

Learning Multi-Modal Nonlinear Embeddings: Performance Bounds and an Algorithm

Semih Kaya and Elif Vural

Abstract—While many approaches exist in the literature to learn low-dimensional representations for data collections in multiple modalities, the generalizability of multi-modal nonlinear embeddings to previously unseen data is a rather overlooked subject. In this work, we first present a theoretical analysis of learning multi-modal nonlinear embeddings in a supervised setting. Our performance bounds indicate that for successful generalization in multi-modal classification and retrieval problems, the regularity of the interpolation functions extending the embedding to the whole data space is as important as the between-class separation and cross-modal alignment criteria. We then propose a multi-modal nonlinear representation learning algorithm that is motivated by these theoretical findings, where the embeddings of the training samples are optimized jointly with the Lipschitz regularity of the interpolators. Experimental comparison to recent multi-modal and single-modal learning algorithms suggests that the proposed method yields promising performance in multi-modal image classification and cross-modal image-text retrieval applications.

Index Terms—Multi-modal learning, multi-view learning, cross-modal retrieval, nonlinear embeddings, supervised embeddings, RBF interpolators.

I. INTRODUCTION

MANY data analysis applications involve the acquirement or analysis of data collections in multiple modalities. In some problems, the purpose is to fuse the information in different modalities to improve the detection or classification accuracy, while some other applications require the retrieval of data samples in a certain modality that are relevant to a query sample provided in another modality. For instance, in an image-text cross-modal retrieval problem, one might be interested in retrieving image samples from the same category as a query text sample. In this paper, we study the problem of learning supervised nonlinear representations for multi-modal classification and cross-modal retrieval applications.

Multi-modal learning algorithms often aim to compute joint representations in a common domain, where the main challenge is to efficiently align different modalities without damaging their inherent geometry. Subspace learning methods such as CCA [1] align different modalities via linear projections or transformations. Supervised linear embedding methods such as GMLDA [2] and its various extensions aim to enhance the separation between different data classes in addition to the alignment of different modalities. However, when different modalities have significantly dissimilar geometric structures, linear methods may fall short of learning effective joint

representations since they are mostly restricted by the original geometry of the individual modalities. Kernel extensions of linear methods such as Kernel CCA [1], Kernel GMLDA [2] and its variants provide nonlinear representations that may improve some of these shortcomings; however, the resulting algorithms might still lack in flexibility in certain problems. In particular, the suitability of the selected kernel type might vary largely depending on the data set and the embedding may generalize poorly to test data. In the recent years, impressive performance has been attained in retrieval and classification problems with deep learning algorithms based on cross-modal CNNs and autoencoders [3], [4], [5]. While these methods compute powerful nonlinear representations, they typically require much larger data sets and their training complexity is significantly higher.

While different multi-modal learning approaches might be preferable to each other depending on the setting, their capacity to generalize to novel test samples is a questionable issue in general. A multi-modal learning method may yield promising performance figures on training data, while its performance may be much lower on previously unseen test data, especially if it involves complex and rich models. In fact, the theoretical characterization of the generalization capability of multi-modal embedding algorithms is a somewhat overlooked problem in the literature. Some previous studies have focused on the analysis of co-training [6], [7] or co-regularized RKHS problems [8], [9], which however do not tell what geometric properties a nonlinear multi-modal embedding should have for successful generalization.

In this paper, we consider the problem of learning supervised nonlinear embeddings for multi-modal classification and cross-modal retrieval applications that can generalize well to new test data. Our main purpose in preferring nonlinear embeddings as opposed to subspace methods is to achieve a relatively high model capacity that can adapt to challenging data geometries. On the other hand, we adhere to a shallow data representation model with a single-stage embedding as opposed to deep methods, in order to achieve applicability to settings with restricted availability of training data or limited computation budget. Our study has two main contributions. We first propose a theoretical analysis of learning supervised multi-modal embeddings. We consider a nonlinear embedding model where the training samples from different modalities are jointly mapped to a common lower-dimensional domain, and the training embeddings are extended to the whole data space via Lipschitz-continuous interpolation functions. Our theoretical bounds suggest that for good generalization performance, the multi-modal embedding of training samples should satisfy

S. Kaya and E. Vural are with the Department of Electrical and Electronics Engineering, METU, Ankara.
E-mail: kaya.semih@metu.edu.tr, velif@metu.edu.tr

three conditions: (1) Different modalities should be aligned sufficiently well; (2) Different classes should be sufficiently well-separated from each other; (3) The geometric structure of each modality (captured through nearest neighborhoods) should be preserved. Then, under these conditions, we show that the embedding generalizes well to test data, provided that the Lipschitz constants of the interpolation functions are sufficiently low. This points to an important trade-off: Multi-modal methods may fail to generalize to test data when the nonlinear interpolation functions are too irregular, even if the embeddings of training samples exhibit good cross-modal alignment and between-class separation properties.

Our next contribution is to propose a new supervised nonlinear multi-modal learning algorithm. Motivated by the above theoretical findings, we formulate an optimization problem where a cross-modal alignment term and a between-class separation term are jointly optimized with the Lipschitz constants of the interpolation functions generalizing the embeddings. The resulting objective function is minimized iteratively, by jointly learning the nonlinear embedding coordinates with the interpolator parameters. Our method has the advantage of providing more flexible representations than subspace methods thanks to the employed nonlinear models, while it entails a relatively lightweight training phase compared to elaborate approaches such as deep learning methods. The proposed method is suitable for multi-modal problems with significantly different data types in different modalities, as well as multi-view problems with closely related or same data types across different views. We test the proposed algorithm in multi-view image classification and image-text cross-modal retrieval applications. Experimental results show that the proposed method yields quite satisfactory performance in comparison with recent multi-modal learning approaches.

The rest of the paper is organized as follows. In Section II, we overview the related literature. In Section III, we present a theoretical analysis of the multi-modal representation learning problem. In Section IV, we describe our supervised nonlinear multi-modal representation learning algorithm. In Section V, we experimentally evaluate the performance of the proposed method, and in Section VI, we conclude.

II. RELATED WORK

The multi-modal learning approaches in the literature can be mainly grouped as co-training methods, subspace learning approaches, kernel methods and deep learning methods. Co-training methods learn separate models in different modalities by encouraging their predictions to be similar [6]. A probabilistic model for Support Vector Machine (SVM) is constructed in [10] based on the Co-EM approach. There also exist co-regression algorithms employing the co-training idea [11]. The co-training technique is also used in graph-based methods such as [12], where a Gaussian process model is used on an undirected Bayesian graph. Co-training algorithms have been used in various data analysis applications [13], [14].

Subspace learning methods are based on computing linear projections or transformations that suitably align samples from different modalities. The well-known unsupervised subspace

learning algorithm CCA (Canonical Correlation Analysis) maximizes the correlation between different modalities [1]. Alternative versions of CCA such as cluster CCA [15], multi-label CCA [16] and three-view CCA [17] have been proposed to improve the performance of CCA in various supervised tasks, all of which employ linear projections. In the recent years, many supervised subspace methods have been proposed, which aim to enhance the between-class separation and cross-modal alignment when learning linear projections of data. The GMLDA (Generalized Multiview Analysis) method proposes a multi-modal extension of the LDA algorithm within this framework [2]. Projection directions for different modalities are learnt by optimizing a quadratic objective function that contains within-class scatter, between-class scatter, and cross-modal correlation terms [2]. A kernel extension of the GMLDA algorithm for learning nonlinear mappings is also presented in [2]. Following the approach in [2], many extensions of this work have been proposed in succeeding studies. In [18], a view consistency term is added to the objective function of GMLDA so as to impose the similarity of the linear projection functions of different views. The cross-media retrieval method in [19] addresses a supervised linear projection learning problem as in [2]; however, applies regularization on projection matrices. Several works have focused on kernel extensions of the problem [20], [21]. The KMvMDA (Kernel Multi-View Modular Discriminant Analysis) method proposed in [21] learns kernel representations by imposing within-class and between-class correlation constraints across different modalities. The kernel method in [22] uncorrelates the feature vectors of the individual modalities as an additional consideration. Another body of methods model data by constructing label-aware data graphs and include a graph-based regularization term in the objective in order to preserve the geometry of the data set [23], [24]. The JFSSL (Joint Feature Selection and Subspace Learning) method uses a joint graphical model for calculating projections with relevant and irrelevant features [25]. Various methods are based on combining kernels in different modalities, such as convex combinations of multiple Laplacian kernels [26], or the power mean of multilayer graph kernels [27]. Among all these methods, our method bears similarities to especially supervised kernel methods in that it learns nonlinear data representations in view of class separation and cross-modality alignment objectives. On the other hand, it has two major differences from these approaches: (1) Our nonlinear representation model is particularly flexible and effective as the pointwise embeddings of training samples are optimized individually by respecting the data geometry. (2) We explicitly incorporate the generalization performance of the algorithm in the objective function, which is a unique and distinctive feature of our method. While there are multi-view algorithms learning pointwise nonlinear mappings as in our method, these often address unsupervised problems such as spectral embedding [28] and multi-modal clustering [29].

In the recent years, deep learning methods have provided quite effective solutions for analyzing large multi-modal data sets. Deep multi-view autoencoders can learn shared representations [3], [5], or cross-weights [30] across different data modalities. Convolutional neural networks are widely used

in multi-modal problems as well, where CNN structures for visual modalities can be combined with other modalities at the feature level [4], [31], or the classifier level [32]. GAN-type architectures adversarially train feature generators and domain discriminators across different modalities [33]. The method in [34] learns a common latent representation for different modalities via a deep matrix factorization scheme.

Some previous studies proposing a theoretical analysis of multi-view learning are the following. The study in [6] analyses the learnability of joint models in the two-view co-training problem, assuming the conditional independence of two views. Several other PAC-style bounds are proposed in [7], [35], mainly stating that the agreement of the classifiers of the two views on training data guarantees a good estimate of the expected test error. Several studies have proposed generalization bounds for co-regularized RKHS methods [8], [9], [36] in terms of the Rademacher complexities of the involved function classes. These previous analyses differ from ours in that they all aim to bound the difference between the training loss and the expected loss in multi-view classification, while our analysis addresses the particular problem of nonlinear dimensionality reduction in multi-modal learning. Our distinctive contribution is that we explicitly characterize the geometric properties and the regularity conditions of the nonlinear embedding to achieve successful generalization.

Finally, some other previous works related to our study are the following. The theoretical analysis in [37] provides performance bounds for supervised nonlinear embeddings in a single modality. The idea in [37] is developed in this paper to perform a theoretical analysis for multi-modal embeddings. The previous work [38] proposes a supervised nonlinear dimensionality reduction algorithm via smooth representations like in our work; however, it treats the embedding problem in a single modality. Lastly, a preliminary version of our work was presented in [39]. The current paper builds on [39] by including a theoretical analysis of the multi-modal learning problem and significantly extending the experimental results.

III. PERFORMANCE BOUNDS FOR MULTI-MODAL LEARNING WITH SUPERVISED EMBEDDINGS

In this section, we first describe the multi-modal representation learning setting considered in this study and then present a theoretical analysis of multi-modal classification and retrieval with supervised embeddings.

A. Notation and Setting

We consider a setting with M data classes and V modalities (also called *views*) such that a data sample x has an observation $x^{(v)}$ in each modality (or view) $v = 1, \dots, V$. Let the data samples from each class $m = 1, \dots, M$ in each modality $v = 1, \dots, V$ be drawn from a probability measure $\nu_m^{(v)}$ on a Hilbert space $H^{(v)}$. We assume that the probability measure $\nu_m^{(v)}$ has a bounded support $\mathcal{M}_m^{(v)} \subset H^{(v)}$ for each v , and that the probability measures $\{\nu_m^{(v)}\}$ in different modalities v are independent for each class m .

Let $\mathcal{X} = \{x_i\}$ be a set of training samples such that each i -th training sample x_i belongs to one of the classes

$m = 1, \dots, M$. In each modality v , the observations of the training samples $\{x_i^{(v)}\}$ from each class m are independent and identically distributed, drawn from the probability measure $\nu_m^{(v)}$. In this paper, we study a setting where the training samples from all modalities are embedded as $\mathcal{Y} = \{y_i^{(v)}\}$ into a common Euclidean domain \mathbb{R}^d , such that each training sample $x_i^{(v)} \in H^{(v)}$ from modality v is mapped to a vector $y_i^{(v)} \in \mathbb{R}^d$. Although we do not impose any conditions on the dimension d of the embedding, d is typically small in many methods.

Focusing mainly on a scenario where the embedding is nonlinear in this work, we assume that the embedding of the training samples is extended to the whole data space through interpolation functions $f^{(v)} : H^{(v)} \rightarrow \mathbb{R}^d$, for $v = 1, \dots, V$, such that each training sample in a modality v is mapped to its embedding as $f^{(v)}(x_i^{(v)}) = y_i^{(v)}$. We characterize the regularity of the interpolation functions $f^{(v)}$ with their Lipschitz continuity, which is defined as follows.

