

WHEN METRIC LEARNING GOES HAND IN HAND WITH DIMENSIONALITY REDUCTION

Masoud Faraki, Mehrtash T. Harandi & Fatih Porikli

Research School of Engineering
Australian National University
Canberra, ACT 0200, Australia

Data61, CSIRO
Locked Bag 8001
Canberra, ACT 2601, Australia

{masoud.faraki, mehrtash.harandi, fatih.porikli}@data61.csiro.au

ABSTRACT

Despite its attractive properties, the recently introduced Keep It Simple and Straightforward METric learning (KISSME) method is very reliant on PCA to prepare neat input data. This dependency can lead to difficulties, e.g., when the dimensionality is not meticulously set. To address this issue, we devise a unified formulation for joint dimensionality reduction and metric learning based on the KISSME algorithm. Our joint formulation is expressed as an optimization problem on the Grassmann manifold, hence enjoys properties of Riemannian optimization techniques. We also devise end-to-end learning of a generic deep network for metric learning using our derivation.

1 INTRODUCTION

Metric learning algorithms are of practical interest when learning from large number of categories (with limited training samples per category) is deemed, if the machinery is meant to deal with unseen classes (e.g., retrieval), or if weaker forms of supervision are considered. In such scenarios, conventional classification approaches are either not applicable or may fail badly.

The Keep It Simple and Straightforward METric learning (KISSME) Koestinger et al. (2012) algorithm is agnostic to the class labels and learns a metric purely from a set of equivalence constraints (similar/dissimilar pairs). Furthermore, the algorithm scales beautifully to large scale problems, makes it a suitable -if not perfect- match for the aforementioned problems.

Some recent studies (e.g., Xiong et al. (2014)) show that the KISSME algorithm is successful only when its input is carefully processed and denoised using PCA. This in turn results in high degree of sensitivity to the dimensionality of the PCA step as shown in Fig. 1. In this paper, we propose to learn a low-dimensional subspace along its metric in the spirit of the KISSME algorithm in a unified fashion. Furthermore, based on our derivation, we show end-to-end learning of a generic deep Convolutional Neural Network (CNN) for metric learning.

Contributions.

1. We propose the Joint Dimensionality Reduction for KISSME formulation (JDR-KISSME) that learns a low-dimensional space along its metric in the spirit of the KISSME algorithm. In doing so, we benefit from the optimization techniques on Riemannian manifolds Absil et al. (2009) and in particular the geometry of Grassmann manifolds.
2. Upon our development, we propose a few simple, yet effective steps, to train a deep network for metric learning using KISSME verification signal as supervision.

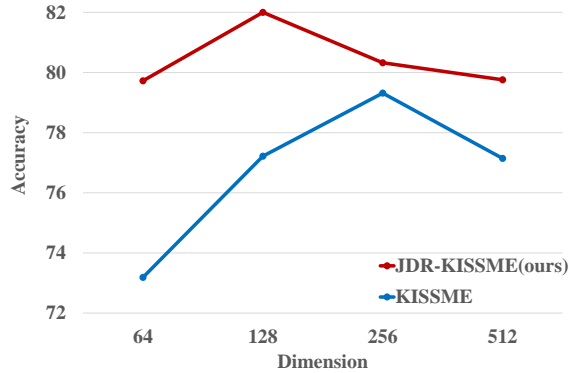


Figure 1: Verification accuracy against dimensionality of the space for an experiment using the surveillance nature images of the CompCars dataset Yang et al. (2015) (see §8 for more details).

2 RELATED WORK

We review some notable examples of conventional and deep metric learning techniques here. We start by two early studies in restricted and unrestricted metric learning¹, namely Mahalanobis Metric for Clustering (MMC) Xing et al. (2003) and Large Margin Nearest Neighbor (LMNN) Weinberger & Saul (2009).

MMC aims to minimize sum of distances over similar pairs while ensuring dissimilar pairs are far apart. The problem is formulated as an iterative gradient descent algorithm where at each iteration the obtained solution is projected back to the set of Positive Semi-Definite (PSD) matrices to ensure the metric is proper. Projection onto the PSD cone requires eigenvalue decomposition, making MMC computationally expensive when dealing with high-dimensional data.

LMNN learns a local linear transformation of labeled input data to improve kNN classification accuracy. Ideally, the learned transformation is expected to unite the k-nearest neighbors of each point sharing the same label (called target neighbors) while minimizing the number of impostors. This is implemented as a semi-definite programming problem and solved by iterating between a gradient descent step followed by projecting the solution onto the PSD cone.

While convexity is an attractive property, recent studies make use of non-convex objective functions to yield more accurate solutions. Two notable examples are the Linear Discriminant Metric Learning (LDML) Guillaumin et al. (2009; 2012) and Pairwise Constrained Component Analysis (PCCA) Mignon & Jurie (2012). The LDML algorithm measures the probability of a pair being similar by a logistic function. The probabilities are then used to determine a projection that maximizes a log-likelihood function. PCCA learns a transformation to project similar pairs inside a ball (with a given radius) while dissimilar pairs are pushed away. Optimization in both algorithms is performed by making use of the gradient descent method.

