# Unsupervised Neural Domain Adaptation for Document Image Binarization

Francisco J. Castellanos[a,*], Antonio-Javier Gallego[a], Jorge Calvo-Zaragoza[a]

[a]*Department of Software and Computing Systems, University of Alicante, Carretera San Vicente del Raspeig s/n, 03690 Alicante, Spain*

## Abstract

Binarization is a well-known image processing task, whose objective is to separate the foreground of an image from the background. One of the many tasks for which it is useful is that of preprocessing document images in order to identify relevant information, such as text or symbols. The wide variety of document types, alphabets, and formats makes binarization challenging. There are multiple proposals with which to solve this problem, from classical manually-adjusted methods, to more recent approaches based on machine learning. The latter techniques require a large amount of training data in order to obtain good results; however, labeling a portion of each existing collection of documents is not feasible in practice. This is a common problem in supervised learning, which can be addressed by using the so-called Domain Adaptation (DA) techniques. These techniques take advantage of the knowledge learned in one domain, for which labeled data are available, to apply it to other domains for which there are no labeled data. This paper proposes a method that combines neural networks and DA in order to carry out unsupervised document binarization. However, when both the source and target domains are very similar, this adaptation could be detrimental. Our methodology, therefore, first measures the similarity between domains in an innovative manner in order to determine whether or not it is appropriate to apply the adaptation process. The results reported in the experimentation, when evaluating up to 20 possible combinations among five different

*Corresponding author.

*Email addresses:* fcastellanos@dlsi.ua.es (Francisco J. Castellanos),
jgallego@dlsi.ua.es (Antonio-Javier Gallego), jcalvo@dlsi.ua.es (Jorge
Calvo-Zaragoza)

domains, show that our proposal successfully deals with the binarization of new document domains without the need for labeled data.

## 1. Introduction

In the context of documents, digital transcription is the process of exporting the information that is physically present on pages into a format that can be processed by a computer. Its objective is to preserve, and often disseminate, the content of these documents [1]. This process could be performed manually, but this would be a tedious and error-prone task. One solution is the development of systems that are capable of automatically extracting the content [2, 3] and subsequently encoding it into a structured digital format [4, 5, 6]. The wide variety of document types, such as old books, medical reports or even handwritten music scores, makes automating this transcription task very complex, since it is necessary for the system to identify the relevant information to be extracted from each type of document.

One of the most common steps in document image processing is binarization. This process reduces the image to a binary representation (that is, a black and white image) by segmenting the relevant content—such as text, ornaments, or other types of symbols, which will depend on the document—and separating it from the background and other artifacts, such as ink stains, shadow-through, bleed-through, or other types of degradation. This process is often a key aspect in document image analysis workflows [7, 8, 9]. The relevant literature contains multiple heuristic methods based on establishing either a local or global threshold with which to perform the process [10, 11].

Binarization can also be formulated as a supervised machine learning task [12, 13], in which a model is trained to classify each pixel of the image. This strategy, in addition to providing competitive results [14], comprises a more generalizable approach that automatically adjusts the model to each new domain by learning from labeled data. This formulation, however, entails the need for labeled data, which are not always available. Moreover, the models learned with these data are highly dependent on them, so they perform well only for the same domain, or at best, for similar ones. This means that, for each new type of document, it is necessary to retrain the method by using new labeled data from that domain, resulting in an inefficient approach in

2

practice.

The issue mentioned above is a common problem in supervised learning, and there is consequently an open line of research, referred to as Domain Adaptation (DA), that studies how to mitigate it [15]. DA algorithms propose mechanisms with which to apply the knowledge learned from a domain, for which labeled data are available, to other unlabeled domains. It is, therefore, necessary only to have a labeled source domain and perform the adaptation process to the new target domains without having to label more data.

This paper proposes a neural network approach that uses unsupervised DA in order to binarize documents. This approach specifically modifies the Selectional Auto-Encoder (SAE) network [14], previously used for binarization, so as to integrate it into a Gradient Reversal Layer (GRL) [16], thus allowing it to learn domain-invariant features and binarize documents from new domains without using new labeled data.

Furthermore, in document binarization, it is likely that a model trained for one domain works well in many other (target) domains. Although this assumption may occur in other contexts related to semantic segmentation, this is more common in binarization due to the presence of the two same categories in all cases: ink and background. Both categories depict some common characteristics across all types of documents. This makes the situation mentioned above, with a source model behaving well in a new domain, more likely to happen for document image binarization. We will demonstrate that a high similarity between the source and the target domains may be detrimental to the adaptive learning process. Given this, and depending on the features of the new domain, it may not be appropriate to carry out the DA process. An innovative mechanism with which to measure the similarity between binarization domains is also proposed in order to evaluate the new data before adapting the features learned by the network.

Our proposal is evaluated with five datasets that contain different domains, such as old handwritten text documents, musical scores or Balinese palm leaf manuscripts. The results obtained are compared with those obtained by the learning-driven state-of-the-art architecture proposed in [14].

The present work is organized as follows: a literature review is provided in Section 2, while the proposed approach is described in Section 3 and the details of the experimental setup and the analysis of results are shown in Section 4. Finally, our conclusions and future work are discussed in Section 5.

