

# Several Tunable GMM Kernels

Ping Li  
Baidu Research USA  
pingli98@gmail.com

## ABSTRACT

While tree methods have been popular in practice, researchers and practitioners are also looking for simple algorithms which can reach similar accuracy of trees. In 2010, [20] developed the method of “abc-robust-logitboost” and compared it with other supervised learning methods on datasets used by the deep learning literature. In this study, we propose a series of “tunable GMM kernels” which are simple and perform largely comparably to tree methods on the same datasets. Note that “abc-robust-logitboost” [20] substantially improved the original “GDBT” in that (a) it developed a tree-split formula based on second-order information of the derivatives of the loss function; (b) it developed a new set of derivatives for multi-class classification formulation.

In the prior study [23], the “generalized min-max” (GMM) kernel was shown to have good performance compared to the “radial-basis function” (RBF) kernel. However, as demonstrated in this paper, the original GMM kernel is often not as competitive as tree methods on the datasets used in the deep learning literature. Since the original GMM kernel has no parameters, we propose tunable GMM kernels by adding tuning parameters in various ways. Three basic (i.e., with only one parameter) GMM kernels are the “ $e$ GMM kernel”, “ $p$ GMM kernel”, and “ $\gamma$ GMM kernel”, respectively. Extensive experiments show that they are able to produce good results for a large number of classification tasks. Furthermore, the basic kernels can be combined to boost the performance.

For large-scale machine learning tasks, it is crucial that learning methods should be able to scale up with the size of the training samples. It has been known that the original GMM kernel can be efficiently linearized (i.e., achieving the result of a nonlinear kernel at the cost of a linear kernel). As demonstrated in this paper, several tunable GMM kernels also inherit this nice property in that they can also be efficiently linearized.

## 1 INTRODUCTION

Kernel methods [30] are an important part of machine learning. Among many types of kernels, the linear kernel and the “radial basis function” (RBF) kernel are probably the most well-known. Recently, the “generalized min-max” (GMM) kernel [22] was introduced for large-scale search and machine learning, owing to its efficient linearization via either hashing or the Nystrom method [28].

For defining the GMM kernel, the first step is a simple transformation on the original data. Consider, for example, the original data vector  $u_i$ ,  $i = 1$  to  $D$ . We define the following transformation, depending on whether an entry  $u_i$  is positive or negative:

$$\begin{cases} \tilde{u}_{2i-1} = u_i, & \tilde{u}_{2i} = 0 & \text{if } u_i > 0 \\ \tilde{u}_{2i-1} = 0, & \tilde{u}_{2i} = -u_i & \text{if } u_i \leq 0 \end{cases} \quad (1)$$

For example, when  $D = 2$  and  $u = [-4 \ 6]$ , the transformed data vector becomes  $\tilde{u} = [0 \ 4 \ 6 \ 0]$ . The GMM kernel is defined [22] as follows:

$$GMM(u, v) = \frac{\sum_{i=1}^{2D} \min\{\tilde{u}_i, \tilde{v}_i\}}{\sum_{i=1}^{2D} \max\{\tilde{u}_i, \tilde{v}_i\}} \quad (2)$$

Even though the GMM kernel has no tuning parameter, it performs surprisingly well for classification tasks as empirically demonstrated in [22] (also see Table 1 and Table 2), when compared to the linear kernel and best-tuned radial basis function (RBF) kernel:

$$RBF(u, v; e) = e^{-\lambda_e \left( 1 - \frac{\sum_{i=1}^D u_i v_i}{\sqrt{\sum_{i=1}^D u_i^2 \sum_{i=1}^D v_i^2}} \right)} \quad (3)$$

where  $\lambda_e > 0$  is the tuning parameter.

Furthermore, the (nonlinear) GMM kernel can be efficiently linearized via hashing [15, 21, 26] (or the Nystrom method [28]). This means we can use the linearized GMM kernel for large-scale machine learning tasks essentially at the cost of linear learning.

Given the deceiving simplicity of the GMM kernel and its surprising performance compared to the RBF kernel, researchers and practitioners might be seriously interested in asking two questions:

- (1) Does the GMM kernel perform comparably to more sophisticated learning methods such as trees, at least in the context of supervised learning when features are already available?
- (2) Can one improve the accuracy of the original (tuning-free) GMM kernel, for example, by adding tuning parameters?

This paper aims at addressing these two questions. It turns out that, in many datasets, the original (tuning-free) GMM kernel can be substantially improved, by adding tuning parameters. Furthermore, we report a comparison study using a set of public datasets in the deep learning literature [17]. This set of datasets were used by an empirical study [20] to compare several boosting & tree methods with deep nets. This paper will show that tunable GMM kernels can achieve comparably accuracy as trees.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

2018, ,

© 2018 Copyright held by the owner/author(s).

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

[https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

### 1.1 Tunable GMM Kernels

In order to improve the performance of the original (tuning-free) GMM kernels, we propose three basic tunable GMM kernels:

$$e\text{GMM}(u, v; \lambda_e) = e^{-\lambda_e \left( 1 - \frac{\sum_{i=1}^{2D} \min(\tilde{u}_i, \tilde{v}_i)}{\sum_{i=1}^{2D} \max(\tilde{u}_i, \tilde{v}_i)} \right)} \quad (4)$$

$$p\text{GMM}(u, v; p) = \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})^p}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})^p} \quad (5)$$

$$\gamma\text{GMM}(u, v; \gamma) = \left( \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})} \right)^\gamma \quad (6)$$

and the combinations of the basic tunable GMM kernels:

$$p\gamma\text{GMM}(u, v; p, \gamma) = \left( \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})^p}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})^p} \right)^\gamma \quad (7)$$

$$ep\text{GMM}(u, v; \lambda_e, p) = e^{-\lambda_e \left( 1 - \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})^p}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})^p} \right)} \quad (8)$$

$$e\gamma\text{GMM}(u, v; \lambda_e, \gamma) = e^{-\lambda_e \left( 1 - \left( \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})} \right)^\gamma \right)} \quad (9)$$

$$ep\gamma\text{GMM}(u, v; \lambda_e, p, \gamma) = e^{-\lambda_e \left( 1 - \left( \frac{\sum_{i=1}^{2D} (\min\{\tilde{u}_i, \tilde{v}_i\})^p}{\sum_{i=1}^{2D} (\max\{\tilde{u}_i, \tilde{v}_i\})^p} \right)^\gamma \right)} \quad (10)$$

In this study, we will provide an empirical study on kernel SVMs based on the tunable GMM kernels. Perhaps not surprisingly, the improvements can be substantial on many datasets. In particular, we will also compare them with tree methods on 11 datasets used by the deep learning literature [17] and later by [20].