METRIC LEARNING AND DEEP NETS

Similarity/distance metric learning using deep nets originates from the advent of Siamese networks Bromley et al. (1993); Chopra et al. (2005). In recent years, deep metric learning has received growing attention, following the trend of deep CNNs in solving large-scale classification problems Krizhevsky et al. (2012). For example, in the spirit of PCCA, a two layer discriminative network for face verification is proposed in Hu et al. (2014). Sun et al. (2014), use an ensemble of networks, each operating on a different face patch for face verification. The networks are trained by making use of a combination of classification and verification cost functions. The cross-entropy loss is used for the classification loss. As for the verification loss, cost function encodes an l_2 -margin between face images.

¹Algorithms in restricted metric learning scenario do not have access to the class labels of the samples. These algorithms can also work in the unrestricted scenario while the opposite is not often the case. This makes the restricted algorithms more appealing as they can address a broader range of problems.

Very recent works in deep metric learning include the work of Hoffer & Ailon (2015) and Schroff et al. (2015) in which an LMNN based triplet loss layer is used to direct CNN parameter learning. Song et al. (2016) show careful construction of batches such as including hard triplets during training, leads to better clustering and retrieval qualities.

3 PROPOSED METHOD

In this part, we present our idea to jointly learn a low-dimensional space and its metric. Throughout the paper, we use bold lower-case letters (e.g., \mathbf{x}) to denote vectors and bold upper-case letters (e.g., \mathbf{X}) to show matrices. \mathbf{I}_d is the $d \times d$ identity matrix. The Frobenius norm of a matrix is $\|\mathbf{X}\|_F = \sqrt{\text{Tr}(\mathbf{X}^T \mathbf{X})}$, where $\text{Tr}(\cdot)$ indicates the matrix trace. $\det(\cdot)$ shows the matrix determinant. $[\cdot]_+$ indicates the hinge loss function, i.e., $\max(0, \cdot)$. The set of d -dimensional subspaces in \mathbb{R}^D forms a smooth manifold known as the Grassmann manifold \mathcal{G}_D^d . The Riemannian metric for \mathcal{G}_D^d is defined as the normal metric through the quotient reconstruction (see Edelman et al. (1998)).

3.1 PROBLEM STATEMENT

Let $\{\mathbf{x}_i, \tilde{\mathbf{x}}_i\}_{i=1}^N$, $\mathbf{x}_i, \tilde{\mathbf{x}}_i \in \mathbb{R}^D$ be a set of N training pairs. Furthermore, let $l_i \in \{1, -1\}$ denote the label of the i -th pair with $l_i = 1$ indicating that the corresponding pair is similar and $l_i = -1$ otherwise. Koestinger et al. propose to measure the probability of a pair being (dis)similar by multivariate Gaussian distributions Koestinger et al. (2012). In particular and making use of the symmetry in decisions², the function

$$\delta(\mathbf{x}_i, \tilde{\mathbf{x}}_i) = \frac{\mathcal{N}(\mathbf{x}_i - \tilde{\mathbf{x}}_i | \Sigma_D, \mathbf{0})}{\mathcal{N}(\mathbf{x}_i - \tilde{\mathbf{x}}_i | \Sigma_S, \mathbf{0})} \quad (1)$$

models the dissimilarity value of the pair $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ with $\mathcal{N}(\mathbf{x} | \Sigma, \mu)$ denoting the multivariate Gaussian with mean μ and covariance Σ evaluated at \mathbf{x} . Here,

$$\Sigma_D = \frac{1}{\#(l_i = -1)} \sum_{i, l_i = -1} (\mathbf{x}_i - \tilde{\mathbf{x}}_i)(\mathbf{x}_i - \tilde{\mathbf{x}}_i)^T. \quad (2)$$

$$\Sigma_S = \frac{1}{\#(l_i = +1)} \sum_{i, l_i = +1} (\mathbf{x}_i - \tilde{\mathbf{x}}_i)(\mathbf{x}_i - \tilde{\mathbf{x}}_i)^T. \quad (3)$$

Intuitively, if the pair $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ is dissimilar, one expects $\mathcal{N}(\mathbf{x}_i - \tilde{\mathbf{x}}_i | \Sigma_D, \mathbf{0})$ attaining a high value while $\mathcal{N}(\mathbf{x}_i - \tilde{\mathbf{x}}_i | \Sigma_S, \mathbf{0})$ being small. This makes $\delta(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ large while the opposite happens if the pair $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ is similar. The lemma below provides the basis of metric construction using $\delta(\cdot, \cdot)$ as suggested in Koestinger et al. (2012) and is central to our developments [presented in the following section](#).

Lemma 1. *Let $\delta : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}^+$ be the function defined in Eq. (1). The $\log(\delta)$ approximates a form of Mahalanobis metric, upto a constant term in \mathbb{R}^D . The approximated Mahalanobis metric is recognized as $\mathbf{M} = \text{Proj}(\Sigma_S^{-1} - \Sigma_D^{-1})$ with $\text{Proj}(\cdot) : \text{Sym}(D) \rightarrow S_{++}^D$ being the projection from the set of $D \times D$ symmetric matrices to the cone of Symmetric Positive Definite (SPD) matrices of the same size.*

Proof. We note that

$$\log(\delta(\mathbf{x}_i, \tilde{\mathbf{x}}_i)) = \frac{1}{2} \left(\log \det(\Sigma_S) - \log \det(\Sigma_D) + (\mathbf{x}_i - \tilde{\mathbf{x}}_i)^T (\Sigma_S^{-1} - \Sigma_D^{-1}) (\mathbf{x}_i - \tilde{\mathbf{x}}_i) \right). \quad (4)$$

