

# Feature Combination beyond Basic Arithmetics

Hao Fu<sup>1</sup>

hxf@cs.nott.ac.uk

Guoping Qiu<sup>1</sup>

qiu@cs.nott.ac.uk

Hangen He<sup>2</sup>

hehangen@yahoo.com

<sup>1</sup> School of Computer Science

University of Nottingham

Nottingham, UK

<sup>2</sup> College of Mechatronic Engineering and

Automation, National University of De-

fense Technology, Changsha, Hunan,

P.R.China, 410073

---

## Abstract

Kernel-based feature combination techniques such as Multiple Kernel Learning use arithmetical operations to linearly combine different kernels. We have observed that the kernel distributions of different features are usually very different. We argue that the similarity distributions amongst the data points for a given dataset should not change with their representation features and propose the concept of relative kernel distribution invariance (RKDI). We have developed a very simple histogram matching based technique to achieve RKDI by transforming the kernels to a canonical distribution. We have performed extensive experiments on various computer vision and machine learning datasets and show that calibrating the kernels to an empirically chosen canonical space before they are combined can always achieve a performance gain over state-of-art methods. As histogram matching is a remarkably simple and robust technique, the new method is universally applicable to kernel-based feature combination.

## 1 Introduction

The importance of feature combination has long been recognized by the computer vision community. Different features, such as local, global, color, texture, etc, capture different characteristics of an image. It is often helpful and sometimes necessary to combine various features together in order to gain a comprehensive understanding of an image.

Kernel based feature combination is an effective method [1, 2, 3], where different types of kernels for the same type of feature or the same type of kernel for different types of features are combined together. In some methods such as the prominent Multiple Kernel Learning (MKL) technique [4], the weights of different kernels are learned adaptively together with the parameters of the final classifier, and these methods can be referred to as adaptive kernel combination (AKC); whilst in other methods, the weights of different kernels are predefined [5] and these methods can be referred to as nonadaptive kernel combination (NAKC).

In this paper, we make an important observation of the distribution of different kernels that are routinely used in the literature. We discovered that the histograms of the kernel values of different features are usually quite different from each other (see Fig.2 for example).

Some histograms may be narrow and occupy only a short range, while others may span a wide range; some histograms may look like a gaussian distribution, while others may look like an exponential distribution. As these histograms differ so much, it means that their units of measure are not the same. In other words, for the same similarity/difference value, it may represent a ‘huge’ difference in one feature channel, but only a ‘tiny’ difference in the other channel. Therefore, it is necessary to standardize each feature channel before they are combined together.

We argue that there may exist some invariant properties that are intrinsic to the data itself and will not change with different feature representations and the use of different forms of kernels. We propose an intuition that the similarity distributions amongst the data points for a given dataset should not change with their representation features. As kernels measure the similarities between samples, we call this intuition the relative kernel distribution invariance (RKDI) property. To achieve RKDI, we propose a simple but very effective method to standardize different kernels through histogram matching and show that this surprisingly simple operation can reliably boost the performances of AKC and NAKC methods for a variety of applications. The new method is very simple, easy to implement and robust; it can therefore be considered as a new baseline for feature combination in addition to simple arithmetics such as average and product [9].

This paper is structured as follows: in section 2, we review related works in feature combination. We present our histogram matching based kernel combination technique in section 3. In section 4, we present experiments on various datasets to show the effectiveness of our algorithm. Discussions and conclusions are given in section 5.

## 2 Related work

The idea of combination appears in every aspect of computer vision. If we consider a classifier as a one-input-one-output black box, combination can happen both in the input level and the output level. In the input level, we can simply concatenate different kinds of features, while on the output level, we can fuse the outputs of different classifiers [6]. These classifiers can be based on different features [6] or even the same feature [10]. Previous works have always shown a performance gain when combination is used. Besides combination in the input or output level, we can consider using the kernel as a middle level fusion stage. The reason why kernel methods is more suitable for combination lies in two aspects: firstly, kernels can directly model the similarities of samples in different feature channels [10]; secondly, in the kernel space, a linear classifier can have sufficient capability of classifying samples.