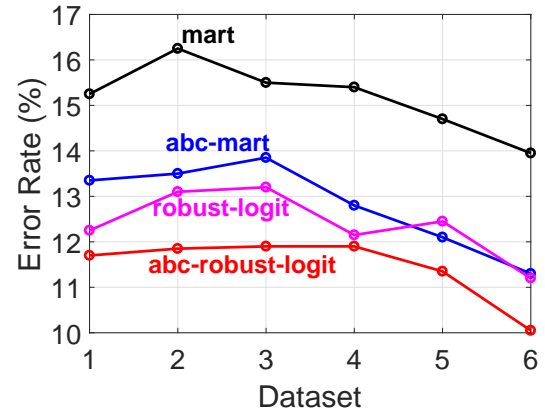
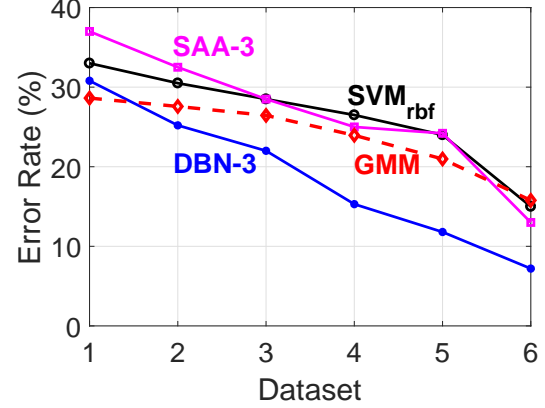
### 1.2 The GMM Kernels versus Tree Methods

[19, 20] developed several boosting & tree methods including “abc-mart”, “robust logitboost”, and “abc-robust-logitboost” and demonstrated their performance on 11 datasets used by the deep learning literature [17]. The good accuracy was achieved by establishing the second-order tree-split formula and new derivatives for multi-class logistic loss. See Table 2 for more information on those datasets.

Figure 1 presents the classification accuracy results on 6 datasets, suggesting that the GMM kernel (upper panel) does not perform as well as tree methods (bottom panel). This observation has motivated us to develop tunable GMM kernels. Later in the paper, Figure 4 will show that the tunable GMM kernels are able to produce roughly comparable results as trees.

### 1.3 Connection to the Resemblance Kernel

The GMM kernel is related to several similarity measures widely used in data mining and web search. When the data are nonnegative, GMM becomes the “min-max” kernel, which has been studied in the literature [5, 15, 16, 21, 26]. When the data are binary (0/1), GMM becomes the well-known “resemblance” similarity. The min-wise hashing algorithm for approximating resemblance has been a highly successful tool in web search for numerous applications [1, 4, 6–8, 10, 11, 14, 24, 25, 27]. Note that for the  $p$ GMM kernel and nonnegative data, when  $p \rightarrow 0$ , it also approaches the resemblance.



**Figure 1: Classification test error rates (lower the better) on 6 datasets (M-Noise1, M-Noise2, ..., M-Noise6) listed in Table 2 which were used by the deep learning literature [17] and later by [20] for comparing tree methods. On these datasets, the results of the GMM kernel (dashed curve in upper panel) are not as accurate as the results produced by the four tree & boosting algorithms (bottom panel).**

### 1.4 Linearization of Nonlinear Kernels

It is common in practice to use linear learning algorithms such as logistic regression or linear SVM. It is also known that one can often improve the performance of linear methods by using nonlinear algorithms such as kernel SVMs, if the computational/storage burden can be resolved. A straightforward implementation of a nonlinear kernel, however, can be difficult for large-scale datasets [3]. For example, for a small dataset with merely 60,000 data points, the  $60,000 \times 60,000$  kernel matrix has  $3.6 \times 10^9$  entries. In practice, being able to linearize nonlinear kernels becomes highly beneficial. Randomization (hashing) is a popular tool for kernel linearization. After data linearization, we can then apply our favorite linear learning packages such as LIBLINEAR [9] or SGD (stochastic gradient descent) [2]. In this study, we focus on linearizing the  $p$ GMM kernels via hashing and we will also discuss how to linearize the  $e$ GMM and  $\gamma$ GMM kernels.

Next, we present an experimental study on the a large number of classification tasks using the variety of kernels we have discussed.

**Table 1: Public classification datasets and  $l_2$ -regularized kernel SVM results. We report the test classification accuracies for the linear kernel, the best-tuned RBF kernel, the original (tuning-free) GMM kernel, the best-tuned  $e$ GMM,  $p$ GMM, and  $\gamma$ GMM kernels, at their individually-best SVM regularization  $C$  values. All datasets are from the UCI repository except for the last 11 datasets, which were used by [17, 20] for testing deep learning algorithms and tree methods.**