Definition 1. A function $f : H \rightarrow \mathbb{R}^d$ defined on a Hilbert space H is Lipschitz continuous with constant $L > 0$ if for any $x_1, x_2 \in H$, the function satisfies $\|f(x_1) - f(x_2)\| \leq L \|x_1 - x_2\|$.

The notation $\|\cdot\|$ will denote the usual norm in the space of interest (e.g. L^2 -norm, or ℓ^2 -norm), unless stated otherwise. Now, for each modality v , let $B_\delta(x^{(v)}) \subset H^{(v)}$ be an open ball of radius δ around the point $x^{(v)}$

$$B_\delta(x^{(v)}) = \{z^{(v)} \in H^{(v)} : \|x^{(v)} - z^{(v)}\| < \delta\}.$$

Then, for each class m , we define a parameter $\eta_{m,\delta}$, which is a lower bound on the measure of the open ball $B_\delta(x^{(v)})$ around any point from class m in any modality

$$\eta_{m,\delta} := \min_{v=1,\dots,V} \inf_{x^{(v)} \in \mathcal{M}_m^{(v)}} \nu_m^{(v)}(B_\delta(x^{(v)})).$$

In the following, $C(\cdot)$ denotes the class label of a sample, $|\cdot|$ refers to the cardinality of a set, the notation $z \sim \nu$ means that the sample z is drawn from the distribution ν , $P(\cdot)$ denotes the probability of an event, and $\|\cdot\|_F$ denotes the Frobenius norm. The notation $\text{tr}(\cdot)$ stands for the trace of a matrix, and $(\cdot)_{ij}$ indicates the entry of a matrix in the i -th row and the j -th column.

B. Theoretical Analysis of Classification and Retrieval Performance

We now present performance bounds for the multi-modal classification problem and the cross-modal retrieval problem.

1) *Multi-Modal Classification Performance:* Let x be a test sample with an observation $x^{(v)}$ available in a specific modality v . Denoting the true class of x by m , we assume that the observation $x^{(v)}$ of the test sample is drawn from the probability measure $\nu_m^{(v)}$ independently of the training samples.

We consider a classification setting where the class label of $x^{(v)}$ is estimated by first embedding $x^{(v)}$ into \mathbb{R}^d as $f^{(v)}(x^{(v)})$ through the interpolator $f^{(v)}$ learnt using the training samples.

Then the estimate $\hat{C}(x)$ of the class label $C(x)$ of x is found via nearest-neighbor classification in \mathbb{R}^d over the embeddings $y_i^{(u)}$ of the training samples $x_i^{(u)}$ from all modalities $u = 1, \dots, V$. Hence, the class label of the test sample x is estimated as $\hat{C}(x) = C(x_{i^*})$, where ¹

$$i^* = \arg \min_i \min_{u=1, \dots, V} \|y_i^{(u)} - f^{(v)}(x^{(v)})\|. \quad (1)$$

In the following theorem, we present our main result for multi-modal classification with supervised embeddings.

Theorem 1. *Let the training sample set \mathcal{X} contain at least N_m training samples $\{x_i\}_{i=1}^{N_m}$ from class m , whose observations $\{x_i^{(u)}\}$ with $x_i^{(u)} \sim \nu_m^{(u)}$ are available in all modalities $u = 1, \dots, V$. Let \mathcal{Y} be an embedding of \mathcal{X} in \mathbb{R}^d with the following properties*

- (P1) $\|y_i^{(v)} - y_i^{(u)}\| \leq \eta$ for all training samples x_i and for all $v, u \in \{1, \dots, V\}$
- (P2) $\|y_i^{(u)} - y_j^{(u)}\| \leq R_\delta$ for all $u \in \{1, \dots, V\}$,
if $\|x_i^{(u)} - x_j^{(u)}\| \leq 2\delta$ and $C(x_i) = C(x_j)$
- (P3) $\|y_i^{(v)} - y_j^{(v)}\| > \gamma$ for all $v, u, \in \{1, \dots, V\}$
if $C(x_i) \neq C(x_j)$

where η and γ are some constants and R_δ is a δ -dependent constant. Assume that the interpolation function $f^{(u)} : H^{(u)} \rightarrow \mathbb{R}^d$ in each modality u is a Lipschitz continuous function with constant L such that for some parameters $\epsilon > 0$ and $\delta > 0$, the following inequality is satisfied

$$6L\delta + 2\sqrt{d}\epsilon + 2R_\delta + 2\eta \leq \gamma. \quad (2)$$

Then for some $Q \geq 1$, if the number of training samples is such that

$$N_m > \frac{Q}{\eta_{m,\delta}}, \quad (3)$$

the probability of correctly classifying a test sample x from class m observed as $x^{(v)}$ in modality v via the nearest neighbor classification rule in (1) is lower bounded as

$$P(\hat{C}(x) = m) \geq 1 - \left[\exp\left(-\frac{2(N_m \eta_{m,\delta} - Q)^2}{N_m}\right) + 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) + (1 - \eta_{m,\delta})^Q \right]^V. \quad (4)$$

The proof of Theorem 1 is given in the Appendix. The theorem intuitively states the following: First, (P1), (P2), and (P3) define the properties that the embedding should have, which are illustrated in Figure 1. (P1) requires the observations $x_i^{(v)}$, $x_i^{(u)}$ of the same training sample x_i in two different modalities to be mapped to nearby points in the common domain \mathbb{R}^d of embedding, so that the distance between their embeddings does not exceed some threshold $\eta > 0$. This property imposes that different modalities be well aligned through the learnt embedding. The property (P2) indicates that two nearby samples from the same modality and the same class

¹We adopt the notation $C(x)$ instead of $C(x^{(v)})$ for class labels as the observation $x^{(v)}$ of a sample x in any modality v has the same class label.

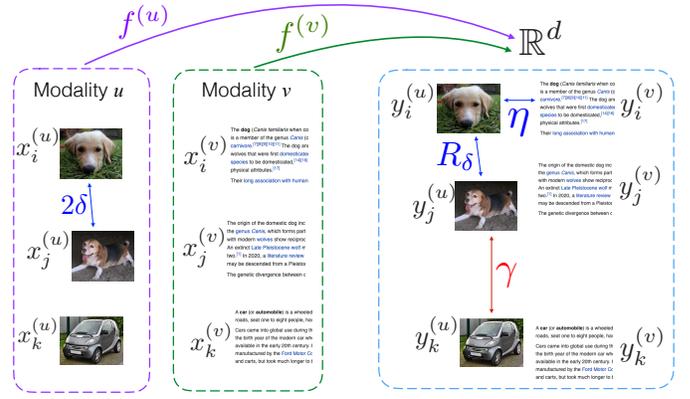


Fig. 1. Illustration of the studied multi-modal embedding setting. Modalities u (image) and v (text) are mapped to the common domain \mathbb{R}^d via interpolators $f^{(u)}$ and $f^{(v)}$. The parameters η , R_δ , and γ respectively measure the alignment between different modalities, the within-class compactness, and the separation between different classes. (Images: wikipedia.org)

should be mapped to nearby points, so that a distance of 2δ in the original domain is mapped to a distance of at most R_δ in the domain of embedding, where R_δ is a constant depending on δ . This can be seen as a condition for the preservation of the local geometry of each modality within the same class. Lastly, the property (P3) imposes samples from different classes to be separated by a distance of at least γ in the domain of embedding, regardless of their modality. Here, the parameter $\gamma > 0$ can be seen as a separation margin between different classes in the learnt embedding.

If the embedding of the training samples has these properties, supposing that the condition in (2) is satisfied, Theorem 1 guarantees that the probability of correctly classifying a test sample from some class m approaches 1 at an exponential rate as the number of training samples N_m from that class increases. This can be verified by observing that N_m should be chosen proportionally to the parameter Q as seen in (3), in which case the correct classification probability in (4) improves at rate $1 - e^{-O(VN_m)}$. Here, an important observation is that as the number of modalities V increases, the correct classification probability improves at an exponential rate. This confirms that the multi-modal learning algorithm can successfully fuse the information obtained from different modalities for improving the classification performance.

Finally, a crucial implication of Theorem 1 is that the condition in (2) must be satisfied in order to achieve high classification accuracy. The condition (2) is quite central to our study and it will be of importance when proposing an algorithm in Section IV. It states that a certain compromise must be sought between the Lipschitz regularity of the interpolator and the separation between different classes: When learning nonlinear embeddings, the separation γ between training samples from different classes should be adjusted in a way to allow the existence of a sufficiently regular interpolator, so that L remains sufficiently small. While an embedding with a too small γ value would fail to satisfy the condition (2), increasing γ too much would result in a highly irregular warping of the training samples,

which typically leads to an increase in the magnitude of the interpolator parameters. This results in an interpolator with poor Lipschitz regularity with a large L value where the condition (2) would fail again. Hence, the condition (2) points to how the separation margin and the interpolator regularity should be jointly taken into account when learning an embedding with good generalization properties.

2) *Cross-Modal Retrieval Performance*: Next, we analyze the performance of cross-modal retrieval via supervised embeddings. Given the multi-modal data set $\mathcal{X} = \{x_i\}$, where each data sample x_i belongs to one of the classes $m = 1, \dots, M$, we formally define the retrieval problem as follows. Let $x^{(v)}$ be a query test sample observed in modality v . We study a cross-modal retrieval setting where the purpose is to retrieve samples from a certain modality u that are ‘‘relevant’’ to the query sample $x^{(v)}$ from modality v . We consider two samples to be relevant if they belong to the same class.

Denoting the modality of the query sample by v and the modality of the retrieved samples by u , we consider a retrieval strategy that returns the most relevant K samples to the query sample, based on the distance of the samples in the domain of embedding. Hence, given the query sample $x^{(v)}$, it is first embedded into \mathbb{R}^d as $f^{(v)}(x^{(v)})$ via the interpolator $f^{(v)}$; and then the K training samples $\{x_i^{(u)}\}$ from modality u whose embeddings $\{f^{(u)}(x_i^{(u)})\}$ have the smallest distance to $f^{(v)}(x^{(v)})$ are retrieved as the most relevant samples, thus returning the set $\{x_{i_k}^{(u)}\}_{k=1}^K$, where

$$\begin{aligned} i_1 &= \arg \min_i \|f^{(u)}(x_i^{(u)}) - f^{(v)}(x^{(v)})\| \\ i_k &= \arg \min_{i \notin \{i_1, \dots, i_{k-1}\}} \|f^{(u)}(x_i^{(u)}) - f^{(v)}(x^{(v)})\|, k = 2, \dots, K. \end{aligned} \quad (5)$$

The precision rate P and the recall rate R of the retrieval algorithm are then given by

$$P = \frac{TP}{TP + FP}, \quad R = \frac{TP}{TP + FN} \quad (6)$$

where TP , FP , and FN respectively denote the number of true positive, false positive, and false negative samples depending on whether the retrieved and unretrieved samples are relevant or not.

We present the following main result regarding the performance of cross-modal retrieval with supervised embeddings.

Theorem 2. *Let the training sample set \mathcal{X} contain N_m training samples $\{x_i\}_{i=1}^{N_m}$ from class m , with observations $\{x_i^{(v)}\}$ and $\{x_i^{(u)}\}$ available in the modalities v and u . Let \mathcal{Y} be an embedding of \mathcal{X} in \mathbb{R}^d with the following properties:*

- (P1) $\|y_i^{(v)} - y_i^{(u)}\| \leq \eta$ for all training samples x_i
- (P2) For two samples x_i and x_j with $C(x_i) = C(x_j)$

$$\|y_i^{(v)} - y_j^{(v)}\| \leq R_\delta \text{ if } \|x_i^{(v)} - x_j^{(v)}\| \leq 2\delta;$$

$$\|y_i^{(u)} - y_j^{(u)}\| \leq R_\delta \text{ if } \|x_i^{(u)} - x_j^{(u)}\| \leq 2\delta$$
- (P3) $\|y_i^{(v)} - y_j^{(u)}\| > \gamma$ if $C(x_i) \neq C(x_j)$,

where η and γ are some constants and R_δ is a δ -dependent constant. Assume that the interpolation functions $f^{(v)}$:

$H^{(v)} \rightarrow \mathbb{R}^d$ and $f^{(u)} : H^{(u)} \rightarrow \mathbb{R}^d$ in modalities v and u are Lipschitz continuous with constant L such that for some parameters $\epsilon > 0$ and $\delta > 0$, the following inequality holds

$$6L\delta + 2\sqrt{d}\epsilon + 2R_\delta + 2\eta \leq \gamma. \quad (7)$$

For some $Q \geq 1$, let the number of training samples from class m be such that

$$N_m > \frac{Q}{\eta_{m,\delta}}.$$

Let $x^{(v)} \sim \nu_m^{(v)}$ be a query sample from class m observed in modality v , the relevant samples to which are sought in modality u . Then, with probability at least

$$1 - \exp\left(-\frac{2(N_m\eta_{m,\delta} - Q)^2}{N_m}\right) - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) - (1 - \eta_{m,\delta})^Q$$

the precision rate P of the retrieval algorithm in (5) satisfies

$$\begin{aligned} P &= 1, & \text{if } K \leq Q \\ P &\geq \frac{Q}{K}, & \text{if } K > Q \end{aligned} \quad (8)$$

and the recall rate R of the retrieval algorithm satisfies

$$\begin{aligned} R &= \frac{K}{N_m}, & \text{if } K \leq Q \\ R &\geq \frac{Q}{N_m}, & \text{if } K > Q. \end{aligned} \quad (9)$$

The proof of Theorem 2 is given in the Appendix. Theorem 2 can be interpreted similarly to Theorem 1. The properties (P1), (P2) and (P3) ensure that the learnt embedding aligns modalities v and u sufficiently well, while mapping nearby samples from the same classes to nearby points, and increasing the distance between samples from different classes. Assuming that the condition (7) is satisfied, the precision and recall rates given in (8) and (9) are attained with probability approaching 1 at an exponential rate as the number of training samples increases. In the proof of the theorem, the precision and recall rates in (8) and (9) are obtained by identifying the conditions under which at least Q samples out of the K samples returned by the retrieval algorithm are relevant to the query sample.

