within each group and used to construct an AKM of significantly reduced size. This reduced size AKM leads to the same results as the original AKM. We demonstrate that the proposed two stage approach is memory efficient and leads to the better performance than the original AKM and is robust to noise. The results are compared with other state-of-the art MKL techniques, and show improvement on challenging object recognition benchmarks.

1 Introduction

Object and image recognition has undergone a rapid progress in the last decade due to advances in both features design and kernel methods [1], [2], [3] in machine learning. In particular, the recent introduction of multiple kernel learning methods set a new direction of research in computer vision and object recognition. The state-of-the-art object and image recognition algorithms use multiple kernel learning based methods for classification, dimensionality reduction and clustering in a wide range of applications [1], [2], [3], [4]. Due to importance of complementary information in MKL, much research was done in the field of feature design [5], [6] to diversify kernels, leading to large number of kernels in typical visual classification tasks. Kernels are often computed independently of each others thus may be highly informative, noisy or redundant. Proper selection and fusion of kernels is therefore crucial to maximize the performance and to address the efficiency issues in large scale visual recognition applications.

The key idea of Multiple Kernel Learning (MKL), in case of SVM, is to learn a linear combination of given base kernels by maximizing the soft margin between two classes. The MKL was first proposed by Lancriet et al. [7] using semi-definite programming, where the kernel weights were learned by ℓ_1 -norm

regularization. Since the algorithm proposed in [7] was limited to small kernel sizes and low number of kernels, a number of other methods were proposed to address these problems [8], [9]. Different formulation for MKL primal are compared in [10], which also extend MKL to multiclass. All these MKL methods focus on linear combination of kernels, in which a single kernel corresponding to a particular feature space is attributed a single weight. This is a strong constraint as it does not exploit the information from individual samples in different feature spaces, e.g., in the context of object recognition, some samples can carry more shape information while others may carry more texture information for the same object category. To address this problem augmented kernel matrix (AKM) was proposed [11] in which different features extracted from the same sample are treated as different samples of the same class. Despite the improvement in classification performance the fundamental problem with AKM is its large augmented matrix which requires a lot of memory and makes it inapplicable to large datasets. In this paper we derive the primal and dual of the AKM, discuss its empirical feature space and address its issues with a two stage architecture. In the first stage, groups are formed from a set of base kernels based on the similarity between kernels. Next, a representative kernel for each group is learned by a linear combination of within group kernels. These representative kernels are highly informative containing most of the information from each group. Our grouping approach is also useful for methods proposed in [12], [13], which assumed that the kernel groups are available. We further reduce the complexity of AKM by exploiting the independence of empirical feature spaces of representative kernels in the augmented kernel matrix. Due to the independence, only the most influential training examples from the representative kernels can be used to build an AKM of a reduced size without compromising its performance. In the second stage, the AKM scheme is used to include the contribution of the most influential samples from all the representative kernels in the final classifier. Our experiments show that the proposed strategy of grouping kernels and selecting subsets of training examples makes the approach efficient and improves the classifier performance. The AKM results are compared to other MKL techniques, using different regularization, ℓ_1 , ℓ_2 , and ℓ_∞ norms. We demonstrate significant improvement on challenging object recognition benchmark Pascal VOC 2007 [14] and multiclass Oxford flower datasets [15], [16]. Moreover, the proposed memory efficient learning strategy is also applicable in other MKL techniques which is particularly important in large scale data scenario.

The rest of paper is organized as follows. In section 2 we discuss the structure of AKM matrix and derive its primal and dual for SVM. We then compare empirical feature spaces of a linear combination MKL and AKM schemes. Our proposed two stage multiple kernel learning for AKM is presented in section 3. In section 4 we present the result and compare with other state-of-art MKL methods for object recognition.

2 Linear Combination vs Augmented Kernel Matrix

In this section, we first present the structure of AKM and give primal formulations for a binary classification. We then present the concept of empirical feature space for the AKM scheme. In the next section, we illustrate feature spaces for linear combination and AKM schemes with a toy example.