Dataset	# train	# test	# dim	linear	RBF	GMM	$e$ GMM ( $\lambda_e$ )	$p$ GMM ( $p$ )	$\gamma$ GMM ( $\gamma$ )
Car	864	864	6	71.53	94.91	98.96	99.31 (2)	<b>99.54</b> (2)	99.31 (6)
Coverttype25k	25000	25000	54	62.64	82.66	82.65	88.32 (20)	83.25 (0.6)	<b>88.34</b> (20)
CTG	1063	1063	35	60.59	89.75	88.81	88.81 (0.01)	<b>100.00</b> (0.1)	90.78 (0.3)
DailySports	4560	4560	5625	77.70	97.61	99.61	99.61 (0.2)	99.61 (0.6)	<b>99.63</b> (0.8)
DailySports2k	2000	7120	5625	72.16	93.71	98.99	99.00 (0.1)	99.07 (.75)	<b>99.16</b> (0.3)
Dexter	300	300	19999	92.67	93.00	94.00	94.00 (17)	94.67 (0.5)	<b>95.67</b> (0.2)
EEGEye	7490	7490	14	61.46	86.82	78.54	95.51 (1000)	87.65 (15)	91.20 (60)
Gesture	4937	4936	32	37.22	61.06	65.50	66.67 (1.9)	66.33 (0.6)	<b>67.16</b> (2.6)
ImageSeg	210	2100	19	83.81	91.38	95.05	95.38 (1.2)	<b>95.57</b> (0.6)	95.38 (1.8)
Isolet2k	2000	5797	617	93.95	95.55	95.53	95.55 (0.2)	95.53 (1.0)	<b>95.60</b> (0.8)
MHealth20k	20000	20000	23	72.62	82.65	85.28	85.33 (0.5)	<b>86.69</b> (0.5)	85.31 (1.3)
MiniBooNE20k	20000	20000	50	88.42	93.06	93.00	93.01 (0.2)	<b>93.69</b> (0.6)	93.01 (1.1)
MSD20k	20000	20000	90	66.72	68.07	71.05	71.18 (0.2)	<b>71.84</b> (0.5)	71.44 (0.6)
Magic	9150	9150	10	78.04	84.43	87.02	86.93 (0.3)	<b>87.57</b> (0.5)	87.09 (0.8)
Musk	3299	3299	166	95.09	<b>99.33</b>	99.24	99.24 (0.3)	99.24 (1.0)	99.24 (1.0)
Musk2k	2000	4598	166	94.80	97.63	98.02	98.02 (0.01)	<b>98.06</b> (1.25)	<b>98.06</b> (0.5)
PageBlocks	2737	2726	10	95.87	97.08	96.56	96.56 (1.4)	<b>97.30</b> (0.1)	96.64(0.8)
Parkinson	520	520	26	61.15	66.73	69.81	<b>70.19</b> (0.6)	69.81 (1.0)	<b>70.19</b> (1.7)
PAMAP101	20000	20000	51	76.86	96.68	98.91	98.91 (0.1)	<b>99.00</b> (1.5)	98.92 (1.1)
PAMAP102	20000	20000	51	81.22	95.67	<b>98.78</b>	98.77 (0.01)	<b>98.78</b> (2)	<b>98.78</b> (1.7)
PAMAP103	20000	20000	51	85.54	97.89	99.69	<b>99.70</b> (0.01)	99.69 (1.0)	<b>99.70</b> (0.8)
PAMAP104	20000	20000	51	84.03	97.32	99.30	<b>99.31</b> (0.6)	99.30 (1.0)	<b>99.31</b> (1.3)
PAMAP105	20000	20000	51	79.43	97.34	99.22	99.24 (1.1)	99.22 (0.75)	<b>99.26</b> (1.8)
PIMA	384	384	8	66.67	71.35	76.30	77.08 (12)	76.56 (0.75)	76.82 (9.5)
RobotNavi	2728	2728	24	69.83	90.69	96.85	96.77 (0.1)	<b>98.20</b> (0.1)	97.65 (0.3)
Satimage	4435	2000	36	72.45	85.20	90.40	<b>91.85</b> (35)	90.95 (5)	91.35(9.5)
SEMG1	900	900	3000	26.00	<b>43.56</b>	41.00	41.22 (0.1)	42.89 (0.25)	42.11 (1.7)
SEMG2	1800	1800	2500	19.28	29.00	54.00	54.00 (0.3)	<b>56.11</b> (2)	55.22 (0.6)
Sensorless	29255	29254	48	61.53	93.01	99.39	99.38 (0.01)	<b>99.76</b> (0.5)	99.62 (0.5)
Shuttle500	500	14500	9	91.81	99.52	99.65	99.65 (0.1)	99.66 (0.5)	<b>99.68</b> (0.4)
SkinSeg10k	10000	10000	3	93.36	99.74	99.81	<b>99.90</b> (20)	99.85 (5)	99.87 (8.5)
SpamBase	2301	2300	57	85.91	92.57	94.17	94.13 (0.6)	<b>95.78</b> (0.25)	94.17 (1.0)
Splice	1000	2175	60	85.10	90.02	95.22	<b>96.46</b> (5)	95.26 (1.25)	<b>96.46</b> (5)
Theorem	3059	3059	51	67.83	70.48	71.53	<b>71.69</b> (1.6)	71.53 (1.0)	<b>71.76</b> (2.1)
Thyroid	3772	3428	21	95.48	97.67	98.31	98.34 (0.3)	<b>99.10</b> (0.1)	98.63 (0.6)
Thyroid2k	2000	5200	21	94.90	97.00	98.40	98.40 (0.01)	<b>98.96</b> (0.1)	98.62 (0.6)
Urban	168	507	147	62.52	51.48	66.08	65.68 (0.5)	<b>83.04</b> (0.1)	67.26 (0.2)
Vertebral	155	155	6	80.65	83.23	89.04	<b>89.68</b> (1.4)	89.04 (1.0)	<b>89.68</b> (1.1)
Vowel	264	264	10	39.39	94.70	96.97	<b>98.11</b> (5)	96.97 (1.0)	<b>98.11</b> (4)
Wholesale	220	220	6	89.55	90.91	93.18	93.18 (0.6)	<b>93.64</b> (1.25)	<b>93.64</b> (0.2)
Wilt	4339	500	5	62.60	83.20	87.20	<b>87.60</b> (1.1)	87.40 (0.75)	<b>87.60</b> (1.7)
YoutubeAudio10k	10000	11930	2000	41.35	48.63	50.59	50.60 (0.01)	<b>51.84</b> (0.6)	50.84 (0.9)
YoutubeHOG10k	10000	11930	647	62.77	66.20	68.63	68.65 (0.01)	<b>72.06</b> (0.5)	68.63 (1.1)
YoutubeMotion10k	10000	11930	64	26.24	28.81	31.95	<b>33.05</b> (4)	32.65 (0.6)	32.98 (4.5)
YoutubeSaiBoxes10k	10000	11930	7168	46.97	49.31	51.28	51.22 (0.001)	<b>52.15</b> (0.6)	51.39 (0.8)
YoutubeSpectrum10k	10000	11930	1024	26.81	33.54	39.23	39.27 (0.1)	<b>41.23</b> (0.5)	39.28 (1.1)
M-Basic	12000	50000	784	89.98	<b>97.21</b>	96.34	96.47 (1.2)	96.40 (0.5)	96.84 (2.3)
M-Image	12000	50000	784	70.71	77.84	80.85	81.20 (1.5)	<b>89.53</b> (50)	81.32 (2.1)
M-Noise1	10000	4000	784	60.28	66.83	71.38	71.70 (0.5)	<b>85.20</b> (80)	71.90 (2.8)
M-Noise2	10000	4000	784	62.05	69.15	72.43	72.80 (3)	<b>85.40</b> (70)	72.95 (2.8)
M-Noise3	10000	4000	784	65.15	71.68	73.55	74.70 (3)	<b>86.55</b> (50)	74.83 (3)
M-Noise4	10000	4000	784	68.38	75.33	76.05	76.80 (2.5)	<b>86.88</b> (60)	77.03 (2.8)
M-Noise5	10000	4000	784	72.25	78.70	79.03	79.48 (3)	<b>87.33</b> (30)	79.70 (3.5)
M-Noise6	10000	4000	784	78.73	85.33	84.23	84.58 (2)	<b>88.15</b> (20)	84.68 (4)
M-Rand	12000	50000	784	78.90	85.39	84.22	84.95 (4)	<b>89.09</b> (40)	85.17 (3.5)
M-Rotate	12000	50000	784	47.99	<b>89.68</b>	84.76	86.02 (1.6)	86.52 (0.25)	87.33 (2.1)
M-RotImg	12000	50000	784	31.44	45.84	40.98	42.88 (4)	<b>54.58</b> (20)	43.22 (3.5)