After representing different kinds of features using different kernels, we can design algorithms to fuse those kernels. The most prominent methods for combining different kernels should be Multiple Kernel Learning (MKL), in which the algorithm tries to learn an optimal weight for each different kernel. These weights and the parameters of the final classifier are learned jointly in a principled framework. The seminal work of Multiple Kernel Learning dates back to [11], where the authors proposed an efficient algorithm to solve this optimization problem. After MKL was proposed, many variants of it have been proposed [12, 13, 14], and have been quickly adopted to deal with various computer vision problems [15, 16].

Despite its huge success, the formulation of MKL is still being questioned by researchers. In essence, MKL is simply a linear combination of different kernels. It implies that the contribution of each kernel is fixed for all the training samples [17]. This seems to be an unnecessary too strong constraint. In [18], the authors propose to learn augmented coeffi-

cients for each sample in each feature channel. They achieve this by augmenting the kernel matrixes. However, as the augmented kernel they used is still a block diagonal matrix, the coefficients they learned are equivalent to learning different kernels separately and adding an appropriate bias term for all the kernel classifiers. The authors of [4] introduced a non-stationary approach and allowed the relative weights of the kernels to be varied with the input samples. In [5], instead of learning different kernels simultaneously, the author takes a boosting-like two stage strategy. At the first stage, a classifier is learned for each kernel separately, then these learned classifiers are treated as weak learners and assembled together at the second stage. Although this method shows good performance in [5], it shows limitations elsewhere [18].

### 3 Standardizing kernel values through piecewise linear histogram matching

Let  $(x_i, y_i), i = 1, 2, \dots, N$  be  $N$  instances consisting of images  $x_i \in X$  and class labels  $y_i \in \{1, 2, \dots, C\}$ ;  $f_m \in R^{d_m}, m = 1, 2, \dots, F$ , represent a given set of features, where  $d_m$  denotes the dimensionality of the  $m$ -th feature. Feature combination is to use all these  $F$  features together to learn a classifier to classify  $X$  into  $Y$ . Kernel methods make use of kernel functions to define a measure of similarity between pairs of instances. Let  $k$  be a kernel function, the similarity between two images based on their  $m$ -th feature,  $f_m$ , is defined as:

$$k_m(x, x') = k(f_m(x), f_m(x')) \quad (1)$$

Kernel based feature combination is about combining different  $k_m$  into a single kernel  $k^*$  and can be done with various arithmetical operations [4] including baseline average (2) and MKL (3).

The baseline average kernel:

$$k^*(x, x') = \frac{1}{F} \sum_{m=1}^F k_m(x, x') \quad (2)$$

In the case of MKL, the combined kernel  $k^*$  is a linear combination of different kernels weighted by a set of adaptive parameters  $\{\beta_m\}$  to be learned by the MKL algorithms.

$$k^*(x, x') = \sum_{m=1}^F \beta_m k_m(x, x') \quad (3)$$

An inspection of the distributions of  $k_m(x, x')$  for different features (see Fig.2) shows that they are very different for different features. Linear combination of the kernels as (2) and (3) can be seen as combining ‘things’ measured with different units directly without converting them to the same standard. We argue that before they are linearly combined, the kernel values should be calibrated to a canonical feature space (CFS). Although the exact form of the CFS is unknown, it is reasonable to assume that in this CFS, there are some invariant properties that are intrinsic to the data itself and will not change with different feature representations or the use of different forms of kernels.