## 2. Background

*2.1. Binarization*

A binarization process traditionally uses a hand-set threshold to separate the relevant information regarding the image from the rest of the content. However, the pixel values of what is relevant may change from one type of document to another, signifying that the threshold also depends on the type of document and its features (color, lighting, and so on). This approach, therefore, lacks the flexibility required for it to be applied directly to all types of documents without having to be adjusted manually each time.

On the one hand, many heuristic methods that estimate an optimal threshold from the features of the image have been developed. There are, for example, methods that calculate a global threshold or that subdivide the image into regions and calculate local thresholds for each of them. One such method is that of Otsu [17], which employs the histogram of the grayscale image in order to compute a global threshold; Perez and Gonzalez [18], meanwhile, proposed a method based on the expansion of Taylor's series to deal with images with non-uniform illumination. It is also worth mentioning the recent adaptive version of Otsu's method proposed by Moghaddam and Cheriet [19], which binarizes images using a grid-based strategy to estimate a background map. There are many other methods that calculate a local rather than a global threshold. This concept was first proposed by Niblack [20], and computes one threshold per pixel on the basis of its neighborhood. The most common extensions to this are that of Sauvola [10] and that proposed by Wolf *et al.* [21], both of which employ more complex equations in order to define the local threshold calculation. There is also a binarization method based on Markov random fields [22], or the method proposed by Jia *et al.* [23], which is used to binarize degraded documents and which computes local thresholds on the basis of structural symmetric pixels. The reader is referred to the available surveys on document image binarization for further literature on heuristic methods [24, 25, 26].

On the other hand, there are methods that formulate binarization as a supervised learning problem. For example, Kita and Wakahara [12] proposed an approach that binarizes images using a combination of $k$-means clustering and Support Vector Machines (SVM) [27]; Chou *et al.* [28] binarize photographed documents by subdividing the image into different regions and calculating local thresholds by using SVM, and Xiong *et al.* [29] applied a similar method in order to binarize degraded scanned documents.

Of the solutions based on machine learning, those that use deep neural networks [30] are especially relevant owing to their good performance. For instance, Convolution Neural Networks (CNN) have been analyzed in order to classify each pixel of an image into a binary value [13]; Afzal *et al.* [31] proposed a Long Short-Term Memory (LSTM) [32] based method that achieves excellent results, even with degraded documents; Milletari *et al.* [33] proposed a new Fully Convolutional Neural Network (FCNN) model with which to segment medical images; the PDNet approach [34] formulates binarization as an energy minimization problem, in which an unrolled primal-dual network based on FCNN is trained to minimize the labeling cost for each pixel; He and Schomaker [11] combined a heuristic method and a deep neural network for binarization: the neural network first processes the image in order to reduce the degradation level, after which Otsu's method is applied in order to binarize the document. Another work is that of Calvo and Gallego [14], who proposed a method that uses a FCNN-based SAE architecture to efficiently binarize images of different collections, such as handwritten text documents or music manuscripts.

However, despite the good results reported for the solutions based on machine learning, and especially on deep learning, they all have the same major drawback: the need for sufficient annotated samples from each document collection. As stated above, this information is not always available, and manually labeling documents can be a costly task, which is not a scalable solution in practice. The development of techniques that allow these models to be used in new domains without requiring labeled data is, therefore, of particular interest.

*2.2. Domain Adaptation*

DA techniques attempt to use the knowledge learned from a domain for which labeled data are available—referred to as the *source domain*—in a different (but related) domain for which labeled data are not available—referred to as the *target domain*. This can be achieved by using several strategies, which can be grouped into three main categories [15]:

- *Divergence-based DA*: A domain-invariant feature representation is obtained by minimizing a measure of divergence between two corpora [35, 36, 37]. For example, Rozantsev *et al.* [38] proposed a pair of parallel neural architectures that regularize the loss function by using a Maximum Mean Discrepancy metric [35]. Another example is

5

the work of Sun and Saenko [39], which uses the Correlation Alignment metric [36] to minimize correlations between source and target domains. The DeepJDOT approach by Damodaran *et al.* [40] learns both a classifier and a common representation of the source and target domains by using loss functions based on the Optimal Transport theory [41].

- *Reconstruction-based DA*: The objective of this strategy is to obtain a common representation of the data in order to use the same classifier for both source and target domains. One example of this is the Deep Reconstruction Classification Network [42], which uses a multi-task learning approach to learn an intermediate feature representation so as to classify both domains. Another is the proposal of Isola *et al.* [43], which employs a conditional Generative Adversarial Network (GAN) to transform one domain into another using an encoder-decoder or a U-Net architecture.

- *Adversarial-based DA*: This strategy trains neural networks, or parts of them, by means of adversarial learning in order to carry out the adaptation process between domains. In this category, GAN [44] can be highlighted as a generative approach that is composed of a generator and a discriminator. The generator attempts to convert source images into target images, and the discriminator attempts to differentiate whether the image is part of the target or a fake. This tunes the weights of the network, signifying that the transformed source images cannot be differentiated from the target ones. Another relevant example is the Domain-Adversarial Neural Network [16], a classification network based on the use of a Gradient Reversal Layer (GRL). This network is divided into two parts that compete during training: one that learns to classify and the other that learns to distinguish the domain. GRL is used to penalize the features that allow a distinction to be made between domains, such that only common features (called *domain-invariant* features) are learned. This strategy was recently extended by Gallego *et al.* [45] in order to incorporate a self-labeling incremental learning process that improves the results obtained.