Consider we are given m training samples $(x_i, y_i)_{i=1, \dots, m}$, where x_i is the sample in the input space and $y_i \in \pm 1$ is its label. Feature extraction results in

n training kernels $(K_p)_{p=1,\dots,n}$ of size $m \times m$ and corresponding n test kernels $(\tilde{K}_p)_{p=1,\dots,n}$ of size $m \times l$. Each kernel $K_p = \langle \Phi_p(x_i), \Phi_p(x_j) \rangle$ implicitly maps samples x_i from the input space to the feature space with mapping function $(\Phi_p(x_i))_{p=1,\dots,n}$. In MKL the aim is to find linear combination $\sum_{p=1}^n \beta_p K_p$, normal vector \mathbf{w} and bias b of the separating hyperplane simultaneously such that the soft margin between two classes is maximized. The primal and its corresponding dual for a linear combination of kernels are derived for various formulations in [17], [7], [8], [9] and compared in [10]. The dual problem can be solved by several existing MKL approaches, e.g., using SDP [7], SMO [8], SILP [9]. The decision function is then $f(x) = \text{sign}(\sum_{i=1}^m \alpha_i y_i k(x_i, x) + b)$, where $k(x_i, x)$ is the dot product of test sample x with the i^{th} training sample in the feature space. The Lagrange multiplier $\alpha \in \mathbb{R}^m$, and b are learnt by maximizing the margin. The contribution of a given feature channel is fixed by β_p , which may be suboptimal, as in a particular feature channel one example can carry more shape information than texture or vice versa. In contrast, in AKM [11], given the set of base training kernels $(K_p)_{p=1,\dots,n}$ the augmented kernel is defined as follows:

$$K = K_1 \oplus \dots \oplus K_n = \begin{bmatrix} K_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & K_n \end{bmatrix} \quad (1)$$

where the base kernels are on the diagonal. The zeros on the off diagonal reflect that there is no cross terms between different kernel matrices. Note that all the base kernels are of size $m \times m$ while the AKM is of size $(n \times m) \times (n \times m)$, thus it uses $n \times m$ training samples instead of m . The SVM primal of AKM scheme is then given:

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{p=1}^n \langle \mathbf{w}_p, \mathbf{w}_p \rangle + C \sum_{i=1}^{n \times m} \xi_i \\ \text{w.r.t.} \quad & w \in \mathbb{R}^{r_1 + \dots + r_n}, \xi \in \mathbb{R}^{n \times m}, b \in \mathbb{R} \\ \text{s.t.} \quad & y_i \left(\sum_{p=1}^n \langle \mathbf{w}_p, \Phi_p(x_i) \rangle + b \right) \geq 1 - \xi_{pi}, \quad \xi_{pi} \geq 0, \quad i = 1, \dots, m, \quad p = 1, \dots, n \end{aligned} \quad (2)$$

The dual optimization problem of equation 2 can be derived using Lagrange multiplier techniques:

$$\begin{aligned} \max \quad & \sum_{p=1}^n \sum_{i=1}^m \alpha_{pi} - \frac{1}{2} \sum_{p=1}^n \sum_{i,j=1}^m \alpha_{pi} \alpha_{pj} y_i y_j k_p(x_i, x_j) \\ \text{w.r.t.} \quad & \alpha \in \mathbb{R}^{n \times m} \\ \text{s.t.} \quad & \sum_{p=1}^n \sum_{i=1}^m \alpha_{pi} y_i = 0, \quad 0 \leq \alpha \leq C, \end{aligned} \quad (3)$$

The decision function of AKM is $f(x) = \text{sign}(\sum_{p=1}^n \sum_{i=1}^m \alpha_{pi} y_i k_p(x_i, x) + b)$, where α_{pi} are Lagrange multipliers and x is the test sample. Note that the same samples from different feature channels are added as separate examples of the

same class, therefore one Lagrange multiplier α_{pi} is learnt for each sample from each feature channel.

The concept of empirical feature space is crucial to analyze the spread and shape of the data. Kernel matrices consist of dot products between samples in some feature spaces. These feature spaces are usually very high or even infinite dimensional. However, in [18] it is shown that there exists an empirical feature space in which the intrinsic geometry of the data is identical to true feature space, thus, in many problems it is sufficient to study the empirical feature space. Empirical feature spaces X and \tilde{X} for training kernel K of size $m \times m$ and test kernel \tilde{K} of size $m \times l$ can be derived by eigen value decomposition as shown in [11].

