

REGULARIZED CONDITIONAL OPTIMAL TRANSPORT FOR FEATURE LEARNING AND GENERALIZATION BOUNDS

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper develops the regularized conditional optimal transport for feature learning in an embedding space. Instead of using joint distributions of data, we introduce conditional distributions to some reference conditional distributions in terms of the Kullback-Leibler (KL) divergence. Using conditional distributions provides the flexibility in controlling the transferring range of given data points. When the alternating optimization technique is employed to solve our model, it is interesting to find that conditional and marginal distributions have closed-form solutions. Moreover, the use of conditional distributions facilitates the derivation of the generalization bound of our model via the Rademacher complexity, which characterizes its convergence speed in terms of the number of samples. By optimizing the anchors (centroids) defined in the model, we also employ optimal transport and autoencoders to explore an embedding space of samples in the clustering problem. In the experimental part, we demonstrate that the proposed model achieves promising performance on some learning tasks. Moreover, we construct a conditional Wasserstein classifier to classify set-valued objects.

1 INTRODUCTION

Achieving effective features of data (Boroujeni et al., 2018; Wang et al., 2019; Qian et al., 2019) is a fundamental task in data analysis, and feature learning has been explored in some fields such as machine learning and computer vision. Feature learning aims at exploring a linear or nonlinear transformation to map the original features into an embedding space by optimizing the defined objective function. In the latent representation space, data can be explored, thereby providing some benefits from various learning tasks (Su & Hua, 2019).

Earlier feature learning algorithms focus on how to develop effective handcrafted extractors for visualizing high-dimensional data and reducing the effect of the curse of dimensionality. Marginal Fisher analysis (Yan et al., 2007) adopts a graph embedding framework to provide an intrinsic graph with intra-class compactness and a penalty graph with inter-class separability. Max-min distance analysis (Bian & Tao, 2011) achieves the low-dimensional data by maximizing the minimum pairwise distance. A robust linear discriminant analysis method based on the $L_{2,1}$ norm (Nie et al., 2021) is developed to obtain robust projection features, and an effective iterative optimization algorithm is derived to solve a general ratio minimization problem. In Flamary et al. (2018), Wasserstein discriminant analysis from optimal transport (Li et al., 2021; Serrurier et al., 2021) is implemented by employing the regularized Wasserstein distance to capture the global and local interactions between classes.

Kernel-based methods that capture the nonlinear features of data have been developed to search for an effective feature space by selecting proper kernel functions. Unlike classical dimensionality reduction methods, the embedding space of data may be an infinite-dimensional feature space since data may be well separated in high-dimensional spaces. Kernel principal component analysis (PCA) and kernel linear discriminant analysis(LDA) are two effective methods for achieving effective features of data. To address the outliers of data, L_1 norm kernel LDA (Zheng et al., 2014) is developed to achieve the nonlinear discriminant features of data. In unsupervised learning, finding effective features contributes to the improvements in the performance of clustering. The classical k-means

method is extended to the kernel k-means method in terms of the kernel trick. To capture multiple feature representations of data, an effective strategy in multiple k-means clustering problems (Yao et al., 2021) is adopted to select the optimal kernel from the prespecified kernels, and an alternating minimization method is used to update the coefficients of the kernels and the cluster membership alternatively. Multiple kernel k-means clustering methods with incomplete kernel matrices (MKCIK) (Liu et al., 2020) embed imputation and clustering into a unified learning framework. One remarkable characteristic of MKCIK is that a complete base kernel matrix over all the samples is not required.

Exploring the local and relevant information of data points is helpful for achieving discriminant features of data (Nie et al., 2022; 2023). For each data point, the conditional distribution of the data point can characterize its local and relevant information. Figure 1 shows that there are three data points in the X space and nine data points in the Y space, where each data point in the X space is relevant to four data points in the Y space in terms of an appropriate structure such as proximity and topology. In supervised learning, data points with the same color in the Y space belong to the same class. When the labels of samples are available, in the Y space, there are four data points whose labels are the same as the label of x_1 , two data points whose labels are the same as the label of x_2 , and three data points whose labels are the same as the label of x_3 . It is clear that the conditional distributions constructed by considering the label information of data points are different from those in unsupervised learning. For data points in the X space, we can obtain their Dirac measures. Thus, we can explore the Wasserstein distance between Dirac measures and conditional measures¹ on two spaces². The Wasserstein distance from optimal transport (Liu et al., 2023; Fatras et al., 2021) can be used to describe the relationship between two probability measures, and autoencoders can explore the latent space of data. Hence, we employ the optimal transport and autoencoders to show how to transport information in an embedding space, which gives a novel framework for learning effective features of data via optimal transport. For each data point, we employ the conditional probability to constrain its transferring range. The merit of using the conditional probability is that varying neighbors of different data points can be explored. To reserve the information of data, we impose the reconstruction error of data on the objective function. In addition, we discuss the properties of our model and extend our model to the clustering problem. Finally, we perform the experiments on a series of data sets. The main contributions of this paper are listed as follows.

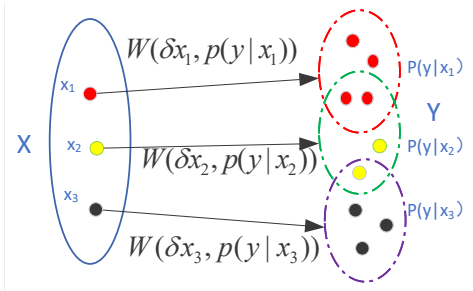


Figure 1: A simple example of conditional optimal transport where the conditional probability is used to characterize the local information of given data points. Each data point in X can be modeled as a Dirac(point) measure. Data points in the same color in Y are taken from the same class, and different conditional distributions can be constructed for unsupervised and supervised learning. $W(\delta_{x_1}, p(y|x_1))$ denotes the Wasserstein distance between δ_{x_1} and $p(y|x_1)$.

- We propose a regularized conditional optimal transport framework for extracting the effective and useful features of data. In this framework, we employ conditional distributions to capture the local behaviors of given data points and use the Kullback-Leibler divergence for conditional distributions, which can consider prior knowledge of conditional distributions.
- We apply the alternating optimization technique to tackle the proposed model. It is noted that marginal and conditional distributions have closed-form solutions. Moreover, we derive the generalization bound of our model in terms of the Rademacher complexity and generalize our model to find anchors in the embedding space, which is available for the clustering problem.
- We perform a series of experiments on some classification and clustering problems to demonstrate the effectiveness of our model. Moreover, we discuss how to modify our mod-

¹The conditional probability of data points outside the transferring range is zero.

²These two spaces may be the same.

el to make it a classifier that can be used to classify set-valued objects, and this classifier degenerates into the deep nearest-neighbor classifier.

2 REGULARIZED CONDITIONAL OPTIMAL TRANSPORT FOR FEATURE LEARNING

2.1 PRELIMINARIES

Let two random vectors $X \in R^m$ and $Y \in R^m$ be taken from two probability spaces (\mathbb{X}, μ) and (\mathbb{Y}, ν) . Assume that Z is the latent space. The L_r norm of a vector $a = (a_1, \dots, a_m)$ is denoted by $\|a\|_r = \sqrt[r]{\sum_{i=1}^m |a_i|^r}$. For measures μ and ν corresponding to X and Y , the Wasserstein distance with the order \bar{r} is defined in the following (Courty et al., 2017; Lin & Chan, 2023):

$$W^{\bar{r}}(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{X} \times \mathbb{Y}} \rho(x, y)^{\bar{r}} d\pi(x, y) \quad (1)$$

where $\Pi(\mu, \nu)$ is the set of probability measures on \mathbb{X} and \mathbb{Y} with marginal measures μ and ν , and $\rho(x, y)$ denotes the distance between $x \in R^m$ and $y \in R^m$. $W^{\bar{r}}(\mu, \nu)$ is the potential cost of moving mass from μ and ν , and the optimal solution provides the optimal transport plan. In real applications, we usually obtain some sampled points in terms of probability measures μ and ν . That is, μ and ν are two discrete measures with a finite number of support points. Thus, $\mu = \sum_{i=1}^n a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^k b_j \delta_{y_j}$, where δ_{x_i} denotes the Dirac measure at the point x_i , and $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_k)$ are vectors in the probability simplex. Assume that $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_k\}$ are sampled data points from μ and ν , respectively. To effectively solve (1), a regularization term is introduced and its discrete version is formulated:

$$\inf_{p_{ij} \in U(a, b)} \sum_{i=1}^n \sum_{j=1}^k \rho(x_i, y_j)^{\bar{r}} p_{ij} - \lambda H(p), \quad (2)$$

where $H(p) = -\sum_{i=1}^n \sum_{j=1}^k p_{ij} (\ln p_{ij} - 1)$ is the information entropy of p , $U(a, b) = \{p_{ij} | p1_k = a, p^T 1_n = b\}$, and 1_k is a k -dimensional vector whose elements are 1. The famous Sinkhorn algorithm (Cuturi, 2013) can be employed to achieve the optimal transport plan with a faster computation.

2.2 PROBLEM FORMULATION

As shown in Figure 1, each data point in a space may locally or semantically correlate with many data points in a space, and conditional distributions can characterize the information of given data points. The theory of optimal transport provides a possible scheme for the movement of data points. Autoencoders facilitate feature formations in an embedding space. For autoencoders, let $f_\theta(x) \in R^d$ be an encoder with parameter θ and its decoder be $f_{\bar{\theta}}(z)$ with parameter $\bar{\theta}$. The functions $g_\phi(y)$ and $\bar{g}_{\bar{\phi}}(z)$ consist of another autoencoder. In encoded spaces, we obtain the Dirac measure at $f_\theta(x_i)$, denoted by $\delta_{f_\theta(x_i)}$. We employ the conditional distribution $p(g_\phi(y)|x_i)$ to characterize the information of y relating to x_i , and $p(g_\phi(y)|x_i)$ can be considered as a push-forward measure induced by ϕ . Data points in encoded spaces can be transported even if X and Y belong to different spaces. To facilitate the learning of the conditional distribution, we adopt a convex combination of Dirac measures to construct $p(g_\phi(y)|x_i)$. That is, y takes k values and $p(g_\phi(y)|x_i) = \sum_{j=1}^k p_{j|i} \delta_{g_\phi(y_{j|i})}$, where $p_{j|i}$ is a nonnegative coefficient that satisfies $\sum_{j=1}^k p_{j|i} = 1$. $y_{j|i}$ may be semantically relevant to x_i . This also ensures that $p(g_\phi(y)|x_i)$ belongs to the Wasserstein space. Thus, the r th-order Wasserstein distance between $\delta_{f_\theta(x_i)}$ and $p(g_\phi(y)|x_i)$ can be achieved, denoted by $\bar{W}_i^{\bar{r}} = \sum_{j=1}^k \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_{j|i}$. To effectively learn $p_{j|i}$ in the conditional distributions, we define the regularized conditional optimal transport for the data point $f_\theta(x_i)$, denoted by

$$\min_{p_{j|i}} L_i := \sum_{j=1}^k \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_{j|i} + \lambda_1 \text{KL}(p_{\cdot|i} || q_{\cdot|i}), \quad (3)$$

where $\text{KL}(p_{\cdot|i} || q_{\cdot|i}) = \sum_{j=1}^k p_{j|i} \ln \frac{p_{j|i}}{q_{j|i}}$, $\sum_{j=1}^k q_{j|i} = 1$, $q_{j|i}$ is the prior probability of transferring x_i to $y_{j|i}$ in the original space, and $y_{j|i}$ is the j th data point determined by x_i . The Kullback-Leibler

(KL) divergence (Bishop, 2007; Zhang et al., 2024) is employed to measure the difference between $\{p_{j|i}\}$ and $\{q_{j|i}\}$. Since we introduce the conditional probability of x_i , we can control the transferring range of x_i . If x_i is not allowed to be moved to $y_{j|i}$, then we set $q_{j|i} = 0$. The nonnegative hyperparameter λ_1 controls the tradeoff between the transport cost and the KL divergence. The variables $p_{j|i}$ ($j = 1, \dots, k$) need to be optimized, and they implicitly depend on the embedding spaces. $q_{j|i}$ is the prior conditional probability independent of embedding spaces. Interestingly, optimizing (3) can also be regarded as a proximal algorithm to obtain the proximal operator (Li et al., 2023; Gu et al., 2024; Zhang et al., 2024). Unlike those proximal operators we explore the conditional distribution in a discrete form. In fact, one may replace the KL term in (3) with the f-divergence (Zhang et al., 2024) between two distributions, which increases the flexibility of the model. For computational convenience, we apply the KL divergence in (3). To explore the transport cost of all data points, we define the following model:

$$\min_{(\theta, \phi, p_{j|i}, p_i)} \hat{L}_w := \sum_{i=1}^n L_i p_i + \lambda_2 \text{KL}(p||q), \quad (4)$$

where $p_i = p(x_i)$, $\sum_{i=1}^n p_i = 1$, $\text{KL}(p||q) = \sum_{i=1}^n p_i \ln \frac{p_i}{q_i}$, q_i is the prior probability of x_i independent of embedding spaces, and λ_2 is a nonnegative hyperparameter.