Most of the DA methods described above are focused principally on classification and are not directly applicable to binarization problems. In binarization, the method has to distinguish the relevant information from the

6

rest of the content and has to return a response for each pixel of the image, rather than a single category.

A task related to binarization is that of image segmentation, which aims to recognize relevant elements from a given image by classifying them according to a set of categories. Some works have applied DA to image segmentation. For example, Danbing *et al.* [46] segments medical images using a combination of reconstruction and adversarial strategies based on the so-called CycleGAN [47]. Despite the good results reported, this solution was evaluated only in very similar domains (in fact, they belonged to the same dataset). Several approaches based on GAN can also be found in the literature, such as the one proposed by Hoffman *et al.* [48], which aims to obtain a common latent representation for the two domains involved, or the work by Yunsheng *et al.* [49], whose proposal performs an image translation from source to target for then applying a segmentation image model. The GAN and the segmentation model are trained through bidirectional learning, by connecting the outputs of each model to the input of the other. Also, Haq and Huang [50] combined a GAN model with an autoencoder to transform the target images so that a discriminator was not able to differentiate them, thus allowing both domains to be segmented with the same model. However, the goodness of these works was proved with domains in which the elements to be recognized keep their shape, thereby representing a very different situation to the binarization case.

As regards the particularities of binarization, images only contain pixels categorized as background or ink; however, the ink class may represent different elements (such as text, musical notes, decorations, etc.) with different colors and shapes and with high-detail labeling. These characteristics depend on the nature of the documents, the engraving mechanisms or the different degradation levels associated with the course of time. These factors are hardly found in segmentation tasks, so they represent a significant challenge for methods such as GAN, in which it does not make sense to transform, for example, musical notes in handwritten text to binarize them, but it is rather necessary to look for other types of common features.

Unlike existing approaches, we propose an unsupervised DA method for document image binarization based on the state-of-the-art SAE architecture and the GRL. We extend the use of this layer to binarization tasks in order to process new domains without using labeled data. As introduced in Section 1, when the involved domains are alike, it could happen that using a model trained only with the source data performs better than with DA. Given

7

that the similarity between domains is, therefore, an important factor to be considered, we will propose a measure to evaluate the new target data before applying DA.

### 2.3. Domain Similarity

Given its relevance in the present work, we now discuss the existing attempts to measure domain similarity. In general, DA techniques are based on the search for similarities or differences between domains. All the methods reviewed in the previous section include some form of comparison that seeks to eliminate the differences or bring together the similarities between the source and target domains. For example, the GRL [16] mentioned in the previous section is connected to a part of the network that performs the comparison of domains, which is used for learning the domain-invariant features.

Many of these methods propose loss functions that calculate the similarity between domains [39, 35, 40], suggest strategies to obtain a common representation of the data in order to use the same classifier for both source and target domains [42, 43], or propose generative networks that analyze the similarities and differences of the domains in order to transform the images from source to target and vice versa [46]. We can also find some proposals that include specific networks for calculating the similarity between domains, such as the Domain2Vec architecture proposed for classification tasks by Peng *et al.* [51], which makes use of a Siamese network to compare the domains, or the work by Osumi *et al.* [52] which weights the contribution of the domain samples in the adaptation process by using their similarity.

In addition, it should be noted that in all these strategies, domain adaptation is always carried out, regardless of whether this process is adequate or not. However, if the domains are similar, there is also the possibility of using the model learned with the source domain to process the samples from the target domain. In this paper, we propose an external process that is carried out before the DA task with the aim of determining whether the adaptation process is necessary. This initial check, in addition to achieving better results (as will be shown in the experimentation section), also means a general reduction in training time.

Our proposal to measure domain similarity (described in detail in Section 3.4) is inspired by the Inception Score (IS) metric [53] used to assess the quality of images created by GAN. IS uses a pre-trained network to classify the images generated and calculates a series of statistics from the probability

distribution obtained, assuming that the images that represent objects in a clear and realistic way will activate certain classes with a higher probability, and that the images without objects or with unclear or unrealistic figures, will obtain a uniform distribution, in which none of the classes will stand out.

Aligned to the idea of IS, in this paper we propose a mechanism to determine the domain similarity by relying on the probability distribution provided by a neural network trained to binarize source images. Our hypothesis is that the activation obtained by this model for a new domain will allow us to compare its similarity with the domain used to train the network. That is, the network will provide similar responses for similar samples, and, therefore, when the domains are not alike, the distributions will not be alike either.

It is also important to clarify that the proposed metric, rather than measuring the similarity between domains in a generic way, it compares the response of the network for a new domain with that given for the domain it was trained with. Hence, with all this, and supported by the idea of IS, this metric is especially suitable for estimating whether the network will be successful with the new domain.