The condition (7) required for successful cross-modal retrieval is the same as the condition (2) for accurate multi-modal classification. Hence, similarly to the findings of our multi-modal classification analysis, the results of our retrieval analysis also suggest that it is necessary to find a good compromise between the Lipschitz continuity of the interpolators and the separation between different classes when learning nonlinear embeddings for cross-modal retrieval applications.

IV. PROPOSED MULTI-MODAL SUPERVISED EMBEDDING METHOD

In this section, we propose a multi-modal nonlinear dimensionality reduction algorithm that relies on the theoretical findings of Section III. We formulate the nonlinear embedding problem in Section IV-A and then discuss its solution in Section IV-B.

A. Problem Formulation

Let $X^{(v)} \in \mathbb{R}^{N^{(v)} \times n^{(v)}}$ denote the training data matrix of modality v , each row of which is the observation $x_i^{(v)}$ of some training sample x_i in the v -th modality. Here $N^{(v)}$ is the total number of observations² from all classes in modality v , and $n^{(v)}$ is the dimension of the Hilbert space $H^{(v)}$ of modality v , assumed to be finite in a practical setting. Given the training samples $X^{(v)}$ from modalities $v = 1, \dots, V$, we would like to compute embeddings $Y^{(v)} \in \mathbb{R}^{N^{(v)} \times d}$ of the training samples into the common domain \mathbb{R}^d , such that each $x_i^{(v)} \in \mathbb{R}^{n^{(v)}}$ is mapped to a vector $y_i^{(v)} \in \mathbb{R}^d$. The embedding is extended to the whole data space through interpolation functions $f^{(v)} : \mathbb{R}^{n^{(v)}} \rightarrow \mathbb{R}^d$ such that each training sample is mapped to its embedding as $f^{(v)}(x_i^{(v)}) = y_i^{(v)}$.

Our main purpose is to find an embedding that can be successfully generalized to initially unavailable test samples. We recall from our theoretical analysis that for successful generalization in multi-modal classification and retrieval, the embedding must have the properties (P1), (P2) and (P3) given in Theorems 1 and 2, while the Lipschitz constant of the interpolators must be kept sufficiently small as imposed by the conditions (2) and (7). We now formulate our multi-modal learning problem in the light of these results.

Lipschitz regularity of the interpolators. For the extension of the embedding, we choose to use RBF interpolation functions, which are analytical functions with well-studied properties. Hence, the interpolator of each modality $v = 1, \dots, V$ has the form $f^{(v)}(x^{(v)}) = [f_1^{(v)}(x^{(v)}) \dots f_d^{(v)}(x^{(v)})]$, where

$$f_k^{(v)}(x^{(v)}) = \sum_{i=1}^{N^{(v)}} C_{ik}^{(v)} \phi^{(v)}(\|x^{(v)} - x_i^{(v)}\|) \quad (10)$$

is the k -th component of $f^{(v)}(x^{(v)})$. Here

$$\phi^{(v)}(r) = e^{-r^2/(\sigma^{(v)})^2}$$

is a Gaussian RBF kernel with scale parameter $\sigma^{(v)}$ and $C_{ik}^{(v)}$ are the interpolator coefficients.

The Lipschitz continuity of Gaussian RBF interpolators has been studied in [38], from which it follows that $f^{(v)}(x^{(v)})$ is Lipschitz-continuous with constant

$$L^{(v)} = \sqrt{2}e^{-\frac{1}{2}} \sqrt{N^{(v)}}(\sigma^{(v)})^{-1} \|C^{(v)}\|_F. \quad (11)$$

Here $C^{(v)}$ is the coefficient matrix with entries $C_{ik}^{(v)}$. The interpolator coefficients can be easily obtained as

$$C^{(v)} = (\Psi^{(v)})^{-1} Y^{(v)}$$

by fitting the embedding coordinates $Y^{(v)}$ to the training data $X^{(v)}$, where $\Psi^{(v)} \in \mathbb{R}^{N^{(v)} \times N^{(v)}}$ is the RBF kernel matrix with entries $\Psi_{ij}^{(v)} = \phi^{(v)}(\|x_i^{(v)} - x_j^{(v)}\|)$.

The conditions (2) and (7) suggest that the Lipschitz constants of the interpolators should be sufficiently small for successful generalization of the embedding to test data. In

²Although the observations of all training samples were assumed to be available in all modalities for the simplicity of the theoretical analysis in Section III, here we remove this assumption and allow some observations to be missing in some modalities. Hence $N^{(v)}$ may be different for different v .

view of these results, when learning a nonlinear embedding, we propose to minimize the kernel scale of each modality v through the term

$$\sum_{v=1}^V (\sigma^{(v)})^{-2}$$

as well as the interpolator coefficients of all modalities through

$$\sum_{v=1}^V \|C^{(v)}\|_F^2 = \sum_{v=1}^V \|(\Psi^{(v)})^{-1} Y^{(v)}\|_F^2 = \text{tr}(\tilde{Y}^T \tilde{\Psi}^{-2} \tilde{Y})$$

so that the Lipschitz constant $L^{(v)}$ in (11) is minimized for each modality v . Here

$$\tilde{Y} = [(Y^{(1)})^T (Y^{(2)})^T \dots (Y^{(V)})^T]^T \in \mathbb{R}^{N \times d}$$

denotes the matrix containing the embeddings from all modalities (with $N = \sum_v N^{(v)}$) and $\tilde{\Psi} \in \mathbb{R}^{N \times N}$ is a block-diagonal matrix containing the kernel matrix $\Psi^{(v)}$ in its v -th block.

Within-class compactness. Theorems 1 and 2 suggest that the constant R_δ in (P2) should be kept small, so that the conditions (2) and (7) are more likely to be met. Although it is not easy to analytically formulate the minimization of R_δ , in practice if nearby samples from the same modality and same class are embedded into nearby points, R_δ will be small. This problem is well-studied in the manifold learning literature. The total weighted distance between the embeddings of same-class samples can be formulated as

$$\sum_{v=1}^V \sum_{i,j=1}^{N^{(v)}} (W_w^{(v)})_{ij} \|y_i^{(v)} - y_j^{(v)}\|^2 = \text{tr}(\tilde{Y}^T \tilde{L}_w \tilde{Y}). \quad (12)$$

Here $W_w^{(v)} \in \mathbb{R}^{N^{(v)} \times N^{(v)}}$ is chosen as a weight matrix whose entries $(W_w^{(v)})_{ij} = \exp(-\|x_i^{(v)} - x_j^{(v)}\|^2/(\theta^{(v)})^2)$ represent the affinity between the data samples when $x_i^{(v)}$ and $x_j^{(v)}$ are from the same class (for a scale parameter $\theta^{(v)}$), and $(W_w^{(v)})_{ij} = 0$ otherwise. In the equality, the block-diagonal matrix $\tilde{L}_w \in \mathbb{R}^{N \times N}$ contains the within-class Laplacian $L_w^{(v)} = D_w^{(v)} - W_w^{(v)}$ in its v -th block, where $D_w^{(v)}$ is the diagonal degree matrix with i -th diagonal entry given by $\sum_j (W_w^{(v)})_{ij}$. The term in (12) hence imposes nearby samples $x_i^{(v)}$, $x_j^{(v)}$ from the same class and the same modality to be mapped to nearby coordinates.

Between-class separation. In Theorems 1 and 2, the between-class margin γ in (P3) must be sufficiently large for conditions (2) and (7) to be satisfied. Since it is difficult to formulate the maximization of the exact value of γ , we relax this problem to the maximization of

$$\sum_{v=1}^V \sum_{i,j=1}^{N^{(v)}} (W_b^{(v)})_{ij} \|y_i^{(v)} - y_j^{(v)}\|^2 = \text{tr}(\tilde{Y}^T \tilde{L}_b \tilde{Y})$$

which aims to increase the separation between the samples from different classes within each modality v . Here the matrix $W_b^{(v)} \in \mathbb{R}^{N^{(v)} \times N^{(v)}}$ has entries $(W_b^{(v)})_{ij} = 1$ when $x_i^{(v)}$ and $x_j^{(v)}$ are from different classes; and $(W_b^{(v)})_{ij} = 0$, otherwise. The block-diagonal matrix $\tilde{L}_b \in \mathbb{R}^{N \times N}$ contains the between-class Laplacian $L_b^{(v)} = D_b^{(v)} - W_b^{(v)}$ in its v -th block, where

$D_b^{(v)}$ is the diagonal between-class degree matrix with i -th diagonal entry given by $\sum_j (W_b^{(v)})_{ij}$.

Cross-modal alignment. Finally, the constant η in property (P1) in Theorems 1 and 2 should be sufficiently small for conditions (2) and (7) to be met. The parameter η represents the distance between the embeddings of the observations of the same sample in different modalities. We relax the minimization of η to the minimization of the following term, which aims to embed samples of high affinity from different modalities v, u into nearby points

$$\sum_{v=1}^V \sum_{u \neq v} \sum_{i=1}^{N^{(v)}} \sum_{j=1}^{N^{(u)}} (W_w^{(vu)})_{ij} \left\| y_i^{(v)} - y_j^{(u)} \right\|^2 = \text{tr}(\tilde{Y}^T \tilde{L}_{cw} \tilde{Y}).$$

Here, the matrix $W_w^{(vu)} \in \mathbb{R}^{N^{(v)} \times N^{(u)}}$ encodes the affinities between sample pairs from different modalities. $(W_w^{(vu)})_{ij}$ is nonzero only if $x_i^{(v)}$ and $x_j^{(u)}$ are from the same class, in which case it is computed with the Gaussian kernel based on the distance between $x_i^{(v)}$ and $x_j^{(u)}$ when transferred to a common modality (i.e., using $\|x_i^{(v)} - x_j^{(u)}\|$ or $\|x_i^{(u)} - x_j^{(u)}\|$, otherwise $\|x_i^{(r)} - x_j^{(r)}\|$ in some other modality r if the former ones are not possible). Denoting by $\tilde{W}_{cw} \in \mathbb{R}^{N \times N}$ the cross-modal within-class weight matrix containing $W_w^{(vu)}$ in its (v, u) -th block, the corresponding Laplacian matrix $\tilde{L}_{cw} \in \mathbb{R}^{N \times N}$ is computed as $\tilde{L}_{cw} = \tilde{D}_{cw} - \tilde{W}_{cw}$, where \tilde{D}_{cw} is the diagonal degree matrix with i -th diagonal entry given by $\sum_j (\tilde{W}_{cw})_{ij}$.

Meanwhile, the property (P3) in Theorems 1 and 2 suggests that two samples from modalities v, u should be separated if they are from different classes. We thus propose to maximize

$$\sum_{v=1}^V \sum_{u \neq v} \sum_{i=1}^{N^{(v)}} \sum_{j=1}^{N^{(u)}} (W_b^{(vu)})_{ij} \left\| y_i^{(v)} - y_j^{(u)} \right\|^2 = \text{tr}(\tilde{Y}^T \tilde{L}_{cb} \tilde{Y})$$

where the matrix $W_b^{(vu)} \in \mathbb{R}^{N^{(v)} \times N^{(u)}}$ is formed by setting $(W_b^{(vu)})_{ij} = 1$ if $x_i^{(v)}$ and $x_j^{(u)}$ are from different classes, and 0 otherwise. The cross-modal between-class weight matrix $\tilde{W}_{cb} \in \mathbb{R}^{N \times N}$ contains the matrix $W_b^{(vu)}$ in its (v, u) -th block, while $\tilde{L}_{cb} \in \mathbb{R}^{N \times N}$ is the corresponding Laplacian matrix given by $\tilde{L}_{cb} = \tilde{D}_{cb} - \tilde{W}_{cb}$, with \tilde{D}_{cb} denoting the diagonal degree matrix with i -th diagonal entry given by $\sum_j (\tilde{W}_{cb})_{ij}$.