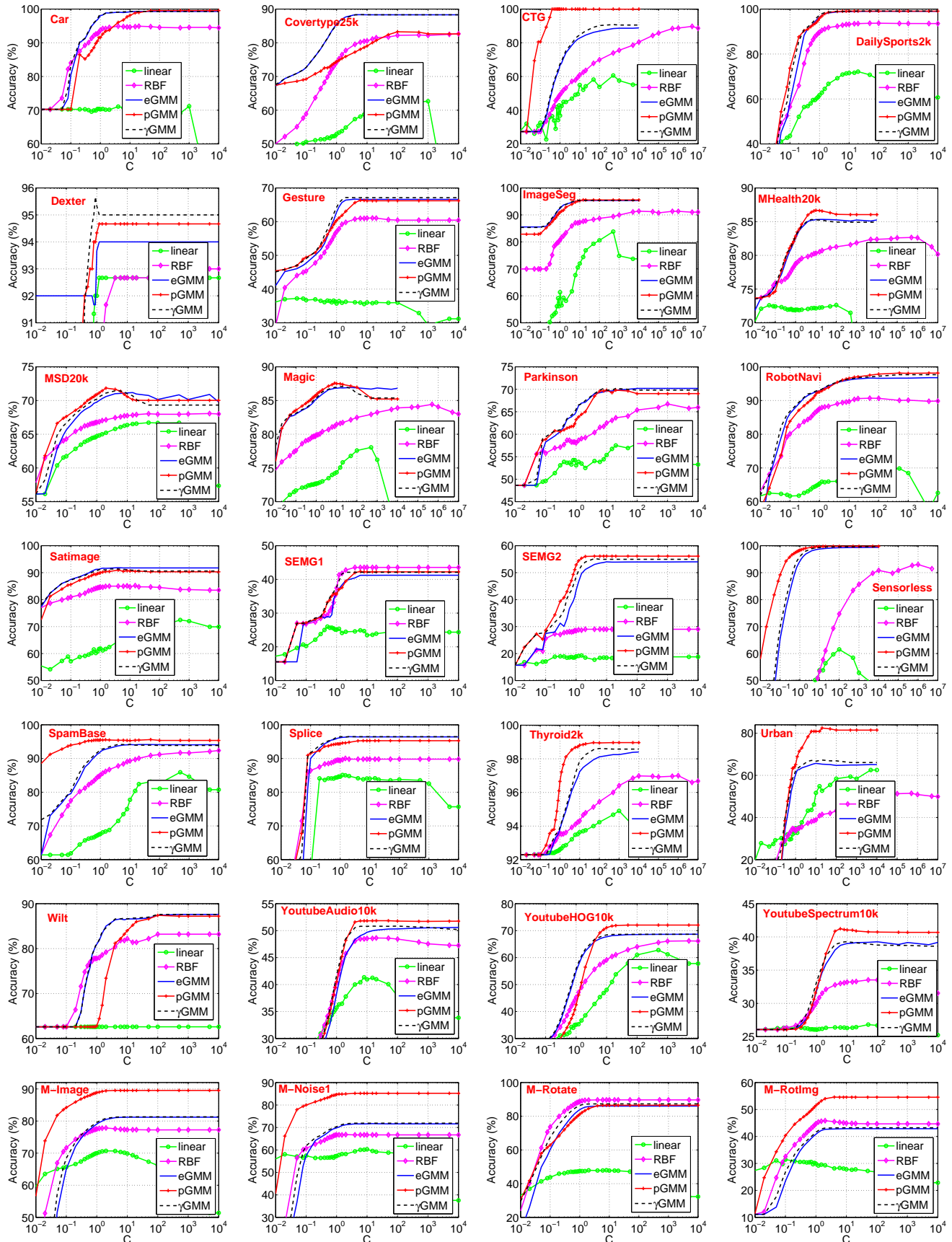


Figure 2: Test classification accuracies of various kernels using LIBSVM pre-computed kernel functionality. The results are presented with respect to  $C$ , which is the  $l_2$ -regularized kernel SVM parameter. For RBF, eGMM, pGMM, and  $\gamma$ GMM, at each  $C$ , we report the best test accuracies from a wide range of kernel parameter values.

## 2 EXPERIMENTAL STUDY ON KERNEL SVMs

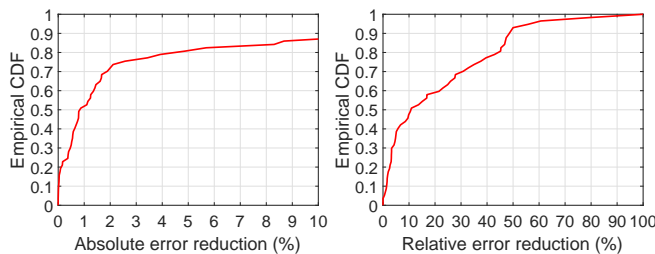
### 2.1 Basic Kernels: $e$ GMM, $p$ GMM, $\gamma$ GMM

Table 1 lists a large number of publicly available datasets from the UCI repository plus the 11 datasets (the last 11 datasets whose names start with “M-”) used by the deep learning literature [17]. In this table, we report the kernel SVM test classification results for a variety of kernels: linear, RBF, GMM,  $e$ GMM,  $p$ GMM,  $\gamma$ GMM.

In all the experiments, we adopt the  $l_2$ -regularization (with a regularization parameter  $C$ ) and report the test classification accuracies at the best  $C$  values in Table 1. More detailed results for a wide range of  $C$  values are reported in Figures 2. To ensure repeatability, we use the LIBSVM pre-computed kernel functionality, at the significant cost of disk space. For the RBF kernel, we exhaustively experiment with 58 different values of  $\lambda_e \in \{0.001, 0.01, 0.1:0.1:2, 2.5, 3:1:20, 25:5:50, 60:10:100, 120, 150, 200, 300, 500, 1000\}$ . Basically, Table 1 reports the best results among all  $C$  and  $\lambda_e$  values in our experiments. Here, 3:1:20 is the matlab notation, meaning that the iterations start at 3 and terminate at 20, at a space of 1.

For the  $e$ GMM kernel, we experiment with the same set of (58)  $\lambda_e$  values as for the RBF kernel. For the  $p$ GMM kernel, however, because we have to materialize (store) a kernel matrix for each  $\gamma$ , disk space becomes a serious concern. Therefore, for the  $p$ GMM kernel, we only search in the range of  $p \in \{0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.4, 0.5, 0.6, 0.75, 1, 1.25, 1.5, 2, 5, 10, 15, 20, 25, 30 : 10 : 100\}$ . For the  $\gamma$ GMM kernel, we experiment with  $\gamma \in \{0.05, 0.1 : 0.1 : 3, 3.5 : 0.5 : 10, 11 : 1 : 20, 30 : 10 : 100\}$ .

The classification results in Table 1 and Figures 2 confirm that the  $e$ GMM,  $p$ GMM, and  $\gamma$ GMM kernels typically improve the original GMM kernel. On a good fraction of datasets, the improvements can be very substantial. Figure 3 quantifies the improvements by plotting the empirical CDF (cumulative distribution function) of the absolute (left panel) and relative (right panel) error reduction (in %), obtained from using one of the  $e$ GMM,  $p$ GMM, or  $\gamma$ GMM kernels, compared to using the original GMM kernel.



**Figure 3: Empirical CDFs of the absolute (left panel) and relative (right panel) error reductions (%) obtained by using one of the  $e$ GMM,  $p$ GMM, or  $\gamma$ GMM kernels, compared to using the original GMM kernel. There are in total 57 datasets.**

For example, the left panel of Figure 3 says that, out of a total of 57 datasets (in Table 1), about 50% of the datasets exhibit an improvement of  $> 1\%$  absolute error reduction, and about 20% of the datasets exhibit an improvement of  $> 5\%$  absolute error reduction. The right panel says that about 50% of the datasets exhibit an improvement of  $> 10\%$  relative error reduction, and about 20% of the data exhibit an improvement of  $> 45\%$  relative error reduction.

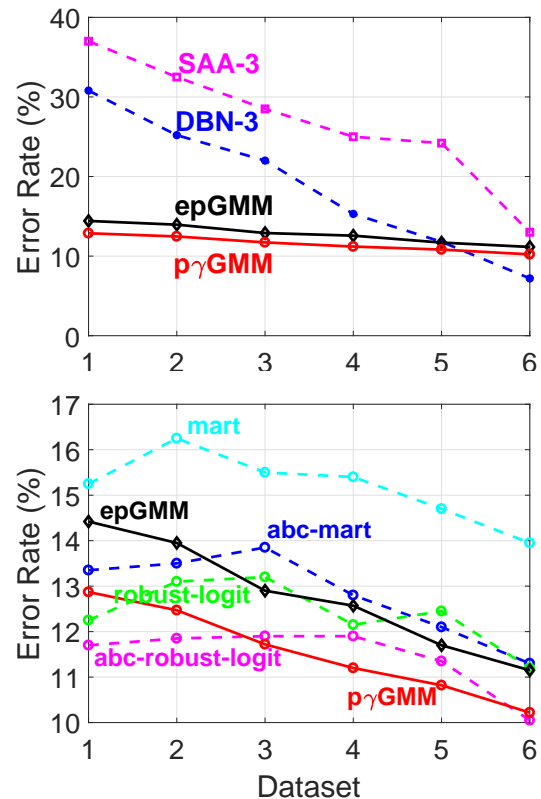
### 2.2 $ep$ GMM and $p\gamma$ GMM

Combining the three basic tunable GMM kernels produces four more tunable GMM kernels. Naturally we expect that, in many datasets, adding more tuning parameters would further improve the accuracies. In Table 2, we report the kernel SVM experimental results on the 11 datasets used by the deep learning literature [17], for the  $ep$ GMM and  $p\gamma$ GMM kernels. Since there are two parameters, the results are obtained by a two-dim grid search.

