

Unified Multi-Domain Learning and Data Imputation using Adversarial Autoencoder

Andre Mendes
New York University
New York, USA
andre.mendes@nyu.edu

Julian Togelius
New York University
New York, USA
julian.togelius@nyu.edu

Leandro dos Santos Coelho
Pontifical Catholic University of Parana
Federal University of Parana
Curitiba, Brazil
leandro.coelho@pucpr.br

Abstract—We present a novel framework that can combine multi-domain learning (MDL), data imputation (DI) and multi-task learning (MTL) to improve performance for classification and regression tasks in different domains. The core of our method is an adversarial autoencoder that can: (1) learn to produce domain-invariant embeddings to reduce the difference between domains; (2) learn the data distribution for each domain and correctly perform data imputation on missing data. For MDL, we use the Maximum Mean Discrepancy (MMD) measure to align the domain distributions. For DI, we use an adversarial approach where a generator fill in information for missing data and a discriminator tries to distinguish between real and imputed values. Finally, using the universal feature representation in the embeddings, we train a classifier using MTL that given input from any domain, can predict labels for all domains. We demonstrate the superior performance of our approach compared to other state-of-art methods in three distinct settings, DG-DI in image recognition with unstructured data, MTL-DI in grade estimation with structured data and MDMTL-DI in a selection process using mixed data.

Index Terms—multi-task, multi-domain, data imputation

I. INTRODUCTION

Many real-world problems in machine learning (ML) require enough labeled data for training classifiers to address a task at hand. Whenever this data is either not available, or hard to collect, an alternative is to use available datasets that are related to the task. However, although related, these other tasks often belong to a different domain and training classifiers on them without addressing their differences result in sub-optimal performance. To address this case, multi-domain learning (MDL) [1] approaches have been proposed to reduce the gap among domains and create more general classifiers.

MDL has different applications, such as Domain Generalization (DG) [2] when the goal is to use source datasets to learn a classifier for an unknown target domain and Domain Adaptation (DA) [3] when the target domain is known but the labels are not available or their sample size is small. In this work, we focus on a case of MDL where datasets and labels for different domains are available. However, all datasets have relatively small sample sizes. Therefore, all datasets are used as source and target and the knowledge has to be shared in all directions. We propose a framework based on an autoencoder [4] that can produce domain-invariant encoded spaces. We adopt the maximum mean discrepancy (MMD) [5]

to measure differences between generated encoded spaces, and we reduce MMD during training.

Another important challenge in ML is dealing with incomplete datasets where a significant part of the data is missing. One solution would be to remove the damaged samples, but this is often impractical because of the reduction in sample size. To address this case, Data Imputation (DI) [4], [6], [7] methods try to learn the data distribution and generate values to fill in for missing data. In our approach, we use an adversarial strategy to incorporate DI in our autoencoder, combining MDL and DI to jointly address missing data and domain differences.

Additionally to this unsupervised approach, we also proposed a supervised method to incorporate the labels available in the training process. To accomplish that, we adapt a multi-task learning (MTL) [8], [9] framework in which a single classifier is trained to make predictions for all tasks/domains. By learning the tasks together, the classifier is exposed to more samples in different domains, allowing it to improve generalization and perform better than classifiers trained individually. Therefore, the main contributions of this work are:

- We present an adversarial autoencoder to perform DI and MDL in the same training process. The process is unsupervised and the encoded space and features generated can be used in different downstream tasks.
- We present a supervised framework that combines the autoencoder with an MTL approach to train a classifier for all domains/tasks together.
- We address an important case of MDL where there is no significant source/target dataset and all datasets have to share knowledge among them.
- By performing experiments in different settings, we show that our method can be applied to DG, MTL, and MDL using the same structure. We also show empirically that it outperforms other state-of-art methods in all settings.

This paper is organized as follows: Section II shows the related work. Section III introduces our notation notation and problem statement. Section IV explains the method and Section V describes the experiments for validation. Section VI presents the results and a conclusion and future work are discussed in Section VII.

II. RELATED WORK

Here we describe related work in MDL, DI and MTL.

A. Multi-Domain Learning (MDL)

While traditional ML methods train classifiers to generalize well for a specific domain, MDL trains classifiers that can predict outputs across domains with different data distribution [1]. For example, [10] uses metadata (semantic descriptors), to explore information on task and domain relatedness. They construct a two-sided network to learn representations from the original input vector and the metadata by minimizing the empirical risk for all domains/tasks. MDL has vast applicability in different contexts [10], [11].

DG uses data from seen source domains to create models that can generalize well to unseen target domains. The work in [2] uses a domain-invariant component analysis method with a kernel-based optimization algorithm to minimize the dissimilarity across domains. Other approaches, such as [12], perform DG and DA using a contrastive loss to guide samples from the same class being embedded nearby in latent space across data sources. Similar to our method, the work in [13] uses an autoencoder, expanding it to a multi-task approach to extract invariant features that are robust to domain variations. In [14], adversarial autoencoders are extended to DG by imposing MMD [5] metric to align multi-domain distributions.

DA [3] is similar to DG, however, in this case the target data is known. DA methods attempt to minimize the shift between source and target distributions in input space [15], feature space [16], [17], or output space [18], using MMD [19] or adversarial learning (GAN) [20]. DA has been applied to computer vision [15], [17], structured data [21] and text [17].

B. Data Imputation (DI)

Methods for DI range from simple column average to complex imputations, and they can be categorized as discriminative or generative. Discriminative methods include MICE [6] and matrix completion [22]. Generative methods include algorithms based on Expectation Maximization [23] and on deep learning (DL) [4], [24], [25].

DL methods can capture the latent structure of high-dimensional data by efficiently emulating complex distributions [26]. Many of the algorithms proposed in this setting use autoencoder based structures [4], [27] or adversarial generative [20] based methods [24], [25]. While most of these approaches focus on computer vision tasks, the work in [28] focuses on solving DI for arbitrary datasets, including tabular data, which are often incomplete and heterogeneous.

In our approach, we deal with structured and unstructured data and our method is similar to [29], which explores the GAN approach to create: a generator to accurately impute missing data; a discriminator to distinguish between observed and imputed components; and a hint mechanism to help generate samples according to the true underlying data distribution.