The first term in (4) denotes the transport cost, the second term is the KL divergence between $\{p_i\}$ and $\{q_i\}$. If $\{x_i\}$ are sampled from the uniform distribution, i.e. $p_i = 1/n$, we let $\lambda_2 = 0$ since the KL term is constant. The introduction of $q_{j|i}$ and q_i helps us use prior knowledge of data from the original space. If no prior knowledge of data is available, $q_{j|i}$ and q_i may take the uniform distribution. Here, we take student's t distribution as the probability of moving from x_i to $y_{j|i}$ (Xie et al., 2016) in the original space, denoted by

$$q_{j|i} = \frac{(1 + \rho(x_i, y_{j|i})^{\bar{r}})^{-1}}{\sum_{j=1}^k (1 + \rho(x_i, y_{j|i})^{\bar{r}})^{-1}}. \quad (5)$$

Note that $q_{j|i}$ depends on the original features instead of the embedding features. Generally, $q_{j|i}$ reflects information in the original space, but $p_{j|i}$ can be learned in the embedding space. The KL divergence does not meet the triangle inequality, so it is not a true distance measure. The KL divergence is not symmetric since $\text{KL}(p|q) \neq \text{KL}(q|p)$.

Unlike the Wasserstein distance from the optimal transport theory, we decompose the joint distribution into the product of two distributions, i.e. $p_{i,j} = p_i p_{j|i}$. Moreover, we utilize the conditional KL divergence as the regularization term by introducing prior conditional probabilities of data points. Note that trivial solutions of θ and ϕ may be obtained if we do not impose additional constraints on encoders. In order to address this problem, we add the reconstruction error of data to the objective function by using decoders. Thus, we define the following model:

$$\min_{(\theta, \bar{\theta}, \phi, \bar{\phi}, p_{j|i}, p_i)} \hat{L} := \hat{L}_w + \lambda_3 \sum_{i=1}^n \|x_i - \bar{f}_{\bar{\theta}} f_{\theta}(x_i)\|_2 p_i + \lambda_4 \sum_{i,j=1}^{n,k} p_i p_{j|i} \|y_{j|i} - \bar{g}_{\bar{\phi}} g_{\phi}(y_{j|i})\|_2, \quad (6)$$

where λ_i ($i = 3, 4$) are nonnegative hyperparameters. The last two terms in (6) involve the reconstruction errors of x_i and $y_{j|i}$. The continuous version of (6) can be found in appendixes. From (6), we find that the loss function in the proposed model consists of the transport cost, reconstruction errors of data and additional regularization terms. The framework is generic since we do not give specific autoencoders and any transport cost can be used to replace $\rho(\cdot)$. Note that in the above model, we assume that $\{x_i\}$ and $\{y_{j|i}\}$ adopt different encoders and decoders. In fact, when $\{x_i\}$ and $\{y_{j|i}\}$ are sampled from the same data source, we can take the same encoders and decoders. **In this paper, we only consider that $\{x_i\}$ and $\{y_{j|i}\}$ take the same encoders and decoders, but we reserve more general notations for future extensions of our framework for different dimensions of features from two data sources.** Since we consider the conditional distribution of x_i , we use it to describe the local information of x_i . That is, $y_{j|i}$ ($j = 1, \dots, k$) are taken from the k neighbors of x_i . In supervised learning, we allow $y_{j|i}$ to be taken from the samples whose labels are the same as the label of x_i . If $\{x_i\}$ and $\{y_{j|i}\}$ are taken from the same data source, we let $f_{\theta} = g_{\phi}$, $\bar{f}_{\bar{\theta}} = \bar{g}_{\bar{\phi}}$, and $\lambda_4 = 0$.

2.3 OPTIMIZATION

Note that there are several groups of parameters to be optimized in our model. Moreover, some parameters such as the conditional probability have additional constraints. Thus, the model of (6) is a constrained and non-convex optimization problem. To solve our model, we resort to the alternating optimization technique. Specifically, we alternatively optimize a group of variables by fixing other groups of optimization variables. In the following, we will demonstrate how to divide these variables into several groups and how to optimize them.

(a): Update $p_{j|i}$ by fixing other variables. In this step, when we fix $\theta, \bar{\theta}, \phi, \bar{\phi}$, and p_i , we solve the following model:

$$\min_{p_{j|i}} \sum_{i,j=1}^{n,k} \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_i p_{j|i} + \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} \| q_{\cdot|i}) + \lambda_4 \sum_{i,j=1}^{n,k} p_i p_{j|i} \|y_{j|i} - \bar{g}_{\bar{\phi}} g_\phi(y_{j|i})\|_2. \quad (7)$$

It is noted that (7) is a strongly convex optimization problem. Hence, it has a unique solution. It is of interest to note that it has a closed-form solution, denoted by

$$p_{j|i} = \frac{q_{j|i} \exp(-(L_{j|i}^{op} + L_{j|i}^{re})/\lambda_1)}{\sum_{j=1}^k q_{j|i} \exp(-(L_{j|i}^{op} + L_{j|i}^{re})/\lambda_1)}, \quad (8)$$

where $L_{j|i}^{op} = \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}}$ and $L_{j|i}^{re} = \lambda_4 \|y_{j|i} - \bar{g}_{\bar{\phi}} g_\phi(y_{j|i})\|_2$.

(b): Update p_i by fixing other variables. Given $\theta, \bar{\theta}, \phi, \bar{\phi}$, and $p_{j|i}$, we achieve p_i by solving the following problem:

$$\min_{p_i} \sum_{i,j=1}^{n,k} \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_i p_{j|i} + \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} \| q_{\cdot|i}) + \lambda_2 \text{KL}(p \| q) + \lambda_3 \sum_{i=1}^n \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 p_i + \lambda_4 \sum_{i,j=1}^{n,k} p_i p_{j|i} \|y_{j|i} - \bar{g}_{\bar{\phi}} g_\phi(y_{j|i})\|_2. \quad (9)$$

It is observed that the objective function in (9) is strongly convex. Thus there exists a unique solution of p_i . The closed-form solution is denoted by

$$p_i = \frac{q_i \exp(-(L_i^{op} + L_i^{enre})/\lambda_2)}{\sum_{i=1}^n q_i \exp(-(L_i^{op} + L_i^{enre})/\lambda_2)}, \quad (10)$$

where $L_i^{op} = \sum_{j=1}^k \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_{j|i}$ and $L_i^{enre} = \lambda_4 \sum_{j=1}^k p_{j|i} \|y_{j|i} - \bar{g}_{\bar{\phi}} g_\phi(y_{j|i})\|_2 + \lambda_3 \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 + \text{KL}(p_{\cdot|i} \| q_{\cdot|i})$.

(c): Update $\theta, \bar{\theta}, \phi, \bar{\phi}$ by fixing other variables. In this step, we try to learn the parameters of autoencoders. Specifically, we solve the following optimization problem:

$$\min_{(\theta, \bar{\theta}, \phi, \bar{\phi})} \sum_{i,j=1}^{n,k} \rho(f_\theta(x_i), g_\phi(y_{j|i}))^{\bar{r}} p_i p_{j|i} + \lambda_3 \sum_{i=1}^n p_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 + \lambda_4 \sum_{i,j=1}^{n,k} p_i p_{j|i} \|y_{j|i} - \bar{g}_{\bar{\phi}} g_\phi(y_{j|i})\|_2. \quad (11)$$

Note that the objective function in (11) is nonconvex. We cannot obtain the global optimal solution. We generally update these parameters of models through the chain rule in the framework of neural networks. In this work, we resort to automatic differentiation to learn these parameters.

For completeness, we summarize the main steps of solving the proposed model in Algorithm 1. It is found that step 2.1 involves the computational complexity of $O(H_1^2 H_2 n)$ in each iteration, step 2.2 involves $O(nk(m+d))$ and step 2.3 is $O(n(m+d))$, where H_1 is the maximum number of hidden units of layers and H_2 is the number of layers. In addition, the convergence of Algorithm 1 comes from the fact that it belongs to the block coordinate descend method (Razaviyayn et al., 2012).

Algorithm 1: Optimization algorithm to (6)

- 1: Given $\lambda_i, q_{j|i}, q_i$, and initialize $p_i = q_i, p_{j|i} = q_{j|i}$
 - 2: **For** $t=1$ to T **do**
 - 2.1: solve (11) to achieve its parameters $(\theta, \bar{\theta}, \phi, \bar{\phi})$;
 - 2.2: solve (7) to achieve $p_{j|i}$;
 - 2.3: solve (9) to achieve p_i ;
 - 3: **Output**: the encoders and decoders.
-

2.4 THEORETICAL ANALYSIS OF OUR MODEL

In this subsection, we theoretically analyze some properties of our model. There are several parameters in our models. We observe that $\lim_{\lambda_1 \rightarrow +\infty} p_{j|i} = q_{j|i}$ and $\lim_{\lambda_2 \rightarrow +\infty} p_i = q_i$ if $p_{j|i}$ and p_i are defined in (8) and (10). This indicates that if parameters λ_1 and λ_2 approach the positive infinity, $p_{j|i}$ and p_i will have the same distributions as prior distributions. If prior distributions are uniform distributions, the optimal transport plan will be uniform distributions. In such a case, the objective function of our model makes the trade-off between the reconstruction error and the transport cost. Note that when deriving the generalization bound of our model, we do not consider the expectation with respect to the random variable Y . Here we assume that Y has the support consisting of k data points. For given x_i , we need to find k data points $y_{1|i}, \dots, y_{k|i}$. These k data points are varying for different x_i . Evidently, it is different from the fixed sampled points y_1, \dots, y_k in the optimal transport theory. To explore the effect of the parameters of networks, we study the generalization bound of our model based on the assumption that $S = \{x_1, \dots, x_n\}$ are independent and identically distributed samples, i.e., $p_i = \frac{1}{n}$. First we define the empirical loss as done in Maurer & Pontil (2010) when λ_4 takes the zero value.

$$\hat{L}_S(\theta, \bar{\theta}) := \min_{p_{j|i}} \frac{1}{n} \left\{ \sum_{i,j=1}^{n,k} p_{j|i} \rho(f_\theta(x_i), f_{\bar{\theta}}(y_{j|i}))^{\bar{r}} + \lambda_1 \sum_{i=1}^n \text{KL}(p_{\cdot|i} \| q_{\cdot|i}) + \sum_{i=1}^n \lambda_3 \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 \right\}. \quad (12)$$

Note that there are several differences between Equations (6) and (12). Here, we let $f_\theta = g_\phi$ and $\bar{f}_{\bar{\theta}} = \bar{g}_{\bar{\phi}}$. Moreover, we do not consider the reconstruction error of $y_{1|i}, \dots, y_{k|i}$ since they are taken from the space that x_i belongs to. In addition, we assume that p_i in (12) is taken from the uniform distribution. Let $L(\theta, \bar{\theta})$ be the expected loss corresponding to (12). We make the following assumptions.

A0: the distance measure has the form of $\rho(x, y) = \varphi(x - y)$ and $\varphi(x)$ has the Lipschitz constant ℓ ;

A1: x_i and $y_{j|i}$ are bounded, i.e., $\exists M$ such that $\|x_i\|_2 \leq M$ and $\|y_{j|i}\|_2 \leq M$;

A2: $\|\bar{f}_{\bar{\theta}}\|_2 \leq M$ and $\|f_\theta\|_2 \leq M$ hold for parameters $\bar{\theta}$ and θ in a parameter space;

A3: if $q_{j|i} = 0$, $p_{j|i} = 0$.

The assumption A0 holds if the metric is induced by the norm in a normed space and the data are taken from a compact space. For example, $\rho(x, y)$ takes the form of the L_r norm. The assumptions A1 and A2 are reasonable since the data we deal with are bounded. The assumption A3 ensures that the KL divergence is well defined. Now we show the uniform deviation bound of the objective function in (12) by using the following theorem.