Consider a linear combination of two training kernels K_1, K_2 with the sample points in r_1, r_2 dimensional empirical feature space given by matrices X_1, X_2 of sizes $r_1 \times m$ and $r_2 \times m$, respectively. By the definition of a dot product, computing the weighted sum of base kernels is equivalent to computing the cartesian product of the associated empirical feature spaces, after scaling them with $\sqrt{\beta_p}$, $p = 1, ..n$. An illustration of the empirical feature space is given in figure 1. K_1, K_2 are two base kernels with rank $r_1 = r_2 = 1$ i.e., the samples live in one dimensional empirical feature space as shown in figure 1(a) and (b). Note, this toy example is for illustration purpose, whereas, in practice the empirical feature spaces can be up to m dimensional. Figure 1(c) shows the empirical feature space of a sum of two kernels. Note that the number of samples in figure 1(c) is equal to m which is the same as the number of samples in K_1 and K_2 .

Consider the augmented kernel matrix K of two training kernels K_1, K_2 . The matrix X of training vectors in the empirical feature space associated with K can be computed by eigen value decomposition [11]. However, by exploiting the property of block diagonal augmented matrix \tilde{K} , its associated matrix X is directly given by:

$$X = \begin{pmatrix} X_1 & 0 \\ 0 & X_2 \end{pmatrix} \quad (4)$$

where X is a block diagonal matrix of size $(r_1 + r_2) \times 2m$, with matrix X_1 and X_2 on its diagonal. The empirical feature space for augmented kernel matrix from two one-dimensional kernels K_1 and K_2 is shown in figure 1(d). Note that there are now total of $2m$ training examples in the empirical feature space of AKM.

3 Two-Stage Multiple Kernel Learning

In this section we present a two stage architecture for multiple kernel learning which combines the MKL and AKM schemes. Kernel matrix of AKM needs large amount of memory and is very slow in training of classifier. For example, the extra memory required by cross terms in a large augmented kernel matrix of n base kernels is $n(n - 1)$ times larger than linear combination of these base kernels. This makes AKM less inapplicable to large datasets especially when n is large. We address this problem by introducing grouping of base kernels followed by a selection of training samples. Two stage approach serves two goals. It addresses the memory problems of AKM but also filters out noisy and redundant feature channels. Adding redundant feature channels as separate examples increases the memory requirements in AKM and adding noisy feature channel as

separate examples leads to a significant performance loss. These two problems are alleviated by applying the grouping stage.

3.1 Kernel Grouping

We define multiple groups of base kernels using a similarity criterion. One such grouping criterion can be based on the modality of features or their extraction technique. For example, feature channels based on colour can belong to one group, texture based feature channels to another group and shape based ones to yet another group. However, this kind of grouping is not automatic and needs prior information about input spaces of kernel which may not be available. We exploit Kernel Alignment [19] as a measure of similarity between kernels to group them in an unsupervised manner. Given an unlabeled sample set $S = \{x_i\}_{i=1}^m$, we use the Frobenius inner product between kernel matrices i.e., $\langle K_1, K_2 \rangle_F = \sum_{i,j=1}^m K_1(x_i, x_j) K_2(x_i, x_j)$. The empirical alignment between kernels with respect to the set S is defined as:

$$\hat{A}(S, K_1, K_2) = \frac{\langle K_1, K_2 \rangle_F}{\sqrt{\langle K_1, K_1 \rangle_F \langle K_2, K_2 \rangle_F}} \quad (5)$$

where K_i is the kernel matrix for the sample S . In [19] concentration and generalization of kernel alignment was introduced and proved. Concentration means that the probability of an empirical estimate deviating from its mean can be bounded as an exponentially decaying function of that deviation. In other words, the alignment is little dependent on the training set S as shown by theorem 3 in [19]. Generalization (test error) of a simple classification function is related to the value of the alignment as shown by theorem 4 in [19].

Using kernel alignment $\hat{A}(S, k_1, k_2)$ defined in equation 5 as a similarity measure we perform agglomerative clustering to find g groups of kernels. We initialize all kernels as clusters. We then merge two most similar clusters at a time. The similarity between two clusters is defined as the largest distance between all possible pairs of the clusters members. This continues until g groups are obtained. We used agglomerative as opposed to k-means to make it independent to initialization. Kullback-Leibler divergence can also be used as a similarity criterion between kernels [20], [21].