The results in Table 2 illustrate that the additional improvements can be also quite substantial. For example, on the M-RotImg dataset, the accuracy of the original GMM kernel is 40.98%, and the accuracy of the  $p$ GMM kernel is 54.58%. However, the accuracy of  $p\gamma$ GMM kernel becomes 59.92%.

We choose to show experiment on these 11 datasets because the prior work [20] in 2010 already conducted a thorough empirical study of a series of tree & boosting methods on the same datasets.

### 2.3 Comparisons with Trees



**Figure 4: Classification test error rates on M-Noise1, M-Noise2, ..., M-Noise6 datasets, for evaluating  $ep$ GMM and  $p\gamma$ GMM kernels. The upper panel includes the results of two deep learning algorithms (SAA-3 and DBN-3) as reported in [17]. The bottom panel compares the  $ep$ GMM and  $p\gamma$ GMM kernels with four boosted tree methods as reported in [20].**

Following [17, 20], we report the results on these 11 datasets in terms of the test error rates instead of accuracies, in Figure 4 (for M-Noise1, M-Noise2, M-Noise3, M-Noise4, M-Noise5, M-Noise6) and Table 3 (for M-Basic, M-Rotate, M-Image, M-Rand, M-RotImg).

**Table 2: We add the results (test classification accuracies) for the  $ep$ GMM and  $p\gamma$ GMM kernels as the last two columns, for the 11 datasets used in the deep learning literature [17] and later by [20] for comparing tree methods.**

Dataset	# train	# test	# dim	linear	RBF	GMM	$e$ GMM	$p$ GMM	$\gamma$ GMM	$ep$ GMM	$p\gamma$ GMM
M-Basic	12000	50000	784	89.98	<b>97.21</b>	96.34	96.47	96.40	96.84	96.71	97.00
M-Image	12000	50000	784	70.71	77.84	80.85	81.20	89.53	81.32	89.96	<b>90.96</b>
M-Noise1	10000	4000	784	60.28	66.83	71.38	71.70	85.20	71.90	85.58	<b>87.13</b>
M-Noise2	10000	4000	784	62.05	69.15	72.43	72.80	85.40	72.95	86.05	<b>87.53</b>
M-Noise3	10000	4000	784	65.15	71.68	73.55	74.70	86.55	74.83	87.10	<b>88.28</b>
M-Noise4	10000	4000	784	68.38	75.33	76.05	76.80	86.88	77.03	87.43	<b>88.80</b>
M-Noise5	10000	4000	784	72.25	78.70	79.03	79.48	87.33	79.70	88.30	<b>89.18</b>
M-Noise6	10000	4000	784	78.73	85.33	84.23	84.58	88.15	84.68	88.85	<b>89.78</b>
M-Rand	12000	50000	784	78.90	85.39	84.22	84.95	89.09	85.17	89.43	<b>90.63</b>
M-Rotate	12000	50000	784	47.99	<b>89.68</b>	84.76	86.02	86.56	87.33	88.36	89.06
M-RotImg	12000	50000	784	31.44	45.84	40.98	42.88	54.58	43.22	55.73	<b>59.92</b>

The results presented in Figure 4 and Table 3 are quite exciting, because, at this point, we merely use kernel SVM with single kernels. The performances of tunable GMM kernels are already comparable to four boosting & tree methods: mart, abc-mart, robust logitboost, and abc-robust-logitboost, whose training procedures are time-consuming with large model sizes (up to 10000 boosting iterations). It is reasonable to expect that additional improvements might be achieved in near future.

The “mart” tree algorithm [12] has been popular in industry practice, especially in search. At each boosting step, it uses the first derivative of the logistic loss function as the residual response to fit regression trees, to achieve excellent robustness and fairly good accuracy. The earlier work on “logitboost” [13] were believed to exhibit numerical issues (which in part motivated the development of mart). It turns out that the numerical issue does not actually exist after [20] derived the tree-split formula using both the first and second order derivatives of the logistic loss function. [20] showed the “robust logitboost” in general improves “mart”, as can be seen from Figure 4 and Table 3.

[18–20] made an interesting observation that the derivatives (as in text books) of the classical logistic loss function can be written in a different form for the multi-class case, by enforcing the “sum-to-zero” constraints. At each boosting step, they identify a “base class” either by the “worst-class” criterion [18] or the exhaustive search method as reported in [19, 20]. This “adaptive base class (abc)” strategy can be combined with either mart or robust logitboost; hence the names “abc-mart” and “abc-robust-logitboost”. The improvements due to the use of “abc” strategy can also be substantial. In all the tree implementations, they [18–20] always used the adaptive-binning strategy for simplifying the implementation and speeding up training. Also, they followed the “best-first” criterion whereas many tree implementations used balanced trees.

Table 3 reports the test error rates on five datasets: M-Basic, M-Rotate, M-Image, M-Rand, and M-RotImg. In group 1 (as reported in [17]), the results show that (i) the kernel SVM with RBF kernel outperforms the kernel SVM with polynomial kernel; (ii) deep learning algorithms usually beat kernel SVM and neural nets. Group 2 presents the same results as in Table 2 (in terms of error rates as

opposed to accuracies). In group 3, overall the tree methods especially abc-robust-logitboost achieve very good accuracies. The results of tunable GMM kernels are largely comparable.

The training of boosted trees is typically slow (especially in high-dimensional data) because a large number of trees are usually needed in order to achieve good accuracies. Consequently, the model sizes of tree methods are usually large. Therefore, it would be exciting to have methods which are much simpler than trees and achieve comparable accuracies.

### 3 HASHING THE $p$ GMM KERNEL

It is now well-understood that it is highly beneficial to be able to linearize nonlinear kernels so that learning algorithms can be easily scaled to massive data. Linearization can be done either through hashing [15, 21, 26] or the Nystrom method [28].

It turns out that developing a hashing method for the  $p$ GMM kernel is quite straightforward, by modifying the prior algorithms. Algorithm 1 summarizes the modified GCWS (generalized consistent weighted sampling).

**Algorithm 1** Modified generalized consistent weighted sampling (GCWS) for hashing the  $p$ GMM kernel with a tuning parameter  $p$ .

**Input:** Data vector  $u_i$  ( $i = 1$  to  $D$ )

Generate vector  $\tilde{u}$  in  $2D$ -dim by (1).