The two terms $\log \det(\Sigma_S)$ and $\log \det(\Sigma_D)$ are constant and hence can be removed without loss of generality. From $(\mathbf{x}_i - \tilde{\mathbf{x}}_i)^T (\Sigma_S^{-1} - \Sigma_D^{-1}) (\mathbf{x}_i - \tilde{\mathbf{x}}_i)$ and by noting that $\Sigma_S^{-1} - \Sigma_D^{-1}$ is not necessarily an SPD matrix, we conclude that an approximated Mahalanobis metric associated to $\log(\delta)$ has the form $\mathbf{M} = \text{Proj}(\Sigma_S^{-1} - \Sigma_D^{-1})$. \square

²The discrimination value between the pair $(\mathbf{x}_i, \tilde{\mathbf{x}}_i)$ should match that of $(\tilde{\mathbf{x}}_i, \mathbf{x}_i)$.

3.2 JOINT DIMENSIONALITY REDUCTION AND METRIC LEARNING

As discussed in § 1, the metric M , while scaling well to the number of (dis)similar pairs, is sensitive to the dimensionality of the space (see for example Xiong et al. (2014)). In practice, finding the optimal dimension (or equivalently the most discriminative subspace) is done through PCA. Obviously, finding the subspace along its metric is more appealing and promises better discriminatory power.

As such, our goal is to find a lower dimensional space and its Mahalanobis metric by making use of lemma 1. Formally, we seek a linear mapping $h : \mathbb{R}^D \rightarrow \mathbb{R}^d$ and an SPD matrix M such that

$$d_M^2(h(\mathbf{x}_i), h(\tilde{\mathbf{x}}_i)) = (h(\mathbf{x}_i) - h(\tilde{\mathbf{x}}_i))^T M (h(\mathbf{x}_i) - h(\tilde{\mathbf{x}}_i)) \quad (5)$$

reflects the dissimilarity function in Eq. (1) better with $h(\mathbf{x}) = \mathbf{W}^T \mathbf{x}$, $\mathbf{W} \in \mathbb{R}^{D \times d}$. Using lemma 1 if \mathbf{W} is known, then the metric in the latent space is defined as

$$M = \text{Proj} \left((\mathbf{W}^T \Sigma_S \mathbf{W})^{-1} - (\mathbf{W}^T \Sigma_D \mathbf{W})^{-1} \right), \quad (6)$$

noting that the covariance matrix Σ in the space defined by \mathbf{W} has the form of $\mathbf{W}^T \Sigma \mathbf{W}$. While being valid, this form of the metric is obtained by keeping \mathbf{W} fixed. A more general and theoretically attractive alternative is to find \mathbf{W} along M . In doing so, we keep an eye on Koestinger et al. (2012) and rewrite Eq.(4) in the space defined by $h(\cdot)$

$$\begin{aligned} \log(\delta(h(\mathbf{x}_i), h(\tilde{\mathbf{x}}_i))) &= \frac{1}{2} \left(\log \det(\mathbf{W}^T \Sigma_S \mathbf{W}) - \log \det(\mathbf{W}^T \Sigma_D \mathbf{W}) \right. \\ &\quad \left. + (\mathbf{x}_i - \tilde{\mathbf{x}}_i)^T \mathbf{W} M \mathbf{W}^T (\mathbf{x}_i - \tilde{\mathbf{x}}_i) \right), \end{aligned} \quad (7)$$

This enables us to define the loss over all training pairs as

$$\mathcal{E}(\mathbf{W}) \triangleq \sum_{i=1}^N l_i \log(\delta(h(\mathbf{x}_i), h(\tilde{\mathbf{x}}_i))). \quad (8)$$

Minimizing Eq. (8) indeed makes the distances between similar pairs smaller while simultaneously increases the distances between dissimilar pairs. To obtain \mathbf{W} and M , we opt for an alternating optimization scheme. That is, we keep \mathbf{W} fixed to update M , followed by updating \mathbf{W} by keeping M fixed. To avoid degeneracy in the solution and inline with general practice in dimensionality reduction Weinberger & Saul (2009), we add an orthogonality constraint to \mathbf{W} , thus casting the problem of updating \mathbf{W} into the following constrained optimization problem

$$\begin{aligned} \min_{\mathbf{W}} \quad & \mathcal{E}(\mathbf{W}) \\ \text{s.t.} \quad & \mathbf{W}^T \mathbf{W} = \mathbf{I}_d \end{aligned} \quad (9)$$

To minimize (9), we make use of the recent advances in optimization over the matrix manifolds Absil et al. (2009). In particular, the constrained optimization problem in (9) can be understood as a minimization problem on the Grassmannian which we solve using Riemannian Conjugate Gradient Descent (RCGD). The Grassmannian manifold \mathcal{G}_D^d consists of the set of all linear d -dimensional subspaces of \mathbb{R}^D which lets us handle constraints of the form $\mathbf{W}^T \mathbf{W} = \mathbf{I}_d$.