Learning a linear combination of kernels within a group can discard or down-weight redundant or noisy kernels thus result in a better kernel. Moreover, linear combination leads to more compact representation without loss of information. Therefore, for each group, MKL-SVM methods using ℓ_1 , ℓ_2 and ℓ_∞ norms are applied to obtain the representative kernels. The kernel that obtains the highest score on the validation data is used as group representant. Thus, the grouping and within group combination results in a set of representative kernels containing most of the information from various feature channels.

3.2 Selection of Training Samples

Kernel grouping partially addresses the issue of large AKM matrix. However, the matrix can be further reduced without compromising the performance by selecting only the samples from representative kernels which are crucial for classification. The decision function of SVM is determined by the α_i , one for each training sample. The α_i are non-zero for the support vectors only. Hence, for a

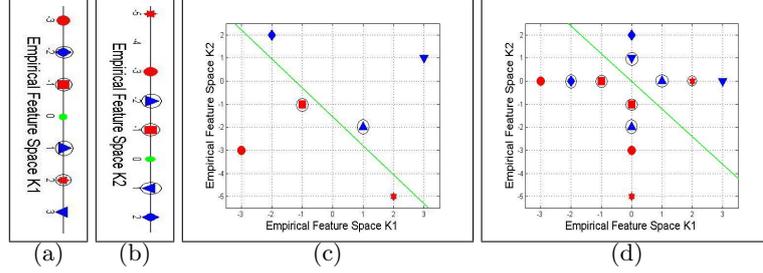


Fig. 1. Empirical feature spaces for Multiple kernels: (a) empirical feature space for K_1 ; (b) empirical feature space for K_2 ; (c) empirical feature space for $K_1 + K_2$; (d) empirical feature space for $K_1 \oplus K_2$.

single kernel, support vectors are sufficient for classification and all other samples can be discarded without performance degradation. This is supported by the fact that the feature spaces do not interfere with each other due to the structure of the augmented kernel matrix (cf. equation 1). It can be proved by considering the dual of AKM, equation 3, which can be rewritten as follows:

$$\begin{aligned}
 \max \quad & \sum_{i=1}^m \alpha_{1i} - \frac{1}{2} \sum_{i,j=1}^m \alpha_{1i} \alpha_{1j} y_i y_j \langle \Phi_1(x_i), \Phi_1(x_j) \rangle + \dots + \quad (6) \\
 & \sum_{i=1}^m \alpha_{ni} - \frac{1}{2} \sum_{i,j=1}^m \alpha_{ni} \alpha_{nj} y_i y_j \langle \Phi_n(x_i), \Phi_n(x_j) \rangle \\
 \text{w.r.t. } & \alpha \in \mathbb{R}^{n \times m} \\
 \text{s.t. } & \sum_{i=1}^m \alpha_{1i} y_i + \dots + \sum_{i=1}^m \alpha_{ni} y_i = 0, \quad 0 \leq \alpha \leq C,
 \end{aligned}$$

The first constraint in equation 6 is the sum of constraints for n kernels. The support vectors for all individual kernels together satisfy this constraint and thus lie in the feasible set of the optimization problem in equation 6. This is also illustrated by a toy example of a binary classification in figure 1. All the support vectors in empirical feature space for base kernels K_1 and K_2 are shown by the enclosing black circles and the hyperplane is represented in green at origin in figure 1(a) and (b), respectively. Figure 1(c) shows the empirical feature space of unweighted linear combination of base kernel. There are only two support vectors in figure 1(c), and the classes are separated by the hyperplane. However, the separability of the training set does not necessarily guarantee better performance as it depends upon the generalization to the test set [1]. Figure 1(d) is the empirical feature space of AKM combination of base kernels. Feature spaces of two base kernels are orthogonal to each other. There are $2m$ training samples and all the support vectors of kernel K_1 and K_2 are support vectors of AKM due to the orthogonality of their feature space. It is clear from equation 6 and figure 1, that the support vectors of representative kernels from each group are sufficient to construct the AKM matrix as the Lagrange multipliers of support vectors lie in the feasible set of equation 6. The use of support vectors only for different combinations of kernels is validated empirically in section 4.