## 3. Methodology

### 3.1. Problem formulation

Let $\mathcal{S}$ be an annotated or *source* dataset for document image binarization, composed of pairs in the form $(\mathcal{X}_\mathrm{S}, \mathcal{Y}_\mathrm{S})$, where $\mathcal{X}_\mathrm{S} = [0, 255]^{h_\mathrm{s} \times w_\mathrm{s} \times c}$ is a document image of size $h_\mathrm{s} \times w_\mathrm{s}$ px and $c$ channels (for instance, 3 in color and 1 in grayscale), and $\mathcal{Y}_\mathrm{S} = \{0, 1\}^{h_\mathrm{s} \times w_\mathrm{s}}$ is its corresponding pixelwise binary annotation.

Let $\mathcal{T}$ be a non-annotated or *target* dataset, which consists solely of a list of images $\mathcal{X}_\mathrm{T} = [0, 255]^{h_\mathrm{t} \times w_\mathrm{t} \times c}$, with a size of $h_\mathrm{t} \times w_\mathrm{t}$ px and $c$ channels.

The task addressed in this paper is that of learning a model from $\mathcal{S}$ and $\mathcal{T}$, with the aim of correctly binarizing images belonging to $\mathcal{T}$. Note that the overall problem is semi-supervised because it employs both labeled and unlabeled data; however, we refer to it as *unsupervised DA* because we assume that there are no labeled data for the target set.

### 3.2. Selectional Auto-Encoder

The backbone of the method presented in this paper is an SAE, which has been successfully used for document binarization in literature [14]. Given

9

its importance in the present work, we introduce this architecture before integrating it into a DA process.

An SAE is a kind of neural network architecture that receives an image as input and provides another image with the same size as output but with values in the range of $[0, 1]$. This network usually consists of two parts: an encoder, which processes the image through the use of consecutive pooling operators in order to extract meaningful features, and a decoder, which contains as many oversampling operators as pooling layers in the encoder with the objective of retrieving the original size of the image.

However, the SAE does not binarize directly, but obtains a probabilistic map on which the value of each pixel represents the probability of being foreground. Given an image $\mathcal{X}_S$, the SAE, therefore, computes the map of probabilities $\mathcal{P}_S$, in order to subsequently apply a decision to each pixel, typically based on a probability threshold $\text{th}_s$.

Once the SAE has been trained with a set of images, it can be used to process the remaining images from the same collection. In this case, since all the documents involved belong to the same domain, the performance is expected to be successful.

It should be kept in mind that the above approach splits the images into multiple patches of a fixed size, $h \times w$ px, during both training and inference. In the latter case, the binarized version of each patch is retrieved and then combined with the others in order to assemble the full binarized image.

### 3.3. Domain adaptation for binarization

One weak point of the SAE is that of dealing with images belonging to different corpora with respect to that used in the training process. Since there is no labeled information for $\mathcal{T}$, the model can be trained only with $\mathcal{S}$. Given the knowledge provided by $\mathcal{S}$, the SAE model learns the features that allow any image $\mathcal{X}_S$ to be binarized. However, it will be able to successfully binarize images from $\mathcal{T}$ if $\mathcal{S}$ and $\mathcal{T}$ are similar domains, and the performance will be severely affected otherwise [14].

In order to alleviate this problem, we propose an adversarial scheme based on the well-known GRL [16], which is a layer that penalizes those features that make it possible to distinguish the domains involved. This is useful as regards obtaining a domain-invariant representation that enables a successful binarization regardless of the domain.

Given $\mathcal{S}$ and $\mathcal{T}$, defined previously, our approach must binarize pages from $\mathcal{T}$ without any ground-truth from that domain. A graphic representation of

the method that we propose is shown in Figure 1. It consists of an SAE architecture with which to binarize images, in which a GRL connects one of its layers to a domain classifier that can differentiate the domain of the input image. GRL assumes a hyper-parameter $\lambda$ in order to adjust the contribution of the domain classifier when training, which should be set empirically.
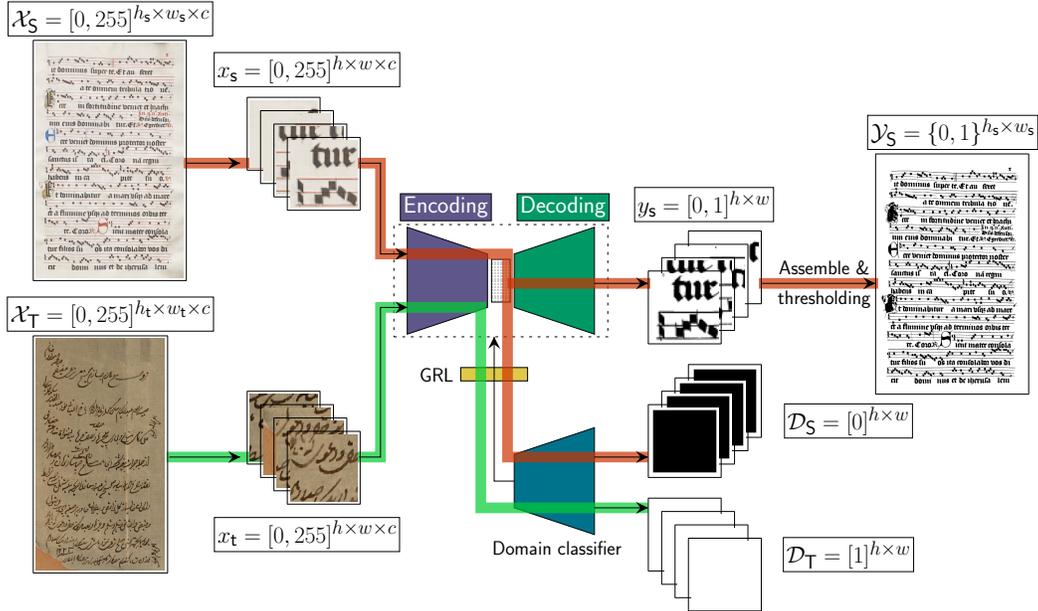