Remark 1. The geometrically correct setting to minimize a problem with the orthogonality constraint is by making use of the geometry of the Stiefel manifold $St(d, D)$. However, since the loss defined in Eq. (8) is invariant to the right action of the orthogonal group \mathcal{O}_d , i.e., $\mathcal{E}(\mathbf{W}) = \mathcal{E}(\mathbf{W}\mathbf{R})$, $\mathbf{R} \in \mathcal{O}_d$, the optimization should be performed using the geometry of Grassmannian \mathcal{G}_D^d .

To perform RCGD on \mathcal{G}_D^d , we need to compute the Riemannian gradient of the loss $\mathcal{E}(\mathbf{W})$. For a smooth function $f : \mathcal{G}_D^d \rightarrow \mathbb{R}$ defined over the Grassmannian, the Riemannian gradient at \mathbf{X} denoted by $\text{grad}_{\mathbf{X}} f$ is an element of the tangent space $T_{\mathbf{X}} \mathcal{G}_D^d$ and is given by

$$\text{grad}_{\mathbf{X}} f = (\mathbf{I}_d - \mathbf{X} \mathbf{X}^T) \nabla_{\mathbf{X}} f, \quad (10)$$

where $\nabla_{\mathbf{X}} f$ is a $D \times d$ matrix of partial derivatives of f with respect to the elements of \mathbf{X} , i.e., $[\nabla_{\mathbf{X}} f]_{i,j} = \frac{\partial f}{\partial X_{i,j}}$.

Algorithm 1 The proposed JDR-KISSME algorithm

Input: $\{x_i, \tilde{x}_i, l_i\}_{i=1}^N$ a set of training pairs in \mathbb{R}^D with their similarity labels, the target dimensionality d
Output: Projection \mathbf{W} and metric \mathbf{M}
1: Compute Σ_D and Σ_S using Eq. (2) and Eq. (3)
2: Initialize \mathbf{W} to an orthonormal matrix (e.g., truncated identity)
3: Compute \mathbf{M} using Eq. (6)
4: **repeat**
5: $\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{E}(\mathbf{W})$ using RCGD on Grassmannian
6: $\mathbf{W} \leftarrow \mathbf{W}^*$
7: Update \mathbf{M} using Eq. (6)
8: **until** convergence

Below, we derive $\nabla_{\mathbf{W}} \mathcal{E}$. For a symmetric matrix Σ we have Petersen et al. (2008)

$$\nabla_{\mathbf{W}} \log \det(\mathbf{W}^T \Sigma \mathbf{W}) = 2 \Sigma \mathbf{W} (\mathbf{W}^T \Sigma \mathbf{W})^{-1}.$$

Also,

$$\nabla_{\mathbf{W}} (x_i - \tilde{x}_i)^T \mathbf{W} \mathbf{M} \mathbf{W}^T (x_i - \tilde{x}_i) = 2(x_i - \tilde{x}_i)(x_i - \tilde{x}_i)^T \mathbf{W} \mathbf{M}.$$

Therefore,

$$\nabla_{\mathbf{W}} \mathcal{E}(\mathbf{W}) = \sum_{i=1}^N l_i \left(\Sigma_S \mathbf{W} (\mathbf{W}^T \Sigma_S \mathbf{W})^{-1} - \Sigma_D \mathbf{W} (\mathbf{W}^T \Sigma_D \mathbf{W})^{-1} + (x_i - \tilde{x}_i)(x_i - \tilde{x}_i)^T \mathbf{W} \mathbf{M} \right). \quad (11)$$

Putting everything together, the algorithm to learn \mathbf{W} and \mathbf{M} is depicted in Alg. 1. In our experiments, we observed that the algorithm typically converges in less than 30 iterations.

3.3 INCORPORATING THE SOLUTION INTO DEEP NETS

In this part, we describe the proposed CNN architecture for metric learning. We denote the functionality of the CNN on an input image x by $f(x)$ below. Deep CNNs are successful techniques to directly learn a compact Euclidean space from the input image space such that distances correspond to a notion of semantics between images. Generally speaking, training a deep CNN for metric learning is cast as one of the following forms (see Song et al. (2016) and Schroff et al. (2015) for more details)

- **Pairwise:** training data consists of pairs of similar and dissimilar training images. Given a predefined margin τ , training is guided by a loss function which seeks to learn an embedding such that distances between similar images are smaller than τ while those between dissimilar images are greater than τ . In this manner, the cost function for a batch with N pairs $\{x_i, \tilde{x}_i\}$ and their corresponding similarity label $l_i \in \{1, -1\}$ can be written as

$$\sum_{i=1}^N \left[(\|f(x_i) - f(\tilde{x}_i)\|_2^2 - \tau) l_i \right]_+ \quad (12)$$

- **Triplewise:** training data consists of triplets of samples: one anchor image x , an image in the same class x^+ , and one differently labelled image x^- , and a predefined margin τ . Then, a loss function supervises training such that for each triplet, the distance between x and x^- becomes greater than the distance between x and x^+ plus τ . Thus, the cost function for a batch with N triplets is

$$\sum_{i=1}^N \left[\|f(x_i) - f(x_i^+)\|_2^2 - \|f(x_i) - f(x_i^-)\|_2^2 + \tau \right]_+ \quad (13)$$

An important difference between the two categories is the fact that only methods in the first group can work in the restricted learning scenario. We present our extension to deep networks utilizing the first protocol, making our method applicable to wider set of problems. As for the distances and similarity learning in deep networks, assume a generic CNN provides us with an embedding