Intuitively, the similarity distributions amongst the data points for a given dataset should not change with their representation features. As kernels measure the similarities between

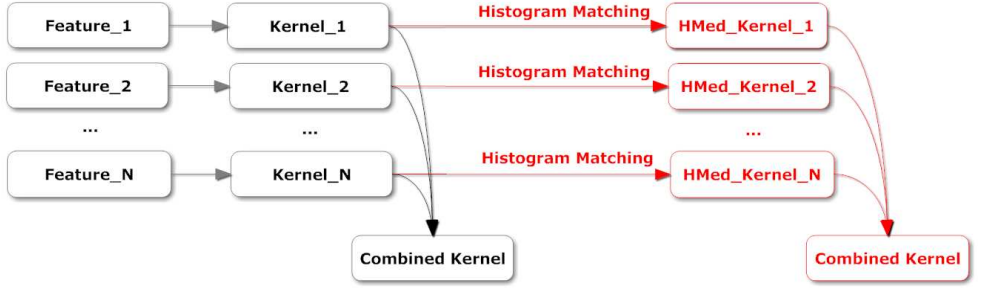


Figure 1: Typical feature combination methods always represent features into their kernel forms. These kernels are then combined. Traditionally, the kernels are combined directly through one of the methods in (2) or (3). In this paper, we proposed to add a histogram matching module before these kernels are combined by one of the methods in (2) or (3)

samples, we call this intuition the relative kernel distribution invariance (RKDI) property. In the following, we will try to make this intuition concrete.

Defining the inverse cumulative density function (ICDF) of kernel  $m$  as:

$$ICDF_m(u) = \inf_{v \in \mathfrak{R}} \left( \int_{-\infty}^v p_m(k_m(x, x') = w) dw \geq u \right) \quad (4)$$

where  $p_m(k_m(x, x'))$  is the probability density function of the  $m$ -th feature channel, then RKDI can be defined as:

$$\int_{-\infty}^{ICDF_m(u)} p_m(k_m(x, x') = w) dw = \int_{-\infty}^{ICDF_{CFS}(u)} p_{CFS}(k_{CFS}(x, x') = w) dw \quad \forall m, u \quad (5)$$

where  $p_{CFS}(k_{CFS}(x, x'))$  and  $ICDF_{CFS}(u)$  represent the probability density function and the inverse cumulative density function of the canonical feature space respectively.

Clearly, (5) states that the percentiles of the relative similarities of the given data should be the same in any feature space and should be calibrated to the canonical feature space. Although there is no formal proof known to us at this stage, we believe it is a reasonable assumption and will show experimentally that maintaining such invariance can help improve performances. In the absence of a known CFS, we use cross-validation to select one of the kernels as the CFS and calibrate all other kernels to this empirical CFS.

The problem of (5) is the well-known histogram matching problem and our new feature combination framework is illustrated in Fig.1. Let  $HM(k_m(x, x'))$  represent the Histogram Matching operator that perform canonical histogram matching on the  $m$ -th kernel, then AKC (MKL) and NAKC (average) are represented as follows.

The new NAKC (average)  $k^*$  kernel is formed as:

$$k^*(x, x') = \frac{1}{F} \sum_{m=1}^F HM(k_m(x, x')) \quad (6)$$

In the case of MKL, the combined kernel  $k^*$  is formed as:

$$k^*(x, x') = \sum_{m=1}^F \beta_m HM(k_m(x, x')) \quad (7)$$

Our histogram matching algorithm is summarized in **Algorithm 1**. It differs from typical histogram matching methods<sup>1</sup> in that the elements in the kernel matrixes are continuous instead of discrete values. Therefore, we need to quantize the kernel values into discrete bins. To reduce the quantization error and maintain the original order, the values are piecewise linear interpolated for each bin. Note that in all our experiments, we use 1500 bins.