C. Multi-Task Learning (MTL)

The goal in MTL [8], [9] is to learn multiple related tasks simultaneously so that knowledge obtained from each task can be re-used by the others. This can increase the sample size for each task, making MTL beneficial for tasks with small

training sample. MTL has been applied in different approaches including neural nets (NN) [9], kernel methods [30], and deep neural nets (DNN) [31].

Regularization is important in MTL since it controls the relationship among tasks and features and prevents overfitting. The work in [32] presents an overview of regularization in MTL, such as group-sparse based [33] and robust feature based [30] regularization. Other methods have explored tensor factorization [34] and evolutionary algorithms [35] to find optimal architectures for MTL networks.

III. PROBLEM STATEMENT

Let's define a domain \mathcal{D} with feature space \mathcal{X} and a marginal probability distribution $P(\mathbf{X})$. Each $\mathbf{x}_i \in \mathbb{R}^d$ is the i -th feature vector (sample), n is the number of samples and d is the number of dimensions in $\mathbf{X} \in \mathbb{R}^{n \times d}$. The i -th one-hot encoding label sample is represented by $\mathbf{y}_i \in \mathbb{R}^{d_y}$ with d_y as the number of classes. We also define a predictive function $f(\cdot)$. Therefore, we define a domain as $\mathcal{D} = \{\mathcal{X}, P(\mathbf{X})\}$ and a task as $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}$.

A. Multiple Domains

Let's define S as the number of domains, where the index of a domain is represented by $s = \{1, \dots, S\}$. In our MDL approach, we use data from multiple domains $\{\mathcal{D}^s\}_{s=1}^S$ to create a function $p(\cdot)$ that produces a domain-invariant space with encoded features $\mathbf{E} \in \mathbb{R}^{n \times d_e}$, where d_e is the dimension for the feature space. Ideally, classifiers trained on this domain-invariant space can generalize to all domains.

B. Missing Values

In many cases, a feature space from a source domain is incomplete. Given an incomplete feature space \mathbf{X} with missing values, we want to find a function $g(\cdot)$ that can learn the data distribution in the domain and impute values that are as close as possible to the original ones.

More formally, we can define a new feature space $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times d}$ by replacing the missing values with a random noise z , not previously observed in \mathbf{X} . We also define $\mathbf{M} \in \mathbb{R}^{n \times d}$ as a mask vector taking values in $\{0, 1\}^d$ that indicates the components of $\tilde{\mathbf{X}}$ that are observed. This process to create $\tilde{\mathbf{X}}$ and \mathbf{M} is represented by the symbol \tilde{I} in Fig. 1.

For imputation, we want to generate samples according to $P(\mathbf{X}|\tilde{\mathbf{X}} = \tilde{\mathbf{x}}^i)$, for each i , to fill in the missing data. Our goal is to approximate $g(\cdot)$ and $p(\cdot)$ in the same autoencoder structure, achieving MDL and DI in the same framework.

C. Multiple Tasks and Small Datasets

An additional problem we face on real-world data is datasets with a relatively small sample size n .

Consider that features \mathbf{X}^s , and labels, \mathbf{Y}^s for S domains, $s = \{1, \dots, S\}$, are available and that the sample size, n^s for each domain is relatively too small. This can cause classifiers, $\{f^s(\cdot)\}_{s=1}^S$, to overfit the data when trained individually.

We aim to train a classifier $f(\cdot)$ that can learn simultaneously from all datasets. Ideally, given a feature input \mathbf{X}^s from

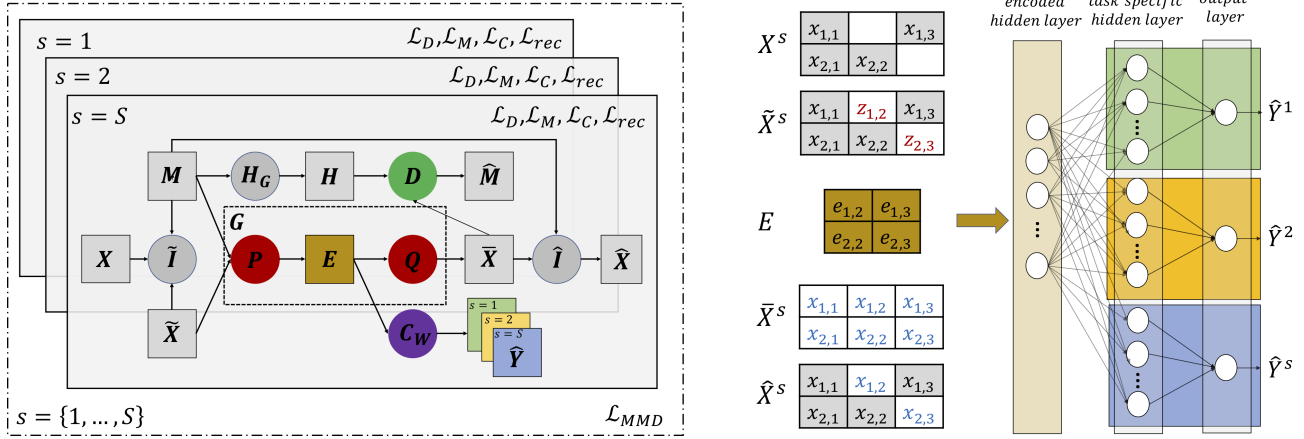


Fig. 1. MD2I framework. Left: Overview of the process. Each square represents a dataset. Each circle represents a function that receives and processes input datasets and produces modified output datasets. Components in gray are specific to the current domain, while color components are trained and updated across all domains. Given a domain s , a dataset \mathbf{X}^s is fed to the adversarial autoencoder. A generator $G = (P, Q)$ and a discriminator D are trained to learn the representation of the data and produce values for data imputation. At the same time, the encoder component P is trained to generate encoded values \mathbf{E} that are domain-invariant. Finally, the parameters \mathbf{W} of the classifier C are trained using MTL to predicted labels for all domains. For each domain, we calculate the reconstruction, classification and adversarial losses. The MMD loss is calculated after all encoded values are produced. Right: Representation of different datasets along the process and MTL structure taking encoded values as input and producing labels, $\{\hat{\mathbf{Y}}^s\}_{s=0}^S$, as output.

a domain s , $f(\cdot)$ can correctly predict labels in all domains, $\{\mathbf{Y}^s\}_{s=1}^S$, and the performance for $f(\cdot)$ in each domain is better than using classifiers trained separately. This can be achieved using an MTL framework that seeks to improve the generalization performance of multiple related tasks by learning them jointly. Hence, we aim to combine $p(\cdot)$, $g(\cdot)$ and $f(\cdot)$, so that the encoded feature space and a final classifier are learned simultaneously.