Figure 2: Incorporating JDR-KISSME into deep nets. The dimensionality reduction can be seen as a fully connected layer immediately before a loss layer. One can either have the metric M in computing the loss (the model depicted in the left panel) or since $M = LL^T$, combine it with the dimensionality reduction layer (the model depicted in the right panel). Empirically, we found the first solution to be more stable.

from image space to \mathbb{R}^d . Since an SPD matrix M can be decomposed as $M = L^T L$, the distance between two images x_i and \tilde{x}_i (of a batch) passing through the CNN can be written as

$$d_M^2(f(x_i), f(\tilde{x}_i)) = (f(x_i) - f(\tilde{x}_i))^T M (f(x_i) - f(\tilde{x}_i)) = \|L(f(x_i) - f(\tilde{x}_i))\|_2^2. \quad (14)$$

Therefore, the metric M can be incorporated into a deep net in the form of a Fully Connected (FC) layer immediately before a loss layer. The whole setup can be trained via BackPropagation (BP). In our case, we start with an initial orthonormal W and compute M using Eq. (6), relying on the CNN to provide features in the low-dimensional space³. To tune W and M via BP, two possibilities are available

- Aim to learn M in a closed-form manner and tune W as developed in § 3. We refer to this solution as “pairwise+KISSME”. To be precise in pairwise+KISSME,
 - We perform Stochastic Gradient Descent (SGD) while M is kept fixed.
 - We update M after a number of SGD iterations or when the network reaches to a reasonably good representation.
- Initialize the FC layer to be WL , followed by training the network using BP. When the training is done, the metric is identified by decomposing the weights of the FC layer using the QR decomposition (or any other spectral decompositions). We refer to this solution as “pairwise+KISSME-Compact”.

Empirically, we observed that pairwise+KISSME solution is more stable and works better. We conjecture that the separation of learning W and M is the reason here. In § 7, we compare the two scenarios in more details. Before concluding this part, we would like to mention that placing a Local Response Normalization (LRN) block (see Fig. 2) before the dimensionality reduction block helps the convergence in our solution. This is inline with other deep metric learning models Liu et al. (2016); Schroff et al. (2015).

4 EXPERIMENTS

In this section, we assess and contrast the performance of our proposal against the KISSME baseline and several state-of-the-art methods. We begin by evaluating JDR-KISSME using the Comprehensive Cars (CompCars) dataset Yang et al. (2015). We then demonstrate the strength of the solution when incorporated into deep networks using the CUB200-2011 dataset Wah et al. (2011) and Cifar100 dataset Krizhevsky (2009)⁴.

4.1 CAR VERIFICATION

As our first experiment, we tackle the task of car verification using the recently released CompCars dataset (see some examples in Fig. 3). The dataset is one of the largest benchmarks for image verification containing 214,345 images of 1,687 car models from two significantly different scenarios: web nature and surveillance nature. Web nature data is split into three subsets without overlap. Related to verification is part II and part III of the dataset. Part II contains 4,454 images in 111 models (classes) while there are 22,236 images spanning 1,145 models in part III.

³Recent works in vision extensively use successful architectures such as AlexNet, GoogleNet, variants of VGGNet or ResNets to address the problem in hand Song et al. (2016); Schroff et al. (2015); Liu et al. (2016). This strategy provides better controls over the training process.

⁴The experiment on Cifar100 dataset Krizhevsky (2009) is presented in § 6 due to space limitations.



Figure 3: From left to right three sample images of the web nature and surveillance nature data of the CompCars dataset Yang et al. (2015) are shown, respectively.



Figure 4: Example images of the CUB200-2011 dataset Wah et al. (2011).

Table 1: Verification accuracy on different protocols of the CompCars dataset Yang et al. (2015).

Method	Easy	Medium	Hard
MixedDiff+CCL Liu et al. (2016) on VGG-CNN-M-1024	83.3%	78.8%	70.3%
BoxCars Sochor et al. (2016)	85.0%	82.7%	76.8%
Joint Bayesian Yang et al. (2015) on GoogLeNet	90.7%	85.2%	78.8%
KISSME Koestinger et al. (2012) on GoogLeNet	88.9%	83.3%	75.4%
JDR-KISSME(ours) on GoogLeNet	90.1%	86.0%	79.5%
KISSME Koestinger et al. (2012) on VGG-CNN-M-1024	77.6%	76.9%	69.4%
JDR-KISSME(ours) on VGG-CNN-M-1024	81.7%	79.7%	75.5%

We followed Yang et al. (2015), the standard verification protocol on this dataset, which splits part III to three sets with different levels of difficulty, namely easy, medium, and hard. Each set contains 20,000 pairs including equal number of similar and dissimilar pairs. Each image in the easy pairs is chosen from the same viewpoint, while each pair in the medium pairs is selected from a random viewpoint. Each dissimilar pair in the hard subset is selected from the same car make. As for feature extraction, we again followed the setup in Yang et al. (2015) to have a fair comparison. More specifically, we utilized the publicly available GoogLeNet Szegedy et al. (2015) pretrained on the ImageNet dataset Deng et al. (2009) and fine tuned on part II of the CompCars.