For  $i$  from 1 to  $2D$

$r_i \sim \text{Gamma}(2, 1)$ ,  $c_i \sim \text{Gamma}(2, 1)$ ,  $\beta_i \sim \text{Uniform}(0, 1)$

$t_i \leftarrow \lfloor p \frac{\log \tilde{u}_i}{r_i} + \beta_i \rfloor$ ,  $a_i \leftarrow \log(c_i) - r_i(t_i + 1 - \beta_i)$

End For

**Output:**  $i^* \leftarrow \arg \min_i a_i$ ,  $t^* \leftarrow t_{i^*}$

With  $k$  samples, we can estimate  $p$ GMM( $u, v$ ) according to the following collision probability:

$$\Pr \left\{ i_{\tilde{u},j}^* = i_{\tilde{v},j}^* \text{ and } t_{\tilde{u},j}^* = t_{\tilde{v},j}^* \right\} = p\text{GMM}(u, v), \quad (11)$$

or, for implementation convenience, the approximate collision probability [21]:

$$\Pr \left\{ i_{\tilde{u},j}^* = i_{\tilde{v},j}^* \right\} \approx p\text{GMM}(u, v) \quad (12)$$

**Table 3: Test error rates of five datasets reported in [17, 20]. The results in group 1 are from [17], where they compared kernel SVM, neural nets, and deep nets. The results in group 3 are from [20], which compared four boosted tree methods.**

Group	Method	M-Basic	M-Rotate	M-Image	M-Rand	M-RotImg
1	SVM-RBF	3.05%	11.11%	22.61%	14.58%	55.18%
	SVM-POLY	3.69%	15.42%	24.01%	16.62%	56.41%
	NNET	4.69%	18.11%	27.41%	20.04%	62.16%
	DBN-3	3.11%	<b>10.30%</b>	16.31%	<b>6.73%</b>	47.39%
	SAA-3	3.46%	<b>10.30%</b>	23.00%	11.28%	51.93%
	DBN-1	3.94%	14.69%	16.15%	9.80%	52.21%
2	Linear	10.02%	52.01%	29.29%	21.10%	68.56%
	RBF	<b>2.79%</b>	<b>10.32%</b>	22.16%	14.61%	54.16%
	GMM	3.64%	15.24%	19.15%	15.78%	59.02%
	eGMM	3.53%	13.98%	18.80%	15.05%	57.12%
	pGMM	3.60%	13.44%	10.47%	10.91%	45.42%
	$\gamma$ GMM	3.16%	12.67%	18.68%	14.83%	56.78%
	epGMM	3.29%	11.64%	10.04%	10.57%	44.27%
	p $\gamma$ GMM	3.00%	10.94%	9.04%	9.53%	<b>40.18%</b>
3	mart	4.12%	15.35%	11.64%	13.15%	49.82%
	abc-mart	3.69%	13.27%	9.45%	10.60%	46.14%
	robust logitboost	3.45%	13.63%	9.41%	10.04%	45.92%
	abc-robust-logitboost	3.20%	11.92%	<b>8.54%</b>	9.45%	44.69%

For each vector  $u$ , we obtain  $k$  random samples  $i_{u,j}^*$ ,  $j = 1$  to  $k$ . We store only the lowest  $b$  bits of  $i^*$ . We need to view those  $k$  integers as locations (of the nonzeros). For example, when  $b = 2$ , we should view  $i^*$  as a binary vector of length  $2^b = 4$ . We concatenate all  $k$  such vectors into a binary vector of length  $2^b \times k$  and then feed the new data vectors to a linear classifier if the task is classification. The storage and computational cost is largely determined by the number of nonzeros in each data vector, i.e., the  $k$  in our case. This scheme can also be used for other tasks including clustering, regression, and near neighbor search.

Figure 5 presents the experimental results on hashing for M-Rotate. For this dataset,  $p = 0.25$  is the best choice (among the range of  $p$  values we have searched). Figure 5 plots the results for both  $p = 0.25$  (left panels) and  $p = 1$  (right panels), for  $b \in \{12, 8, 4, 2\}$ . Recall here  $b$  is the number of bits for representing each hashed value. The results demonstrate that: (i) hashing using  $p = 0.25$  produces better results than hashing using  $p = 1$ ; (ii) It is preferable to use a fairly large  $b$  value, for example,  $b \geq 4$  or 8. Using smaller  $b$  values (e.g.,  $b = 2$ ) hurts the accuracy; (iii) With merely a small number of hashes (e.g.,  $k = 128$ ), the linearized pGMM kernel can significantly outperform the original linear kernel. Note that the original dimensionality is 784. This example illustrates the significant advantage of nonlinear kernel and hashing.

Figure 6 (for CTG dataset) and Figure 7 (for SpamBase dataset) are somewhat different from the previous figures. For both datasets, using  $\gamma = 0.05$  achieves the best accuracy. We plot the results for  $\gamma = 0.05, 0.25, 0.5, 0.75$ , and  $b = 8, 4, 2$ , to visualize the trend.

We have conducted significantly more experiments than we have presented here, but we hope they are convincing enough.

#### 4 DISCUSSION: HASHING $\gamma$ GMM AND eGMM

At least for  $\gamma$  being integers, it is possible to modify the Algorithm 1 to develop a hashing method for the  $\gamma$ GMM kernel. Basically, for each hash value, we just need to generate  $\gamma$  independent samples. For the original data vectors  $u$  and  $v$ , we require all  $\gamma$  samples of  $u$  to match all  $\gamma$  samples of  $v$ . The collision probability will be exactly equal to the  $\gamma$ GMM kernel.

Developing a hashing method for the eGMM is more challenging. A more straightforward approach is a two-stage hashing scheme. We first generate hashed values using Algorithm 1 (with  $p = 1$ ), then we apply random Fourier features (RFF) [29] on the hashed values. Based on the analysis in [23], the RFF method needs a very large number of samples in order to reach a satisfactory accuracy. Therefore, we do not expect this two-stage scheme would be practical for hashing the eGMM kernel.

#### 5 CONCLUSION

It is commonly believed that deep learning algorithms and tree methods can produce the state-of-the-art results in many statistical machine learning tasks. In 2010, [20] reported a set of surprising experiments on the datasets used by the deep learning community [17], to show that tree methods can outperform deep nets on a majority (but not all) of those datasets and the improvements can be substantial on a good portion of datasets. [20] introduced several ideas including the second-order tree-split formula and the new derivatives for multi-class logistic loss function. Nevertheless, tree methods are slow and their model sizes are typically large.