IV. METHODS

To achieve MDL and DI, we propose an autoencoder using a generator $G = (Q, P)$. Q is an encoder that maps values from an input $\mathbf{X} \in \mathbb{R}^{n \times d}$ to a encoded state $\mathbf{E} \in \mathbb{R}^{d \times d_e}$. P is a decoder that reconstructs the values from \mathbf{E} back to $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times d}$. A successful autoencoder is able to produce values for $\tilde{\mathbf{X}}$ that are as close as possible to \mathbf{X} .

Given S different, but related domains, our goal is to create an autoencoder that can generate the input features for all domains. This can be achieved if \mathbf{E} produced by G is a domain-invariant encoded state (See Fig 1).

A. Multi-Domain Learning

To make hidden codes \mathbf{E} invariant underlying different domains, we first need to define a metric to measure domain similarity. We adopt MMD [5], which finds the largest difference in expectations over functions in the unit ball of a reproducing kernel Hilbert space (RKHS) [36]. Hence, we introduce an MMD-based loss term $\mathcal{L}_{mmd}(\mathbf{E}^1, \dots, \mathbf{E}^S)$ on \mathbf{E}^s 's among different domains. Similar to [14], given two

domains \mathbf{E}^s and \mathbf{E}^t , with number of samples d^s and d^t respectively, the MMD loss is given by:

$$\begin{aligned} MMD(\mathbf{E}^s, \mathbf{E}^t)^2 &= \left\| \left(\frac{1}{d^s} \sum_{i=1}^{d^s} \phi(\mathbf{E}_i^s) - \frac{1}{d^t} \sum_{i=1}^{d^t} \phi(\mathbf{E}_i^t) \right) \right\|_{\mathcal{H}}^2 \\ &= \frac{1}{(d^s)^2} \sum_{i=1}^{d^s} \sum_{j=1}^{d^s} k(\mathbf{E}_i^s, \mathbf{E}_j^s) + \frac{1}{(d^t)^2} \sum_{i=1}^{d^t} \sum_{j=1}^{d^t} k(\mathbf{E}_i^t, \mathbf{E}_j^t) \\ &\quad - \frac{2}{d^s d^t} \sum_{i=1}^{d^s} \sum_{j=1}^{d^t} k(\mathbf{E}_i^s, \mathbf{E}_j^t). \end{aligned} \quad (1)$$

We use a map operation $\phi(\cdot)$ to project the domain into an RKHS \mathcal{H} . The work in [37] shows that if an arbitrary distribution of the features, represented by a kernel embedding technique $k(\cdot, \cdot)$ [38] is characteristic, then the mapping to the \mathcal{H} is injective. This injectivity indicates that the probability distribution is uniquely represented by an element in RKHS.

Therefore, we have a kernel function $k(\mathbf{e}_i^s, \mathbf{e}_j^t) = \phi(\mathbf{e}_i^s) \phi(\mathbf{e}_j^t)^\top$ induced by $\phi(\cdot)$. For this function, we use the RBF kernel with bandwidth σ given by

$$k(\mathbf{e}_i^s, \mathbf{e}_j^t) = \exp\left(-\frac{1}{2\sigma} \|\mathbf{e}_i^s - \mathbf{e}_j^t\|^2\right). \quad (2)$$

Finally, [14] shows it is possible to minimize the upper bound of the distribution variance among domains using

$$\mathcal{L}_{mmd}(\mathbf{E}^1, \dots, \mathbf{E}^S) = \frac{1}{S^2} \sum_{1 \leq i, j \leq S} MMD(\mathbf{E}^i, \mathbf{E}^j). \quad (3)$$

B. Adversarial Data Imputation

To also perform data imputation, we modify our generator G with an adversarial component based on [29]. More

specifically, \mathbf{X} is the original input with missing values. In order to indicate which positions are observed and which are missing, we define a variable $\mathbf{M} \in \mathbb{R}^{n \times d}$, with $m_{i,j} = 1$ if $x_{i,j}$ is observed and $m_{i,j} = 0$, otherwise. We then create the random variable $\tilde{\mathbf{X}}$, where missing values, indicated by \mathbf{M} , are replaced by noise from a noise variable \mathbf{Z}_m .

Given realizations of $\tilde{\mathbf{X}}$, \mathbf{M} and \mathbf{Z}_m as input, G produces a vector of imputations $\bar{\mathbf{X}}$,

$$\bar{\mathbf{X}} = G(\tilde{\mathbf{X}}, \mathbf{M}, (\mathbf{1} - \mathbf{M}) \odot \mathbf{Z}_m). \quad (4)$$

To create the final output $\hat{\mathbf{X}}$, we copy the original \mathbf{X} and replace the positions for missing values with values from $\bar{\mathbf{X}}$,

$$\hat{\mathbf{X}} = \mathbf{M} \odot \tilde{\mathbf{X}} + (\mathbf{1} - \mathbf{M}) \odot \bar{\mathbf{X}}. \quad (5)$$

The process in Eq. 5 is represented by \hat{T} in Fig. 1.

Inspired in [29], we introduce a discriminator, D , that will be used as an adversary to train G . However, instead of predicting if the entire dataset produced by G is fake or real, we use D to predict which values in $\hat{\mathbf{X}}$ are observed and which are imputed. Therefore, a successful D produces values in $\hat{\mathbf{M}}$ that are as close as possible to the real values in \mathbf{M} .