Overall problem. We now combine all these objectives in the following overall optimization problem

$$\begin{aligned} & \text{minimize}_{\tilde{Y}, \{\sigma^{(v)}\}} \text{tr}(\tilde{Y}^T \tilde{L}_w \tilde{Y}) - \mu_1 \text{tr}(\tilde{Y}^T \tilde{L}_b \tilde{Y}) + \mu_2 \text{tr}(\tilde{Y}^T \tilde{\Psi}^{-2} \tilde{Y}) \\ & + \mu_3 \sum_{v=1}^V (\sigma^{(v)})^{-2} + \mu_4 \text{tr}(\tilde{Y}^T \tilde{L}_{cw} \tilde{Y}) - \mu_5 \text{tr}(\tilde{Y}^T \tilde{L}_{cb} \tilde{Y}) \end{aligned} \quad (13)$$

subject to $\tilde{Y}^T \tilde{Y} = I$, where μ_1, \dots, μ_5 are positive weight parameters, $I \in \mathbb{R}^{d \times d}$ is the identity matrix, and the optimization constraint $\tilde{Y}^T \tilde{Y} = I$ is for the normalization of the learnt coordinates.

B. Solution of the Optimization Problem

Defining

$$A = \tilde{L}_w - \mu_1 \tilde{L}_b + \mu_2 \tilde{\Psi}^{-2} + \mu_4 \tilde{L}_{cw} - \mu_5 \tilde{L}_{cb} \quad (14)$$

the problem in (13) can be rewritten as

$$\text{minimize}_{\tilde{Y}, \{\sigma^{(v)}\}} \text{tr}(\tilde{Y}^T A \tilde{Y}) + \mu_3 \sum_{v=1}^V (\sigma^{(v)})^{-2}, \text{ subject to } \tilde{Y}^T \tilde{Y} = I. \quad (15)$$

The above problem is not jointly convex in \tilde{Y} and $\{\sigma^{(v)}\}$, hence it is not easy to find its global optimum. We minimize the objective function with an iterative alternating optimization scheme, where we first optimize \tilde{Y} by fixing $\{\sigma^{(v)}\}$, and then optimize $\{\sigma^{(v)}\}$ by fixing \tilde{Y} in each iteration as follows.

Optimization of \tilde{Y} : When $\{\sigma^{(v)}\}$ are fixed, the optimization problem in (15) becomes

$$\text{minimize}_{\tilde{Y}} \text{tr}(\tilde{Y}^T A \tilde{Y}) \text{ subject to } \tilde{Y}^T \tilde{Y} = I. \quad (16)$$

The solution to this problem is given by the d eigenvectors of the matrix A corresponding to its smallest d eigenvalues.

Optimization of $\{\sigma^{(v)}\}$: Fixing \tilde{Y} , the problem (15) becomes

$$\text{minimize}_{\{\sigma^{(v)}\}} \mu_2 \text{tr}(\tilde{Y}^T \tilde{\Psi}^{-2} \tilde{Y}) + \mu_3 \sum_{v=1}^V (\sigma^{(v)})^{-2}. \quad (17)$$

Note that the first term in the objective depends on the kernel scale parameters $\{\sigma^{(v)}\}$ through the entries of the kernel matrix $\tilde{\Psi}$. Due to the block diagonal structure of $\tilde{\Psi}$ and the separability of the second term, the objective (17) can be decomposed into V individual objectives, each one of which is a function of only one scale parameter $\sigma^{(v)}$. We minimize these objective functions one by one, by optimizing one scale parameter $\sigma^{(v)}$ at a time through exhaustive search.

If μ_1 and μ_5 are sufficiently small, the matrix A becomes positive semi-definite. In this case, the objective function is guaranteed to converge since it is nonnegative, and both updates on \tilde{Y} and $\{\sigma^{(v)}\}$ reduce it. We continue the iterations until the convergence of the objective. We call the proposed algorithm Multi-modal Nonlinear Supervised Embedding (MNSE), which is summarized in Algorithm 1.

Algorithm 1 Multi-modal Nonlinear Supervised Embedding (MNSE)

Input: Training data matrices $\{X^{(v)}\}$ and training data labels

Initialization:

Obtain the graph Laplacian matrices $\tilde{L}_w, \tilde{L}_b, \tilde{L}_{cw}$, and \tilde{L}_{cb} .

Assign weight parameters $\{\mu_1, \mu_2, \dots, \mu_5\}$, and initial kernel scales $\sigma^{(v)}$

repeat

 Compute the nonlinear embeddings \tilde{Y} through (16) by fixing $\{\sigma^{(v)}\}$

 Compute the kernel scale parameters $\{\sigma^{(v)}\}$ through (17) by fixing \tilde{Y}

until the maximum number of iterations or the convergence of the objective

Output:

Kernel scale parameters $\sigma^{(v)}$ and projected training data $Y^{(v)}$

Kernel coefficients $C^{(v)} = (\Psi^{(v)})^{-1} Y^{(v)}$

C. Complexity Analysis

The complexity of the proposed MNSE method is mainly determined by those of the problems (16) and (17) repeated in the main loop of the algorithm. When computing the matrix A in (14), the matrices $\tilde{L}_w, \tilde{L}_b, \tilde{\Psi}, \tilde{L}_{cw}$, and \tilde{L}_{cb} can be constructed with complexity not exceeding $O(N^2)$, where $N = \sum_v N^{(v)}$ is the total number of observations from all modalities. The eigenvalue decomposition step in (16) is of

complexity $O(N^3)$. In the optimization problem (17), the evaluation of the objective for each $\sigma^{(v)}$ value requires $O((N^{(v)})^3)$ operations in modality v ; hence, the total complexity of finding all $\{\sigma^{(v)}\}$ is smaller than $O(N^3)$. Therefore, the overall complexity of the algorithm is determined as $O(N^3)$.

V. EXPERIMENTAL RESULTS

A. Data sets

The following data sets are used in the experiments.

The *MIT-CBCL multi-view face data set* [40] contains face images of 10 participants captured under 36 illumination conditions and 9 different pose angles. Images with frontal and profile poses are used as Modality 1 and Modality 2, respectively. Images are converted to greyscale and downsampled to a resolution of 30×30 pixels. The experiments are repeated 10 times by randomly dividing the data set into training and test images.

The *Multi-PIE multi-view face data set* [41] consists of face images of many participants under varying camera angles and facial expressions. We conduct our experiments on a cropped and reduced version of this data set [42], where the images of 120 participants captured under 6 camera angles, 20 lighting conditions, and 2 facial expressions (neutral and smiling) are used. Greyscale images of resolution 32×32 pixels are used. Modality 1 and Modality 2 are respectively chosen as the frontal camera angle (0°) and the five other camera angles (15° , 30° , 45° , 60° , 75°). Four experimentation settings are prepared. In each setting, the images of the participants under either the first 10 or the last 10 lighting conditions, and either the neutral or the smiling facial expression are used for training; and the rest of the images are used for test. The results are averaged over these four experimentation settings.

In the classification experiments with the MIT-CBCL and the Multi-PIE data sets, embedding parameters are learnt using the training images from both modalities. Test images are assumed to be available in only one modality and mapped to the common domain with the learnt embeddings. The class labels of test images are estimated via NN classification on the embeddings of the training images of their own modality.

The *Wikipedia image-text data set* [43] contains 2866 image-text pairs describing the contents of Wikipedia articles, which are categorized into 10 classes. The image-text pairs are randomly divided into 1300 training and 1566 test pairs in each trial of the experiments and the results are averaged over 10 trials. 128-dimensional SIFT histogram features are used in the image modality, and 10-dimensional text features obtained with a latent Dirichlet allocation model are used in the text modality [44], [45].

The *Pascal VOC2007 image-text data set* [46] contains image-text pairs from 20 different object classes, where the experiments are done on 2808 training and 2841 test pairs whose images contain only one object. GIST feature vectors in the image modality and word count feature vectors in the text modality are used in the experiments.

In the retrieval experiments with the Wikipedia and the Pascal VOC2007 data sets, embedding functions are learnt using the training set, and then the relevant matches of an

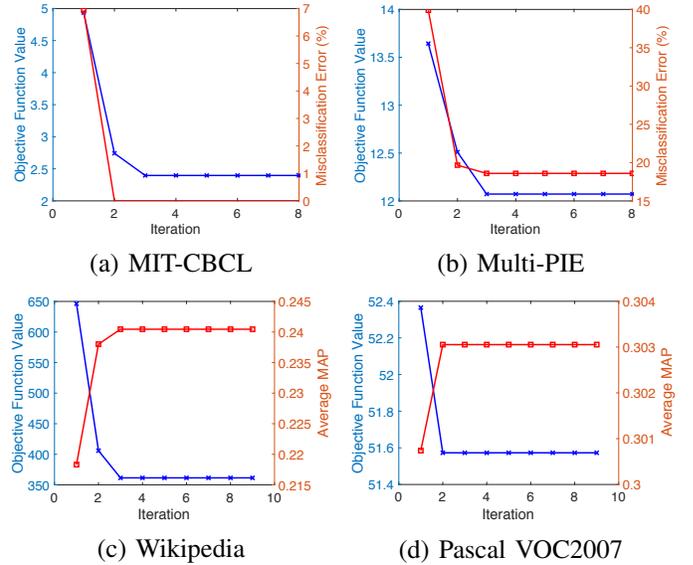


Fig. 2. The evolution of the objective function and the algorithm performance during the iterations for all four data sets

image (text) query are searched in the embedding of the text (image) database based on cosine similarity. The precision and recall rates are computed as in (6) by considering a retrieved item relevant if it is from the same class as the query sample. The Mean Average Precision (MAP) scores of the methods are computed by averaging all average precision values over all query samples. When computing the MAP scores, for each query sample the number of retrieved items is set so as to retrieve all training samples relevant to it.

B. Stabilization and Sensitivity Analysis of MNSE

We first study the stabilization of the proposed MNSE algorithm and its sensitivity to the algorithm parameters. The image classification performance of MNSE is analyzed on the MIT-CBCL and the Multi-PIE data sets. 100 training and 260 test images are used for the MIT-CBCL data set. The average of the misclassification errors of Modalities 1 and 2 is reported. The retrieval performance of the algorithm is studied on the Wikipedia and the Pascal VOC2007 data sets. The MAP scores are averaged over the image and the text queries.

We first study in Figure 2 the evolution of the objective function (13) along with the classification and the retrieval performance of MNSE throughout the optimization iterations on all four data sets. The objective function is seen to steadily decrease throughout the iterations as expected. The updates on both the embeddings $\{Y^{(v)}\}$ and the kernel scale parameters $\{\sigma^{(v)}\}$ ensure that the objective function is non-increasing. The improvement in the misclassification errors or MAP scores follows the decrease in the objective during the iterations. This suggests that the proposed objective function is indeed well-representative of the performance of the algorithm.

Next, the effect of the weight parameters $\mu_1, \mu_2, \mu_3, \mu_4, \mu_5$ on the algorithm performance is examined in Tables I-IV for all four data sets. The performance on the Pascal VOC2007 and Wikipedia data sets is measured on half of the training

TABLE I

THE VARIATION OF THE MISCLASSIFICATION ERROR WITH THE WEIGHT PARAMETERS FOR THE MIT-CBCL DATA SET. FIXED PARAMETERS ARE CHOSEN AS $\mu_1 = \mu_4 = \mu_5 = 10^2$ IN UPPER TABLE AND $\mu_2 = 10^{-3}$, $\mu_3 = 1$ IN LOWER TABLE.