Theorem 1. Under the above assumptions, with probability at least $1 - \tau$, the following inequality holds for θ and $\bar{\theta}$ in proper parameter spaces:

$$\hat{L}_S(\theta, \bar{\theta}) \leq L(\theta, \bar{\theta}) + 4\sqrt{2}M_1R_1 + 2\sqrt{2}R_2 + \chi_1 \sqrt{\frac{-\log \tau}{2n}} \quad (13)$$

where $M_1 = \bar{r}(2M\ell)^{\bar{r}-1}\ell$, $\chi_1 = \frac{2(2M)^{\bar{r}} + 4\lambda_3 M}{n}$, $R_1 = E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d |\sum_{i=1}^n \sigma_{it}(f_\theta(x_i))_t|$, $R_2 = E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m |\sum_{i=1}^n \sigma_{it}(f_{\bar{\theta}} f_\theta(x_i))_t|$, σ_{it} denotes the Rademacher random variable, E_S denotes the expectation with respect to S , and $f_\theta(x_i)_t$ denotes the t -th element of the vector $f_\theta(x_i)$.

R_1 denotes the Rademacher complexity of the encoder $f_\theta(\cdot)$, and R_2 denotes the Rademacher complexity of the encoder-decoder $\bar{f}_{\bar{\theta}} f_\theta(\cdot)$. In the case of a single-layer linear network, if the parameters of the network satisfy $\theta^T \theta = I_d$ and $\bar{\theta} = \theta^T$, then we have $R_1 \leq kdM/\sqrt{n}$ and $R_2 \leq dM/\sqrt{n}$. It has been proved in Truong (2019) that the Rademacher complexity of deep learning models is of order $O(1/\sqrt{n})$ under proper conditions. Thus, R_1 and R_2 have the order of $O(1/\sqrt{n})$.

2.5 EXTENSIONS TO THE CLUSTERING PROBLEM

In the above section, we assume that x_i is transported to data points $y_{1|i}, \dots, y_{k|i}$. These data points are taken from the class of x_i or from k -neighbors of x_i . This implicitly uses prior knowledge from the original data. **Without using prior knowledge, are they learned from**

the data via some optimization methods? This may be a trivial thing since the number of $\{y_{j|i} | i = 1, \dots, n, j = 1, \dots, k\}$ is much bigger than that of $\{x_i | i = 1, \dots, n\}$ as shown in Figure 1. To avoid triviality, we can impose additional constraints on $\{y_{j|i}\}$ to reduce the number of $\{y_{j|i}\}$. In supervised learning, we may consider that the data points in the same class are transported to unknown data points (anchors). That is, $y_{j|s} = y_{j|t}$ if x_s and x_t are from the same class. In unsupervised learning where the labels of samples are not available, we may consider the case where all the data points $\{x_i | i = 1, \dots, n\}$ are transported to unknown data points $\{y_j | j = 1, \dots, k\}$, i.e. $y_j = y_{j|1} = y_{j|2} = \dots = y_{j|n}$. Thus, the conditional distribution is denoted by $p(g_\phi(y)|x_i) = \sum_{j=1}^k p_{j|i} \delta_{g_\phi(y_j)}$, where $p_{j|i}$ and y_j need to be learned. Instead of finding $\{y_j | j = 1, \dots, k\}$ in the original space, we explore unknown data points in an embedding space and let $z_j = g_\phi(y_j)$ ($j = 1, \dots, k$). Since we directly look for $\{z_j\}$ in the embedding space, we do not need to consider the encoder g_ϕ and the decoder \bar{g}_ϕ . Thus, the following model is formulated to learn $\{z_j\}$ in an embedding space of data in an unsupervised way.

$$\begin{aligned} \min_{(\theta, \bar{\theta}, p_{j|i}, p_i, z_j)} \hat{L} := & \sum_{i,j=1}^{n,k} p_i \rho(f_\theta(x_i), z_j)^{\bar{r}} p_{j|i} + \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} || q_{\cdot|i}) + \\ & \sum_{i=1}^n \lambda_3 p_i \|x_i - \bar{f}_\theta f_\theta(x_i)\|_2 + \lambda_2 \text{KL}(p || q). \end{aligned} \quad (14)$$

Note that z_1, \dots, z_k are optimization variables in an embedding space. We refer to z_1, \dots, z_k as anchors. These anchors can also be taken as the cluster centroids of data in the embedding space if k is equal to the number of clusters. In such a case, the conditional probability $p_{j|i}$ can be regarded as the probability of x_i closing to z_j . We also employ the alternating optimization method to solve (14), which can be found in appendixes. The conditional probability $p_{j|i}$ and marginal probability p_i have closed-form solutions in each step. In such a case, we can learn the anchors (centroids) in the embedding space by using autoencoders. The main aim of designing our model of (14) is to obtain features in the embedding space in an unsupervised way. Here, we employ (14) to learn the embedding space of data and perform the possible clustering tasks in the embedding space. In fact, pretrained autoencoders may be employed to initialize the weights of autoencoders. When data points are independent and identically distributed, we can explore the generalization bound of our model of (14). To this end, we define the following empirical loss.

$$\begin{aligned} \hat{L}_S^c(\theta, \bar{\theta}, z_j) := & \min_{p_{j|i}} \sum_{i,j=1}^{n,k} \frac{1}{n} \rho(f_\theta(x_i) - z_j)^{\bar{r}} p_{j|i} \\ & + \frac{\lambda_1}{n} \sum_{i=1}^n \text{KL}(p_{\cdot|i} || q_{\cdot|i}) + \sum_{i=1}^n \frac{\lambda_3}{n} \|x_i - \bar{f}_\theta f_\theta(x_i)\|_2. \end{aligned} \quad (15)$$

Let $L^c(\theta, \bar{\theta}, z_j)$ be the expected loss corresponding to $\hat{L}_S^c(\theta, \bar{\theta}, z_j)$. We give the following Theorem 2 to characterize the generalization bound of (15).

Theorem 2. *As with the assumptions in Theorem 1, with probability at least $1 - \tau$, the following inequality holds for $\theta, \bar{\theta}, z_j$ in proper parameter spaces:*

$$\hat{L}_S^c(\theta, \bar{\theta}, z_j) \leq L^c(\theta, \bar{\theta}, z_j) + 2\sqrt{2}M_1R_1 + 2\sqrt{2}R_2 + \frac{\chi_1 + \chi_2}{\sqrt{n}} \quad (16)$$

where $M_1 = \bar{r}(2M\ell)^{\bar{r}-1}\ell$, $\chi_1 = \frac{2(2M)^{\bar{r}} + 4\lambda_3M}{n} \sqrt{\frac{-\ln \tau}{2}}$, $R_1 = E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} (f_\theta(x_i))_t \right|$, $\chi_2 = 2\sqrt{2}M_1Mdk$, $R_2 = E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sigma_{it} (\bar{f}_\theta f_\theta(x_i))_t \right|$, σ_{it} denotes the Rademacher random variable, E_S denotes the expectation with respect to S , and $f_\theta(x_i)_t$ denotes the t -th element of the vector $f_\theta(x_i)$.

Compared to Theorem 1, an additional term χ_2 appears in Theorem 2 due to the optimization of anchors (centroids). It is found that the upper bound of the empirical loss depends on the number of anchors. Our generalization bound has a similar form to the bound in the k-means (Maurer & Pontil, 2010). A tighter upper bound of the kernel k-means can be found in Yin et al. (2022) by using the infinity vector contraction (Foster & Rakhlin, 2019).

3 EXPERIMENTAL RESULTS

3.1 EXPERIMENTS ON SUPERVISED LEARNING

We perform the experiments on some data sets to obtain effective representations of features for classification tasks. In experiments, x_1, \dots, x_n consist of the training set and $y_{j|i}$ ($j = 1, \dots, k$) are taken from the samples that have the same label as x_i . It is found that some face data sets belong to the small-sample-size problem since the number of each class in the training set is much smaller than the dimension of the samples. When our model adopts one linear layer, we refer to our model as linear conditional Wasserstein supervised learning (LCWSL). When our model contains several linear layers and ReLU functions, we refer to our model as deep conditional Wasserstein supervised learning (DCWSL). The dimension of embedding spaces in our model is equal to the number of classes. We compare our model with several kernel-based methods including kernel discriminant analysis (KDA)(Zheng et al., 2014), kernel discriminant analysis based on the L_1 norm (KDAL1)(Zheng et al., 2014) and regularized kernel discriminant analysis (RKDA) (Diaz Vico & Dorrnsoro, 2020). In addition, deep Fisher discriminant analysis (DFDA) (Diaz Vico & Dorrnsoro, 2020) and deep Wasserstein discriminant analysis (DWDA) (Su et al., 2022) are tested. Since our model is employed to explore the latent space in supervised learning, we adopt the nearest neighbor (NN) classifier with the Euclidean norm. Experimental results on the data sets are shown in Table 1 and experimental details are in appendixes.

Table 1: Error rates (%) of various methods and their standard deviations on data sets

| data sets | KDA | KDAL1 | RKFDA | DFDA | DWDA | LCWSL | DCWSL |
|-----------|------------|------------|------------|------------|------------------|------------------|-------------------|
| Dna | 10.18±2.37 | 9.74±2.46 | 9.56±2.35 | 9.41±3.05 | 9.49±2.27 | 9.58±2.47 | 9.21±2.35 |
| Pendigits | 7.36±1.29 | 6.25±1.04 | 6.17±3.24 | 6.27±2.08 | 6.32±2.38 | 6.20±1.45 | 6.05±1.38 |
| Iris | 4.00±2.28 | 3.33±2.04 | 3.33±2.04 | 4.00±2.28 | 3.33±2.04 | 3.33±2.04 | 2.56±1.78 |
| Satimage | 24.57±2.26 | 24.38±2.67 | 24.69±3.05 | 16.77±2.59 | 16.86±2.62 | 23.46±2.77 | 16.57±2.61 |
| Waveform | 22.26±1.72 | 20.34±1.51 | 20.11±1.47 | 20.19±1.52 | 19.87±2.24 | 20.21±1.85 | 19.02±1.76 |
| ORL | 8.76±2.12 | 8.53±2.09 | 8.36±2.24 | 10.46±2.37 | 10.55±2.16 | 8.21±2.45 | 10.38±1.92 |
| Yale | 7.52±3.50 | 7.44±3.95 | 7.26±3.41 | 11.47±3.09 | 11.90±3.51 | 7.20±1.05 | 11.93±1.55 |
| UMIST | 8.97±2.25 | 8.76±2.34 | 8.45±3.02 | 10.56±3.50 | 10.78±3.05 | 8.21±2.92 | 10.33±2.06 |
| COIL | 8.45±2.21 | 9.43±1.65 | 8.22±1.69 | 8.13±2.02 | 8.06±1.72 | 8.19±1.98 | 8.08±1.23 |
| MSRA | 10.12±1.05 | 9.56±0.98 | 9.35±0.97 | 9.43±1.02 | 9.46±1.15 | 9.21±1.98 | 9.72±1.09 |

From Table 1, we can see that deep learning models such as DFDA, DWDA and DCWSL perform poorly on ORL, Yale and UMIST data sets. This comes from the fact that overfitting occurs since there are not enough training samples to learn the parameters of deep learning models. However, LCWSL obtains better performance than other methods on these face data sets. It is found that KDAL1 is superior to KDA on these data sets since KDAL1 is robust to outliers. DFDA and DWDA do not explore the reconstruction error of samples, whereas DCWSL makes use of the reconstruction error of the samples. Overall, it is more reasonable to use conditional distributions to transport data points in an embedding space.

3.2 CLUSTERING EXPERIMENTS

We verify the proposed model on some data sets in terms of clustering tasks. We use the normalized mutual information (NMI) to show the performance of the clustering methods. We also implement kernel k-means (KKM)(Paul et al., 2022), kernel fuzzy k-means (KFKM)(Paul et al., 2022), kernel power k-means(KPKM) (Paul et al., 2022), the deep clustering model based on the t distribution (DEC) (Xie et al., 2016), the improved DEC(IDEC) based on autoencoders (Guo et al., 2017), and the deep fuzzy k-means method (DFKM) (Zhang et al., 2020). Since the aim of our framework of (14) is to search for the embedding space of data in terms of autoencoders, we can use any clustering method after the embedding space of data is obtained. Here, we perform the spectral clustering on obtained features, where the number of neighbors is 5. In such a case, we refer to our model as deep conditional Wasserstein plus spectral clustering (DCWSC). Table 2 shows the NMI of various methods where we list the best result of each method. From Table 2, we note that our model is superior to other models since we optimize anchors to learn features in an embedding space. Note that deep-learning models such as DEC, IDEC and DFKM jointly learn data embedding and clustering. KKM, KFKM, and KPKM make use of kernel functions to learn the embedding space. It is found that the features based on deep learning models are better than those from kernel functions. The experimental results show that feature learning via optimal transport and autoencoders is effective for unsupervised learning.