4 Experiments and Discussion

This section presents the experimental results on challenging binary and multi-class object recognition datasets Pascal VOC 2007 [14], Oxford Flower 17 [15] and Oxford Flower 102 [16].

Pascal VOC 2007 [14] consists of 20 object classes with 9963 image examples. The images include indoor and outdoor scenes, truncated and occluded objects at various scales and different lighting conditions. The classification of 20 object categories is handled as 20 independent binary classification problems (as recommended by organizers of Pascal challenge). We present the results using average precision (AP) [14], which is proportional to the area under precision recall curve. Mean average precision (MAP) is computed by averaging scores for all 20 classes.

We compute 20 kernels by combining features introduced in [6] with 2 sampling strategies (dense, interest points) and spatial location grids [22]: whole image (1x1), horizontal bars (1x3), vertical bars (3x1) and image quarters (2x2). To form a codebook of 4000 visual words the descriptors are clustered using k-means. Each spatial grid is then represented by histograms of codebook occurrences and a separate kernel matrix is computed based on the χ^2 distance between histograms.

In the experiments we use SVM to compare several kernel combination schemes and the two stage AKM scheme proposed in this paper. The multiple kernel SVM (MK-SVM) schemes differ by the regularization norms used during learning, which include ℓ_1 [9], ℓ_2 [17], and ℓ_∞ (equal weights). We divide the 20 kernels into 4 groups as discussed in section 3.1. For each group, MKL-SVM methods using ℓ_1 , ℓ_2 and ℓ_∞ norms are applied to obtain the representative kernels. The results for various learning techniques are presented in table 1.

The consistently lower performance of ℓ_1 -norm, which typically leads to sparsely selected kernels, indicates that most of the base kernels carry complementary information. Therefore, the non-sparse multiple kernel methods, ℓ_2 -norm and ℓ_∞ -norm, give better results. The proposed two stage AKM scheme outperforms the other MKL combination schemes. In case of ℓ_2 within group and AKM between groups, (AKM, ℓ_2), we obtain an improvement of 0.6%, and in case of ℓ_∞ within group and AKM between groups, an improvement of 0.7% over all linear combinations of MKL-SVM. In case of informative kernels, the use of kernel grouping achieves comparable performance to the corresponding non-grouping schemes. The best performance of the state-of-the-art multiple kernel learning for these kernels is 62.1% , as shown in table 1 while the performance of win method for this challenge is 59.4% [14]. We beat the wining method by 3.4%, moreover, the 0.7% improvement by the proposed two stage AKM over state-of-the-art MKL is still significant given that all the kernels are highly informative due to carefully designed features. For example, the leading methods in PASCAL VOC often differ by a fraction of a percent in MAP. It is important to note that AKM on its own is giving 61.0%, however, when it is used together with grouping stage it is performing 1.8% better. It is because the linear combination within grouping stage gives good representative kernel with less noisy or redundant data. These highly informative representative kernel should be combined with AKM scheme so that information in each example of these kernels is exploited. We expect the grouping scheme to show better performance if there are noisy or redundant kernels in the set as shown by the noisy feature channels experiment in next section.

Table 1. MAP of PASCAL VOC 2007 with various MKL and AKM approaches.

	within group	no grouping	linear ℓ_1	linear ℓ_2	linear ℓ_∞
between groups					
linear ℓ_1		56.0	55.3	56.5	56.6
linear ℓ_2		61.4	60.8	61.3	56.5
linear ℓ_∞		62.1	61.1	62.1	62.0
AKM		61.0	60.8	62.7	62.8

We have also validated empirically the selection of support vectors for AKM on 20 binary classification problems of the Pascal 2007 [14] dataset. Only 0.3% to 0.5% of the support vectors of the AKM differs from the union of individual support vectors of representative kernels, while the MAP results are the same up to sixth decimal place. However, due to the use of the significant examples only, we are using 3 to 4 times less samples per base kernel. Hence, size of the AKM matrix is 60% to 70% less than the original size without compromising the performance. It is important note that it is not possible to apply AKM without the selection of significant examples in this benchmark due to memory requirements. We have used 4 groups of kernels thus the augmented kernel matrix is even smaller than the original training matrix of size 5011×5011 . Note that for each group a classifiers has to be trained i.e. 4 in this experiment. This is however done efficiently on small kernels and acceptable considering the performance gain achieved over other multiple kernel learning methods. Moreover, in α -step of alternative MKL techniques [9], [17] we have to train the linear combination of base kernels for different regularization norms several times before obtaining optimal weights values β for base kernels. All the results presented for AKM in this paper are obtained using “support vectors only scheme”.