In machine learning practice with massive data, it is desirable to develop algorithms which run almost as efficient as linear methods (such as linear logistic regression or linear SVM) and achieve similar accuracies as nonlinear methods. In this study, the tunable linearized GMM kernels are promising tools for achieving those goals. Our extensive experiments on the same datasets used

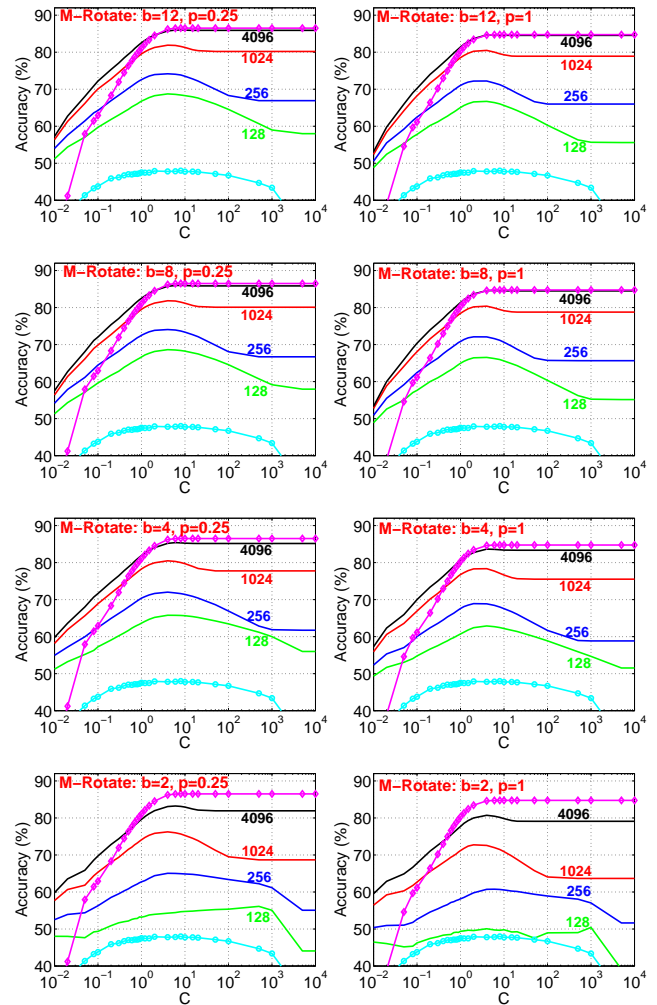


for testing tree methods and deep nets demonstrate that tunable GMM kernels and their linearized versions through hashing can achieve comparable accuracies as trees. On some datasets, “abc-robust-logitboost” achieves better accuracies than the proposed tunable GMM kernels. Also, on some datasets, deep learning methods or RBF kernel SVM outperform tunable GMM kernels. Therefore, there is still room for future improvements.

In this study, we focus on testing tunable GMM kernels and their linearized versions using classification tasks. It is clear that these techniques basically generate new data representations and hence can be applied to a wide variety of statistical learning tasks including clustering and regression. Due to the discrete name of the hashed values, the techniques naturally can also be used for building hash tables for fast near neighbor search.

## REFERENCES

- [1] Michael Bendersky and W. Bruce Croft. 2009. Finding text reuse on the web. In *WSDM*. Barcelona, Spain, 262–271.
- [2] Leon Bottou. <http://leon.bottou.org/projects/sgd>. (????).
- [3] Léon Bottou, Olivier Chapelle, Dennis DeCoste, and Jason Weston (Eds.). 2007. *Large-Scale Kernel Machines*. The MIT Press, Cambridge, MA.
- [4] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. 1997. Syntactic clustering of the Web. In *WWW*. Santa Clara, CA, 1157 – 1166.
- [5] Moses S. Charikar. 2002. Similarity estimation techniques from rounding algorithms. In *STOC*. Montreal, Canada, 380–388.
- [6] Ludmila Cherkasova, Kave Eshghi, Charles B. Morrey III, Joseph Tucek, and Alstair C. Veitch. 2009. Applying syntactic similarity algorithms for enterprise information management. In *KDD*. Paris, France, 1087–1096.
- [7] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, Michael Mitzenmacher, Alessandro Panconesi, and Prabhakar Raghavan. 2009. On compressing social networks. In *KDD*. Paris, France, 219–228.
- [8] Yon Dourisboure, Filippo Geraci, and Marco Pellegrini. 2009. Extraction and classification of dense implicit communities in the Web graph. *ACM Trans. Web* 3, 2 (2009), 1–36.
- [9] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. LIBLINEAR: A Library for Large Linear Classification. *Journal of Machine Learning Research* 9 (2008), 1871–1874.
- [10] Dennis Fetterly, Mark Manasse, Marc Najork, and Janet L. Wiener. 2003. A large-scale study of the evolution of web pages. In *WWW*. Budapest, Hungary, 669–678.
- [11] George Forman, Kave Eshghi, and Jaap Suermondt. 2009. Efficient detection of large-scale redundancy in enterprise file systems. *SIGOPS Oper. Syst. Rev.* 43, 1 (2009), 84–91.
- [12] Jerome H. Friedman. 2001. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics* 29, 5 (2001), 1189–1232.
- [13] Jerome H. Friedman, Trevor J. Hastie, and Robert Tibshirani. 2000. Additive logistic regression: a statistical view of boosting. *The Annals of Statistics* 28, 2 (2000), 337–407.
- [14] Sreenivas Gollapudi and Aneesh Sharma. 2009. An axiomatic approach for result diversification. In *WWW*. Madrid, Spain, 381–390.
- [15] Sergey Ioffe. 2010. Improved Consistent Sampling, Weighted Minhash and L1 Sketching. In *ICDM*. Sydney, AU, 246–255.
- [16] Jon Kleinberg and Eva Tardos. 1999. Approximation Algorithms for Classification Problems with Pairwise Relationships: Metric Labeling and Markov Random Fields. In *FOCS*. New York, 14–23.
- [17] Hugo Larochelle, Dumitru Erhan, Aaron C. Courville, James Bergstra, and Yoshua Bengio. 2007. An empirical evaluation of deep architectures on problems with many factors of variation. In *ICML*. Corvallis, Oregon, 473–480.
- [18] Ping Li. 2008. Adaptive Base Class Boost for Multi-class Classification. *CoRR* abs/0811.1250 (2008).
- [19] Ping Li. 2009. ABC-Boost: Adaptive Base Class Boost for Multi-Class Classification. In *ICML*. Montreal, Canada, 625–632.
- [20] Ping Li. 2010. Robust LogitBoost and Adaptive Base Class (ABC) LogitBoost. In *UAI*. Catalina Island, CA.
- [21] Ping Li. 2015. 0-Bit Consistent Weighted Sampling. In *KDD*. Sydney, Australia, 665–674.
- [22] Ping Li. 2016. *Linearized GMM Kernels and Normalized Random Fourier Features*. Technical Report. arXiv:1605.05721.
- [23] Ping Li. 2017. Linearized GMM Kernels and Normalized Random Fourier Features. In *KDD*. 315–324.



**Figure 5: Test classification accuracies for using linear classifiers combined with hashing in Algorithm 1 on M-Rotate dataset, for  $p = 0.25$  (left panels) and  $p = 1$  (right panels), and  $b \in \{12, 8, 4, 2\}$ . In each panel, the four solid curves correspond to results with  $k$  hashes for  $k \in \{128, 256, 1024, 4096\}$ . For comparisons, each panel also plots the results of linear classifiers on the original data (lower marked curve) and the results of  $p$ GMM kernel SVMs (higher marked curve).**

- [24] Ping Li and Arnd Christian König. 2010. b-Bit Minwise Hashing. In *WWW*. Raleigh, NC, 671–680.
- [25] Ping Li, Anshumali Shrivastava, Joshua Moore, and Arnd Christian König. 2011. Hashing Algorithms for Large-Scale Learning. In *NIPS*. Granada, Spain, 2672–2680.
- [26] Mark Manasse, Frank McSherry, and Kunal Talwar. 2010. *Consistent Weighted Sampling*. Technical Report MSR-TR-2010-73. Microsoft Research.
- [27] Marc Najork, Sreenivas Gollapudi, and Rina Panigrahy. 2009. Less is more: sampling the neighborhood graph makes SALSA better and faster. In *WSDM*. Barcelona, Spain, 242–251.
- [28] E. J. Nyström. 1930. Über Die Praktische Auflösung von Integralgleichungen mit Anwendungen auf Randwertaufgaben. *Acta Mathematica* 54, 1 (1930), 185–204.
- [29] A. Rahimi and B. Recht. 2007. Random features for large-scale kernel machines. In *NIPS*. Vancouver, Canada, 1177–1184.
- [30] Bernhard Schölkopf and Alexander J. Smola. 2002. *Learning with Kernels*. The MIT Press, Cambridge, MA.