We also add a hint mechanism to limit the distributions produced by G that would be optimal for D . This hint mechanism is a random variable, \mathbf{H} that depends on \mathbf{M} . For each (imputed) sample $(\hat{\mathbf{x}}, \mathbf{m})$, we draw \mathbf{h} according to the distribution $\mathbf{H}|\mathbf{M} = \mathbf{m}$. We pass \mathbf{h} as an additional input to the discriminator so that the i -th component of $D(\hat{\mathbf{x}}, \mathbf{h})$ corresponds to the probability that the i -th component of $\hat{\mathbf{x}}$ was observed conditional on $\hat{\mathbf{X}} = \hat{\mathbf{x}}$ and $\mathbf{H} = \mathbf{h}$.

To generate \mathbf{H} , we use a random variable \mathbf{Z}_h defined by first sampling k from $1, \dots, d$ uniformly at random, such as

$$\mathbf{H} = \mathbf{Z}_h \odot \mathbf{M} + 0.5(\mathbf{1} - \mathbf{Z}_h). \quad (6)$$

We can define a minimax problem where we train D to maximize the probability of correctly predicting \mathbf{M} and train G to minimize the probability of D predicting \mathbf{M} .

$$\min_G \max_D \mathbb{E}_{\hat{\mathbf{X}}, \mathbf{M}, \mathbf{H}} \left[\mathbf{M}^T \log D(\hat{\mathbf{X}}, \mathbf{H}) + (\mathbf{1} - \mathbf{M})^T \log (1 - D(\hat{\mathbf{X}}, \mathbf{H})) \right],$$

where \log is element-wise logarithm and dependence on G is through $\hat{\mathbf{X}}$. Writing $\hat{\mathbf{M}} = D(\hat{\mathbf{X}}, \mathbf{H})$, we can define

$$\min_G \max_D \mathbb{E} [\mathcal{L}(\mathbf{M}, \hat{\mathbf{M}})], \quad (7)$$

where \mathcal{L} is the cross-entropy function defined by

$$\mathcal{L}(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^d \left[a_i \log(b_i) + (1 - a_i) \log(1 - b_i) \right]. \quad (8)$$

C. Multi-Task Classification

To incorporate label information into the training process, we define a classifier C with parameters $\mathbf{W} \in \mathbb{R}^{d_e \times m}$ that is trained with the encoded values \mathbf{E} . By using MTL, we train C to predict the labels $\{\mathbf{Y}^s\}_{s=1}^S$ in each domain in parallel and we use the loss for each task to update C and G .

To introduce sparsity into the model to reduce model complexity, we use a simplification of Lasso in MTL [32]. Such regularization controls the sparsity shared among all tasks, assuming that different tasks share the same sparsity parameter. The equation to optimize C is given by

$$\min_{\mathbf{W}} \sum_{s=1}^S L_C(\mathbf{W}^T \mathbf{E}, \mathbf{Y}^s) + \rho_0 \|\mathbf{W}\| + \rho_{L2} \|\mathbf{W}\|_F^2, \quad (9)$$

where L_C is a standard cross-entropy function, \mathbf{E} are the encoded values produced by Q and \mathbf{Y}^s denotes its corresponding labels in domain s . The regularization parameter ρ_0 controls sparsity, and ρ_{L2} controls the l_2 -norm penalty.

D. Training Procedure

Let's recap that our generator G generates values $\bar{\mathbf{X}}$ for all the positions in \mathbf{X} . Therefore, we can use the real and predicted values to calculate the reconstruction error as

$$\mathcal{L}_{rec} = \sum_{i=1}^n \sum_{j=1}^d m_{ij} L(\hat{x}_{ij}, x_{ij}), \quad (10)$$

where

$$L_{rec}(x_{ij}, \hat{x}_{ij}) = \begin{cases} (\hat{x}_{ij} - x_{ij})^2, & \text{if } x_{ij} \text{ is continuous,} \\ -x_{ij} \log(\hat{x}_{ij}), & \text{if } x_{ij} \text{ is binary.} \end{cases}$$

The loss in equation 10 ensures that generated values for observed components ($m_{ij} = 1$) are close to those observed.

We also optimize the generator G to ensure that the imputed values for missing components ($m_i = 0$) successfully fool the discriminator using

$$\mathcal{L}_M = - \sum_{i:m_i=0} (1 - m_i) \log(\hat{m}_i), \quad (11)$$

For the discriminator D , we want to ensure that it does not overfit to the hint mechanism. Therefore we only train D to give the outputs that depend on G , i.e. $m_i = 0$, using

$$\mathcal{L}_D = \sum_{i:m_i=0} \left[m_i \log(\hat{m}_i) + (1 - m_i) \log(1 - \hat{m}_i) \right] \quad (12)$$

\mathcal{L}_M is minimized when D fails to correctly reproduce \mathbf{M} and \mathcal{L}_{rec} is minimized when G can produce values close to the observed values.

The final loss for parameter optimization is given by

$$\min_G \max_D \sum_{s=1}^S \left[\mathcal{L}_{rec} + \lambda_0 \mathcal{L}_M + \lambda_1 \mathcal{L}_D + \lambda_2 \mathcal{L}_C \right] + \lambda_3 \mathcal{L}_{mmd}, \quad (13)$$

where hyperparameters $\{\lambda_i\}_{i=0}^3$ can be defined using cross-validation (See Section V-E for details on hyperparameters)

The overall algorithm is described in Algorithm 2. In lines 5, 10 and 14, Algorithm 2 calls the procedure shown in Algorithm 1 to generate the necessary variables. In Algorithm 2, lines 3 to 7 are related to the optimization of D . Lines 8 to 12 are the steps to optimize G using the optimized D . In lines 13 to 15, the parameters \mathbf{W} for the classifier C are updated using the optimized G . Lines 17 and 18, guarantee that G is

Algorithm 1 AAE

Input: $\mathbf{X}, \mathbf{Z}_m, \mathbf{Z}_h, G$ **Output:** $\tilde{\mathbf{X}}, \hat{\mathbf{X}}, \mathbf{H}, \mathbf{M}, \mathbf{E}$