We briefly describe existing state-of-the-art methods on the web nature part of the CompCars here. Yang et al. (2015) utilize Joint Bayesian algorithm (Chen et al., 2012) which works similar in spirit to the KISSME. More specifically, the inference is based on a symmetry hypothesis comparable to Eq. (1) but with the Gaussian distributions built based on the between-class and within-class scatter matrices. Hence, it takes the class labels of the samples into its inference. Liu et al. (2016) propose Mixed Difference (MixedDiff) network by making use of VGG-CNN-M-1024 model Chatfield et al. (2014). Given triplets, a new loss function, namely Coupled Clusters Loss (CCL), is used for fine tuning the network. Lastly, Sochor et al. (2016) obtain recognition improvements by augmenting CNN inputs with complementary 3D vehicle information, such as information about 3D orientation.

In Table 1, we compare our JDR-KISSME method against the conventional KISSME algorithm and state-of-the-art deep methods. First, we note that the JDR-KISSME shows consistent improvements over the KISSME algorithm. For example, the accuracy gap between the JDR-KISSME and the KISSME method over the hard subset reaches 4.9%. Interestingly, while being shallow, the JDR-KISSME achieves the state-of-the-art verification accuracies on the medium and hard protocols while working on par with Yang et al. (2015) on the easy test.

We note that the work of Yang et al. (2015) makes use of class labels for training (and hence not applicable to the restricted metric learning scenarios) while our method does not rely on the availability of such information. As long as the other deep solutions are considered, the MixedDiff+CCL Liu et al. (2016) uses the VGG-CNN-M-1024 towards verification. As such, we also report the verification accuracies of JDR-KISSME and KISSME by using VGG-CNN-M-1024 to extract features instead of GoogLeNet. These results are also shown in Table 1. Again, we observe that our JDR-KISSME method is superior to the conventional KISSME algorithm when the VGG-CNN-M-1024

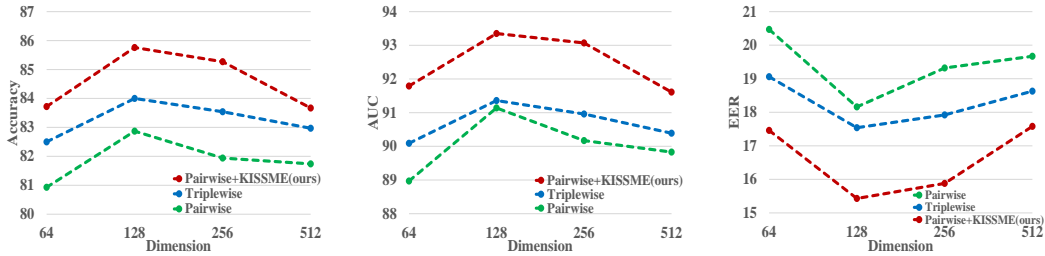


Figure 5: Verification accuracy, AUC, and EER score metrics on the CUB200-2011 dataset Wah et al. (2011).

is considered. Compared to the MixedDiff+CCL, the JDR-KISSME achieves better verification accuracy values on the medium and hard subsets.

4.2 BIRD VERIFICATION

In this part, we evaluate the incorporation of our JDR-KISSME solution into CNNs using CUB200-2011 dataset Wah et al. (2011). CUB200-2011 has 200 classes of birds with 11,788 images (some examples are shown in Fig. 4). We used images of the first 100 classes as training and validation sets and the remaining classes for testing. To measure the performances, we randomly generated 30,000 similar and 30,000 dissimilar pairs from our test data. We report verification accuracy, Area Under the ROC Curve (AUC), and the Equal Error Rate (EER)⁵ values for our algorithm and of the baselines for comparisons. We used Matconvnet package Vedaldi & Lenc (2015) to implement our method. As the CNN, we used the VGG-CNN-M-1024 pretrained on the ImageNet. More details about this experiment is provided in § 8.

To show the effectiveness of pairwise+KISSME, our deep metric learning method, we perform comparisons with the two most common ways of training a CNN for metric learning, i.e., pairwise metric learning (i.e., Eq. (12)) and triplet solution (i.e., Eq. (13)).

We recall from § 3.3 that the training starts by initializing two matrices, \mathbf{W} (or equivalently the last FC layer) and \mathbf{M} . To initialize \mathbf{W} , we rely on the initial CNN to provide embeddings of training images up to the last FC layer. Then, the FC layer is initialized with PCA (of a certain dimensionality). This is consistent with the original KISSME. Next, the metric \mathbf{M} is initialized using the output of the FC layer (Eq. (6)). This is required in the loss layer to perform BP.

Fig. 5 summarizes the three score metrics of our deep metric learning technique and of our two baselines for various embedding sizes (or equivalently subspace dimensionality). Similar to the previous experiment, our solution is consistently superior to the baseline techniques for all embedding sizes. For example, the difference in the verification accuracy between our method and its closest competitor, i.e., triplet solution, is about 2% for the size 128.

Before concluding this section, we note that our aim here was to present a better, yet general way of metric learning for deep networks. In doing so, we were not chiefly concerned about best mining practices as suggested in Song et al. (2016); Schroff et al. (2015).

5 CONCLUSIONS

In this paper, we introduced a joint dimensionality reduction technique for the KISSME algorithm, namely JDR-KISSME. Our motivation stems from the fact that the KISSME fails badly when its input is not meticulously denoised using PCA. Along the way, we formulated the solution as a Riemannian optimization problem. Furthermore, based on our proposal, we showed an end-to-end learning of a generic deep network for metric learning. Our experiments demonstrate consistent improvements of the JDR-KISSME and its deep extension over the original KISSME and state-of-the-art methods. In the future, we plan to benefit from more advanced mining techniques (e.g., multi-class N-pair loss Sohn (2016)) to further improve the accuracy of our proposed method.