Table 2: NMI values (%) of various methods on some data sets

| data sets | KKM | KFKM | KPKM | DEC | IDEC | DFKM | DCWSC |
|-----------|-------------|------------|------------|------------|-------------------|------------|-------------------|
| Dna | 40.25±2.46 | 42.55±3.17 | 46.34±2.81 | 49.22±2.70 | 49.46±3.10 | 48.26±2.92 | 50.21±2.86 |
| Pendigits | 72.46±2.23 | 73.82±1.90 | 73.25±2.24 | 72.35±2.72 | 72.36±2.67 | 77.85±1.90 | 80.55±1.79 |
| Iris | 76.74±2.50 | 77.83±2.89 | 78.35±3.01 | 80.79±2.68 | 81.44±2.71 | 80.25±3.05 | 88.46±2.96 |
| Satimage | 62.33 ±3.05 | 62.25±3.23 | 63.35±3.19 | 65.56±3.52 | 65.26±3.43 | 68.33±3.66 | 70.05±3.25 |
| Waveform | 27.22±2.04 | 28.46±2.17 | 30.51±2.51 | 30.11±2.62 | 34.61±2.73 | 35.14±2.29 | 38.16±2.14 |
| ORL | 62.45±2.78 | 63.64±3.07 | 65.57±3.12 | 70.65±2.46 | 70.24±2.71 | 69.33±2.66 | 80.63±3.41 |
| Yale | 60.22 ±4.52 | 61.25±4.23 | 63.33±4.62 | 65.12±4.02 | 64.18±4.19 | 65.26±4.31 | 73.22±4.51 |
| UMIST | 72.35±3.12 | 72.33±3.25 | 76.59±3.28 | 81.26±3.42 | 80.36±3.30 | 82.35±3.66 | 87.18±3.20 |
| COIL | 78.69±2.21 | 80.42±2.32 | 82.62±2.44 | 90.12±2.50 | 90.25±2.36 | 86.26±2.68 | 92.35±2.23 |
| MSRA | 56.12±2.02 | 58.24±2.12 | 57.22±2.29 | 60.22±2.30 | 62.21±2.19 | 59.23±2.26 | 61.26±2.25 |

3.3 EXPERIMENTS ON TWO LARGE-SCALE DATA SETS

We find that on small-scale data sets, using a one-layer network sometimes obtains much better performance than using multiple-layer networks. Does this phenomenon occur on large-scale data sets? In the following experiments, we find that this phenomenon does not occur. Here we select two large-scale image data sets (MNIST and FashionMNIST) to evaluate the proposed model. The aim of using these two data sets is that we do not need to employ complex networks to achieve relatively good performance. Unlike the deep learning models based on data augmentation, we only use our autoencoders to achieve the embedding features. The training samples are employed to select the parameters of models and test samples are used to measure the performance of models. In our experiments, we adopt a large batch size of 2000. Since there are a large number of samples in the training set, we employ the class-mean classifier in the classification task. In the kernel-based methods, 100 anchors taken from the k-means algorithms are employed to compute kernel matrices since computing kernel matrices for all the samples is impossible. Table 3 lists the experimental results from classification and clustering tasks. From Table 3, we note that the performance of DCWSL in the classification experiments is much better than that of LCWS. It shows that using multiple-layer networks is beneficial for large-scale data sets. It is clear that our method is superior to other methods since we explore the transferring range of data in the embedding space via conditional optimal transport. In the clustering experiments, we observe that our model outperforms other models since we employ conditional distributions to learn the optimal anchors.

Table 3: Classification (error rate) and clustering (NMI) on two large-scale data sets

| classification | KDA | KDALI | RKFDA | DFDA | DWDA | LCWSL | DCWSL |
|----------------|------------|------------|------------|------------|------------|------------|-------------------|
| MNIST | 10.30±2.12 | 12.15±2.57 | 9.55±1.86 | 9.39±2.32 | 8.86±2.73 | 9.19±2.26 | 8.21±2.31 |
| Fashion | 12.30±3.49 | 14.36±3.89 | 11.79±4.01 | 11.22±3.37 | 10.35±3.58 | 10.41±3.3 | 9.21±3.17 |
| Clustering | KKM | KFKM | KPKM | DEC | IDEC | DFKM | DCWSC |
| MNIST | 54.33±2.53 | 59.40±2.67 | 58.37±2.49 | 67.46±2.56 | 79.21±2.73 | 70.23±2.12 | 81.24±2.51 |
| Fashion | 46.62±3.72 | 47.39±3.69 | 50.28±3.10 | 54.35±3.53 | 56.45±3.44 | 54.37±2.76 | 62.08±3.16 |

3.4 CLASSIFICATION OF SET-VALUED OBJECTS

Here we modify our model to make it capable of handling set-valued classification problems. For set-valued classification problems, each object contains many examples. Unlike previous experiments, we assume that the set $\{x_1, \dots, x_n\}$ is a set-valued object containing n examples in the validation set or test set. For the data point x_i , we can obtain its k neighbours $y_{1|i}, \dots, y_{k|i}$ and these k neighbours are from the training set. Since we know the labels of $y_{1|i}, \dots, y_{k|i}$ in the training set, we assign the label of x_i to the label of $y_{j|i}$ with the largest $p_{j|i}$ ($j = 1, \dots, k$). Thus, we obtain the label of each example in a set-valued object. Finally, the majority voting strategy is employed to achieve the label of the set-valued object. We refer to our model as the deep conditional Wasserstein classifier (DCWC). Our model will degenerate into the deep nearest-neighbor classifier if each object only contains an example and the parameter λ_1 approaches the positive infinity. Here we need to use the validation set to learn the embedding space of data and hyperparameters. In the test stage, we fix the parameters of autoencoders and optimize $p_{j|i}$. We test DCWC on two medical image sets in binary classification problems (Yang et al., 2021). We use 780 images from the breast image set and 4708 images from the pneumonia image set. To evaluate the performance of DCWC, we compare it with several set-valued data classification methods such as the second-order cone programming (SOCP) approach (Shivaswamy et al., 2006), the sparse approximated nearest point (SANP) method (Hu et al., 2011), regularized collaborative representation classification (RCRC) (Zhu et al., 2014),

support measure machines(SMMs) (Muandet et al., 2012), and support function machines (SFMs) (Chen et al., 2017). Figure 2 shows experimental results on two medical image sets.

As can be seen from Figure 2, SFMs are not superior to DCWC since SFMs generally give sparse support vectors. It is found that DCWC yields the best performance on these data sets since DCWC explores the weight of each example in the set-valued objects. Among these methods, SFMs are sampling-based methods. SANP, RCRC and SMMs explore all possible representations of images. If the representations of images contain distorted features, these distorted features will affect the performance of classifiers. Our DCWC makes use of the conditional optimal transport and reconstruction errors of data to achieve effective features. The experimental results indicate that it is reasonable to employ the optimal transport theory to classify set-valued data.

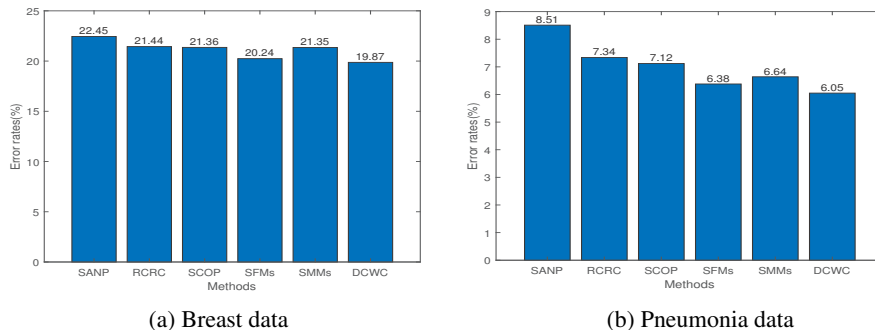


Figure 2: Experimental results on medical image data sets.

4 RELATED WORK

Many feature learning methods based on the deep architectures of neural networks have been developed. Multi-layer learning models (Yuan et al., 2015) have been proposed to deal with the scene recognition problem, and they are available in an unsupervised way. The deep semi-nonnegative matrix factorization (Trigeorgis et al., 2014) can find the latent representation of data in a low-dimensional space, and the new description can improve the clustering performance of data. Deep Fisher discriminant analysis (Diaz Vico & Dorronsoro, 2020) takes advantage of deep neural networks to capture the nonlinear features of data. To deal with sequence data, deep order-preserving Wasserstein discriminant analysis (Su et al., 2022) achieves a nonlinear transformation by maximizing the inter-class distance and minimizing the intra-class distance. The Wasserstein autoencoder (Tolstikhin et al., 2018) was proposed to achieve a generative model of data distributions. However, these feature learning methods do not explore their generalization bounds.

Kernel k-means clustering methods can deal with the nonlinear structure of data in unsupervised learning. For bounded random vectors, the expected excess clustering risk was studied in the work (Maurer & Pontil, 2010). An upper generalization bound of the kernel k-means method in a reduced space (Yin et al., 2022) is derived in terms of the Rademacher complexity. The deep clustering model via the t distribution (DEC) (Xie et al., 2016) has been proposed. The improved DEC (IDEC) (Guo et al., 2017) used autoencoders to enhance the performance of DEC. However, the generalization bounds of DEC and IDEC are not explored.

5 CONCLUSIONS AND FURTHER WORK

In this paper we have introduced a feature learning framework relying on optimal transport and autoencoders. The use of the conditional probability in the proposed model is to make each data point adapt to its neighborhood, and this may well be suitable for the characteristics of data. The experimental results on real data sets demonstrate the feasibility of the proposed model on some classification and clustering tasks. Since the performance of the proposed model is affected by autoencoders, how to select proper autoencoders for data sets is worth exploring. In the future, we will focus on the problem of how to employ advanced autoencoders to improve our model to deal with complex data sets in the real world.