---

**Algorithm 1** Piecewise Linear Histogram Matching
 

---

```

Input: template, orig_kernel, num_of_bins;
Output: HMed_kernel (Histogram Matched kernel)
Normalize template and orig_kernel to (0,1);
sorted_template = sort( template );
for i=1 to num_of_bins do
  cut_point_index = size( find( template < i/num_of_bins ) );
  cut_point_value = sorted_template[cut_point_index];
end for
for i in orig_kernel do
  lower_bound = max( orig_kernel[i] > cut_point_value(:) );
  upper_bound = min( orig_kernel[i] < cut_point_value(:) );
  HMed_kernel[i] = Linear_interpolate( lower_bound,orig_kernel[i],upper_bound );
end for
Normalize HMed_kernel back to the original range
  
```

---

An important question in this method is finding the canonical feature space which is likely to be dataset dependant. In all our experiments, we use cross validation to choose the canonical feature histogram. Note here we should ensure the histogram matched kernels be positive definite. Although we cannot theoretically proved that, we found the histogram matched kernels always satisfy this condition in our experiments if we choose one of the feature histograms as the canonical histogram.

## 4 Experimental Results

### 4.1 Corel5K dataset

In [8], the authors studied the problem of image annotation. They showed that by simply adding the distances of different features, they can achieve superior performance on the corel5K benchmark image annotation dataset. They used features representing color and texture, and the distances of each feature channel are equally weighted. They called their algorithm Joint Equal Contribution (JEC). In [9], the authors proposed to use another 15 kinds of features including global and local features. They reported similar results to JEC. They have also released their features<sup>2</sup> used in the experiments. We did experiments directly based on these features. Different metric measures [9] are adopted to calculate the distances in each feature channel. The histograms corresponding to the distances of those 15 kinds of features are shown in Fig.2. From there we can see an obvious difference between different feature channels.

---

<sup>1</sup>For a brief introduction on histogram matching, please refer to [http://paulbourke.net/texture\\_colour/equalisation/](http://paulbourke.net/texture_colour/equalisation/)

<sup>2</sup><http://lear.inrialpes.fr/people/guillaumin/data.php>

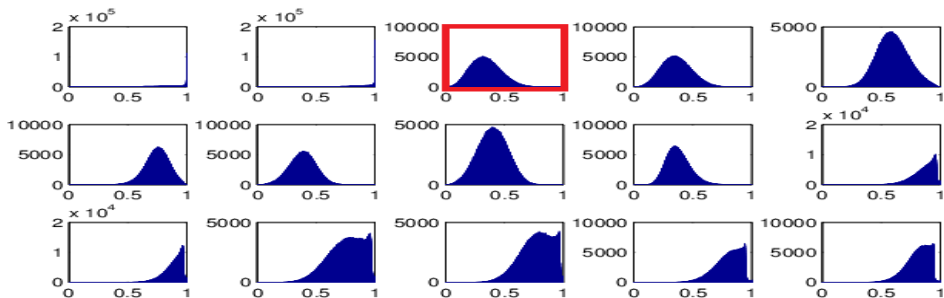


Figure 2: Kernel histograms of the features used in [4]. The histogram in the red box is chosen as the standard histogram, which corresponds to the feature of hue descriptors extracted at Harris-Laplacian interest points

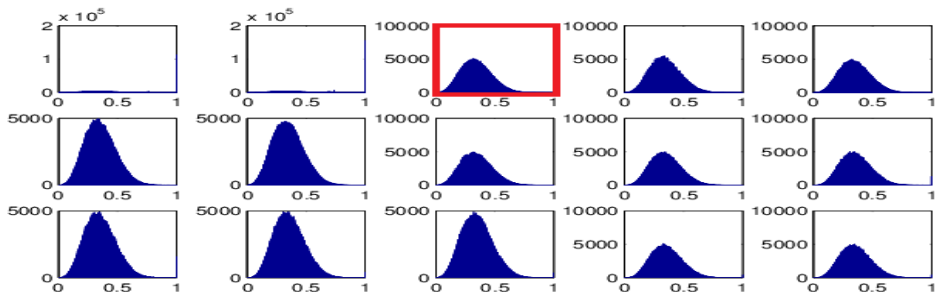


Figure 3: Histograms in Fig.2 after histogram matching.