Figure 1: Scheme of the neural domain adaptation approach proposed for document image binarization. The red line represents the path that the source images follow, while the green line represents the route of the target images. The classifier domain is trained in order to obtain images with all 0s or 1s, according to the origin of the image, i.e., from the source or target, respectively.

With regard to the architecture of the domain classifier, it is important to stress that the different number of weights to be adjusted in both streams—binarization and domain classification—might yield an unbalanced structure during the training, thus causing both parts of the network to contribute in different scales to the tuning of the weights. In preliminary experiments, we observed that employing a simple architecture to discriminate the domain with a number of trainable parameters much lower than the binarization stream barely affects the overall training. This can be controlled by manually tweaking the parameter $\lambda$. However, a great difference between both parts of the network forces the tuning of this parameter to very specific values. We

consequently decided to replicate part of the SAE architecture for the domain classification in order to equal the number of weights in both branches.

### 3.4. Adaptation applicability via domain similarity

When both the source and target domains are different, the supervised training from $\mathcal{S}$ may not be sufficient to binarize images from $\mathcal{T}$ reliably. In this situation, using the DA process to adapt the learning of the target domain is a potential solution by which to binarize target images. However, the DA process might not be appropriate when the $\mathcal{S}$ and $\mathcal{T}$ domains depict similar features. In this case, the GRL would be forced to forget useful features for binarization in both domains, and would pay attention to nuances in order to discriminate domains, which would be detrimental for the overall performance. It is for this reason that, in addition to adding DA to the SAE network, we consider it necessary to employ a strategy with which to determine when it is worth applying DA. We, therefore, propose an unsupervised and innovative mechanism that can be used to assess the similarity between $\mathcal{S}$ and $\mathcal{T}$, in the context of binarization, in order to decide whether or not to eventually employ DA.
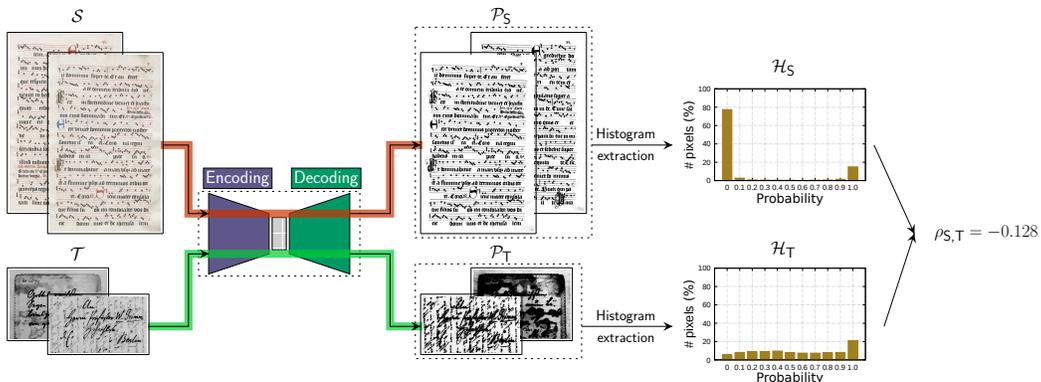


Figure 2: Scheme of the domain similarity computation between two different corpora. The red path represents source processing, while the green path represents the route of the target images. The SAE trained with $\mathcal{S}$ predicts the binarization of both domains. The probabilistic map of each image is then used to obtain a global and normalized histogram of each domain in order to eventually calculate their correlation $\rho_{S,T}$. Note that the images are processed by means of patches, as explained in Section 3.2.

The scheme of our mechanism is shown in Figure 2. Since we are dealing with an unsupervised scenario, the target domain does not include ground-truth data to train the model or to evaluate the result. Instead, we will

resort to the probability map obtained at the network's output, for which it is not necessary to use labels. As explained previously, given any image from $\mathcal{S}$, the SAE model (trained using only source images) obtains a map $\mathcal{P}_{\mathrm{S}}$ showing the probability of each pixel being foreground. It is thus possible to process the images from $\mathcal{S}$ and obtain their probability distribution $\mathcal{P}_{\mathrm{S}}$. This model, when trained with $\mathcal{S}$, can similarly be used to compute the probability distribution $\mathcal{P}_{\mathrm{T}}$ of the images from $\mathcal{T}$.