$\mu_3 \setminus \mu_2$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	0.73	0.35	0.35	0.35	0.35	0.35	1.5	10.7
10^{-3}	0.16	0.27	0.30	0.33	0.35	0.35	1.5	10.7
10^{-2}	0.16	0.23	0.27	0.30	0.33	0.35	1.51	10.6
10^{-1}	0.16	0.20	0.23	0.26	0.29	0.31	1.51	10.6
1	0.16	0.18	0.21	0.24	0.27	0.29	1.34	10.1
10^1	0.16	0.19	0.22	0.23	0.27	0.31	1.52	9.75
10^2	0.16	0.19	0.22	0.24	0.36	0.91	5.00	11.4
10^3	0.16	0.19	0.22	0.27	2.21	7.28	10.5	12.9

$\mu_4 \setminus \mu_1, \mu_5$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	7.94	6.75	6.23	4.68	1.42	0.24	0.21	0.18
10^{-3}	7.61	5.26	5.28	4.63	1.43	0.24	0.21	0.18
10^{-2}	6.88	5.83	5.88	3.38	1.34	0.24	0.21	0.18
10^{-1}	2.92	2.62	1.62	1.90	0.86	0.24	0.21	0.18
1	1.35	1.10	0.71	0.65	0.24	0.23	0.20	0.18
10^1	0.94	0.77	0.59	1.00	0.22	0.22	0.20	0.18
10^2	0.86	0.81	0.75	0.76	0.18	0.18	0.18	0.18
10^3	0.78	0.75	0.66	0.62	0.18	0.18	0.18	0.18

TABLE II

THE VARIATION OF THE MISCLASSIFICATION ERROR WITH THE WEIGHT PARAMETERS FOR THE MULTI-PIE DATA SET. FIXED PARAMETERS ARE CHOSEN AS $\mu_1 = \mu_4 = \mu_5 = 10^2$ IN UPPER TABLE AND $\mu_2 = 10^{-3}$, $\mu_3 = 1$ IN LOWER TABLE.

$\mu_3 \setminus \mu_2$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	70.5	44.0	49.8	45.1	59.9	88.8	86.9	76.0
10^{-3}	17.0	35.6	49.8	45.1	59.9	88.8	86.9	76.0
10^{-2}	17.0	18.5	42.4	48.0	58.9	88.8	86.9	76.0
10^{-1}	17.0	17.4	24.5	39.6	59.3	88.8	86.9	76.0
1	17.0	15.3	19.8	26.8	53.8	88.8	86.9	76.0
10^1	17.0	15.2	16.4	21.4	50.7	88.8	86.9	76.0
10^2	17.0	16.8	16.6	17.1	28.0	80.0	83.8	62.7
10^3	17.0	17.4	18.1	19.4	21.1	43.2	54.4	52.8

$\mu_4 \setminus \mu_1, \mu_5$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	29.9	29.1	41.9	44.9	77.3	89.4	91.8	95.5
10^{-3}	28.2	28.3	40.2	43.6	77.5	89.4	91.8	95.5
10^{-2}	22.1	22.2	20.7	31.5	75.9	89.4	91.8	95.5
10^{-1}	17.3	17.5	17.9	18.3	67.9	89.9	92.3	95.9
1	16.44	16.5	17.6	24.1	43.3	84.7	88.9	96.3
10^1	16.0	16.1	16.2	21.7	21.8	16.7	48.3	63.4
10^2	16.0	16.0	16.2	22.0	22.2	23.1	15.3	39.6
10^3	16.0	16.1	16.1	21.0	21.1	22.6	24.0	15.5

samples assigned for validation. The upper half of each table shows the variation of the average performance with μ_2 and μ_3 , and the lower half of each table shows the variation with μ_4 and $\mu_1 = \mu_5$. The parameters μ_1 and μ_5 are set to be equal, motivated by the similarity in the construction of the between-class separation matrices associated with these parameters. Tables I and II show that setting $\mu_2 = 10^{-3}$ and choosing μ_3 in the interval $[1, 10^3]$ leads to reasonably small misclassification error in both data sets. A suitable choice for μ_1, μ_4 , and μ_5 seems to be $\mu_4 = \mu_1 = \mu_5 \in [10^2, 10^3]$. On the other hand, the optimal parameter ranges are slightly different in the retrieval experiments. Choosing $\mu_2 = 1$ and $\mu_3 \in [1, 10^3]$ maximizes the MAP score for both the Wikipedia and the Pascal VOC2007 data sets, while setting $\mu_4 = 100\mu_1 = 100\mu_5$ and selecting $\mu_4 \in [1, 10]$ gives close to optimal performance. As an overall conclusion, the observation that μ_2 and μ_3

TABLE III

THE VARIATION OF THE MAP WITH THE WEIGHT PARAMETERS FOR THE WIKIPEDIA DATA SET. FIXED PARAMETERS ARE CHOSEN AS $\mu_1 = \mu_4 = \mu_5 = 10^{-3}$ IN UPPER TABLE AND $\mu_2 = \mu_3 = 1$ IN LOWER TABLE.

$\mu_3 \setminus \mu_2$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	0.10	0.19	0.17	0.19	0.22	0.21	0.17	0.17
10^{-3}	0.11	0.18	0.17	0.19	0.22	0.21	0.17	0.17
10^{-2}	0.11	0.17	0.18	0.19	0.23	0.21	0.17	0.17
10^{-1}	0.11	0.17	0.20	0.22	0.23	0.21	0.17	0.17
1	0.11	0.20	0.21	0.22	0.23	0.22	0.18	0.17
10^1	0.11	0.20	0.21	0.22	0.23	0.22	0.17	0.17
10^2	0.11	0.20	0.21	0.22	0.23	0.22	0.19	0.18
10^3	0.11	0.20	0.21	0.22	0.23	0.22	0.20	0.18

$\mu_4 \setminus \mu_1, \mu_5$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	0.13	0.19	0.19	0.18	0.11	0.13	0.13	0.13
10^{-3}	0.22	0.23	0.20	0.18	0.11	0.13	0.13	0.13
10^{-2}	0.22	0.22	0.23	0.20	0.11	0.13	0.13	0.13
10^{-1}	0.22	0.22	0.22	0.23	0.11	0.13	0.13	0.13
1	0.21	0.21	0.21	0.21	0.11	0.12	0.13	0.13
10^1	0.20	0.20	0.20	0.20	0.20	0.11	0.13	0.13
10^2	0.19	0.19	0.19	0.19	0.19	0.18	0.11	0.13
10^3	0.19	0.19	0.19	0.19	0.19	0.18	0.18	0.11

TABLE IV

THE VARIATION OF THE MAP WITH THE WEIGHT PARAMETERS FOR THE PASCAL VOC2007 DATA SET. FIXED PARAMETERS ARE CHOSEN AS $\mu_1 = \mu_5 = 10^{-1}$, $\mu_4 = 10$ IN UPPER TABLE AND $\mu_2 = \mu_3 = 1$ IN LOWER TABLE.

$\mu_3 \setminus \mu_2$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	0.21	0.24	0.23	0.22	0.29	0.28	0.22	0.22
10^{-3}	0.25	0.25	0.24	0.23	0.29	0.28	0.22	0.22
10^{-2}	0.25	0.25	0.24	0.23	0.29	0.28	0.22	0.22
10^{-1}	0.25	0.25	0.24	0.25	0.29	0.28	0.23	0.22
1	0.25	0.25	0.24	0.25	0.30	0.28	0.24	0.22
10^1	0.25	0.25	0.24	0.25	0.30	0.28	0.24	0.24
10^2	0.25	0.25	0.24	0.25	0.30	0.28	0.24	0.24
10^3	0.25	0.25	0.24	0.25	0.30	0.28	0.24	0.23

$\mu_4 \setminus \mu_1, \mu_5$	0	10^{-3}	10^{-2}	10^{-1}	1	10^1	10^2	10^3
0	0.10	0.16	0.16	0.16	0.12	0.12	0.12	0.12
10^{-3}	0.19	0.20	0.19	0.16	0.12	0.12	0.12	0.12
10^{-2}	0.19	0.19	0.20	0.20	0.12	0.12	0.12	0.12
10^{-1}	0.21	0.21	0.21	0.22	0.11	0.11	0.12	0.12
1	0.29	0.29	0.29	0.30	0.10	0.12	0.12	0.12
10^1	0.30	0.30	0.30	0.30	0.25	0.06	0.09	0.09
10^2	0.30	0.30	0.30	0.30	0.25	0.23	0.06	0.08
10^3	0.30	0.30	0.30	0.30	0.25	0.25	0.23	0.06

should have nonzero values in both classification and retrieval experiments confirms that the Lipschitz regularity terms in the objective function are necessary for high performance, which is a central idea in our study. Next, it may be more advantageous to select μ_4 relatively higher than μ_1 and μ_5 in retrieval experiments, unlike in classification where they can be chosen equal. This suggests that the alignment across different modalities may become even more critical in retrieval problems.

Finally, we examine the effect of the embedding dimension d on the algorithm performance in Figure 3. Different curves correspond to the misclassification errors obtained at different training sizes (ratio of training samples) for MIT-CBCL, and at different camera angles for Modality 2 for the Multi-PIE data set. The MAP scores of the image and text queries are reported individually for the Wikipedia and the Pascal VOC2007 data

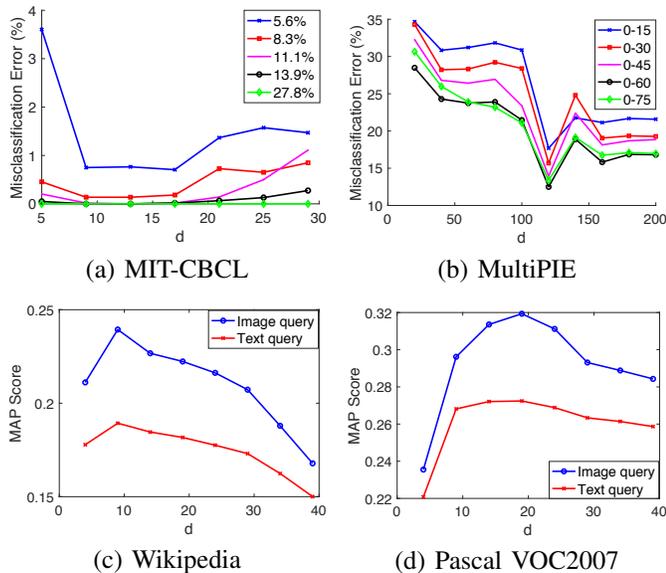


Fig. 3. Variation of algorithm performance with embedding dimension d

sets. The optimal embedding dimension is consistently seen to be close to the number of classes in all four data sets. In the rest of our experiments, the embedding dimension d is chosen as $M - 1$ for each data set, where M is the number of classes.

C. Evaluation of the Algorithm Performance

We now evaluate the performance of the proposed MNSE algorithm with comparative experiments in image classification and image-text retrieval applications. MNSE is compared to the multi-modal representation learning algorithms CCA, Kernel CCA [1], GMLDA [2], Kernel GMLDA [2], KMvMDA [21], JFSSL [25], and DeepMF [34], as well as the baseline single-modal methods PCA, NN classification in the original domain, and NSSE [38]. We use our own implementations for the GMLDA, Kernel GMLDA, KMvMDA, and JFSSL methods in the experiments. The CCA and GMLDA algorithms are applied after a preprocessing step of dimensionality reduction with PCA, which has been seen to improve their performance. The single-modal methods are applied independently in each modality. The parameters of the compared algorithms are optimized for the best performance on the test set. The parameter settings used for the algorithms that require parameter tuning are given in Table V for all data sets (intervals mean that parameters may vary in different repetitions of the experiments.)

1) *Multi-modal image classification*: The multi-modal classification experiments are done on the MIT-CBCL and the Multi-PIE face data sets. The weight parameters of the proposed MNSE method are chosen within the regions suggested in Section V-B as $\mu_1 = \mu_4 = \mu_5 = 10^2$, $\mu_2 = 10^{-3}$, $\mu_3 = 1$ for both data sets. Tables VI and VII show the test misclassification errors in percentage for the MIT-CBCL and the Multi-PIE data sets. The upper and lower tables show the average errors and the error standard deviations over the 10 random trials for MIT-CBCL, and over the 4 experimentation

TABLE V
ALGORITHM PARAMETER VALUES IN THE EXPERIMENTS (DENOTED AS IN THE PAPERS CITED IN THE TABLE)

Algorithm	Parameter	MIT-CBCL	Multi-PIE	Pascal VOC	Wikipedia
PCA	d (dim.)	15	82	-	-
	d	8-20	119	-	-
NSSE [38]	μ_1	42-400	400	-	-
	μ_2	$10^{-4} - 0.1$	0.1	-	-
	μ_3	1-3	1.2	-	-
	μ_4	-	-	-	-
CCA [1]	d (dim.)	15	80	7	8
Kernel CCA [47]	d (dim.)	15	220	19	8
	c'	10^{-5}	0.1	100	1000
	Kernel type	Gauss.	Gauss.	Gauss.	Gauss.
GMLDA [2]	k (dim.)	15	124	9	9
	α	10	10	500	500
	μ	1	2	0.05	0.01
Kernel GMLDA [2]	k (dim.)	15-18	124	9	9
	α	10-100	100	500	500
	μ	0.8-1.5	2	0.05	0.01
	Kernel type	Chi sq.	Chi sq.	Chi sq.	Chi sq.
JFSSL [25]	k	30	40	300	130
	λ_1	0.01	0.1	0.01	0.01
	λ_2	0.01	0.01	0.001	0.001
	β	1	1	1	1
	ϵ	10^{-8}	10^{-8}	10^{-8}	10^{-8}
	ϵ	10^{-8}	10^{-8}	10^{-8}	10^{-8}
DeepMF [34]	γ	0.95	0.95	0.5	0.5
	β	0.01	0.01	0.01	0.01
	$[p_1, p_2]$	[250, 45]	[250, 45]	[100, 50]	[100, 10]
	k	5	5	5	5
KMvMDA [21]	d	38-120	119	19	9
	Kernel type	Gauss.	Gauss.	Gauss.	Gauss.
MNSE	d	9	119	19	9
	$\mu_1 = \mu_5$	100	100	0.1	0.001-0.1
	μ_2	0.001	0.001	1	0.1-1
	μ_3	1	1	1	0.001-10
	μ_4	100	100	10	0.001-0.1

settings considered for the Multi-PIE data set. The error is studied with respect to the training size (ratio of the training samples) for MIT-CBCL, and the camera angle of Modality 2 for Multi-PIE. The errors obtained for Modalities 1 and 2 are given in the top and the bottom rows for each method.