As there is no theoretical guidance on how to choose a standard histogram, we use cross validation to choose one from these 15 features as the canonical feature. The histograms after histogram matching are shown in Fig.3.

Those distances after histogram matching are added together. Based on this added distance, the  $K$  nearest neighbors for each test sample are retrieved from the training set. The tags of each test sample are solely determined by these  $K$  nearest neighbors. In predicting the tags from these  $K$  neighbors, we also adopted the label transfer strategy used in [8]. Precision and recall are used to evaluate the performance and the results are shown in Table 1. From there we can see a performance boost by introducing the histogram matching module. It is important to note that the purpose here is not to compete with the state of the art image tagging performances but rather to demonstrate that by calibrating the kernels using simple histogram matching before combining them can improve performances.

## 4.2 Caltech101 in 39 kernels [3]

In [3], the authors thoroughly studied the problem of feature combination. One of their important findings is that simple average kernel may even outperform sophisticated MKL

Table 1: Image Annotation Performances on Corel5K Dataset. HM is short for **H**istogram **M**atching. Rate+ is the number of tags whose recall is above zero.

models	Prec	Recall	Rate+
HPM [20]	0.25	0.28	136/260
JEC [8]	0.27	0.32	139/260
JEC-15 [9]	0.28	0.33	140/260
JEC-15 + HM	0.30	0.36	150/260

algorithms. They have also released their code and the gram matrixes<sup>3</sup> used in their experiments. The best result they got was based on a combination of 39 kernels. These different kernels are mainly based on 5 different kinds of features: LBP, PHOG, SIFT, Region covariance and Gabor filter banks. Those features are assembled in different layouts, resulting in a total of 39 kernels. In their work, they have already compared their results with typical MKL algorithms, including SILP [12] and SimpleMKL [11]. In some cases, simple average kernel may outperform these complicated MKL methods.

We did experiments directly on these publicly available gram matrixes. Experimental results are shown in Fig.4. Again, we can see a performance boost by introducing the histogram matching module before combining the kernels.

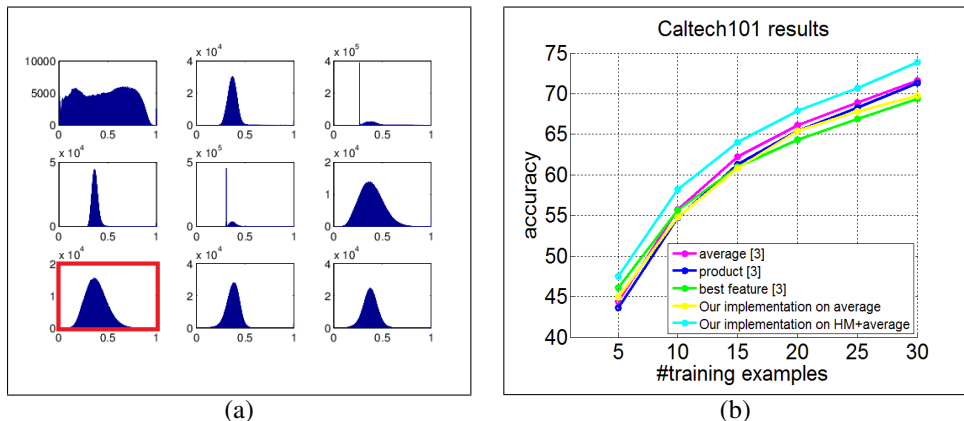


Figure 4: (a) Some representative kernel histograms among the 39 kernels used in [8], the one in the red box is chosen as the standard histogram; (b) The classification results on Caltech101. From the figure, we can see the average of histogram matched kernels can always perform better than averaging original kernels. Note that the author of [8] report results on five random splits of the dataset. However, they have only released their gram matrixes of one split. We did experiments only on this split. This results in the slight difference between our implementation on average and the average accuracy reported in [8].