REFERENCES

- 540
541
542 Wei Bian and Dacheng Tao. Max-min distance analysis by using sequential sdp relaxation for
543 dimension reduction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):
544 1037–1050, 2011.
- 545 Christopher M. Bishop. *Pattern Recognition and Machine Learning*. 2007.
- 546
547 Forough Rezaei Boroujeni, Sen Wang, Zhihui Li, Nicholas West, Bella Stantic, Lina Yao, and
548 Guodong Long. Trace ratio optimization with feature correlation mining for multiclass discrimi-
549 nant analysis. In *AAAI-18 AAAI Conference on Artificial Intelligence*, pp. 2746–2753, 2018.
- 550 Jiqiang Chen, Qinghua Hu, Xiaoping Xue, Minghu Ha, and Litao Ma. Support function machine for
551 set-based classification with application to water quality evaluation. *Inf. Sci.*, 388:48–61, 2017.
552
- 553 Nicolas Courty, Remi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for do-
554 main adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(9):1853–
555 1865, 2017. doi: 10.1109/TPAMI.2016.2615921.
- 556 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Proceedings*
557 *of the 26th International Conference on Neural Information Processing Systems, NIPS’13*, pp.
558 2292–2300, 2013.
559
- 560 David Diaz Vico and Jose R. Dorransoro. Deep least squares fisher discriminant analysis. *IEEE*
561 *Transactions on Neural Networks and Learning Systems*, 31(8):2752–2763, 2020.
- 562 Kilian Fatras, Thibault Sejourne, Rémi Flamary, and Nicolas Courty. Unbalanced minibatch optimal
563 transport; applications to domain adaptation. In *Proceedings of the 38th International Conference*
564 *on Machine Learning*, volume 139, pp. 3186–3197. PMLR, 2021.
565
- 566 Remi Flamary, Marco Cuturi, Nicolas Courty, and Alain Rakotomamonjy. Wasserstein discriminant
567 analysis. *Machine Learning*, 107(12):1923–1945, 2018.
568
- 569 Dylan J. Foster and Alexander Rakhlin. Vector contraction for rademacher complexity. *CoRR*,
570 abs/1911.06468, 2019. URL <http://arxiv.org/abs/1911.06468>.
- 571 Hyemin Gu, Markos A. Katsoulakis, Luc Rey-Bellet, and Benjamin J. Zhang. Combining
572 wasserstein-1 and wasserstein-2 proximals: robust manifold learning via well-posed generative
573 flows. *ArXiv*, abs/2407.11901v1, 2024. URL <https://arxiv.org/abs/2407.11901v1>.
574
- 575 Xifeng Guo, Long Gao, Xinwang Liu, and Jianping Yin. Improved deep embedded clustering
576 with local structure preservation. In *Proceedings of the 26th International Joint Conference on*
577 *Artificial Intelligence, IJCAI’17*, pp. 1753–1759. AAAI Press, 2017. ISBN 9780999241103.
- 578 Yiqun Hu, Ajmal S. Mian, and Robyn A. Owens. Sparse approximated nearest points for image set
579 classification. *CVPR 2011*, pp. 121–128, 2011.
580
- 581 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd*
582 *International Conference on Learning Representations*, volume abs/1412.6980, 2015. URL
583 <https://api.semanticscholar.org/CorpusID:6628106>.
- 584 Yu. G. Kuritsyn. The khinchin inequality. *Journal of Soviet Mathematics*, 35(9):2363–2365, 1986.
585
- 586 Qian Li, Zhichao Wang, Gang Li, Jun Pang, and Guandong Xu. Hilbert sinkhorn divergence for
587 optimal transport. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
588 *(CVPR)*, pp. 3834–3843, 2021. doi: 10.1109/CVPR46437.2021.00383.
- 589 Wuchen Li, Siting Liu, and Stanley Osher. A kernel formula for regularized wasserstein proximal
590 operators. *SIAM J. Optim.*, 10:1–16, 2023.
591
- 592 Wei Lin and Antoni B. Chan. Optimal transport minimization: Crowd localization on density maps
593 for semi-supervised counting. In *2023 IEEE/CVF Conference on Computer Vision and Pattern*
Recognition (CVPR), pp. 21663–21673, 2023. doi: 10.1109/CVPR52729.2023.02075.

- 594 Xinwang Liu, Xinzhong Zhu, Miaomiao Li, Lei Wang, En Zhu, Tongliang Liu, Marius Kloft, D-
595 inggang Shen, Jianping Yin, and Wen Gao. Multiple kernel kk -means with incomplete kernels.
596 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(5):1191–1204, 2020. doi:
597 10.1109/TPAMI.2019.2892416.
- 598 Yang Liu, Zhipeng Zhou, and Baigui Sun. Cot: Unsupervised domain adaptation with clustering and
599 optimal transport. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
600 *(CVPR)*, pp. 19998–20007, 2023. doi: 10.1109/CVPR52729.2023.01915.
- 602 Andreas Maurer and Massimiliano Pontil. k -dimensional coding schemes in hilbert spaces. *IEEE*
603 *Transactions on Information Theory*, 56(11):5839–5846, 2010. doi: 10.1109/TIT.2010.2069250.
- 604 Krikamol Muandet, Kenji Fukumizu, Francesco Dinuzzo, and Bernhard Schölkopf. Learning from
605 distributions via support measure machines. In *NIPS*, 2012.
- 607 Feiping Nie, Z.Wang, R.Wang, Z.Wang, and X.Li. Towards robust discriminant projections learning
608 via non-greedy l_{21} norm minmax. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(6):2086–2100,
609 2021.
- 610 Feiping Nie, Zheng Wang, Rong Wang, and Xuelong Li. Adaptive local embedding learning for
611 semi-supervised dimensionality reduction. *IEEE Transactions on Knowledge and Data Engineer-*
612 *ing*, 34(10):4609–4621, 2022. doi: 10.1109/TKDE.2021.3049371.
- 614 Feiping Nie, Canyu Zhang, Zheng Wang, Rong Wang, and Xuelong Li. Local embedding learning
615 via landmark-based dynamic connections. *IEEE Transactions on Neural Networks and Learning*
616 *Systems*, 34(11):9481–9492, 2023. doi: 10.1109/TNNLS.2022.3203014.
- 617 Debolina Paul, Saptarshi Chakraborty, Swagatam Das, and Jason Xu. Implicit annealing in kernel
618 spaces: A strongly consistent clustering approach. *IEEE Transactions on Pattern Analysis and*
619 *Machine Intelligence*, pp. 1–10, 2022. doi: 10.1109/TPAMI.2022.3217137.
- 621 Q. Qian, L. Shang, B. Sun, J. Hu, T. Tacoma, H. Li, and R. Jin. Softtriple loss: Deep metric
622 learning without triplet sampling. In *2019 IEEE/CVF International Conference on Computer*
623 *Vision (ICCV)*, pp. 6449–6457, Los Alamitos, CA, USA, nov 2019.
- 624 Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Tom Luo. A unified convergence analysis of
625 block successive minimization methods for nonsmooth optimization. *SIAM J. Optim.*, 23:1126–
626 1153, 2012.
- 627 Mathieu Serrurier, Franck Mamalet, Alberto Gonzalez-Sanz, Thibaut Boissin, Jean-Michel Loubes,
628 and Eustasio del Barrio. Achieving robustness in classification using optimal transport with hinge
629 regularization. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
630 *(CVPR)*, pp. 505–514, 2021. doi: 10.1109/CVPR46437.2021.00057.
- 632 Pannagadatta K. Shivaswamy, Chiranjib Bhattacharyya, and Alexander J. Smola. Second order
633 cone programming approaches for handling missing and uncertain data. *J. Mach. Learn. Res.*, 7:
634 1283–1314, 2006.
- 635 Bing Su and Gang Hua. Order-preserving optimal transport for distances between sequences. *IEEE*
636 *Transactions on Pattern Analysis and Machine Intelligence*, 41(12):2961–2974, 2019. doi: 10.
637 1109/TPAMI.2018.2870154.
- 639 Bing Su, Jiahuan Zhou, Ji-Rong Wen, and Ying Wu. Linear and deep order-preserving wasserstein
640 discriminant analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(6):
641 3123–3138, 2022. doi: 10.1109/TPAMI.2021.3050750.
- 642 I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schölkopf. Wasserstein auto-encoders. In
643 *6th International Conference on Learning Representations (ICLR)*, May 2018. URL
644 <https://openreview.net/forum?id=HkL7n1-0b>.
- 645 George Trigeorgis, Konstantinos Bousmalis, Stefanos Zafeiriou, and Björn Schuller. A deep semi-
646 nmf model for learning hidden representations. In *International Conference on Machine Learning*,
647 2014.

- 648 Lan V. Truong. Rademacher complexity-based generalization bounds for deep learning. *CoRR*,
649 abs/2208.04284, 2019. URL <http://arxiv.org/abs/2208.04284>.
650
- 651 Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R. Scott. Multi-similarity loss
652 with general pair weighting for deep metric learning. In *2019 IEEE/CVF Conference on Computer
653 Vision and Pattern Recognition (CVPR)*, pp. 5017–5025, 2019. doi: 10.1109/CVPR.2019.00516.
- 654 Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis.
655 In *Proceedings of the 33rd International Conference on International Conference on Machine
656 Learning - Volume 48, ICML16*, pp. 478–487. JMLR.org, 2016.
657
- 658 Shuicheng Yan, Dong Xu, Benyu Zhang, Hong-jiang Zhang, Qiang Yang, and Stephen Lin. Graph
659 embedding and extensions: A general framework for dimensionality reduction. *IEEE Transac-
660 tions on Pattern Analysis and Machine Intelligence*, 29(1):40–51, 2007. doi: 10.1109/TPAMI.
661 2007.250598.
- 662 Jiancheng Yang, Rui Shi, and Bingbing Ni. Medmnist classification decathlon: A lightweight au-
663 toml benchmark for medical image analysis. In *2021 IEEE 18th International Symposium on
664 Biomedical Imaging (ISBI)*, pp. 191–195, 2021.
- 665 Yaqiang Yao, Yang Li, Bingbing Jiang, and Huanhuan Chen. Multiple kernel k-means clustering by
666 selecting representative kernels. *IEEE Transactions on Neural Networks and Learning Systems*,
667 32(11):4983–4996, 2021. doi: 10.1109/TNNLS.2020.3026532.
668
- 669 Rong Yin, Yong Liu, Weiping Wang, and Dan Meng. Scalable kernel k -means with randomized
670 sketching: From theory to algorithm. *IEEE Transactions on Knowledge and Data Engineering*,
671 pp. 1–14, 2022. doi: 10.1109/TKDE.2022.3222146.
- 672 Yuan Yuan, Lichao Mou, and Xiaoqiang Lu. Scene recognition by manifold regularized deep learn-
673 ing architecture. *IEEE Transactions on Neural Networks and Learning Systems*, 26(10):2222–
674 2233, 2015. doi: 10.1109/TNNLS.2014.2359471.
675
- 676 Benjamin J. Zhang, Siting Liu, Wuchen Li, Markos A. Katsoulakis, and Stanley Osher. Wasserstein
677 proximal operators describe score-based generative models and resolve memorization. *ArXiv*,
678 abs/2402.06162, 2024. URL <https://arxiv.org/abs/2402.06162>.
- 679 Rui Zhang, Xuelong Li, Hongyuan Zhang, and Feiping Nie. Deep fuzzy k-means with adaptive loss
680 and entropy regularization. *IEEE Transactions on Fuzzy Systems*, 28(11):2814–2824, 2020. doi:
681 10.1109/TFUZZ.2019.2945232.
- 682 Wenming Zheng, Zhouchen Lin, and Haixian Wang. L1-norm kernel discriminant analysis via bayes
683 error bound optimization for robust feature extraction. *IEEE Transactions on Neural Networks*,
684 25(4):793–805, 2014.
685
- 686 Peng Fei Zhu, Wangmeng Zuo, Lei Zhang, Simon C. K. Shiu, and David Dian Zhang. Image
687 set-based collaborative representation for face recognition. *IEEE Transactions on Information
688 Forensics and Security*, 9:1120–1132, 2014.
689
690
691
692
693
694
695
696
697
698
699
700
701

A THE CONTINUOUS VERSION OF (6) IN OUR PAPER

Here, we give a continuous version of (6) which provides insights into understanding our discrete version. The following framework is used to extract effective features of data.

$$\begin{aligned} \min L(\theta, \bar{\theta}, \phi, \bar{\phi}, p(g_\phi(y)|x), p(x)) &:= \int \rho(f_\theta(x), g_\phi(y))^\bar{r} p(g_\phi(y)|x) p(x) dx dy + \\ \lambda_1 \int p(x) \text{KL}(p(g_\phi(y)|x) || q(g_\phi(y)|x)) dx &+ \lambda_2 \text{KL}(p(x) || q(x)) + \\ \lambda_3 \int \|x - \bar{f}_\theta f_\theta(x)\|_2 d\mu &+ \lambda_4 \int \|y - \bar{g}_\phi g_\phi(y)\|_2 d\nu. \end{aligned} \quad (17)$$

where $\text{KL}(p(g_\phi(y)|x) || q(y|x)) = \int p(g_\phi(y)|x) \ln \frac{p(g_\phi(y)|x)}{q(y|x)} dy$, $\text{KL}(p(x) || q(x)) = \int p(x) \ln \frac{p(x)}{q(x)} dx$, $q(y|x)$ and $q(x)$ are prior (conditional) probabilities in the original space, and $\lambda_i (i = 1, \dots, 4)$ are nonnegative hyperparameters. $p(g_\phi(y)|x)$ is actually the induced distribution of $q(y|x)$ via the encoder g_ϕ . The first term in (17) denotes the transport cost in the embedding space. The second term is the Kullback-Leibler (KL) divergence to control conditional probabilities between $p(g_\phi(y)|x)$ and $q(y|x)$. The third term is the KL divergence between $p(x)$ and $q(x)$. The last two terms involve the reconstruction errors of data x and y . Trivial solutions of $\theta, \bar{\theta}, \phi, \bar{\phi}$ may be obtained if we do not employ the reconstruction error of data or the regularization terms for these parameters. From (17), we find that the loss function in the proposed model consists of the transport cost, reconstruction errors of data and additional regularization terms. The framework is generic since we do not give specific autoencoders and any transport cost can be used to replace $\rho(\cdot)$. That is, we can employ some existing autoencoders to our framework.