⁵EER is the rate at which both acceptance and rejection errors are equal (the lower is the better).

REFERENCES

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- Jane Bromley, James W Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. Signature verification using a siamese time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(04):669–688, 1993.
- K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *Proc. British Machine Vision Conference (BMVC)*, 2014.
- Dong Chen, Xudong Cao, Liwei Wang, Fang Wen, and Jian Sun. Bayesian face revisited: A joint formulation. In *Proc. European Conference on Computer Vision (ECCV)*, pp. 566–579, 2012.
- Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pp. 539–546, 2005.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255, 2009.
- Alan Edelman, Tomás A Arias, and Steven T Smith. The geometry of algorithms with orthogonality constraints. 20(2):303–353, 1998.
- Matthieu Guillaumin, Jakob Verbeek, and Cordelia Schmid. [Is that you? Metric learning approaches for face identification](#). In *Proc. Int. Conference on Computer Vision (ICCV)*, pp. 498–505. IEEE, 2009.
- Matthieu Guillaumin, Thomas Mensink, Jakob Verbeek, and Cordelia Schmid. [Face recognition from caption-based supervision](#). *Int. Journal of Computer Vision (IJCV)*, 96(1):64–82, 2012.
- Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International Workshop on Similarity-Based Pattern Recognition*, pp. 84–92. Springer, 2015.
- Junlin Hu, Jiwen Lu, and Yap-Peng Tan. Discriminative deep metric learning for face verification in the wild. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1875–1882, 2014.
- Martin Koestinger, Martin Hirzer, Paul Wohlhart, Peter M Roth, and Horst Bischof. Large scale metric learning from equivalence constraints. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2288–2295, 2012.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 1097–1105, 2012.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Hongye Liu, Yonghong Tian, Yaowei Yang, Lu Pang, and Tiejun Huang. Deep relative distance learning: Tell the difference between similar vehicles. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2167–2175, 2016.
- David G Lowe. Distinctive image features from scale-invariant keypoints. *Int. Journal of Computer Vision (IJCV)*, 60(2):91–110, 2004.
- Alexis Mignon and Frédéric Jurie. Pcca: A new approach for distance learning from sparse pairwise constraints. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2666–2672, 2012.

- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7:15, 2008.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 815–823, 2015.
- Jakub Sochor, Adam Herout, and Jiri Havel. Boxcars: 3d boxes as cnn input for improved fine-grained vehicle recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3006–3015, 2016.
- Kihyuk Sohn. [Improved Deep Metric Learning with Multi-class N-pair Loss Objective](#). In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 1849–1857, 2016.
- Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 1988–1996, 2014.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1–9, 2015.
- A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/>, 2008.
- Andrea Vedaldi and Karel Lenc. Matconvnet: Convolutional neural networks for matlab. In *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 689–692. ACM, 2015.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- Kilian Q Weinberger and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research (JMLR)*, 10:207–244, 2009.
- Eric P Xing, Andrew Y Ng, Michael I Jordan, and Stuart Russell. Distance metric learning with application to clustering with side-information. *Proc. Advances in Neural Information Processing Systems (NIPS)*, 15:505–512, 2003.
- Fei Xiong, Mengran Gou, Octavia Camps, and Mario Sznajder. Person re-identification using kernel-based metric learning methods. In *Proc. European Conference on Computer Vision (ECCV)*, pp. 1–16. 2014.
- Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine-grained categorization and verification. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3973–3981, 2015.

6 EXPERIMENT ON THE CIFAR100

As another experiment for deep metric learning, we studied the task of image verification using the Cifar100 dataset Krizhevsky (2009) which has 60,000 images of size 32×32 . To this end, we trained the LeNet-5 LeCun et al. (1998) network on the Cifar10 dataset Krizhevsky (2009) for 22,500 iterations of SGD. We then cropped the pretrained network at the fourth layer and fine tuned it on the Cifar100 similar to the experiment on the CUB200-2011 dataset. We kept the embedding size to 32 for this experiment (i.e., \mathbf{W} was 64×32 and \mathbf{M} was 32×32). All other experimental details (e.g., train/test split, number of test pairs, values of the parameters, etc) were the same as the bird verification experiment.

In Table 2, we compare our method against the so called pairwise and triplet methods as well as the original KISSME on the pretrained network (without fine tuning). Here again, our solution comfortably outperforms the other methods for all the studied metrics over the 60,000 test pairs.

Table 2: Verification accuracy, AUC, and EER score metrics on the Cifar100 dataset Krizhevsky (2009).