The results in Tables VI and VII show that the proposed MNSE method outperforms all methods except for NSSE and JFSSL in all experiments. On the MIT-CBCL data set, the linear supervised JFSSL algorithm performs the best. The approach of aligning the modalities via linear projections in JFSSL is particularly suited to this synthetic and regularly structured face data set. The proposed MNSE method closely follows JFSSL on MIT-CBCL with very small misclassification rates. On the other hand, MNSE outperforms JFSSL and the other multi-modal algorithms on the Multi-PIE data set. The nonlinear MNSE method learns a relatively rich model while explicitly incorporating generalization performance in its objective. These features bring our method an advantage on the Multi-PIE data set, which has a more challenging structure than MIT-CBCL due to the large number of classes and facial expression and illumination variations. MNSE is seen to perform better than NSSE in most experiments, which also performs very well on this data set. NSSE computes nonlinear and smooth embeddings as MNSE; however, in a single modality. The fact that MNSE often outperforms NSSE confirms that it can successfully exploit and combine the information from both modalities when optimizing the embedding parameters. The standard deviations of the errors suggest that the algorithms with small average errors maintain rather stable performance figures over different repetitions of the experiments for both data sets.

TABLE VI
MISCLASSIFICATION ERRORS (%) FOR THE MIT-CBCL DATA SET:
AVERAGE ERRORS (UPPER TABLE) AND ERROR STANDARD DEVIATIONS
(LOWER TABLE) OVER 10 RANDOM TRIALS.

Avg. errors	Training size				
	5.6%	8.3%	11.1%	13.9%	27.8%
NN	22.12	19	10.69	2.97	0.77
	19.68	17.64	6.94	1.71	0
PCA	3.68	0.06	0.34	0.10	0
	4.29	0.54	0.06	0	0
NSSE [38]	1.94	0.03	0.03	0	0
	4.56	1.00	0.09	0.03	0
CCA [1]	3.67	0.06	0.34	0.10	0
	4.29	0.55	0.06	0	0
Kernel CCA [1]	1.50	0.21	0.21	0.09	0
	4.00	0.72	0.66	0.39	0.19
GMLDA [2]	2.56	0	0	0	0
	5.82	0.30	0.06	0.03	0
Kernel GMLDA [2]	7.68	0.70	1.81	0.39	0
	3.15	2.09	0.50	0.19	0
JFSSL [25]	0	0	0	0	0
	0.12	0	0	0	0
DeepMF [34]	7.12	1.58	1.00	0.65	0
	5.97	1.36	0.63	0.45	0
KMvMDA [21]	11.91	0.48	1.72	0.84	0
	17.50	0.70	1.22	0.58	0
MNSE	0.15	0	0	0	0
	1.35	0.27	0.03	0	0
Error st. dev.	Training size				
	5.6%	8.3%	11.1%	13.9%	27.8%
NN	21.13	23.25	9.97	2.32	0.65
	18.39	21.81	7.71	2.73	0
PCA	4.46	0.13	0.88	0.31	0
	9.55	1.07	0.20	0	0
NSSE [38]	3.14	0.10	0.10	0	0
	3.98	2.38	0.30	0.10	0
CCA [1]	4.46	0.13	0.88	0.31	0
	9.55	1.08	0.20	0	0
Kernel CCA [1]	2.88	0.67	0.69	0.31	0
	5.40	0.98	0.76	0.60	0.49
GMLDA [2]	2.16	0	0	0	0
	3.76	0.96	0.20	0.10	0
Kernel GMLDA [2]	7.60	1.51	2.86	1.02	0
	4.71	4.32	1.47	0.61	0
JFSSL [25]	0	0	0	0	0
	0.28	0	0	0	0
DeepMF [34]	5.88	1.70	0.75	1.05	0
	4.93	1.53	0.61	0.88	0
KMvMDA [21]	10.52	1.13	2.97	1.65	0
	10.02	0.89	1.92	1.02	0
MNSE	0.25	0	0	0	0
	1.55	0.76	0.10	0	0

2) *Cross-modal image-text retrieval*: The retrieval experiments are done on the Wikipedia and the Pascal VOC2007 image-text data sets. The weight parameters of MNSE are set to the values giving the highest average MAP score over the validation set of each experiment in the results of Section V-B. Figures 4 and 5 show the precision-recall and precision-scope curves for both types of queries, respectively on the Wikipedia and the Pascal VOC2007 data sets. Table VIII reports the MAP scores of the methods on both data sets.

The proposed MNSE method outperforms all other multi-modal methods on both data sets. The Wikipedia and the Pascal VOC2007 data sets have diverse and irregular structures, with the two modalities bearing little resemblance. This makes the multi-modal representation learning task rather challenging, where the flexibility of the proposed nonlinear supervised embedding approach brings clear advantages over the other multi-modal methods in comparison. The performance gap between the proposed nonlinear MNSE method and the linear JFSSL and GMLDA algorithms can be explained in the way

TABLE VII
MISCLASSIFICATION ERRORS (%) FOR THE MULTI-PIE DATA SET:
AVERAGE ERRORS (UPPER TABLE) AND ERROR STANDARD DEVIATIONS
(LOWER TABLE) OVER 4 EXPERIMENTATION SETTINGS

Avg. errors	Camera angle for Modality 2				
	15°	30°	45°	60°	75°
NN	20.28	20.28	20.28	20.28	20.28
	21.40	18.93	18.21	12.27	17.01
PCA	22.51	22.51	22.51	22.51	22.51
	24.59	22.06	20.03	15.81	17.23
NSSE [38]	17.46	17.46	17.46	17.46	17.46
	19.85	18.12	16.12	10.78	15.51
CCA [1]	22.80	22.80	22.80	22.80	22.80
	24.47	22.03	20.06	15.88	17.15
Kernel CCA [1]	20.88	21.22	21.08	21.67	21.08
	23.31	24.25	29.01	21.53	29.18
GMLDA [2]	26.92	26.18	25.74	25.35	25.26
	28.73	26.08	21.54	16.71	18.42
Kernel GMLDA [2]	45.42	45.43	49.87	53.42	56.36
	42.99	37.63	29.72	34.74	36.73
JFSSL [25]	24.60	24.61	24.62	24.61	24.62
	30.77	25.30	18.93	16.12	19.99
DeepMF [34]	29.28	26.67	26.31	26.01	26.06
	33.15	32.24	32.73	25.48	29.48
KMvMDA [21]	46.88	48.67	44.58	39.40	46.20
	50.97	52.74	47.16	37.68	50.92
MNSE	17.10	17.11	18.51	18.22	17.60
	18.85	14.24	11.26	9.85	10.44
Error st. dev.	Camera angle for Modality 2				
	15°	30°	45°	60°	75°
NN	1.35	1.35	1.35	1.35	1.35
	1.41	1.80	2.16	0.28	1.73
PCA	1.03	1.03	1.03	1.03	1.03
	1.51	2.37	1.50	1.99	3.94
NSSE [38]	0.21	0.21	0.21	0.21	0.21
	2.09	1.89	1.30	0.72	1.84
CCA [1]	1.15	1.15	1.15	1.15	1.15
	1.63	2.16	1.55	2.00	3.86
Kernel CCA [1]	1.86	1.22	1.67	1.24	1.66
	1.79	0.75	3.79	0.98	3.39
GMLDA [2]	2.11	2.14	1.62	2.19	1.66
	3.74	4.44	3.47	3.96	5.34
Kernel GMLDA [2]	16.33	12.87	13.54	22.92	14.99
	13.37	11.74	8.84	19.36	19.74
JFSSL [25]	0.69	0.67	0.67	0.69	0.67
	1.35	1.35	1.73	1.29	1.23
DeepMF [34]	1.93	1.80	1.81	1.82	1.77
	0.64	1.58	0.54	1.96	1.85
KMvMDA [21]	19.53	6.69	16.68	7.07	16.99
	18.33	2.05	19.06	3.20	19.33
MNSE	0.25	0.18	0.72	1.49	0.76
	0.40	1.10	1.80	2.64	1.76

that nonlinear representations capture the intricate geometries of these two data sets better than linear representations. MNSE also performs significantly better than the supervised nonlinear methods Kernel GMLDA and KMvMDA, as well as the unsupervised Kernel CCA and DeepMF methods. These observations confirm the efficacy of the principle idea underlying MNSE: explicitly including a generalizability objective via the Lipschitz regularity of the interpolators improves the performance of nonlinear representation learning in data sets with complex geometries.

VI. CONCLUSION

We have first proposed a theoretical analysis of the performance of multi-modal supervised embedding methods in multi-modal classification and cross-modal retrieval applications. The main finding of our performance bounds is that achieving good between-class separation and cross-modal alignment is not sufficient, and the regularity of the multi-modal interpolation functions is also important for ensuring

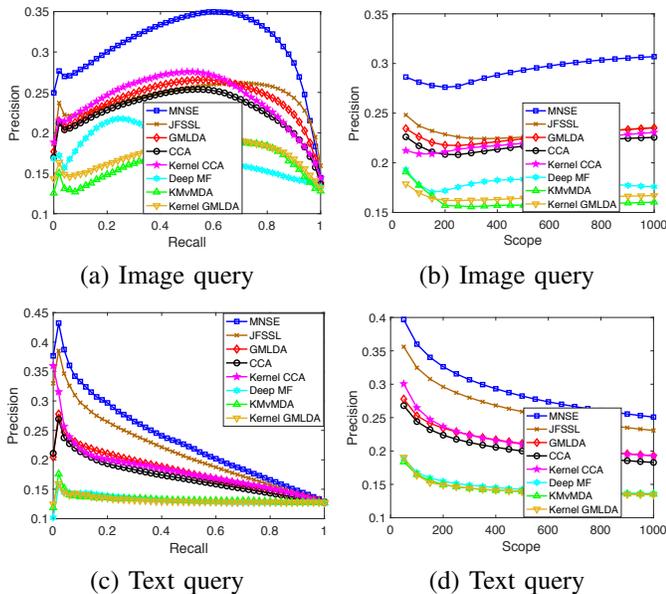


Fig. 4. Retrieval performance of the methods for Wikipedia

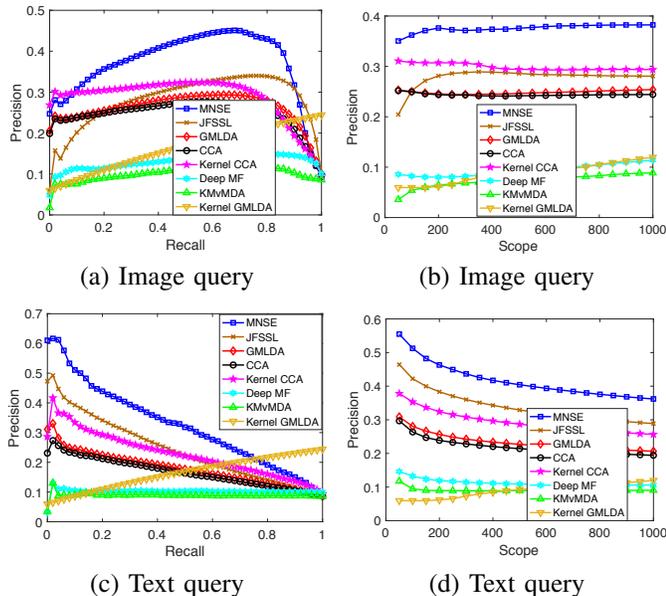


Fig. 5. Retrieval performance of the methods for Pascal VOC2007

good generalization performance. Next, relying on these theoretical findings, we have proposed an algorithm for learning supervised multi-modal nonlinear embeddings, with particular focus on the generalizability of the learnt representations to new test samples. The efficacy of the proposed method has been demonstrated in multi-modal classification and cross-modal retrieval problems, where it has been shown to yield quite satisfactory performance in comparison with recent multi-modal learning algorithms. We hope that our theoretical insights along with our methodological contributions will be useful towards improving the interpretability and the performance of nonlinear representation learning algorithms in multiple domains.