Oxford Flower 17 [15] dataset contain pictures of some of the common flowers in the UK. It consists of 17 categories with 80 images in each category. There is large variation within a category and similarity with other categories. Dataset is split into training (40 images per class), validation (20 images per class) and test (20 images per class) using 3 predefined random splits. For the experiments we have used used seven χ^2 distance matrices provided online. The features used to compute these distance matrices include different types of shape, texture and color based descriptors [15]. We have used SVM as the base classifier and follow one-vs-all setup for multiclass classification [15]. We train an AKM classifier for each category and use the maximum response of the classifiers for each example to obtain the label and score for evaluation. Regularization parameter for the SVM is in the range $C \in \{10^{(-2,-1,\dots,3)}\}$.

The results are given in figure 2(a). For comparison we use recent evaluation results from [23] of state-of-the-art feature fusion techniques including MKL and boosting based classifier fusion. There are two baseline techniques, MKL-prod-SVM and MKL-avg-SVM, which are obtained from element wise product and averaging of base kernels and classifying with SVM. MKL baseline for kernel product gives the highest score of 85.5%. Moreover, it is very simple and fast in comparison to other MKL methods in figure 2(a). Our proposed scheme based on AKM gives 86.7%, which is better than all MKL and Boosting based methods.

we investigate the effect of adding random feature channels on different fusion schemes. In addition to 7 informative kernels of Oxford Flower 17 dataset we have generated 20 RBF kernels from 20 set of random vectors. We started with all informative kernels, i.e., number of noisy kernels is 0, then we added different number of noisy kernel. The mean accuracy of different state-of-the-art methods under noisy channels is presented in figure 2(b). MKL baseline drops

down significantly with the number of noisy kernels while two-stage AKM is robust to noisy feature channels and perform significantly better than MKL or boosting based approaches.

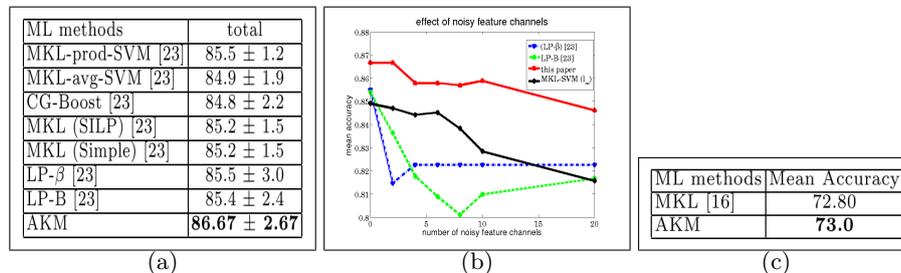


Fig. 2. Empirical feature spaces for Multiple kernels: (a) Mean accuracy on Oxford Flower17 and comparison with different machine learning methods; (b) Oxford Flower17. Mean Accuracy of different fusion methods under noisy feature channels; (c) Mean accuracy on Oxford Flower 102 dataset.

Oxford Flower 102 [16] is an extended multi-class dataset containing 102 flower categories. It consists of 8189 images with 40 to 250 images in each class. The dataset is split into training (10 images per class), validation (10 images per class) and test (with a minimum of 20 images per class) using predefined splits. For the experiments we have used 4 χ^2 distance matrices provided online. The details of the features used to compute these distance matrices can be found in [16]. The experimental setup is the same as for Oxford Flower 17. AKM is performing comparable to MKL as shown in figure 2(c).