Method	Accuracy	AUC	EER
KISSME Koestinger et al. (2012)	63.1%	69.2%	36.4%
Pairwise	64.3%	70.0%	35.4%
Triplewise	65.8%	72.0%	34.2%
Pairwise+KISSME(ours)	68.2%	75.2%	31.7%

7 FURTHER ANALYSIS

In this part, we first contrast the pairwise+KISSME method against the pairwise+KISSME-Compact solution. We recall from § 3.3 that there are two ways for training a network based on the KISSME objective

1. **Pairwise+KISSME:** Initialize the weights of the last FC layer to be \mathbf{W} and engage the metric \mathbf{M} directly in computing the distances in the loss layer. To this end, \mathbf{M} is kept fixed during a number of SGD runs and is updated afterwards using the output of the CNN. In this case, BP updates the CNN according to the KISSME loss (i.e., Eq. (8)) while \mathbf{M} is learned in a closed form manner (see the left panel in Fig. 2 for a conceptual diagram).
2. **Pairwise+KISSME-Compact:** Since \mathbf{M} is an SPD matrix (i.e., $\mathbf{M} = \mathbf{L}^T \mathbf{L}$), it can be absorbed in the last FC layer. Here, we initialize the weights of the last FC layer to be $\mathbf{W} \mathbf{L}$. Then, we train the network using BP. If the explicit form of the metric is required, the weights of the FC layer can be factorized into an orthogonal matrix \mathbf{W} and a full-rank matrix \mathbf{L} using any spectral decomposition such as QR decomposition (see the right panel in Fig. 2 for a conceptual diagram).

To empirically compare the two solutions, we conducted an experiment on the Birds verification (§ 4.2). The accuracy, AUC and EER values for various embedding sizes are depicted for the pairwise+KISSME and pairwise+KISSME-Compact solutions in Fig. 6. From Fig. 6, we conclude that the pairwise+KISSME solution leads to superior performances and is more stable, hence our proposal. This is a consistent extension to the JDR-KISSME, our developments in the shallow mode, where we have an alternating algorithm to find the two matrices \mathbf{W} and \mathbf{M} .

We conjecture that the separation of learning the subspace \mathbf{W} and the metric \mathbf{M} is the reason here. An intuitive explanation is as follows. Assume the ideal projection is \mathbf{W}^* . From Lemma 1, we conclude that the ideal metric is obtained as $\mathbf{M}^* = \Gamma(\mathbf{W})$ where the function $\Gamma(\cdot)$ is a nonlinear function as a result of the projection to the positive definite cone. In the deep extension of our work, the pairwise+KISSME is more aligned with Lemma 1 as the metric is explicitly obtained from the representation. On the other hand, the pairwise+KISSME-Compact removes the dependency of \mathbf{M}^* on \mathbf{W}^* in the hope of learning \mathbf{M}^* , \mathbf{W}^* and the underlying nonlinear projection together.

To have a complete picture, we report the computational burden of our pairwise+KISSME and two baselines, namely the pairwise and triplewise metric learning using the CUB200-2011 dataset (§ 4.2). In average for an embedding of size 128 (i.e., $\mathbf{M} \in \mathcal{S}_{++}^{128}$), performing 10,000 SGD iterations takes 41,050 and 40,803 seconds for the pairwise+KISSME and pairwise, respectively, using a

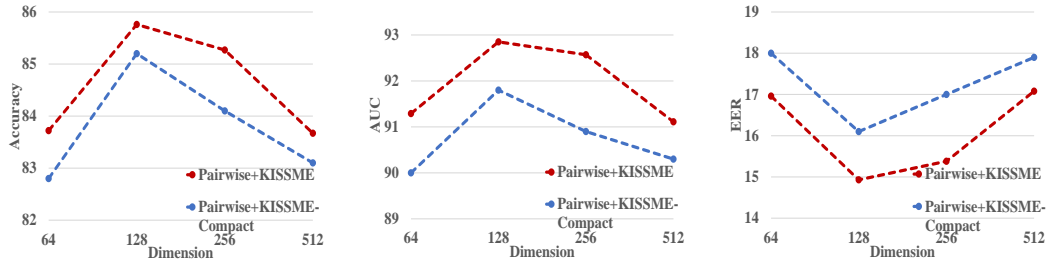


Figure 6: Comparison between the pairwise+KISSME and pairwise+KISSME-Compact deep solutions on the CUB200-2011 dataset Wah et al. (2011). See text for more information.

moderate NVIDIA Quadro M4000 GPU. The slight difference is because computing the metric M is the only extra step required for training in the pairwise+KISSME. For the case of triplewise this time is 42,345 seconds.

8 EXPERIMENTAL SETUP

For fine tuning the CNNs, we randomly generated 300,000 similar pairs and 300,000 dissimilar pairs (and equal number of triplets) and fed them to the CNNs. For all experiments, we set maximum number of SGD iterations to 20,000, margin to $\tau = 1.0$, momentum to $\mu = 0.9$, and learning rate to $\eta = [10^{-4}, 10^{-7}]$ in log-space range. We observed that increasing the learning rate of the fully connected layer by a factor of 10 helps faster convergence. A similar observation is reported in Song et al. (2016). To augment the data, we resorted to only flipping the images at random.

To do complete justice, we also provide details of the experiment used to generate Fig. 1. We used the surveillance images of the CompCars dataset (see Fig. 3 for some examples). There exist 44,481 images in 281 classes (car models). We randomly split the dataset into 140 classes for training and used the remaining 141 classes for testing. This was to ensure that there is no overlap between the training and testing images. We generated 200,000 training pairs and 60,000 testing pairs, randomly from the training and testing sets. To extract image descriptors, we computed SIFT features Lowe (2004) on a dense grid and then computed Bag Of Word representations using a dictionary of size 4096, trained by the k-means algorithm Vedaldi & Fulkerson (2008).