B THE PROOF OF THEOREM 1

Lemma 1. For any $r \geq 1$ and two vectors x and y with proper dimensions, we have

$$\|x + y\|_r \leq \|x\|_r + \|y\|_r. \quad (18)$$

Lemma 2 (Yin et al., 2022). Let x_1, \dots, x_n be n data points, and let F be a class of vector-valued functions $f : X \mapsto \mathbb{R}^d$ and $h_i : \mathbb{R}^d \mapsto \mathbb{R}$ be functions with the Lipschitz constant ℓ . Then we have

$$E \sup_{f \in F} \sum_{i=1}^n \sigma_i h_i(f(x_i)) \leq \sqrt{2} \ell E \sup_{f \in F} \sum_{i=1}^n \sum_{j=1}^d \sigma_{ij} (f_i)_j, \quad (19)$$

where σ_{ij} is an independent doubly indexed Rademacher sequence and $(f_i)_j$ is the j -th component of $f(x_i)$.

Lemma 3 (Kuritsyn, 1986). Let a be a vector containing m elements. The following Khintchine inequality holds

$$A_r \left(\sum_{i=1}^m a_i^2 \right)^{\frac{1}{2}} \leq \left(E \left| \sum_{i=1}^m \sigma_i a_i \right|^r \right)^{\frac{1}{r}} \leq B_r \left(\sum_{i=1}^m a_i^2 \right)^{\frac{1}{2}}, \quad (20)$$

where A_r and B_r are constants depending on r . When $r = 1$, we have $B_p = 1$.

To give the generalization bound of $\hat{L}_S(\theta, \bar{\theta})$, we rewrite $\hat{L}_S(\theta, \bar{\theta})$ by removing $p_{j|i}$. Thus, $\hat{L}_S(\theta, \bar{\theta})$ can be formulated as

$$\hat{L}_S(\theta, \bar{\theta}) = -\frac{\lambda_1}{n} \sum_{i=1}^n \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_i) - f_\theta(y_{j|i}))^\bar{r} / \lambda_1) + \sum_{i=1}^n \frac{\lambda_3}{n} \|x_i - \bar{f}_\theta f_\theta(x_i)\|_2. \quad (21)$$

Let S' be the data set where only a data point is different from the data set S , e.g., \bar{x}_s . Let $\hat{L}_S(\theta, \bar{\theta})$ denote the empirical loss from S' . Let us define the following functions:

$$\psi_S = \sup_{\theta, \bar{\theta}} (L(\theta, \bar{\theta}) - \hat{L}_S(\theta, \bar{\theta})), \quad (22)$$

$$\psi'_S = \sup_{\theta, \bar{\theta}} (L(\theta, \bar{\theta}) - \hat{L}_{S'}(\theta, \bar{\theta})). \quad (23)$$

From (22) and (23), we have

$$\begin{aligned} |\psi_S - \psi'_S| &\stackrel{(1)}{\leq} |\hat{L}_S(\theta, \bar{\theta}) - \hat{L}_{S'}(\theta, \bar{\theta})| = \frac{1}{n} \sup_{\theta, \bar{\theta}} \left| -\lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_s) - f_\theta(y_{j|s})))^{\bar{r}} / \lambda_1 \right. \\ &\quad \left. + \lambda_3 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(\bar{x}_s) - f_\theta(y_{j|s})))^{\bar{r}} / \lambda_1 + \lambda_3 \|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 - \lambda_3 \|\bar{x}_s - \bar{f}_{\bar{\theta}} f_\theta(\bar{x}_s)\|_2 \right| \\ &\leq \frac{1}{n} \left\{ \sup_{\theta} \lambda_1 \left| \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_s) - f_\theta(y_{j|s})))^{\bar{r}} \right| + \sup_{\theta} \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(\bar{x}_s) - f_\theta(y_{j|s})))^{\bar{r}} \right| \right. \\ &\quad \left. + \sup_{\theta, \bar{\theta}} \lambda_3 \|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 + \sup_{\theta, \bar{\theta}} \lambda_3 \|\bar{x}_s - \bar{f}_{\bar{\theta}} f_\theta(\bar{x}_s)\|_2 \right\}. \end{aligned} \quad (24)$$

From assumptions of A0 in our paper, we have

$$\rho(f_\theta(x_s) - f_\theta(y_{j|s}))^{\bar{r}} = \varphi(f_\theta(x_s) - f_\theta(y_{j|s}))^{\bar{r}}. \quad (25)$$

Note that the function $\varphi(\cdot)$ is Lipschitz continuous and its Lipschitz constant is ℓ . Hence, we have

$$\varphi(f_\theta(x_s) - f_\theta(y_{j|s})) \leq \ell \|f_\theta(x_s) - f_\theta(y_{j|s})\|_2. \quad (26)$$

From $\|f_\theta(x_t)\|_2 \leq M$ and $\|f_\theta(y_{j|t})\|_2 \leq M$, we have

$$\|f_\theta(x_s) - f_\theta(y_{j|s})\|_2 \leq 2M. \quad (27)$$

Similarly, we have

$$\|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 \leq 2M. \quad (28)$$

From (24), (26), and (28), we have

$$\begin{aligned} |\psi_S - \psi'_S| &\leq \frac{1}{n} \left\{ \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-(2M\ell)^{\bar{r}} / \lambda_1) \right| \right. \\ &\quad \left. + \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-(2M\ell)^{\bar{r}} / \lambda_1) \right| + 2\lambda_3 2M \right\} \leq \frac{2(2M\ell)^{\bar{r}} + 4\lambda_3 M}{n}. \end{aligned} \quad (29)$$

In the following, we consider the expectation of ψ_S with respect to the data set S , denoted by $E_S(\psi_S)$:

$$\begin{aligned} E_S(\psi_S) &= E_S(\sup_{\theta, \bar{\theta}} (L(\theta, \bar{\theta}) - \hat{L}_S(\theta, \bar{\theta}))) \\ &\stackrel{(1)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \left\{ \sum_{i=1}^n -\lambda_1 \sigma_i \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_i) - f_\theta(y_{j|i})))^{\bar{r}} / \lambda_1 + \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 \right\} \\ &\stackrel{(2)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \left\{ \sum_{i=1}^n \sigma_i \sum_{j=1}^k q_{j|i} \rho(f_\theta(x_i) - f_\theta(y_{j|i}))^{\bar{r}} + \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 \right\} \\ &\stackrel{(3)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sigma_i \sum_{j=1}^k q_{j|i} \rho(f_\theta(x_i) - f_\theta(y_{j|i}))^{\bar{r}} + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2. \end{aligned} \quad (30)$$

In (30), the first inequality comes from the symmetrization of random variables, and the second inequality uses Jensen inequality from the fact that $-\ln x$ is a convex function. Note that the Lipschitz constant of the norm $\|\cdot\|$ is 1. The function $\|x\|_r^r$ is not Lipschitz continuous if the variable

810 x takes the infinite values. However, with the assumptions we provide, there exists the constant
 811 $M_1 = \bar{r}(2M\ell)^{\bar{r}-1}l$ such that $\rho(f_\theta(x_i) - f_\theta(y_{j|i}))^{\bar{r}}$ is Lipschitz continuous. Using Lemma 2, we
 812 have

$$\begin{aligned}
 813 & E_S(\psi_S) \stackrel{(1)}{\leq} 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (f_\theta(x_i) - f_\theta(y_{j|i}))_t \\
 814 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \\
 815 & \stackrel{(2)}{\leq} 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (f_\theta(x_i))_t \\
 816 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (-f_\theta(y_{j|i}))_t \\
 817 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \\
 818 & \stackrel{(3)}{\leq} 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sqrt{2} M_1 (f_\theta(x_i))_t + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (-f_\theta(y_{j|i}))_t \\
 819 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \\
 820 & \stackrel{(4)}{\leq} 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sqrt{2} M_1 (f_\theta(x_i))_t + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (-f_\theta(x_i))_t \\
 821 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t.
 \end{aligned}$$

822 (31)

823 In (31), the first inequality uses Lemma 3, and the second inequality uses the property of *sup*. The
 824 third inequality uses the fact that $\sum_{j=1}^k q_{j|i} = 1$. Since $y_{j|i}$ depends on x_i , we assume that $y_{j|i}$ is an
 825 independent copy of x_i . Hence, the fourth inequality is to replace $y_{j|i}$ with x_i . From (31), we have

$$\begin{aligned}
 826 & E_S(\psi_S) \leq 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} \sqrt{2} M_1 (f_\theta(x_i))_t \right| + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} \sqrt{2} M_1 (-f_\theta(x_i))_t \right| \\
 827 & + 2E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \right| \leq 4\sqrt{2} M_1 E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} (f_\theta(x_i))_t \right| \\
 828 & + 2\sqrt{2} E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \right|.
 \end{aligned}$$

829 (32)

830 From (32), we find that the upper bound of $E_S(\psi_S)$ depends on the Rademacher complexity of
 831 the encoders and decoders. From (29) and (32), we obtain that with probability at least $1 - \tau$, the
 832 following inequality holds for θ and $\bar{\theta}$ in proper parameter spaces by using McDiarmid inequality:

$$833 \hat{L}_S(\theta, \bar{\theta}) \leq L(\theta, \bar{\theta}) + 4\sqrt{2} M_1 R_1 + 2\sqrt{2} R_2 + \chi_1 \sqrt{\frac{-\log \tau}{2n}} \quad (33)$$

where $\chi_1 = \frac{2(2M)^{\bar{r}} + 4\lambda_3 M}{n}$, $R_1 = E_{S,\sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it}(f_{\theta}(x_i))_t \right|$, and $R_2 = E_{S,\sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sigma_{it}(\bar{f}_{\bar{\theta}} f_{\theta}(x_i))_t \right|$.

C OPTIMIZATION OF (14) IN OUR PAPER

In the following, we use the alternative optimization method to solve the optimization model:

$$\begin{aligned} \min_{(\theta, \bar{\theta}, p_{j|i}, p_i, z_j)} \hat{L} := & \sum_{i,j=1}^{n,k} p_i \rho(f_{\theta}(x_i), z_j)^{\bar{r}} p_{j|i} + \\ & \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} \| q_{\cdot|i}) + \sum_{i=1}^n \lambda_3 p_i \|x_i - \bar{f}_{\bar{\theta}} f_{\theta}(x_i)\|_2 + \lambda_2 \text{KL}(p \| q). \end{aligned} \quad (34)$$

(a): Update $p_{j|i}$ by fixing other variables. When we fix $\theta, \bar{\theta}, z_j$ and p_i , we solve the following model:

$$\min_{p_{j|i}} \sum_{i,j=1}^{n,k} \rho(f_{\theta}(x_i), z_j)^{\bar{r}} p_i p_{j|i} + \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} \| q_{\cdot|i}). \quad (35)$$

It is noted that (35) is a strongly convex optimization problem. It is of interest to note that it has a closed-form solution, denoted by

$$p_{j|i} = \frac{q_{j|i} \exp(-L_{j|i}^{op}/\lambda_1)}{\sum_{j=1}^k q_{j|i} \exp(-L_{j|i}^{op}/\lambda_1)}, \quad (36)$$

where $L_{j|i}^{op} = \rho(f_{\theta}(x_i), z_j)^{\bar{r}}$.

(b): Update p_i by fixing other variables. Given $\theta, \bar{\theta}, z_j$, and $p_{j|i}$, we achieve p_i by solving the following problem:

$$\begin{aligned} \min_{p_i} \sum_{i,j=1}^{n,k} \rho(f_{\theta}(x_i), z_j)^{\bar{r}} p_i p_{j|i} + & \lambda_1 \sum_{i=1}^n p_i \text{KL}(p_{\cdot|i} \| q_{\cdot|i}) + \lambda_2 \text{KL}(p \| q) + \\ & \lambda_3 \sum_{i=1}^n \|x_i - \bar{f}_{\bar{\theta}} f_{\theta}(x_i)\|_2 p_i. \end{aligned} \quad (37)$$

It is observed that the objective function in (37) is strongly convex. Thus, there exists a unique solution to p_i . The closed-form solution is denoted by

$$p_i = \frac{q_i \exp(-(L_i^{op} + L_i^{enre})/\lambda_2)}{\sum_{i=1}^n q_i \exp(-(L_i^{op} + L_i^{enre})/\lambda_2)}, \quad (38)$$

where $L_i^{op} = \sum_{j=1}^k \rho(f_{\theta}(x_i), z_j)^{\bar{r}} p_{j|i}$ and $L_i^{enre} = \lambda_3 \|x_i - \bar{f}_{\bar{\theta}} f_{\theta}(x_i)\|_2 + \text{KL}(p_{\cdot|i} \| q_{\cdot|i})$.