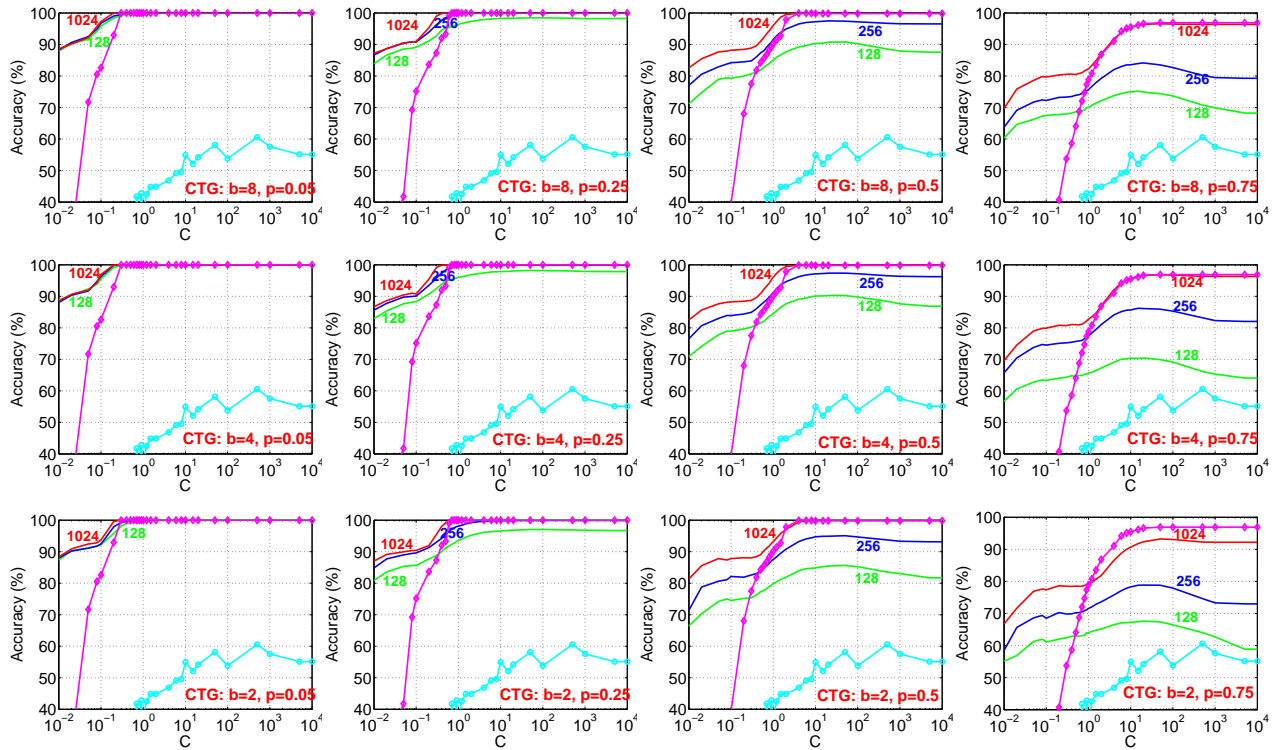


Figure 6: Test classification accuracies for using linear classifiers combined with hashing in Algorithm 1 on CTG dataset, for  $p \in \{0.05, 0.25, 0.5, 0.75\}$  to visualize the trend that, for this dataset, the accuracy decreases with increasing  $p$ .

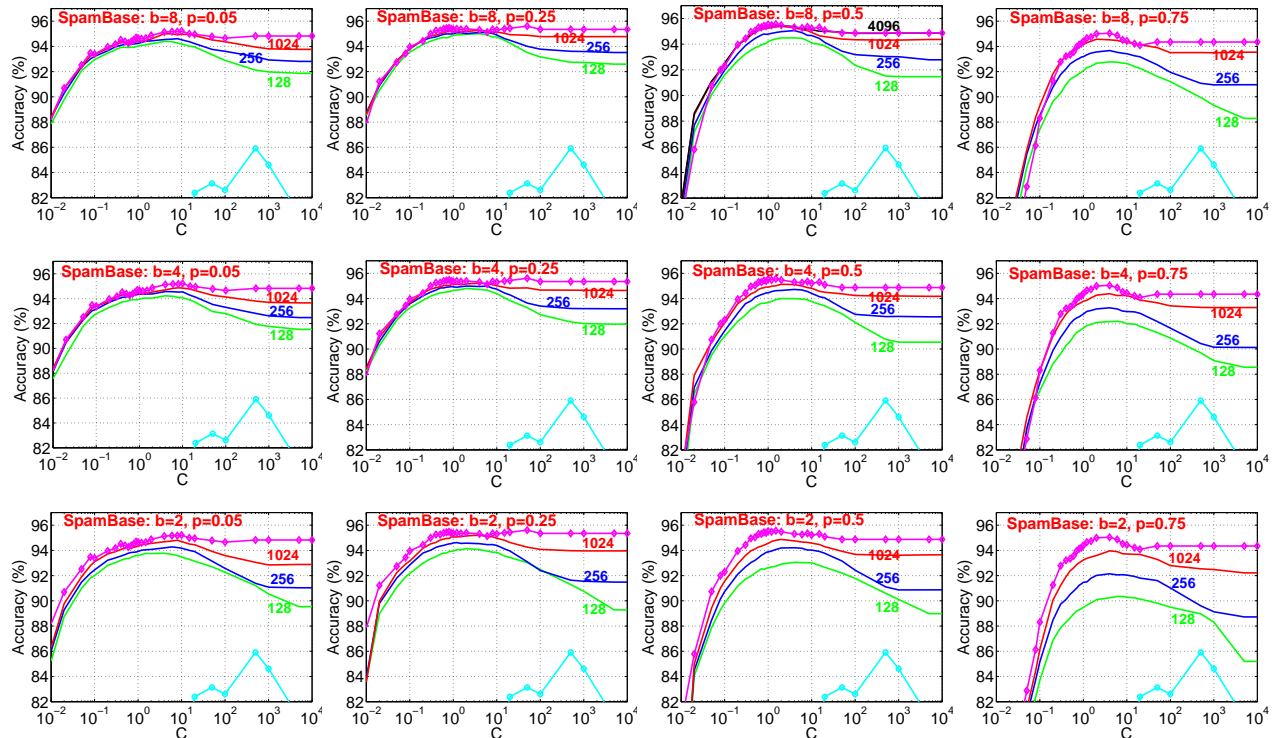


Figure 7: Test classification accuracies for using linear classifiers combined with hashing in Algorithm 1 on SpamBase dataset, for  $p \in \{0.05, 0.25, 0.5, 0.75\}$  to visualize the trend that, for this dataset, the accuracy decreases with increasing  $p$ .