- 1: $\tilde{\mathbf{X}}, \mathbf{M} \leftarrow \tilde{I}(\mathbf{X}, \mathbf{Z}_m)$
 - 2: $\tilde{\mathbf{X}}, \mathbf{E} \leftarrow G(\tilde{\mathbf{X}}, \mathbf{M}, (\mathbf{1} - \mathbf{M}) \odot \mathbf{Z}_m)$
 - 3: $\hat{\mathbf{X}} \leftarrow \mathbf{M} \odot \tilde{\mathbf{X}} + (\mathbf{1} - \mathbf{M}) \odot \tilde{\mathbf{X}}$
 - 4: $\mathbf{H} \leftarrow \mathbf{Z}_h \odot \mathbf{M} + 0.5(\mathbf{1} - \mathbf{Z}^h)$
-

Algorithm 2 MD2I

Input: $\mathbf{X} = \{\mathbf{X}^1, \dots, \mathbf{X}^S\}, \mathbf{Y} = \{\mathbf{Y}^1, \dots, \mathbf{Y}^S\}, \mathbf{Z}_m, \mathbf{Z}_h, \rho_1, \rho_{L2}, \alpha, \{\lambda_1, \dots, \lambda_4\}$ **Parameter:** Initialized parameters $G = (Q, P), D, W$ **Output:** G^*, D^*, W^*

- 1: **while** training loss not converge **do**
 - 2: **for** $s = 1$ to S **do**
 - 3: draw B batches with n_b samples $\{\mathbf{X}_b^s\}$ from \mathbf{X}
 - 4: **for** $b = 1$ to B **do**
 - 5: $\tilde{\mathbf{X}}^s, \hat{\mathbf{X}}^s, \mathbf{H}^s, \mathbf{M}^s, \mathbf{E}^s \leftarrow AAE(\mathbf{X}^s, \mathbf{Z}_m, \mathbf{Z}_h, G)$
 - 6: Update D using SGD and equation 12
 - 7: **end for**
 - 8: draw B batches with n_b samples $\{\mathbf{X}_b^s\}$ from \mathbf{X}
 - 9: **for** $b = 1$ to B **do**
 - 10: $\tilde{\mathbf{X}}^s, \hat{\mathbf{X}}^s, \mathbf{H}^s, \mathbf{M}^s, \mathbf{E}^s \leftarrow AAE(\mathbf{X}^s, \mathbf{Z}_m, \mathbf{Z}_h, G)$
 - 11: Update G using SGD and equations 10 and 11.
 - 12: **end for**
 - 13: draw 1 batch with n_b samples $\{\mathbf{X}_b^s\}$ from \mathbf{X}
 - 14: $\tilde{\mathbf{X}}^s, \hat{\mathbf{X}}^s, \mathbf{H}^s, \mathbf{M}^s, \mathbf{E}^s \leftarrow AAE(\mathbf{X}^s, \mathbf{Z}_m, \mathbf{Z}_h, G)$
 - 15: Update W and G using SGD and equation 9
 - 16: **end for**
 - 17: Calculate loss \mathcal{L}_{mmd} using equation 3
 - 18: Update G using SGD and \mathcal{L}_{mmd}
 - 19: **end while**
-

also optimized to create domain-invariant encoded spaces. All optimizations use stochastic gradient descent (SGD).

One can notice that the method can be used in an unsupervised way by skipping the update in lines 13 to 15. This will produce results that are independently of labels in the dataset.

V. EXPERIMENTS

We perform experiments on real-world problems and compare our framework with other methods in distinct settings.

A. DG-DI in Image Recognition with unstructured data

We first apply our method to image recognition using the MNIST [39] dataset. To create the dataset for DG, we used the method proposed in [13] by producing digit images from six different angles. We use 1,000 digit images of ten classes (with 100 images per class) to represent the basic view. We denote the digit images with 0° by I_0 , and rotate the images in a counter-clockwise direction by $15^\circ, 30^\circ, 45^\circ, 60^\circ$ and 75° .

To create the dataset for DG-DI, we remove a square patch of the size of 13×13 from all images (each image is 28

$\times 28 = 784$ pixels). The location of the patch was uniformly sampled for each image. This creates a missing completely at random (MCAR) case, where there is no dependency on any of the variables. We also perform normalization to the range $[0, 1]$ on pixels and feed the final values to the methods.

For validation, we apply the leave-one-domain-out strategy, using one of the domains as target and training the model in all other domains as a source. We apply the missing mask in both source and target. Finally, we evaluate the methods in terms of classification using the accuracy score. We repeat the experiments for 50 times and report the averaged results.

B. MTL-DI in Regression with structured data

In this setting, we use the school dataset¹, which consists of exam grades from 15,362 students in 139 schools. This is a structured dataset with categorical and numerical features, where categorical attributes have been expressed as binary attributes using one-hot encoding. Our goal is to predict each student's exam grade, hence, the school is the task and exam score is the output, creating an MTL setting.

For MTL-DI, we introduce missingness by appending a uniform random vector \mathbf{r} with values between 0 and 1. Each observation in r corresponds to an observation in the dataset. We then randomly sample two attributes x_1 and x_2 from the dataset and calculate their median mx_1 and mx_2 . Finally, we pick half of the attributes, at random, to have missing values where $r_i \leq t, i \in 1 : n$ and $(x_1 \leq mx_1$ or $x_2 \geq mx_2)$, where n is the number of observations and t is a threshold, that we define as 30%. This creates a missing at random (MAR) case, where the missingness depends on the observed variables.

For validation, we perform 50 experiments with a 50/50% training/test split in each one. Methods are compared in terms of Root Mean Squared Error (RMSE) on the test set.

C. MDMTL-DI in Classification with mixed data

We also evaluate our method in a real-world problem with mixed data, i.e. structured data with continuous and categorical features, and text data. Our information comes from 3 different companies that use selection processes with a similar structure.

For privacy requirements, we refer to the companies with indexes, such that C_1 refers to Company 1. C_1 is an organization that selects students to receive fellowships. C_2 is a big retail company that selects students to work in the company after graduation. C_3 is a governmental agency that selects students for positions in public administration. Table I presents the description for each stage in each process. All processes happen yearly and we have data from two years.

Each stage in the process contains its own set of variables. The data collected in each process is very similar in stages such as *Demographics* (Demo) and *Education* (Edu), both in the form of content and structure. For stages with open questions such as *Video Submission* and *Experience* (Exp), each process has its own specific questions.