5 Conclusions

In this paper we have presented a novel two stage multiple kernel learning approach for augmented kernel matrix. The proposed method addresses the complexity problems of AKM and makes it robust to redundant and noisy kernels. We propose automatic grouping of kernels based on kernel alignment by agglomerative clustering of kernels. Learning representative kernels for each group results in a small set of highly informative kernels. Learning a combination within a group discards or downweights redundant and noisy kernels thus results in an optimal kernel from a set of informative base kernels. The complexity is further reduced by exploiting the property of independence of empirical feature spaces in the AKM scheme. It allows to use only the most influential examples from each representative kernel to construct the AKM matrix. We perform experiments on challenging object recognition datasets and the results validate our technique. The proposed approach makes it possible to use the AKM method for 20 kernels with several thousands of training examples. A performance increase is observed compared to MKL based on a linear combination of all base kernels. This observation is significant as it suggests that the information in the kernels can be exploited more effectively and the classification rate increases without using additional features.

References

1. B. Scholkopf and A. Smola, *Learning with Kernels*. MIT Press, 2002.
2. B. Scholkopf, A. Smola, and K. Muller, "Kernel Principal Component Analysis," in *International Conference on Artificial Neural Networks*, 1997.
3. J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
4. V. Vapnik, *The Nature of Statistical Learning Theory*. Springer Verlag, 2000.
5. K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615–1630, 2005.
6. K. van de Sande, T. Gevers, and C. Snoek, "Evaluation of color descriptors for object and scene recognition," in *Computer Vision and Pattern Recognition*, 2008.
7. G. Lanckriet, N. Cristianini, P. Bartlett, L. Ghaoui, and M. Jordan, "Learning the Kernel Matrix with Semidefinite Programming," *The Journal of Machine Learning Research*, vol. 5, pp. 27–72, 2004.
8. F. Bach, G. Lanckriet, and M. Jordan, "Multiple Kernel Learning, Conic Duality, and the SMO Algorithm," in *International Conference on Machine Learning*, 2004.
9. S. Sonnenburg, G. Ratsch, C. Schafer, and B. Scholkopf, "Large Scale Multiple Kernel Learning," *The Journal of Machine Learning Research*, vol. 7, pp. 1531–1565, 2006.
10. A. Zien and C. Ong, "Multiclass Multiple Kernel Learning," in *International Conference on Machine Learning*, 2007.
11. F. Yan, K. Mikolajczyk, J. Kittler, and A. Tahir, "Combining Multiple Kernels by Augmenting the Kernel Matrix," in *International Workshop on Multiple Classifier Systems*, 2010.
12. M. Szafranski, Y. Grandvalet, and A. Rakotomamonjy, "Composite Kernel Learning," *Machine Learning*, vol. 79, no. 1, pp. 73–103, 2010.
13. J. Nath, G. Dinesh, S. Raman, C. Bhattacharyya, A. Ben-Tal, and K. Ramakrishnan, "On the Algorithmics and Applications of a Mixed-norm Based Kernel Learning Formulation," in *Neural Information Processing Systems*, 2009.
14. M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
15. M. Nilsback and A. Zisserman, "A visual Vocabulary for Flower Classification," in *Computer Vision and Pattern Recognition*, 2006.
16. M.-E. Nilsback and A. Zisserman, "Automated Flower Classification over a Large Number of Classes," in *ICCVGIP*, 2008.
17. M. Kloft, U. Brefeld, P. Laskov, and S. Sonnenburg, "Nonsparse Multiple Kernel Learning," in *Neural Information Processing Systems Workshop on Kernel Learning: Automatic Selection of Optimal Kernels*, 2008.
18. B. Scholkopf, S. Mika, C. Burges, P. Knirsch, K. Muller, G. Ratsch, and A. Smola, "Input Space Versus Feature Space in Kernel-Based Methods," *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 1000–1017, 1999.
19. N. Cristianini, J. Shawe-Taylor, A. Elisseeff, and J. Kandola, "On Kernel-Target Alignment," in *Neural Information Processing Systems*, 2001.
20. J. Davis, B. Kulis, P. Jain, S. Sra, and I. Dhillon, "Information-Theoretic Metric Learning," in *International Conference on Machine Learning*, 2007.
21. N. Lawrence and G. Sanguinetti, "Matching Kernel through Kullback-Leibler Divergence Minimisation," in *Technical Report CS-04-12, Department of Computer Science, University of Sheffield*, 2005.
22. S. Lazebnik, C. Schmid, and J. Ponce, "Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories," in *Computer Vision and Pattern Recognition*, 2006.
23. P. Gehler and S. Nowozin, "On Feature Combination for Multiclass Object Classification," in *International Conference on Computer Vision*, 2009.