Once the probability distributions of all the images from both domains ($\mathcal{P}_{\mathrm{S}}$ and $\mathcal{P}_{\mathrm{T}}$) have been attained, we obtain their normalized histograms ($\mathcal{H}_{\mathrm{S}}$ and $\mathcal{H}_{\mathrm{T}}$, respectively) by quantifying their probability figures with equitable intervals of $\mathcal{H}_{\mathrm{prec}}$. In other words, the SAE model provides the probability of each pixel to be ink or background, with values within the range $[0, 1]$. Our proposal maps this probability information onto a histogram for all the pages that belong to a domain. This histogram is a one-dimensional vector obtained by counting the number of pixels with the same probability values within a small range (precision or granularity $\mathcal{H}_{\mathrm{prec}} = 0.1$). Following this process, we generate two one-dimensional vectors, one for the source domain and one for the target domain, which represent the probability distribution histograms of these domains. Note that both are obtained with the SAE trained with the labeled domain (source). Since the number of samples and their resolution can vary among different domains, these histograms are normalized according to the total number of pixels in each domain.

We then compute the correlation $\rho_{\mathrm{S,T}}$ between these two normalized histograms. For this, we considered the Pearson's correlation coefficient [54]. As we will show in Section 4.4, we experimented with other similar metrics and reached equivalent conclusions. We selected the Pearson's correlation coefficient because it is a linear correlation that can measure the similarity between two distributions symmetrically and with bounded values. The following equation shows its mathematical definition

$$\rho_{S,T} = \frac{\mathrm{cov}(\mathcal{H}_{\mathrm{S}}, \mathcal{H}_{\mathrm{T}})}{\sigma_{\mathrm{S}}\sigma_{\mathrm{T}}}, \tag{1}$$

where $\mathrm{cov}\,(\cdot)$ is the covariance function between two variables, while $\sigma_{\mathrm{S}}$ and $\sigma_{\mathrm{T}}$ represent the standard deviation of the histograms of $\mathcal{S}$ and $\mathcal{T}$, respectively. Note that Eq.1 computes a similarity value between $-1$ and $1$, where $-1$ is employed for opposite distributions and $1$ for similar ones.

The coefficient $\rho_{\mathrm{S,T}}$ provides a single value of source and target similarity. The decision to apply DA is made by using a threshold $\rho_{\mathrm{th}}$, whose importance

will be studied in greater detail in Section 4.3. This threshold divides the range of correlation values into two groups: when $\rho_{S,T} > \rho_{th}$, it is assumed that both domains are similar, and the SAE model with no DA is chosen; and when $\rho_{S,T} \leq \rho_{th}$, the domains involved can be considered as being different, and the adaptation strategy explained in Section 3.3 is used.

The entire process described in this section is detailed in Algorithm 1, in which $b$ and $e$ are the batch size and the number of epochs used to train the involved neural networks, respectively, and $h$ and $w$ are the height and width of the extracted patches of the documents, respectively; `getHistogram`$(\cdot)$ is a function that provides the histogram from a probabilistic map ($\mathcal{P}_S$ or $\mathcal{P}_T$) with a precision $\mathcal{H}_{prec}$, and `correlation`$(\cdot)$ is another function with which to compute the linear correlation between two distributions using Pearson's coefficient.

The algorithm begins training the SAE model with $\mathcal{S}$ (line 3). After this training, it then stores a threshold $th_s$, with the best thresholding applied to the probabilistic map obtained with $\mathcal{S}$. This is computed in each epoch, validating the performance with different equidistant thresholds. That which optimizes the performance is then the threshold selected to be applied in binarization.

The next step is to obtain the probabilistic map $\mathcal{P}_S$ of each image within $\mathcal{S}$ by means of the SAE model (line 6) and to compute an accumulative global histogram for that domain $\mathcal{H}_S$ (line 7). The same operation is performed for $\mathcal{T}$, obtaining the histogram $\mathcal{H}_T$ from each $\mathcal{P}_T$ provided by the SAE (lines 10,11).

The similarity measure $\rho_{S,T}$ is then computed between $\mathcal{H}_S$ and $\mathcal{H}_T$ (line 13) and compared with an input threshold $\rho_{th}$. If the correlation does not surpass the threshold, it is assumed that $\mathcal{S}$ and $\mathcal{T}$ are different domains, and it is, therefore, necessary to apply the DA process. In this case, the domain adaptation model (represented as Bin-DANN) is trained with $\mathcal{S}$ and $\mathcal{T}$, and finally binarizes $\mathcal{T}$ by using $th_s$ (lines 15,16). In the opposite case, when $\mathcal{S}$ and $\mathcal{T}$ obtain high correlation, they are not considered as different domains, and $\mathcal{T}$ is, therefore, binarized with the SAE (line 18).

Note that $\mathcal{S}$ should be split into two partitions, for training and validation. The threshold $th_s$ and the histogram $\mathcal{H}_S$ are adjusted with the validation partition.