¹Available at <http://multilevel.ioe.ac.uk/intro/datasets.html>

TABLE I
STAGES IN THE SELECTION PROCESSES

Stages	C_1	C_2	C_3
Demo	Provide country, city, age.	Same as C_1 .	Same as C_1 .
Edu	Provide university, major, minor and extra activities.	Same as C_1 .	Same as C_1 .
Profile Test	Online tests to measure profile characteristics such as ambition and interests.	Online tests to measure big 5 characteristics	Same as C_1 .
Exp	Write about professional experience (Open format).	Same as C_1	Same as C_1 .
Logic Test	Perform online tests to map levels in problem-solving involving logic puzzles.	Same goal as C_1 but with specific test for C_2	Same goal as C_1 but with specific test for C_3
Video Submission	2-min, explaining why they deserve the fellowship.	5-min, making a case to be selected for the position	1-min, explaining a problem that motivates the applicant
Video Evaluation	Applicants are evaluated based on their entire profile.	Same as C_1 but different criteria.	Same as C_1 but different criteria.

1) *Feature preparation and missing data*: Categorical variables are converted to numerical using a one-hot encoder. In stages such as video, the speech is extracted and the data is used as a text. To convert text data to numerical values, we create word embeddings using Word2Vec [40], where we assign high-dimensional vectors to words in a text corpus but preserving their syntactic and semantic relationships.

Approximately, 30% of the information is missing across all datasets. The source of the missing data is not available, but from different observations, we estimate that it has cases of both MCAR and MAR.

2) *Validation and performance metrics for selection*: Our goal is to predict which applicants will be selected in the *Video Evaluation* stage. Although the real evaluation in each process happens independently, we believe the datasets are related due to the similar structure of the features and the common focus on undergraduate students in their last year.

We perform longitudinal experiments, in which we use a $year_1$ as a training and test set and $year_2$ as a validation set. We also use the standard cross-validation (CV) setting by using data from a specific year and splitting it into train, test, and validation. We repeat this process 3 times, one for each year and one for a dataset with both years combined.

We use a 10-fold CV in each experiment with a total of 50 runs. Furthermore, we compare the models in terms of F1-score for the positive class, which balances precision and recall with a focus on the selected applicants.

D. Benchmark Methods

In our baseline method (B1), we train an individual classifier for each domain. The second baseline (B2) uses the aggregation of all domains into one and trains one individual classifier for all domains. For missing data, we perform mean-imputation, which imputes the mean of each continuous attribute and the mode of each discrete attribute.

For comparison in DI, we repeat the baselines (B1, B2) but use GAIN [29] to perform data imputation. Next, we compare with other MDL methods, using: D-MTAE [13], MMD-AAE [14] and CCSA [12]. At last, we use methods that

TABLE II
BENCHMARK METHODS AND THEIR COMPONENTS.

Method	Group	Noise	Labels	MTL	Domain
B1	Baseline		x		
B2	Baseline		x		
B1-DI	DI	x	x		
B2-DI	DI	x	x		
CCSA	MDL		x		x
D-MTAE	MDL	x			x
MMD-AAE	MDL		x		x
DMTRL	MTL		x	x	
UMDMTL	MDMTL		x	x	x
MD2I-U	Ours	x			x
MD2I-S	Ours	x	x	x	x

address MTL such as DMTRL [34], and UMDMTL (Unified MD and MTL) [10] that addresses both MTL and MDL.

For all methods, we use network specification, hyperparameters, and initialization as reported in the respective papers. For prediction in D-MTAE, Baseline and DI methods, we use an NN with 2 hidden layers with cross-entropy loss for classification and RMSE for regression.

In Table II, we present all methods, the components that they address and their group. *Noise* and *Domain* are marked if the method directly addresses DI and MDL, respectively. In the training process, *Labels* is marked if labels are used and *MTL* is marked if there is an MTL component.

E. Network Structure

We use a single hidden layer for the autoencoder with dimension $d_e = \frac{d}{\log(d)}$. This layer is the input for a prediction layer with 10 neurons for each task with single-output (binary or regression). If the task is categorical with d_y classes, the number of neurons is $\max(10, d_y)$. We use ReLU as the non-linear activation and SGD optimizer with batch size=32, learning rate=0.01 and momentum=0.9. We obtain consistent results across all experiments using $\sigma = 10$ in the RBF kernel, $\{\lambda_i\}_{i=0}^3$ set to 1,0.1,0.1 and 2, respectively, for the final loss and, $\rho_0 = 1$ and $\rho_{L2} = 0.1$ for the regularization in the MTL classifier. We test our approach using two configurations: unsupervised (MD2I-U) and supervised (MD2I-S), which uses labels and MTL in the training process. MD2I-U is used to generate the encoded values and a final 2-layer NN classifier is used for prediction.

VI. RESULTS AND ABLATION STUDY

Here, we report the results for all settings and discuss our main findings. We also highlight the importance of the components shown in Table II and their impact on performance.

A. DG-DI for Image Recognition

Table III show the results obtained for the image recognition task. By comparing B2 and B2-DI, we see that the DI component is important as it improves performance in 31%. For DG methods, all can improve upon baselines. However, CCSA has an inferior performance compared to D-MTAE and MMD-AAE. This can be explained by the autoencoder structure that can generalize well across domains while being robust to incomplete data. In fact, D-MTAE is designed to

TABLE III
RESULTS IN DG-DI SETTING FOR IMAGE RECOGNITION

	I_0	I_{15}	I_{30}	I_{45}	I_{60}	I_{75}	Avg
MD2I-S	80.8	89.5	91.3	82.7	93.6	74.5	85.4
MD2I-U	79.9	90.3	91.2	82.4	92.4	74.7	85.2
D-MTAE	79.6	89.4	91.8	80.2	90.7	77.1	84.8
MMD-AAE	77.2	92.8	87.9	72.2	90.0	75.0	82.5
CCSA	72.1	86.8	85.7	76.5	89.1	70.9	80.2
B2-DI	66.4	84.3	83.8	76.2	84.9	69.2	77.5
B2	48.2	69.1	65.4	55.6	65.2	51.0	59.1

use noise to improve robustness, making its performance better than MMD-AEE, which does not directly handle missing data.