<sup>3</sup><http://people.ee.ethz.ch/~pgehler/projects/iccv09/caltech/>

### 4.3 Oxford flowers dataset

The Oxford flowers dataset [9] contains 17 different kinds of flowers. Each class contains 80 samples, 40 for training, 20 for validation, and the rest 20 for testing. The authors of [9] have also made the distance matrixes they used publicly available<sup>4</sup>. Following [10], these distance matrixes are transformed to kernels using  $k = \exp(-\gamma^{-1} \cdot d)$ , where  $\gamma$  is the mean of the distance matrix, and  $d$  is the distance between samples. The kernel histogram of these 7 features are shown in Fig.5.

Firstly, we use (2) and (6) to combine the kernels. A standard SVM solver<sup>5</sup> is adopted as the classifier. The results are shown in Table 2. As expected, HM+average (6) performs better than average (2). Then we use OBSCURE [10], a state-of-art MKL method to learn the optimal weights of different kernels. We choose OBSCURE as the MKL algorithm mainly because of its efficiency. The results are also shown in Table 2. We also compare our results with some other recently proposed MKL algorithms.

From the table, we can see that OBSCURE shows a similar performance with other MKL algorithms. HM+OBSCURE performs better than all other MKL algorithms. Notice that between the algorithm of OBSCURE and HM+OBSCURE, they use exactly the same feature and the same MKL solver, the only difference lies in whether they use the Histogram matching module to calibrate the kernels. Thus the performance gain should be purely the contribution of our histogram matching module.

Table 2: Experimental results on Oxfordflower

methods	[10]	LP- $\beta$ [3]	average	HM+average	OBSCURE	HM+OBSCURE
accuracy	86.7 $\pm$ 1.2	85.5 $\pm$ 3.0	84.1 $\pm$ 1.0	85.3 $\pm$ 1.4	85.5 $\pm$ 1.5	<b>87.3<math>\pm</math>0.7</b>

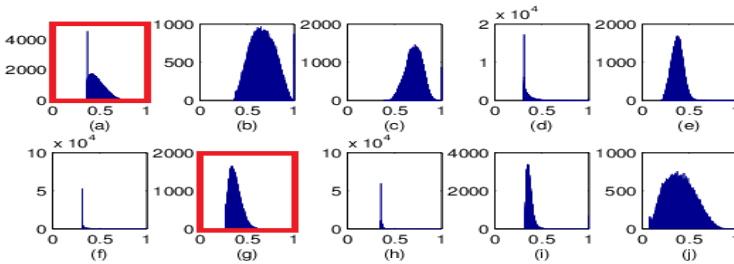


Figure 5: Kernel histograms of three features used in MSRC21 ((a) to (c)) and seven features ((d) to (j)) used in Oxfordflowers. The histogram in the red box is chosen as the standard histogram.

### 4.4 MSRC21 dataset

Next, we consider another example in semantic segmentation area. MSRC21 is a well-known dataset which contains 591 images. Each image has pixel level ground truth labels from 21 semantic classes. Following [13], these 591 images are split into 276 for training, 59 for validation, and the remaining 256 for testing.

<sup>4</sup><http://www.robots.ox.ac.uk/vgg/data/flowers/17/index.html>

<sup>5</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvm/>



Here, our objective is to evaluate feature combination rather than aiming to achieve state of the art semantic labeling performances. Therefore, we simplify the original semantic segmentation problem into a region labeling problem, i.e. we assume the image has already been segmented into regions, and our task is to assign each region a semantic label. To avoid bias, we directly use the ground truth segmentation. We adapted the code from [5], where they used three features to represent each region. These three features are texon histograms, color histograms and pyramid of HOG. Their respective kernel histograms are shown in Fig.5. A typical MKL solver<sup>6</sup> is adopted to learn the weights of these three kernels. Results are shown in Table 3, it is again seen that our histogram matching strategy can help boost performances.