---
**Algorithm 1:** Adaptation applicability via domain similarity.
---

**Input:** $\mathcal{S} \leftarrow \{(\mathcal{X}_S,\ \mathcal{Y}_S)\}$
$\mathcal{T} \leftarrow \{(\mathcal{X}_T)\}$
$\lambda,\ \rho_{\mathrm{th}},\ \mathcal{H}_{\mathrm{prec}},\ h,\ w,\ e,\ b\ \leftarrow\ $ hyper-parameters

**Result:** $\mathcal{B}_T \leftarrow$ Binarized images from $\mathcal{T}$.

**1** $\mathcal{H}_S \leftarrow \varnothing$
**2** $\mathcal{H}_T \leftarrow \varnothing$
**3** $\mathrm{th_s} \leftarrow$ Fit SAE with $\{\mathcal{S}, e, b, h, w\}$
**4** $\mathcal{P}_T = $ SAE prediction with $\{(\mathcal{X}_T, h, w)\}$
**5 foreach** $\mathcal{X}_S$ **in** $\mathcal{S}$ **do**
**6** $\quad$ $\mathcal{P}_S = $ SAE prediction with $\{(\mathcal{X}_S, h, w)\}$
**7** $\quad$ $\mathcal{H}_S \leftarrow \mathcal{H}_S \cup \mathtt{getHistogram}(\mathcal{P}_S, \mathcal{H}_{\mathrm{prec}})$
**8 end**
**9 foreach** $\mathcal{X}_T$ **in** $\mathcal{T}$ **do**
**10** $\quad$ $\mathcal{P}_T = $ SAE prediction with $\{(\mathcal{X}_T, h, w)\}$
**11** $\quad$ $\mathcal{H}_T \leftarrow \mathcal{H}_T \cup \mathtt{getHistogram}(\mathcal{P}_T, \mathcal{H}_{\mathrm{prec}})$
**12 end**
**13** $\rho_{S,T} = \mathtt{correlation}(\mathcal{H}_S, \mathcal{H}_T)$
**14 if** $\rho_{S,T} \leq \rho_{\mathrm{th}}$ **then**
**15** $\quad$ $\mathrm{th_s} \leftarrow$ Fit Bin-DANN with $\{\mathcal{S}, \mathcal{T}, e, b, h, w, \lambda\}$
**16** $\quad$ $\mathcal{B}_T \leftarrow \mathtt{binarize\ with\ Bin\text{-}DANN}(\mathcal{T}, h, w, \mathrm{th_s})$
**17 else**
**18** $\quad$ $\mathcal{B}_T \leftarrow \mathtt{binarize\ with\ SAE}(\mathcal{T}, h, w, \mathrm{th_s})$
**19 end**
**20 return** $\mathcal{B}_T$

---

## 4. Experiments

In this section, we first describe the datasets and the metrics considered for evaluation purposes. We then detail the neural architecture used and the tuning of hyper-parameters, and finally, we present the results of the experiments and analyze them by comparing with those obtained with the state-of-the-art method. The experiments were implemented with the Keras v. 2.3.1 [55] library and TensorFlow v. 1.14 as backend.

## 4.1. Corpora

We evaluated our method by considering several datasets commonly used for document binarization (some examples can be seen in Figure 3):

- DIBCO: the Document Image Binarization Contest has been held from 2009 [56], and its datasets, containing different content, have been published each year. The experiments were carried out using the 2014 edition for the test set and the other editions until 2016 for the training set, as in [14]. This corpus consists of 86 pages with an average size of $659 \times 1560$ px.

- EINSIEDELN: collection of 10 pages of mensural music documents, specifically those of Einsiedeln, Stiftsbibliothek, Codex 611(89)[1]. The images have an average size of $5550 \times 3650$ px.

- SALZINNES: set of 10 pages of music score images of Salzinnes Antiphonal (CDM-Hsmu 2149.14)[2] with $5100 \times 3200$ px., on average.

- PHI: dataset published for the Persian Heritage Image Binarization Competition (PHI) with a collection of Persian documents [57]. It includes 15 pages with an average size of $1022 \times 1158$ px.

- PALM: set of Balinese Palm Leaf manuscripts for the binarization competition organized at the 15th International Conference on Frontiers in Handwriting Recognition [58], which consists of 97 documents with a size of $492 \times 5116$ px., on average, and whose ground truth was built by employing semi-automatic frameworks.

## 4.2. Metrics

The binarization issue is a two-class problem in which pixels are classified into two possible classes: background and foreground. However, the amount of pixels belonging to each class is typically not balanced, thus promoting a possible bias towards the majority class (usually the background).

The metric of the binarization performance considered in this work was the *F-measure* ($F_1$). In a two-class classification task, this metric is defined as

---

[1]http://www.e-codices.unifr.ch/en/sbe/0611/
[2]https://cantus.simssa.ca/manuscript/133/

(a) PHI           (b) DIBCO           (c) PALM

(d) SALZINNES           (e) EINSIEDELN

Figure 3: Some representative image regions from the corpora considered.

$$F_1 = \frac{2 \cdot TP}{2 \cdot TP + FP + FN},\qquad(2)$$

where TP, FP, and FN represent *True Positives* or correctly classified elements, *False Positives* or type I errors, and *False Negatives* or type II errors, respectively. In the following, we consider the foreground as the *positive* class.

