

Table 1 displays the k NN error for $k=3, 7$ and 11 across various datasets (Isolnet, USPS, letters, DSLR, Amazon, Webcam, Caltech) and methods (Euclidean, LMNN, GB-LMNN, MLCR, ITML, NCA, and **ours**). The error is shown as a mean value with standard deviation. The best performing methods are shown in bold. The only non-linear metric learning method in the above is GB-LMNN.

Table 1 (k=3):

Dataset	Isolnet	USPS	letters	DSLR	Amazon	Webcam	Caltech
Euclidean	170	256	76	866	800	800	800
LMNN	779.4	9298	20000	157	958	295	1123
GB-LMNN	26	70	26	70	70	70	70
MLCR	8.66	6.16	4.79	73.26	60.13	56.27	80.5
ITML	4.43	5.48	3.26	24.17	26.72	13.59	46.93
NCA	4.13	5.48	2.92	21.65	26.72	13.56	46.11
ours	6.61	8.27	14.25	36.92	24.61	23.05	46.76

Table 1 (k=7):

Dataset	Isolnet	USPS	letters	DSLR	Amazon	Webcam	Caltech
Euclidean	7.44	6.08	5.46	76.45	62.21	57.29	80.76
LMNN	3.78	4.9	3.58	25.44	29.23	14.58	46.75
GB-LMNN	3.54	4.9	2.66	25.44	29.12	12.45	46.17
MLCR	5.6	8.25	19.92	33.72	23.17	18.98	46.85
ITML	7.57	5.68	5.87	22.32	31.42	10.85	51.74
NCA	6.09	5.83	5.28	36.94	29.22	22.03	45.50
ours	4.61	4.5	2.54	21.61	22.44	11.19	41.61

Table 1 (k=11):

Dataset	Isolnet	USPS	letters	DSLR	Amazon	Webcam	Caltech
Euclidean	8.92	6.86	5.89	73.87	64.61	59.66	81.39
LMNN	3.72	4.76	4.69	23.64	30.12	13.90	49.66
GB-LMNN	3.98	4.78	2.86	23.64	30.07	13.90	49.15
MLCR	5.71	11.11	13.54	36.25	24.32	17.97	44.97
ITML	7.77	6.53	6.52	22.28	30.48	11.86	50.76
NCA	5.90	5.73	6.04	40.06	30.69	26.44	46.48
ours	4.11	4.98	5.05	22.28	24.11	11.19	40.76

Table 1: k NN error, for $k=3, 7$ and 11 . Features were scaled by z-scoring. Mean and standard deviation are shown for data sets on which 5-fold partition was used. Best performing methods are shown in bold. Note that the only non-linear metric learning method in the above is GB-LMNN.

than previously proposed similarity learning methods. Our learning algorithm is simple yet efficient, converging on all the data sets we have experimented upon in reasonable time as compared to the competing methods.

Our choice of Frobenius regularizer is motivated by desire to control model complexity without biasing towards a particular form of the matrix. We have experimented with alternative regularizers, both the trace norm of \mathbf{W} and the shrinkage towards Euclidean distance, $\|\mathbf{W} - \mathbf{I}\|_F^2$, but found both to be inferior to $\|\mathbf{W}\|_F^2$. We suspect that often the optimal \mathbf{W} corresponds to a highly anisotropic scaling of data dimensions, and thus bias towards \mathbf{I} may be unhealthy.

The results in this paper are restricted to Mahalanobis metric, which is an appealing choice for a number of reasons. In particular, learning such metrics is equivalent to learning linear embedding of the data, allowing very efficient methods for metric search. Still, one can consider non-linear embeddings $\mathbf{x} \rightarrow \phi(\mathbf{x}; \mathbf{w})$ and define the distance D in terms of the embeddings, for example, as $D(\mathbf{x}, \mathbf{x}_i) = \|\phi(\mathbf{x}) - \phi(\mathbf{x}_i)\|$ or as $-\phi(\mathbf{x})^T \phi(\mathbf{x}_i)$. Learning S in the latter form can be seen as learning a kernel with discriminative objective of improving k NN performance. Such a model would be more expressive, but also more challenging to optimize. We are investigating this direction.

Acknowledgments

This work was partly supported by NSF award IIS-1409837.