(c): Update z_j by fixing other variables. If ρ takes the Euclidean distance and $\bar{r} = 2$, z_j has the following solution:

$$z_j = \frac{\sum_{i=1}^n q_i q_{j|i} f_{\theta}(x_i)}{\sum_{i=1}^n q_i q_{j|i}}. \quad (39)$$

(d): Update θ and $\bar{\theta}$ by fixing other variables. In this step, we try to learn the parameters of autoencoders. Specifically, we solve the following optimization problem:

$$\min_{(\theta, \bar{\theta})} \sum_{i,j=1}^{n,k} \rho(f_{\theta}(x_i), z_j)^{\bar{r}} p_i p_{j|i} + \lambda_3 \sum_{i=1}^n p_i \|x_i - \bar{f}_{\bar{\theta}} f_{\theta}(x_i)\|_2. \quad (40)$$

Note that the objective function in (40) is nonconvex. We cannot obtain the global optimal solution. We generally update these parameters of models through the chain rule in the framework of neural networks. In this work, we resort to automatic differentiation to learn these parameters. For the sake of completeness, we summarize the main steps of solving (34) in Algorithm 2.

Algorithm 2: Optimization algorithm to (34)

-
- 1: Given $\lambda_i, q_{j|i}, q_i$, and initialize $p_i = q_i, p_{j|i} = q_{j|i}$
2: **For** $t=1$ to T **do**
2.1: solve (40) to achieve the parameters $(\theta, \bar{\theta})$;
2.2: solve (35) to achieve $p_{j|i}$;
2.3: solve (37) to achieve p_i ;
2.4: use (39) to achieve z_j ;
3: Output: the encoders and decoders.
-

D THE PROOF OF THEOREM 2

We obtain $\hat{L}_S^c(\theta, \bar{\theta}, z_j)$ by removing $P_{j|i}$. Thus, $\hat{L}_S^c(\theta, \bar{\theta}, z_j)$ can be formulated as

$$\hat{L}_S^c(\theta, \bar{\theta}, z_j) = \frac{-1}{n} \sum_{i=1}^n \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_i), z_j)^{\bar{r}}/\lambda_1) + \sum_{i=1}^n \frac{\lambda_3}{n} \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2. \quad (41)$$

Let S' be the data set where only a data point is different from the data set S , e.g., \bar{x}_s . Let $\hat{L}_{S'}^c(\theta, \bar{\theta})$ denote the empirical loss from S' . Let us define the following functions:

$$\psi_S = \sup_{\theta, \bar{\theta}, z_j} (L(\theta, \bar{\theta}) - \hat{L}_S^c(\theta, \bar{\theta})), \quad (42)$$

$$\psi_{S'} = \sup_{\theta, \bar{\theta}, z_j} (L(\theta, \bar{\theta}) - \hat{L}_{S'}^c(\theta, \bar{\theta})). \quad (43)$$

From (42) and (43), we have

$$\begin{aligned} |\psi_S - \psi_{S'}| &\stackrel{(1)}{\leq} |\hat{L}_S^c(\theta, \bar{\theta}) - \hat{L}_{S'}^c(\theta, \bar{\theta})| \\ &= \frac{1}{n} \sup_{\theta, \bar{\theta}, z_j} \left| -\lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_s), z_j)^{\bar{r}}/\lambda_1) + \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(\bar{x}_s), z_j)^{\bar{r}}/\lambda_1) \right. \\ &\quad \left. + \lambda_3 \|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 - \lambda_3 \|\bar{x}_s - \bar{f}_{\bar{\theta}} f_\theta(\bar{x}_s)\|_2 \right| \\ &\leq \frac{1}{n} \left\{ \sup_{\theta, y_j} \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_s), z_j)^{\bar{r}}/\lambda_1) \right| + \sup_{\theta, z_j} \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_s), z_j)^{\bar{r}}/\lambda_1) \right| \right. \\ &\quad \left. + \sup_{\theta, \bar{\theta}} \lambda_3 \|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 + \sup_{\theta, \bar{\theta}} \lambda_3 \|\bar{x}_s - \bar{f}_{\bar{\theta}} f_\theta(\bar{x}_s)\|_2 \right\}. \end{aligned} \quad (44)$$

Using the assumption of A0 in our paper, we have

$$\rho(f_\theta(x_s), z_j)^{\bar{r}} = \varphi(f_\theta(x_s) - z_j)^{\bar{r}}. \quad (45)$$

From $\|f_\theta(x_t)\| \leq M$ and $\|z_j\| \leq M$, we have

$$\rho(f_\theta(x_s), z_j)^{\bar{r}} \leq (2M\ell)^{\bar{r}}. \quad (46)$$

Similarly, we have

$$\|x_s - \bar{f}_{\bar{\theta}} f_\theta(x_s)\|_2 \leq 2M. \quad (47)$$

Thus, (44) leads to

$$\begin{aligned} |\psi_S - \psi_{S'}| &\leq \frac{1}{n} \left\{ \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-(2M\ell)^{\bar{r}}/\lambda_1) \right| \right. \\ &\quad \left. + \left| \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-(2M\ell)^{\bar{r}}/\lambda_1) \right| + 2\lambda_3 2M \right\} \\ &\leq \frac{2(2M\ell)^{\bar{r}} + 4\lambda_3 M}{n}. \end{aligned} \quad (48)$$

In the following, we consider the expectation of ψ_S with respect to the data set S , denoted by $E_S(\psi_S)$

$$\begin{aligned}
E_S(\psi_S) &= E_S\left(\sup_{\theta, \bar{\theta}, z_j} (L(\theta, \bar{\theta}, z_j) - \hat{L}_S^c(\theta, \bar{\theta}, z_j))\right) \\
&\stackrel{(1)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}, z_j} \left\{ \sum_{i=1}^n -\sigma_i \lambda_1 \ln \sum_{j=1}^k q_{j|i} \exp(-\rho(f_\theta(x_i), z_j)^{\bar{r}} / \lambda_1) + \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 \right\} \\
&\stackrel{(2)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}, z_j} \left\{ \sum_{i=1}^n \sigma_i \sum_{j=1}^k q_{j|i} \rho(f_\theta(x_i), z_j)^{\bar{r}} + \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2 \right\} \\
&\stackrel{(3)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}, z_j} \sum_{i=1}^n \sigma_i \sum_{j=1}^k q_{j|i} \rho(f_\theta(x_i), z_j)^{\bar{r}} + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sigma_i \|x_i - \bar{f}_{\bar{\theta}} f_\theta(x_i)\|_2.
\end{aligned} \tag{49}$$

In (49), the first inequality comes from the symmetrization of random variables, and the second inequality uses Jensen inequality from the fact that $-\ln x$ is a convex function. Note that the Lipschitz constant of the norm $\|\cdot\|$ is 1, but the function $\|x\|_r^r$ is not Lipschitz if the variable x takes the infinite values. However, with the assumptions we provide, there exists the constant M_1 such that $\rho(f_\theta(x_i), y_j)^{\bar{r}}$ is Lipschitz continuous. Using Lemma 2, we have

$$\begin{aligned}
E_S(\psi_S) &\stackrel{(1)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, z_j} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (f_\theta(x_i) - z_j)_t \\
&\quad + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \\
&\stackrel{(2)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (f_\theta(x_i))_t + 2E_{S, \sigma} \frac{1}{n} \sup_{z_j} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (-z_j)_t \\
&\quad + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \\
&\stackrel{(3)}{\leq} 2E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sqrt{2} M_1 (f_\theta(x_i))_t + 2E_{S, \sigma} \frac{1}{n} \sup_{z_j} \sum_{i=1}^n \sum_{t=1}^d \sigma_{it} \sum_{j=1}^k q_{j|i} \sqrt{2} M_1 (-z_j)_t \\
&\quad + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{i=1}^n \sum_{t=1}^m \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t
\end{aligned} \tag{50}$$

In (50), the first inequality comes from Lemma 2, and the second inequality uses the property of *sup*. The third inequality uses the fact that $\sum_{j=1}^k q_{j|i} = 1$. From (50), we have

$$\begin{aligned}
E_S(\psi_S) &\leq 2E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} \sqrt{2} M_1 (f_\theta(x_i))_t \right| \\
&\quad + 2E_{S, \sigma} \frac{1}{n} \sup_{z_j} \sum_{t=1}^d \sum_{j=1}^k |(z_j)_t| \left| \sum_{i=1}^n \sigma_{it} \sqrt{2} M_1 q_{j|i} \right| + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \right| \\
&\leq 2\sqrt{2} M_1 E_{S, \sigma} \frac{1}{n} \sup_{\theta} \sum_{t=1}^d \left| \sum_{i=1}^n \sigma_{it} (f_\theta(x_i))_t \right| + 2\sqrt{2} M_1 M E_{S, \sigma} \frac{1}{n} \sum_{t=1}^d \sum_{j=1}^k \left| \sum_{i=1}^n \sigma_{it} q_{j|i} \right| \\
&\quad + 2E_{S, \sigma} \frac{1}{n} \sup_{\theta, \bar{\theta}} \lambda_3 \sum_{t=1}^m \left| \sum_{i=1}^n \sqrt{2} \sigma_{it} (\bar{f}_{\bar{\theta}} f_\theta(x_i))_t \right|.
\end{aligned} \tag{51}$$

Further we have

$$\begin{aligned}
 E_S(\psi_S) &\leq 2\sqrt{2}M_1E_{S,\sigma}\frac{1}{n}\sup_{\theta}\sum_{t=1}^d\left|\sum_{i=1}^n\sigma_{it}(f_{\theta}(x_i))_t\right|+\frac{2\sqrt{2}M_1Mdk}{\sqrt{n}} \\
 &+2E_{S,\sigma}\frac{1}{n}\sup_{\theta,\bar{\theta}}\lambda_3\sum_{t=1}^m\left|\sum_{i=1}^n\sqrt{2}\sigma_{it}(\bar{f}_{\bar{\theta}}f_{\theta}(x_i))_t\right|.
 \end{aligned}
 \tag{52}$$

In (52), we use Lemma 3 to obtain the last inequality. From (52), we find that the upper bound of $E_S(\psi_S)$ depends on the number of anchors, and Rademacher complexities of the encoders and decoders. From (48) and (52), we obtain that with probability at least $1 - \tau$, the following inequality holds for $\theta, \bar{\theta}$, and z_j in proper parameter spaces by using McDiarmid inequality:

$$\hat{L}_S^c(\theta, \bar{\theta}, z_j) \leq L(\theta, \bar{\theta}, z_j) + 2\sqrt{2}M_1R_1 + 2\sqrt{2}R_2 + \frac{\chi_1 + \chi_2}{\sqrt{n}} \tag{53}$$

where $\chi_1 = \frac{2(2M)^r + 4\lambda_3M}{n}\sqrt{\frac{-\ln\tau}{2}}$, $\chi_2 = 2\sqrt{2}M_1Mdk$, $R_1 = E_{S,\sigma}\frac{1}{n}\sup_{\theta}\sum_{t=1}^d\left|\sum_{i=1}^n\sigma_{it}(f_{\theta}(x_i))_t\right|$, and $R_2 = E_{S,\sigma}\frac{1}{n}\sup_{\theta,\bar{\theta}}\lambda_3\sum_{t=1}^m\left|\sum_{i=1}^n\sigma_{it}(\bar{f}_{\bar{\theta}}f_{\theta}(x_i))_t\right|$.

E EXPERIMENTAL SETTING AND ADDITIONAL EXPERIMENTS FOR OUR MODEL

All experiments are conducted on a PC with an Intel Core i7- CPU and a RTX 3080 GPU. The structure of encoders we use in this paper is a fully-connected network with the form of $[m, 500, 500, 2000, d]$ and the decoder is a mirror of the encoder, where m is the dimension of input data and d is the dimension of the latent space. The popular ReLU functions are employed in each layer. We employ Adam (Kingma & Ba, 2015) as the backpropagation optimizer. In the experiments, ρ takes the L2 norm, $\bar{r} = 2$, and $\lambda_2 = 1000$. We let $\lambda_4 = 0$ due to the use of the same encoders. The parameters λ_1 and λ_3 are selected from the set $\{10^i, i = -3, -2, \dots, 2, 3\}$. In the classification experiments, we need to determine k in (6) in our paper. Note that x_1, \dots, x_n consist of the training set and $y_{j|i}(j = 1, \dots, k)$ are taken from the samples that have the same label as x_i . We think that the samples in the same class are neighbors. Thus, k will be determined by the number of samples in each class. As a result, k will vary since the number of samples in each class is different in the training set. In the clustering experiments, k in the model of (15) is set to be the number of clusters. We find that good performance can be obtained by this setting since the encoder has strong representations of features.