Table 3: Results on MSRC21. Note that the results of [13] is achieved fully automatically, whilst our results are based on the ground truth segmentations, therefore they are not directly comparable. Our objective is to compare feature combination techniques instead of automatic semantic labeling performances.

	[13]	average	HM+average	MKL	HM+MKL
per-class accuracy	0.58	0.75	0.76	0.81	<b>0.84</b>
global accuracy	0.72	0.86	0.88	0.90	<b>0.93</b>

## 4.5 UCI machine learning repository

To further verify the robustness of our algorithm, we also did experiments on UCI repository datasets. On Sonar and Heart dataset, we did experiments exactly following [2]. Three kernels are used: a quadratic kernel, an RBF kernel and a linear kernel. We report mean test accuracy across ten random replications of three-fold cross validation. The results are shown in Table 4. As those kernels are based on the same features, we did not see big improvements by introducing the histogram matching module. But still, we can get slightly better results.

Next we did experiments on another UCI dataset which contains multiple features for each sample. Mfeat contains 2000 handwritten numerals from 0 to 9. Six kinds of features are provided for each sample. We use RBF kernel for all these features. Following [2], we report results over 50 repeated trials in which we randomly select 20 training and 20 testing for each class. The results are also shown in Table 4. Again, we can see a performance gain by introducing our histogram matching module, no matter directly adding the kernels or using OBSCURE to learn different weights for different kernels.

Table 4: Results on UCI datasets. For sonar, the histogram of RBF kernel is chosen as the standard histogram, while on Heart, quadratic kernel is chosen. On mfeat, the histogram on Zernike moments feature is chosen as the standard histogram.

Dataset	[2]	NSKC[2]	average	OBSCURE	HM+average	HM+OBSCURE
Sonar	-	86.3±0.8	86.9±4.5	89.0±3.7	87.0±3.1	<b>89.2±3.1</b>
Heart	-	60.5±1.6	80.9±2.8	80.1±3.3	<b>82.7±2.7</b>	82.6±2.8
mfeat	94.9±1.7	-	97.2±0.6	97.0±0.4	97.5±0.1	<b>98.3±0.9</b>

<sup>6</sup>Downloadable from <http://www.di.ens.fr/%7Eeobozinski/SKMsmo.tar>

## 5 Discussion and concluding remarks

In this paper, we have proposed a new feature combination method which calibrates the kernels to a canonical feature space before linearly combining them. Experiments on various datasets have shown the effectiveness of this simple strategy. This method can be used in the unsupervised scenario where it consistently performs better than average baseline. In supervised case, it can be seamlessly combined with various kinds of multiple kernel learning algorithms and we have shown it again can consistently boost performances.

Future work will try to understand the histogram matching mechanism thoroughly. It will be very interesting if we can find the shape of the canonical histogram automatically instead of cross-validation used in this paper. A promising way of achieving this would be to treat the canonical histogram as parameters, and directly integrate them into the MKL optimization function.