TABLE VIII
MAP SCORES FOR THE WIKIPEDIA AND PASCAL VOC2007 DATA SETS

Algorithm	Wikipedia Image q.	Wikipedia Text q.	Pascal VOC Image q.	Pascal VOC Text q.
CCA [1]	0.2280	0.1720	0.2470	0.1674
Ker. CCA [1]	0.2419	0.1815	0.2873	0.2282
GMLDA [2]	0.2407	0.1815	0.2609	0.1791
Ker. GMLDA [2]	0.1737	0.1326	0.1640	0.1640
JFSSL [25]	0.2440	0.2143	0.2814	0.2418
DeepMF [34]	0.1760	0.1335	0.1305	0.1038
KMvMDA [21]	0.1661	0.1339	0.1023	0.0894
MNSE	0.3109	0.2332	0.3710	0.3221

REFERENCES

- [1] C. Xu, D. Tao, and C. Xu, "A survey on multi-view learning," *arXiv Preprint*, 2013.
- [2] A. Sharma, A. Kumar, H. Daumé, and D. W. Jacobs, "Generalized multiview analysis: A discriminative latent space," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2160–2167.
- [3] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th International Conference on Machine Learning*, 2011, pp. 689–696.
- [4] Y. Wei, Y. Zhao, C. Lu, S. Wei, L. Liu, Z. Zhu, and S. Yan, "Cross-modal retrieval with CNN visual features: A new baseline," *IEEE Trans. Cybernetics*, vol. 47, no. 2, pp. 449–460, 2017.
- [5] F. Feng, X. Wang, and R. Li, "Cross-modal retrieval with correspondence autoencoder," in *Proc. ACM International Conference on Multimedia*, 2014, pp. 7–16.
- [6] A. Blum and T. Mitchell, "Combining labeled and unlabeled data with co-training," in *Proc. 11th Annual Conference on Computational Learning Theory, COLT*, 1998, pp. 92–100.
- [7] S. Dasgupta, M. L. Littman, and D. A. McAllester, "PAC generalization bounds for co-training," in *Adv. Neural Information Processing Systems*, 2001, pp. 375–382.
- [8] D. S. Rosenberg and P. L. Bartlett, "The Rademacher complexity of co-regularized kernel classes," in *Artificial Intelligence and Statistics*, 2007, vol. 2, pp. 396–403.
- [9] S. Sun, "Multi-view laplacian support vector machines," in *Advanced Data Mining and Applications*, 2011, vol. 7121, pp. 209–222.
- [10] U. Brefeld and T. Scheffer, "Co-EM support vector learning," in *Proc. 21st Int. Conf. Mach. Learn.*, 2004.
- [11] Z. Zhou and M. Li, "Semi-supervised regression with co-training," in *Proc. 19th International Joint Conference on Artificial Intelligence*, 2005, pp. 908–913.
- [12] S. Yu, B. Krishnapuram, R. Rosales, and R. B. Rao, "Bayesian co-training," *Journal of Machine Learning Research*, vol. 12, pp. 2649–2680, 2011.
- [13] N. Poh, J. Kittler, and A. Rattani, "Handling session mismatch by fusion-based co-training: An empirical study using face and speech multimodal biometrics," in *IEEE Symposium on Computational Intelligence in Biometrics and Identity Management*, 2014, pp. 81–86.
- [14] X. Duan, N. B. Thomsen, Z. Tan, B. Lindberg, and S. H. Jensen, "Weighted score based fast converging co-training with application to audio-visual person identification," in *IEEE 29th Int. Conf. Tools with Artif. Intel.*, 2017, pp. 610–617.
- [15] N. Rasiwasia, D. Mahajan, V. Mahadevan, and G. Aggarwal, "Cluster canonical correlation analysis," in *Proc. 17th Int. Conf. Artificial Intelligence and Statistics*, 2014, pp. 823–831.
- [16] V. Ranjan, N. Rasiwasia, and C. V. Jawahar, "Multi-label cross-modal retrieval," in *IEEE Int. Conf. Computer Vision*, 2015, pp. 4094–4102.
- [17] Y. Gong, Q. Ke, M. Isard, and S. Lazebnik, "A multi-view embedding space for modeling internet images, tags, and their semantics," *International Journal of Computer Vision*, vol. 106, no. 2, pp. 210–233, 2014.
- [18] M. Kan, S. Shan, H. Zhang, S. Lao, and X. Chen, "Multi-view discriminant analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 188–194, 2016.
- [19] Y. Wei, Y. Zhao, Z. Zhu, S. Wei, Y. Xiao, J. Feng, and S. Yan, "Modality-dependent cross-media retrieval," *ACM Trans. Intell. Syst. Technol.*, vol. 7, no. 4, pp. 57:1–57:13, 2016.
- [20] X. Jing, R. Hu, Y. Zhu, S. Wu, C. Liang, and J. Yang, "Intra-view and inter-view supervised correlation analysis for multi-view feature learning," in *Proc. 28th AAAI Conference on Artificial Intelligence*, 2014, pp. 1882–1889.

- [21] G. Cao, A. Iosifidis, K. Chen, and M. Gabbouj, “Generalized multi-view embedding for visual recognition and cross-modal retrieval,” *IEEE Trans. Cybern.*, vol. 48, no. 9, pp. 2542–2555, 2018.
- [22] S. Sun, X. Xie, and M. Yang, “Multiview uncorrelated discriminant analysis,” *IEEE Trans. Cybern.*, vol. 46, no. 12, pp. 3272–3284, 2016.
- [23] Q. Yin, S. Wu, and L. Wang, “Unified subspace learning for incomplete and unlabeled multi-view data,” *Pattern Recognit.*, vol. 67, pp. 313–327, 2017.
- [24] L. Zhang, B. Ma, G. Li, Q. Huang, and Q. Tian, “Generalized semi-supervised and structured subspace learning for cross-modal retrieval,” *IEEE Trans. Multimed.*, vol. 20, no. 1, pp. 128–141, 2018.
- [25] K. Wang, R. He, L. Wang, W. Wang, and T. Tan, “Joint feature selection and subspace learning for cross-modal retrieval,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 2010–2023, 2016.
- [26] A. Argyriou, M. Herbster, and M. Pontil, “Combining graph laplacians for semi-supervised learning,” in *Adv. Neural Inf. Proc. Sys 18*, 2005, pp. 67–74.
- [27] P. Mercado, F. Tudisco, and M. Hein, “Generalized matrix means for semi-supervised learning with multilayer graphs,” in *Adv. Neur. Inf. Proc. Sys. 32*, 2019, pp. 14848–14857.
- [28] T. Xia, D. Tao, T. Mei, and Y. Zhang, “Multiview spectral embedding,” *IEEE Trans. Systems, Man, and Cybernetics, Part B*, vol. 40, no. 6, pp. 1438–1446, 2010.
- [29] L. Liu, F. Nie, A. Wiliem, Z. Li, T. Zhang, and B. C. Lovell, “Multi-modal joint clustering with application for unsupervised attribute discovery,” *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4345–4356, 2018.
- [30] S. Rastegar, M. S. Baghshah, H. R. Rabiee, and S. M. Shojaei, “MDL-CW: A multimodal deep learning framework with crossweights,” in *Proc. IEEE Conf. Comp. Vis. Pattern Rec.*, 2016, pp. 2601–2609.
- [31] Y. Niu, Z. Lu, J. Wen, T. Xiang, and S. Chang, “Multi-modal multi-scale deep learning for large-scale image annotation,” *IEEE Trans. Image Processing*, vol. 28, no. 4, pp. 1720–1731, 2019.
- [32] C. Zhang, J. Cheng, and Q. Tian, “Multi-view image classification with visual, semantic and view consistency,” *IEEE Trans. Image Processing*, vol. 29, pp. 617–627, 2020.
- [33] F. Yang, J. Chang, C. Tsai, and Y. F. Wang, “A multi-domain and multi-modal representation disentangler for cross-domain image manipulation and classification,” *IEEE Trans. Image Processing*, vol. 29, pp. 2795–2807, 2020.
- [34] H. Zhao, Z. Ding, and Y. Fu, “Multi-view clustering via deep matrix factorization,” in *Proc. Thirty-First AAAI Conf. Artif. Intel.*, 2017, pp. 2921–2927, AAAI Press.
- [35] S. Sun, J. Shawe-Taylor, and L. Mao, “Pac-bayes analysis of multi-view learning,” *Inf. Fusion*, vol. 35, pp. 117–131, 2017.
- [36] C. Xu, D. Tao, and C. Xu, “Multi-view intact space learning,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 12, pp. 2531–2544, 2015.
- [37] E. Vural and C. Guillemot, “A study of the classification of low-dimensional data with supervised manifold learning,” *Journal of Machine Learning Research*, vol. 18, no. 373, pp. 1–55, 2018.
- [38] C. Örnek and E. Vural, “Nonlinear supervised dimensionality reduction via smooth regular embeddings,” *Pattern Recognition*, vol. 87, pp. 55–66, 2019.
- [39] S. Kaya and E. Vural, “Multi-modal learning with generalizable nonlinear dimensionality reduction,” in *IEEE Int. Conf. Image Proc.*, 2019, pp. 2139–2143.
- [40] “MIT-CBCL face recognition database,” Available: <http://cbcl.mit.edu/software-datasets/heisele/facerecognition-database.html>.
- [41] R. Gross, I. A. Matthews, J. F. Cohn, T. Kanade, and S. Baker, “Multiple,” *Image Vis. Comput.*, vol. 28, no. 5, pp. 807–813, 2010.
- [42] Y. Tian, X. Peng, L. Zhao, S. Zhang, and D. N. Metaxas, “CR-GAN: learning complete representations for multi-view generation,” in *Proc. IJCAI*, 2018, pp. 942–948.
- [43] N. Rasiwasia, J. C. Pereira, E. Coviello, G. Doyle, G. Lanckriet, R. Levy, and N. Vasconcelos, “A New Approach to Cross-Modal Multimedia Retrieval,” in *ACM International Conference on Multimedia*, 2010, pp. 251–260.
- [44] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [45] K. Wang, Q. Yin, W. Wang, S. Wu, and L. Wang, “A comprehensive survey on cross-modal retrieval,” *arXiv Preprint*, 2016.
- [46] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results,” <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [47] S. Van Vaerenbergh, “Kernel methods for nonlinear identification, equalization and separation of signals,” *Ph.D. dissertation, University of Cantabria*, 2010.
- [48] R. Herbrich, “Exact tail bounds for binomial distributed variables,” *Online: Available at <http://research.microsoft.com/apps/pubs/default.aspx?id=66854>*, 1999.

APPENDIX

Proof of Theorem 1

Before we prove Theorem 1, we first present the following lemma, whose proof is given after that of Theorem 1.

Lemma 1. *Let the training sample set \mathcal{X} contain at least N_m training samples $\{x_i\}_{i=1}^{N_m}$ from class m , whose observations $\{x_i^{(u)}\}$ with $x_i^{(u)} \sim \nu_m^{(u)}$ are available in all modalities $u = 1, \dots, V$. Assume that the interpolation function $f^{(u)} : H^{(u)} \rightarrow \mathbb{R}^d$ in each modality u is Lipschitz continuous with constant L .*

Let x be a test sample from class m with an observation $x^{(v)}$ given in modality v , drawn with respect to the probability measure $\nu_m^{(v)}$ independently of the training samples. Let $x^{(u)}$ be the observation of the same sample x in an arbitrary modality u , which need not be available to the learning algorithm. For an arbitrary modality $u \in \{1, \dots, V\}$, define $A^{(u)}$ as the set of the training samples from class m within a δ -neighborhood of $x^{(u)}$ in $H^{(u)}$

$$A^{(u)} = \{x_i^{(u)} : x_i \in \mathcal{X}, C(x_i) = m, x_i^{(u)} \in B_\delta(x^{(u)})\}.$$

Assume that for some $Q \geq 1$ and $\delta > 0$, the number of training samples from class m satisfies

$$N_m > \frac{Q}{\eta_{m,\delta}}.$$

Then for any $\epsilon > 0$, with probability at least

$$1 - \exp\left(-\frac{2(N_m\eta_{m,\delta} - Q)^2}{N_m}\right) - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) - (1 - \eta_{m,\delta})^Q,$$

the set $A^{(u)}$ contains at least Q samples, the distance between $f^{(u)}(x^{(u)})$ and the sample mean of the embeddings of its neighboring training samples is bounded as

$$\left\| f^{(u)}(x^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_i^{(u)} \in A^{(u)}} f^{(u)}(x_i^{(u)}) \right\| \leq L\delta + \sqrt{d}\epsilon, \quad (18)$$

and also there is at least one $x_i^{(u)} \in A^{(u)}$ such that its observation $x_i^{(v)}$ in modality v satisfies $\|x_i^{(v)} - x^{(v)}\| \leq \delta$.