The best results are achieved by our method in both versions (MD2I-S, MD2I-U), showing that it can produce domain-invariant features and better learn the distribution of the data.

In this setting, including the labels during training does not significantly improve the performance in our methods and does not help MMD-AAE to outperform D-MTAE. In general, we see that methods that address the *Noise* and *Domain* components have better performance.

B. MTL-DI for Grade Prediction

For the results in the MTL-DI setting shown in Figure 2, we first observe that B1 significantly outperforms B2, which indicates that using all domains together can cause negative transfer. However, the influence of DI seems to be higher than MDL, as B1-DI and B2-DI outperform MDL methods, such as CCSA and MMD-AAE. By comparing D-MTAE (which handles DI and MDL but no labels) and B1-DI, we see that using labels in training improves performance.

DMTRL and UDMTML alone can't beat the baselines, as they are affected by the missing data. To improve them, we perform experiments with a pipeline. First, we use the DI method to generate a complete dataset and then the MTL methods to predict outputs. These new pipelines achieve better scores, with both DI+UMDMTML and DI+DMTRL outscoring MD2I-U. We also implement MDL methods using a pipeline, but the results were similar to the best baseline.

The best results are achieved by our supervised method (MD2I-S), especially due to the robustness to incomplete datasets and the inclusion of label information in training using an MTL approach. Therefore, for this setting, we see that *Labels* and *Noise* have higher impact while *Domain* has lower impact in performance.

C. MDMTL-DI for Classification in Selection Process

Table IV shows the results for the selection processes. In terms of the rank of the methods, we can see that the bottom 6 (from CCSA down) and the top 3 are consistent across all companies. The positions from 4 to 7 change for each company, showing that none of the methods (MMD-AAE, UDMTML, D-MTAE, and DI+DMTRL) in this middle tier have a general superior performance over the others.

When comparing methods in the average column, the worst ones are the baselines, which shows that individual classifiers for each process achieve inferior results. In this setting, DI

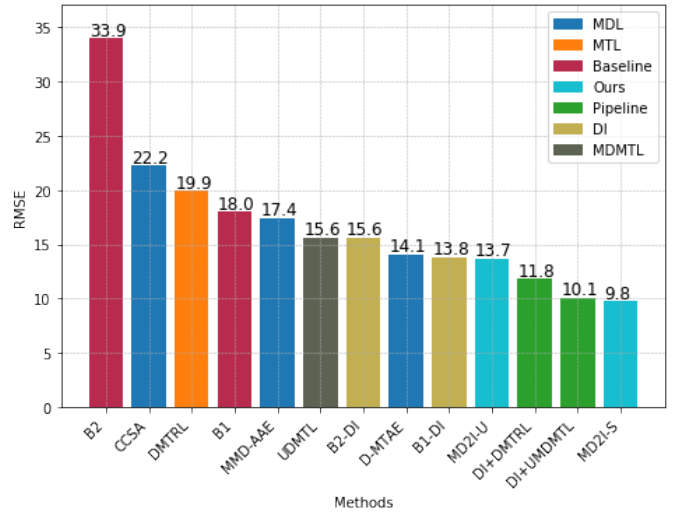


Fig. 2. Results in MTL-DI setting for grade prediction.

TABLE IV
RESULTS IN MDMTL-DI SETTING FOR CLASSIFICATION

Method	Group	C_1	C_2	C_3	Avg
MD2I-S	Ours	83.2	69.7	63.7	72.2
DI+UMDMTML	Pipeline	79.3	67.1	61.2	69.2
MD2I-U	Ours	78.5	60.4	54.7	64.5
D-MTAE	MDL	69.9	55.1	52.4	59.1
DI+DMTRL	Pipeline	67.4	57.6	52.3	59.1
UDMTL	MDMTL	71.6	57.7	45.7	58.3
MMD-AAE	MDL	72.9	52.2	48.6	57.9
CCSA	MDL	53.7	46.5	37.5	45.9
DMTRL	MTL	42.5	33.4	36.4	37.4
B2-DI	DI	43.2	25.0	38.5	35.6
B1-DI	DI	44.6	18.0	37.3	33.3
B2	Baseline	41.3	18.5	33.5	31.1
B1	Baseline	42.4	15.5	31.2	29.7

does not have a high impact as previously, as observed in the marginal gain comparing B1 and B1-DI. In fact, the component *Domain* has a higher impact in this setting, as MDL methods can significantly outperform DI methods.

By having MDL and MTL, UDMTML can beat DI methods and stay in the middle tier. However, when combined with DI using a pipeline, it consistently achieves the second top score, repeating the results from the MTL-DI setting.

MD2I-U is the third-best, beating all other methods but falling short to the top due to not using label information in training. We show that our method can produce domain-invariant encoded spaces and address data imputation jointly. By performing it without labels, we suspect that the features generated from this method can be better generalized to unseen domains where predictions are not available.

Nonetheless, when labels are available, our supervised, MD2I-S, is proven to be effective as it achieves the best score in all settings. These results show the importance of learning all components together in the same framework.

VII. CONCLUSION

We presented a method to combine MDL, DI, and MTL in a single framework. Our loss function optimizes an adversar-

ial autoencoder to jointly produce domain-invariant encoded spaces and learn the data distribution for data imputation. In the unsupervised version (MD2I-U), the encoded and generated features can be used for different downstream tasks. In the supervised version, we add an MTL component for predictions in all domains and to learn from different tasks to improve a single classifier. By combining small datasets from different but related domains, our method can perform better than any classifier trained separately. Through many experiments in different settings including structured, unstructured and mixed data, our method outperformed other state-of-art methods, by successfully combining important components in the same structure. For future work, we want to investigate the features generated and how they can be applied for clustering and classification in unlabelled datasets in different domains.