## References

- [1] Francis R. Bach, Gert R. G. Lanckriet, and Michael I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. *Twenty-first international conference on Machine learning - ICML '04*, page 6, 2004. doi: 10.1145/1015330.1015424.
- [2] Theodoros Damoulas and Mark a. Girolami. Combining feature spaces for classification. *Pattern Recognition*, 42(11):2671–2683, November 2009. ISSN 00313203. doi: 10.1016/j.patcog.2009.04.002.
- [3] Peter Gehler and Sebastian Nowozin. On Feature Combination for Multiclass Object Classification. In *ICCV*, 2009.
- [4] Matthieu Guillaumin, Thomas Mensink, Jakob Verbeek, and Cordelia Schmid. Tag-Prop: Discriminative metric learning in nearest neighbor models for image auto-annotation. *2009 IEEE 12th International Conference on Computer Vision*, pages 309–316, September 2009. doi: 10.1109/ICCV.2009.5459266.
- [5] Yong Jae Lee and Kristen Grauman. Object-Graphs for Context-Aware Category Discovery. *CVPR*, 2010.
- [6] Josef Kirtler, Mohamad Hatefm, Robert P.W.Duin, and Jiri Matas. On combining classifiers. *IEEE Trans on Pattern Analysis and Machine Intelligent*, 20(3):226–239, 1998.
- [7] Darrin P. Lewis, Tony Jebara, and William Stafford Noble. Nonstationary kernel combination. In *Proceedings of the 23rd international conference on Machine learning - ICML '06*, pages 553–560, New York, New York, USA, 2006. ACM Press. ISBN 1595933832. doi: 10.1145/1143844.1143914.
- [8] Ameesh Makadia, Vladimir Pavlovic, and Sanjiv Kumar. A New Baseline for Image Annotation. In *ECCV*, volume 90, May 2008. doi: 10.1007/s11263-010-0338-6.
- [9] Maria-Elena Nilsback and Andrew Zisserman. Automated Flower Classification over a Large Number of Classes. *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729, December 2008. doi: 10.1109/ICVGIP.2008.47.

- [10] Francesco Orabona, Luo Jie, and Barbara Caputo. Online-Batch Strongly Convex Multi Kernel Learning. In *CVPR*, 2010.
- [11] Alain Rakotomamonjy, Francis Bach, Stéphane Canu, and Yves Grandvalet. More efficiency in multiple kernel learning. *Proceedings of the 24th international conference on Machine learning - ICML '07*, pages 775–782, 2007. doi: 10.1145/1273496.1273594.
- [12] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. April 2006.
- [13] Jamie Shotton, John Winn, Carsten Rother, and Antonio Criminisi. TextonBoost: Joint Appearance, Shape and Context Modeling for Multi-Class Object Recognition and Segmentation. *ECCV*, 2006.
- [14] Soren Sonnenburg, Gunnar Ratsch, Christin Schafer, and Bernhard Scholkopf. Large Scale Multiple Kernel Learning. *Journal of Machine Learning Research*, 7:1531–1565, 2006.
- [15] Manik Varma and Debajyoti Ray. Learning The Discriminative Power-Invariance Trade-Off. *2007 IEEE 11th International Conference on Computer Vision*, October 2007. ISSN 1550-5499. doi: 10.1109/ICCV.2007.4408875.
- [16] Andrea Vedaldi, Varun Gulshan, Manik Varma, and Andrew Zisserman. Multiple Kernels for Object Detection. *ICCV*, 2009.
- [17] Stefan Walk, Konrad Schindler, and Bernt Schiele. Disparity Statistics for Pedestrian Detection: Combining Appearance, Motion and Stereo. In *ECCV*, 2010.
- [18] Fei Yan, Krystian Mikolajczyk, Mark Barnard, Hongping Cai, and Josef Kittler. lp Norm Multiple Kernel Fisher Discriminant Analysis for Object and Image Categorisation. In *CVPR*, 2010.
- [19] Fei Yan, Krystian Mikolajczyk, Josef Kittler, and Atif Tahir. Combining Multiple Kernels by Augmenting the Kernel Matrix. In *In Proceedings of the 9th International Workshop on Multiple Classifier Systems*, 2010.
- [20] Ning Zhou, William K. Cheung, Guoping Qiu, and Xiangyang Xue. A Hybrid Probabilistic Model for Unified Collaborative and Content-Based Image Tagging. *PAMI, IEEE Transactions on*, X(99), November 2010. ISSN 0162-8828. doi: 10.1109/TPAMI.2010.204.
- [21] Alexander Zien and Cheng Soon Ong. Multiclass multiple kernel learning. *Proceedings of the 24th international conference on Machine learning - ICML '07*, 1(2):1191–1198, 2007. doi: 10.1145/1273496.1273646.