The purpose of Lemma 1 is to see how much the embedding of a test sample through a Lipschitz-continuous interpolator is expected to deviate from the average embedding of the training samples surrounding it. Lemma 1 provides a probabilistic upper bound on this deviation, which is used in Theorems 1 and 2 for bounding the classification and retrieval errors. Note that the classification algorithm knows the observation $x^{(v)}$ of the test sample x only in modality v , and classifies it through its embedding $f^{(v)}(x^{(v)})$ with respect to the rule in (1). The entity $x^{(u)}$ in the lemma denotes a hypothetical observation of

x in an arbitrary modality u . Although we conceptually refer to $x^{(u)}$ in the derivations, it is not known to the classification algorithm in practice (unless $u = v$).

We can now prove Theorem 1.

Proof. We first recall from Lemma 1 that for a particular modality $u \in \{1, \dots, V\}$, with probability at least

$$1 - \exp\left(-\frac{2(N_m \eta_{m,\delta} - Q)^2}{N_m}\right) - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) - (1 - \eta_{m,\delta})^Q$$

the set $A^{(u)}$ has $|A^{(u)}| \geq Q$ elements, the inequality in (18) holds, and for at least one $x_i^{(u)} \in A^{(u)}$ we have $\|x_i^{(v)} - x^{(v)}\| \leq \delta$. Since there are V modalities and the probability measures $\nu_m^{(u)}$ are independent, with probability at least

$$1 - \left(\exp\left(-\frac{2(N_m \eta_{m,\delta} - Q)^2}{N_m}\right) + 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) + (1 - \eta_{m,\delta})^Q\right)^V$$

there is at least one modality $u \in \{1, \dots, V\}$ such that all of these three events occur. Now, let $x_i^{(u)}, x_j^{(u)} \in B_\delta(x^{(u)})$ denote two training samples from this modality u and class m . As $\|x_i^{(u)} - x_j^{(u)}\| \leq 2\delta$, from the assumption (P2) on the embedding, we have $\|y_i^{(u)} - y_j^{(u)}\| = \|f^{(u)}(x_i^{(u)}) - f^{(u)}(x_j^{(u)})\| \leq R_\delta$. This gives

$$\begin{aligned} & \left\| f^{(u)}(x_i^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} f^{(u)}(x_j^{(u)}) \right\| \\ &= \left\| \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} \left(f^{(u)}(x_i^{(u)}) - f^{(u)}(x_j^{(u)}) \right) \right\| \\ &\leq \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} \|f^{(u)}(x_i^{(u)}) - f^{(u)}(x_j^{(u)})\| \leq R_\delta. \end{aligned} \quad (19)$$

Then, for any $x_i^{(u)} \in B_\delta(x^{(u)})$, we have

$$\begin{aligned} & \|f^{(u)}(x^{(u)}) - f^{(u)}(x_i^{(u)})\| \\ &= \left\| f^{(u)}(x^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} f^{(u)}(x_j^{(u)}) \right\| \\ &\quad + \left\| \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} f^{(u)}(x_j^{(u)}) - f^{(u)}(x_i^{(u)}) \right\| \\ &\leq \left\| f^{(u)}(x^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} f^{(u)}(x_j^{(u)}) \right\| \\ &\quad + \left\| f^{(u)}(x_i^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_j^{(u)} \in A^{(u)}} f^{(u)}(x_j^{(u)}) \right\| \\ &\leq L\delta + \sqrt{d}\epsilon + R_\delta \end{aligned} \quad (20)$$

where the last inequality follows from (18) and (19).

Now, through the training sample $x_i^{(u)} \in B_\delta(x^{(u)})$ whose observation in modality v satisfies $\|x_i^{(v)} - x^{(v)}\| \leq \delta$, and from

property (P1) of the embedding, we observe that the deviation between the embedding $f^{(v)}(x^{(v)})$ of the observation $x^{(v)}$ of the test sample used by the classification algorithm and the unknown embedding $f^{(u)}(x^{(u)})$ of its unavailable observation $x^{(u)}$ is bounded as

$$\begin{aligned} & \|f^{(v)}(x^{(v)}) - f^{(u)}(x^{(u)})\| \leq \|f^{(v)}(x^{(v)}) - f^{(v)}(x_i^{(v)})\| \\ & \quad + \|f^{(v)}(x_i^{(v)}) - f^{(u)}(x_i^{(u)})\| + \|f^{(u)}(x_i^{(u)}) - f^{(u)}(x^{(u)})\| \\ & \leq L\delta + \eta + L\delta = 2L\delta + \eta \end{aligned}$$

where the second inequality follows from the Lipschitz continuity of the interpolators $f^{(v)}$ and $f^{(u)}$, and (P1). Combining this with (20), we get that for the training samples $x_j^{(u)} \in A^{(u)}$

$$\begin{aligned} & \|f^{(v)}(x^{(v)}) - f^{(u)}(x_j^{(u)})\| \\ & \leq \|f^{(v)}(x^{(v)}) - f^{(u)}(x^{(u)})\| + \|f^{(u)}(x^{(u)}) - f^{(u)}(x_j^{(u)})\| \\ & \leq (2L\delta + \eta) + (L\delta + \sqrt{d}\epsilon + R_\delta) \\ & = 3L\delta + \sqrt{d}\epsilon + R_\delta + \eta. \end{aligned} \quad (21)$$

Next, let $x_k^{(r)}$ be a training sample from another class than m , observed in any view $r = 1, \dots, V$. The distance between the embeddings of $x_k^{(r)}$ and the test sample $x^{(v)}$ is lower bounded as

$$\begin{aligned} & \|f^{(v)}(x^{(v)}) - f^{(r)}(x_k^{(r)})\| \geq \\ & \|f^{(u)}(x_j^{(u)}) - f^{(r)}(x_k^{(r)})\| - \|f^{(v)}(x^{(v)}) - f^{(u)}(x_j^{(u)})\| \\ & > \gamma - (3L\delta + \sqrt{d}\epsilon + R_\delta + \eta) \end{aligned} \quad (22)$$

where the last inequality is obtained from the property (P3) of the embedding and the inequality in (21). Using in (22) the assumption (2) on the embedding, we get

$$\|f^{(v)}(x^{(v)}) - f^{(r)}(x_k^{(r)})\| > 3L\delta + \sqrt{d}\epsilon + R_\delta + \eta. \quad (23)$$

We finally observe from the inequalities (21) and (23) that the embedding of any training sample $x_k^{(r)}$ from another class than m has distance larger than $3L\delta + \sqrt{d}\epsilon + R_\delta + \eta$ to the embedding $f^{(v)}(x^{(v)})$ of the test sample $x^{(v)}$, whereas there are at least Q samples from the same class as $x^{(v)}$ within a distance of at most $3L\delta + \sqrt{d}\epsilon + R_\delta + \eta$. We conclude that the test sample $x^{(v)}$ is then correctly classified via nearest neighbor classification through its embedding $f^{(v)}(x^{(v)})$. \square

Proof of Lemma 1

Proof. For an arbitrary modality $u \in \{1, \dots, V\}$, the observation $x_i^{(u)}$ of a training sample x_i from class m drawn independently from the test sample x lies in a δ -neighborhood of $x^{(u)}$ with probability

$$P\left(x_i^{(u)} \in B_\delta(x^{(u)})\right) = \nu_m^{(u)}\left(B_\delta(x^{(u)})\right) \geq \eta_{m,\delta}.$$

Then, the probability that $B_\delta(x^{(u)})$ contains at least Q samples among the N_m training samples drawn from $\nu_m^{(u)}$ is given by

$$\begin{aligned} & P(|A^{(u)}| \geq Q) \\ &= \sum_{k=Q}^{N_m} \binom{N_m}{k} \left(\nu_m^{(u)}(B_\delta(x^{(u)})) \right)^k \left(1 - \nu_m^{(u)}(B_\delta(x^{(u)})) \right)^{N_m-k} \\ &\geq \sum_{k=Q}^{N_m} \binom{N_m}{k} (\eta_{m,\delta})^k (1 - \eta_{m,\delta})^{N_m-k}. \end{aligned}$$

This is obtained by evaluating the probability that at least Q successes occur within N_m independent Bernoulli trials with success probability more than $\eta_{m,\delta}$ in each trial. Following the approach in the proof of [37, Theorem 5], from the assumption $N_m > \frac{Q}{\eta_{m,\delta}}$, we can lower bound this probability using a tail bound for distributions [48]. We thus get

$$P(|A^{(u)}| \geq Q) \geq 1 - \exp\left(-\frac{2(N_m\eta_{m,\delta} - Q)^2}{N_m}\right).$$

Now assume that the event $|A^{(u)}| \geq Q$ has occurred for the modality u , i.e., there are at least Q training samples from class m within a δ -neighborhood of $x^{(u)}$. Then, from [37, Lemma 3], with probability at least

$$1 - 2d \exp\left(-\frac{|A^{(u)}|\epsilon^2}{2L^2\delta^2}\right) \geq 1 - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right)$$

the distance between $f^{(u)}(x^{(u)})$ and the sample average of the embeddings of its neighboring training samples is bounded as

$$\left\| f^{(u)}(x^{(u)}) - \frac{1}{|A^{(u)}|} \sum_{x_i^{(u)} \in A^{(u)}} f^{(u)}(x_i^{(u)}) \right\| \leq L\delta + \sqrt{d}\epsilon. \quad (24)$$

Next, still assuming that the event $|A^{(u)}| \geq Q$ has occurred for the modality u , for each sample $x_i^{(u)} \in A^{(u)}$, the probability that its observation $x_i^{(v)}$ in modality v is outside $B_\delta(x^{(v)})$ is

$$1 - \nu_m^{(v)}(B_\delta(x^{(v)})) \leq 1 - \eta_{m,\delta}.$$

Therefore, with probability at least $1 - (1 - \eta_{m,\delta})^Q$, there is at least one $x_l^{(u)} \in B_\delta(x^{(u)})$ whose observation in modality v satisfies $x_l^{(v)} \in B_\delta(x^{(v)})$, or equivalently, $\|x_l^{(v)} - x^{(v)}\| \leq \delta$. Combining the probability expressions we obtained so far, we conclude that for an arbitrary modality q , with probability at least

$$\begin{aligned} & 1 - \exp\left(-\frac{2(N_m\eta_{m,\delta} - Q)^2}{N_m}\right) - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) \\ & - (1 - \eta_{m,\delta})^Q \end{aligned}$$

we have $|A^{(u)}| \geq Q$, the event in (24) occurs, and there is at least one $x_l^{(u)} \in B_\delta(x^{(u)})$ such that $\|x_l^{(v)} - x^{(v)}\| \leq \delta$. \square

Proof of Theorem 2

Proof. Recall from Lemma 1 that with probability at least

$$\begin{aligned} & 1 - \exp\left(-\frac{2(N_m\eta_{m,\delta} - Q)^2}{N_m}\right) - 2d \exp\left(-\frac{Q\epsilon^2}{2L^2\delta^2}\right) \\ & - (1 - \eta_{m,\delta})^Q \end{aligned}$$

the embedding $f^{(u)}(x^{(u)})$ of the query sample in modality u has $|A^{(u)}| \geq Q$ neighboring training samples from the same class m , the inequality in (18) holds, and for at least one $x_l^{(u)} \in A^{(u)}$ we have $\|x_l^{(v)} - x^{(v)}\| \leq \delta$. Assuming that all these three events have occurred and following the same steps as in the proof of Theorem 1, we conclude that there are at least Q samples $x_j^{(u)} \in B_\delta(x^{(u)})$ such that

$$\|f^{(v)}(x^{(v)}) - f^{(u)}(x_j^{(u)})\| \leq 3L\delta + \sqrt{d}\epsilon + R_\delta + \eta \quad (25)$$

while the distance of the embedding of any training sample $x_k^{(r)}$ from another class to $f^{(v)}(x^{(v)})$ is lower bounded as

$$\|f^{(v)}(x^{(v)}) - f^{(r)}(x_k^{(r)})\| > 3L\delta + \sqrt{d}\epsilon + R_\delta + \eta. \quad (26)$$

This implies that the Q samples of smallest distance to the embedding $f^{(v)}(x^{(v)})$ of the query sample are all from class m . Hence, for $K \leq Q$, the precision rate over the K nearest neighbors is

$$P = K/K = 1.$$

Similarly, when $K > Q$, at least Q of the K nearest neighbors of $f^{(v)}(x^{(v)})$ are from the same class m , hence we get

$$P \geq Q/K.$$

Meanwhile, when $K \leq Q$, the retrieval algorithm returns K samples out of the N_m training samples from the same class m , hence

$$R = \frac{K}{N_m}.$$

Finally, when $K > Q$, since at least Q of the retrieved training samples will be from the same class as the query sample, we have

$$R \geq \frac{Q}{N_m}.$$

\square