REFERENCES

- [1] M. Joshi, W. W. Cohen, M. Dredze, and C. P. Rosé, “Multi-domain learning: when do domains matter?” in *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*. Association for Computational Linguistics, 2012, pp. 1302–1312.
- [2] K. Muandet, D. Balduzzi, and B. Schölkopf, “Domain generalization via invariant feature representation,” in *International Conference on Machine Learning*, 2013, pp. 10–18.
- [3] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.
- [4] X. Chen, D. P. Kingma, T. Salimans, Y. Duan, P. Dhariwal, J. Schulman, I. Sutskever, and P. Abbeel, “Variational lossy autoencoder,” *arXiv preprint arXiv:1611.02731*, 2016.
- [5] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, “A kernel two-sample test,” *Journal of Machine Learning Research*, vol. 13, no. Mar, pp. 723–773, 2012.
- [6] S. v. Buuren and K. Groothuis-Oudshoorn, “mice: Multivariate imputation by chained equations in r,” *Journal of statistical software*, pp. 1–68, 2010.
- [7] D. J. Stekhoven and P. Bühlmann, “Missforest—non-parametric missing value imputation for mixed-type data,” *Bioinformatics*, vol. 28, no. 1, pp. 112–118, 2011.
- [8] Y. Zhang and Q. Yang, “A survey on multi-task learning,” *arXiv preprint arXiv:1707.08114*, 2017.
- [9] R. Caruana, “Multitask learning,” *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [10] Y. Yang and T. M. Hospedales, “A unified perspective on multi-domain and multi-task learning,” *arXiv preprint arXiv:1412.7489*, 2014.
- [11] H. Nam and B. Han, “Learning multi-domain convolutional neural networks for visual tracking,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4293–4302.
- [12] S. Motiian, M. Piccirilli, D. A. Adjeroh, and G. Doretto, “Unified deep supervised domain adaptation and generalization,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5715–5725.
- [13] M. Ghifary, W. Bastiaan Kleijn, M. Zhang, and D. Balduzzi, “Domain generalization for object recognition with multi-task autoencoders,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2551–2559.
- [14] H. Li, S. Jialin Pan, S. Wang, and A. C. Kot, “Domain generalization with adversarial feature learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5400–5409.
- [15] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell, “Cycada: Cycle-consistent adversarial domain adaptation,” *arXiv preprint arXiv:1711.03213*, 2017.
- [16] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, “Domain-adversarial training of neural networks,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [17] H. Liu, M. Long, J. Wang, and M. Jordan, “Transferable adversarial training: A general approach to adapting deep classifiers,” in *International Conference on Machine Learning*, 2019, pp. 4013–4022.
- [18] Y. Luo, L. Zheng, T. Guan, J. Yu, and Y. Yang, “Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2507–2516.
- [19] D. Sejdinovic, B. Sriperumbudur, A. Gretton, K. Fukumizu *et al.*, “Equivalence of distance-based and rkhs-based statistics in hypothesis testing,” *The Annals of Statistics*, vol. 41, no. 5, pp. 2263–2291, 2013.
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [21] M. Long, J. Wang, G. Ding, S. J. Pan, and S. Y. Philip, “Adaptation regularization: A general framework for transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 5, pp. 1076–1089, 2013.
- [22] R. Mazumder, T. Hastie, and R. Tibshirani, “Spectral regularization algorithms for learning large incomplete matrices,” *Journal of machine learning research*, vol. 11, no. Aug, pp. 2287–2322, 2010.
- [23] P. J. García-Laencina, J.-L. Sancho-Gómez, and A. R. Figueiras-Vidal, “Pattern classification with missing data: a review,” *Neural Computing and Applications*, vol. 19, no. 2, pp. 263–282, 2010.
- [24] S. Nowozin, B. Cseke, and R. Tomioka, “f-gan: Training generative neural samplers using variational divergence minimization,” in *Advances in neural information processing systems*, 2016, pp. 271–279.
- [25] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein gan,” *arXiv preprint arXiv:1701.07875*, 2017.
- [26] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [27] J. Tomczak and M. Welling, “Vae with a vampprior,” in *International Conference on Artificial Intelligence and Statistics, AISTATS 2018*, 2018, pp. 1214–1223.
- [28] A. Nazabal, P. M. Olmos, Z. Ghahramani, and I. Valera, “Handling incomplete heterogeneous data using vaes,” *arXiv preprint arXiv:1807.03653*, 2018.
- [29] J. Yoon, J. Jordon, and M. Van Der Schaar, “Gain: Missing data imputation using generative adversarial nets,” *arXiv preprint arXiv:1806.02920*, 2018.
- [30] A. Kumar and H. Daumé III, “Learning task grouping and overlap in multi-task learning,” in *Proceedings of the 29th International Conference on International Conference on Machine Learning*. Omnipress, 2012, pp. 1723–1730.
- [31] S. Ruder, “An overview of multi-task learning in deep neural networks,” *arXiv preprint arXiv:1706.05098*, 2017.
- [32] J. Zhou, J. Chen, and J. Ye, “Malsar: Multi-task learning via structural regularization,” *Arizona State University*, vol. 21, 2011.
- [33] J. Chen, J. Zhou, and J. Ye, “Integrating low-rank and group-sparse structures for robust multi-task learning,” in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2011, pp. 42–50.
- [34] Y. Yang and T. Hospedales, “Deep multi-task representation learning: A tensor factorisation approach,” *arXiv preprint arXiv:1605.06391*, 2016.
- [35] J. Liang, E. Meyerson, and R. Miikkulainen, “Evolutionary architecture search for deep multitask networks,” in *Proceedings of the Genetic and Evolutionary Computation Conference*. ACM, 2018, pp. 466–473.
- [36] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. J. Smola, “A kernel method for the two-sample-problem,” in *Advances in neural information processing systems*, 2007, pp. 513–520.
- [37] K. Fukumizu, A. Gretton, G. R. Lanckriet, B. Schölkopf, and B. K. Sriperumbudur, “Kernel choice and classifiability for rkhs embeddings of probability distributions,” in *Advances in neural information processing systems*, 2009, pp. 1750–1758.
- [38] A. Smola, A. Gretton, L. Song, and B. Schölkopf, “A hilbert space embedding for distributions,” in *International Conference on Algorithmic Learning Theory*. Springer, 2007, pp. 13–31.
- [39] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [40] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” *arXiv preprint arXiv:1301.3781*, 2013.