The outer loop is 10 iterations and the inner loop for autoencoders is run with an Adam optimizer for 100 epochs, with an initial learning rate of 0.001. In the data sets except for two large-scale data sets, all data sets are handled with a full-batch mode. For the large-scale data sets, the batch size is 2000.

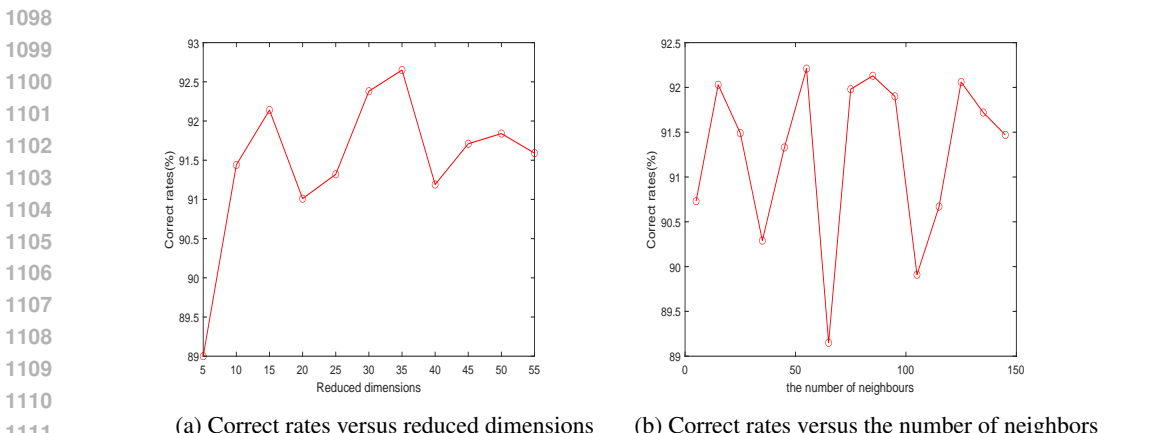
The data sets from the UCI repository are Dna (180 attributes /3 classes /2000 samples), Pendigits (16/10/7494), Satimage (36/6/4435), Iris(4/3/150), and WaveForm (21/3/5000). In addition, we also explore four face image data sets and an object data set. The ORL face database contains 40 distinct persons and each person has taken 10 different images. The UMIST face database contains 564 face images of 20 distinct subjects. The Yale face database contains 165 images of 15 individuals. The COIL database contains 1440 images with black background of 20 objects. The MSRA face data set consists of 1799 images of 12 subjects. All the images are normalized to a resolution of 32×32 pixels for computational efficiency. For each data set from the UCI repository, we randomly choose fifty percent samples to form the training set and the rest is used as the test set. The performance of each model is evaluated over twenty random splits of each data set. The additional five runs are employed to select the parameters of each model.

The MNIST data set contains 60,000 training samples and 10,000 test samples, and the FashionMNIST data set has 60,000 training samples and 10,000 test samples. The dimension of samples in these two data sets is 784.

For set-valued objects, we employ the features extracted from a pre-trained convolutional neural network (CNN), i.e., ReNet18, and the features are taken from the layer of res5b-relu. The extracted

1080 features of each image have a tensor representation of $7 \times 7 \times 512$ dimensions. To reduce the com-
 1081 putational cost, we pre-process the features of each image. That is, we perform the mean operation
 1082 along the first axis and downsample the features with a factor of 4 along the third axis. Thus, we ob-
 1083 tain the features whose dimensions are 7×128 . Namely, each image can be regarded as a set-valued
 1084 object containing seven examples with 128 dimensions. We randomly choose 50% of the samples
 1085 as the training set, 30% of the samples as the validation set, and the other images as the testing set.
 1086 Experimental results are averaged over 10 runs.

1087 For the large-scale FashionMNIST data set, we carry out the experiment to check the effect of
 1088 reduced dimensions and the number of neighbors. Figure 3 (a) denotes the correct rates of DCWSL
 1089 with the change of reduced dimensions, and Figure 3 (b) shows the correct rate of DCWSL as the
 1090 number of neighbors varies. From Figure 3 (a), we observe that the reduced dimensions affect
 1091 the performance of DCWSL. But when the dimensions of the samples exceed 10, our model can
 1092 achieve good performance. From Figure 3 (b), we can see that it is not necessary to employ too
 1093 many neighbors to obtain good better performance since we consider the samples from the same
 1094 class. In addition, we visualize 2000 samples in a two-dimensional space via t-SNE. Figure 4 shows
 1095 the experimental results. Figure 4 (a) denotes the visualization of original images via t-SNE, and
 1096 Figure 4 (b-d) denote the results of DCWSL in the case of different iterations. As can be seen from
 1097 Figure 4, the embedding features in a two-dimensional space from DCWSL are well separated.



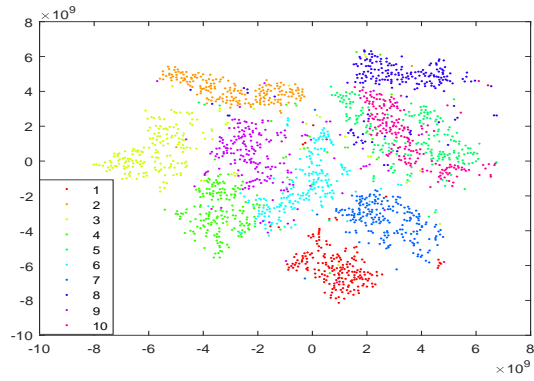
(a) Correct rates versus reduced dimensions (b) Correct rates versus the number of neighbors

Figure 3: Performance of our model on the FashionMNIST data set

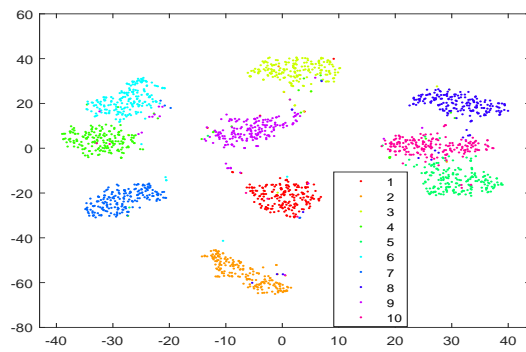
1115 There are several parameters in the proposed model since these parameters determine trade-offs in
 1116 several terms. We let $\lambda_2=1000$ so that $\{p_i|i = 1, \dots, n\}$ approach uniform distributions. If we
 1117 consider that y_1, \dots, y_k are taken from x_1, \dots, x_n , then we set $\lambda_4 = 0$. In such a case, we first
 1118 explore the effect of different λ_1 and λ_3 on supervised learning tasks. To this end, we randomly
 1119 choose half of samples from each person to form the training set and others are used for testing
 1120 on the ORL data set. Assume that the reduced dimension is equal to the number of classes (40)
 1121 and the parameters λ_1 and λ_3 take values from $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$. Thus each
 1122 parameter takes seven values. We also report the experimental results over ten runs. Figure 5 shows
 1123 the experimental results on the classification problem, where the x -axis denotes the parameter λ_3 ,
 1124 the y -axis denotes the parameter λ_1 , and the z -axis denotes the error rate of our model.

1125 As can be seen from Figure 5, the error rates of the proposed model vary with the change of pa-
 1126 rameters. It is found that the error rates of our model are very high when the parameter λ_3 takes
 1127 relatively small values. We observe that λ_3 is more sensitive than λ_1 in affecting the performance of
 1128 the model. From Figure 5, we see that the running time of our model is affected by the parameters.
 1129 Figure 6 shows the experimental results on the clustering problem. From Figure 6, we find that the
 1130 parameters affect the performance of DCWSP in the clustering problem. Overall, the experiments
 1131 indicate that we need to select proper parameters to attain the best performance in real applications.
 1132 In fact, the cross-validation is often employed to select optimal parameters.

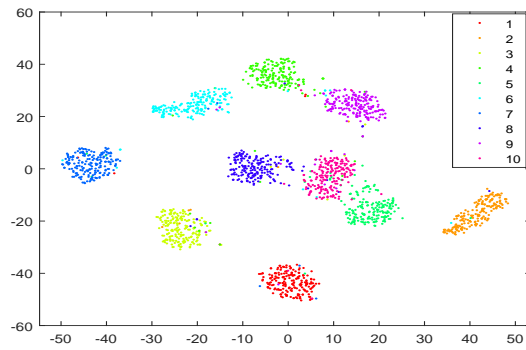
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



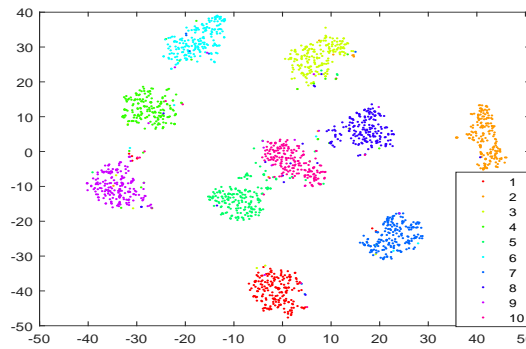
(a) visualization of original images via t-SNE



(b) visualization from our model at the first iteration via t-SNE



(c) visualization from our model at the fifth iteration via t-SNE



(d) visualization from our model at the tenth iteration via t-SNE

Figure 4: Visualization of 2000 images on the FashionMNIST data set

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

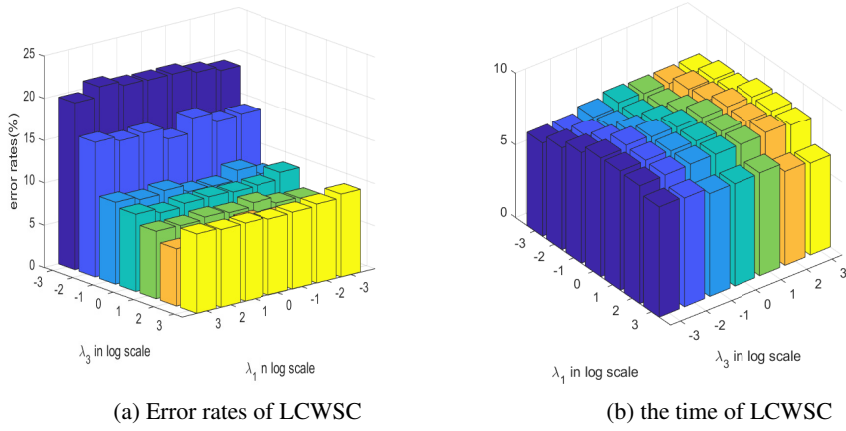


Figure 5: Performance of LCWSC with varying parameters

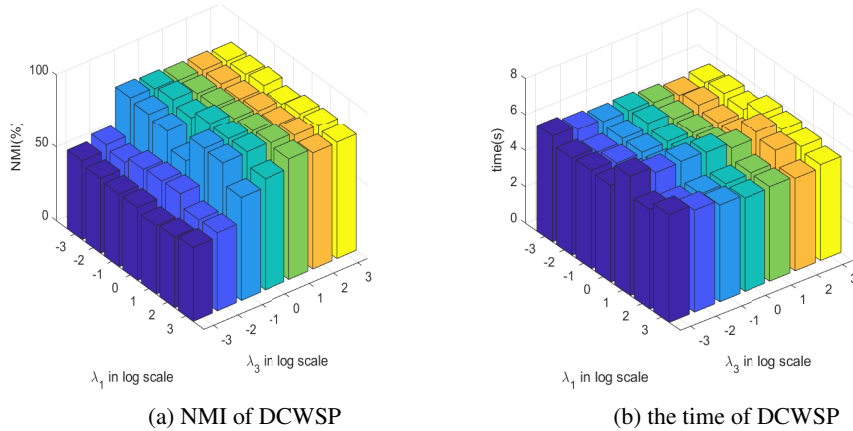


Figure 6: Performance of DCWSP with varying parameters