To further analyze the performance, we also include results in terms of precision (P) and recall (R) metrics, defined mathematically as follows:

$$P = \frac{TP}{TP + FP}\qquad(3)$$

$$R = \frac{TP}{TP + FN}\qquad(4)$$

Note that $F_1$ is the harmonic mean of these two metrics.

17

*4.3. Impact of SAE and DA hyper-parameterization*

Since an invariant-domain binarization method based on SAE and GRL is presented in this paper, it is necessary to study multiple parameters: those involved in the SAE model and those associated with GRL. In order to achieve a correct and robust configuration for our method, we first used three corpora—DIBCO, SALZINNES and PHI—to subsequently assess its performance with the remaining collections. With regard to the SAE proposed, in [14] it was demonstrated that, although the method attains high accuracy in on-domain scenarios, it is not robust in cross-domain ones. Since our proposal is based on that method, in this paper, the experiments were carried out in order to compare them in terms of domain-adaptation ability. Note that the images are processed by the neural network in patches of $256 \times 256$ px, and that we considered 300 epochs for the training of all the models, keeping the best model according to a validation partition of the source dataset.

The SAE considered in our experiments consists of an encoder with six convolution layers, with 64 filters, kernels of $3\times3$, strides of $2\times2$, and Rectifier Linear Unit (ReLU) activations. The decoder performs the inverse operation with six transposed convolutional layers with the same configuration. After each ReLU activation, a 0.2 of dropout is performed for both the encoder and the decoder. We shall denote the set of convolution, ReLU activation, and dropout layers as a *block* of layers. After the last block of the decoder, we shall then apply a non-stride convolution layer with a sigmoid activation in order to obtain the probability of each pixel being foreground. Since residual connections from each encoding layer to its analogous decoding layer were a key factor in the original SAE topology, our implementation also includes them.

With regard to the integration with GRL, the position of this layer within the SAE architecture is relevant, since it affects only the previous layers to which it is connected. For instance, if GRL is connected to the middle of the SAE model, i.e., the so-called latent space of an auto-encoder, it will affect only the encoder section, and not that of the decoder (at least not directly).

As mentioned in Section 3.3, we propose an architecture in which the domain classifier and the binarization model contain a balanced number of trainable parameters. We have, therefore, implemented the classifier by replicating part of the SAE architecture, specifically with the same structure of the SAE from the layer in which GRL is connected. In preliminary experiments, we noticed that the model performed the binarization task better

when the GRL was connected before the last convolutional block of the decoder. As regards the hyper-parameter $\lambda$ of GRL, we considered an incremental function that starts at 0.1 and adds 0.01 per epoch.

Since the model obtains a map of probabilities, thresholding is necessary in order to eventually determine which pixels belong to the foreground class and which belong to the background. However, it is supposed that only $\mathcal{S}$ has available ground truth. For each epoch in training, we, therefore, calculate the best threshold, denoted as $\text{th}_\text{s}$, using a validation partition of $\mathcal{S}$. This threshold is then used to evaluate the images from $\mathcal{T}$ in the experiments.

In the experiments, we first compare the baseline model, which we shall denote as SAE, with our GRL approach, henceforth Bin-DANN—, without, as yet, calculating the similarity domains. In these experiments, we considered all the possible combinations of pairs of datasets from the three selected. The pairs of datasets in the text were, for convenience, denoted as $\mathcal{S} \rightarrow \mathcal{T}$ to refer to the combination of the labeled source dataset $\mathcal{S}$ with the unlabeled target dataset $\mathcal{T}$.

The first set of results is reported in Table 1. This table shows that, despite the clear increase in accuracy achieved in several cases with respect to the state-of-the-art method, with an average performance of between 51.6% and 66.9% of $F_1$, not all the pairs of datasets considered in these preliminary experiments improve the baseline. This is particularly the case of the pairs DIBCO $\rightarrow$ SALZINNES, PHI $\rightarrow$ SALZINNES and PHI $\rightarrow$ DIBCO, for which there is no such improvement. These results show that, for certain pairs of domains, it is not always better to perform the adaptation process, probably because of their similarity in the neural feature space. As we shall show later, thanks to the method proposed to measure this similarity, we shall improve the final result.

Besides, we observe some differences between the P and R figures. The reason behind this is that the binarization method is based on a threshold set to the probability map provided by SAE, which is calculated using only the images from $\mathcal{S}$. For example, for SALZINNES $\rightarrow$ PHI and DIBCO $\rightarrow$ PHI, SAE presents a severe reduction in the P obtained. However, the adaptation process carried out by Bin-DANN manages to improve this result. In other cases, we observe a high P but a low R. This is also associated with the threshold set in an unsupervised manner. In these cases, binarization provides few false positives, but some false negatives, mainly because the method classifies as ink only when the probability of being ink is very high. These results are also a symptom that the SAE is not working well with the

Table 1: Preliminary experiments carried out for the three datasets selected in terms of precision, recall and $F_1$ (%). The average row includes only the figures concerning the unsupervised scenarios. The SAE columns represent the state-of-the-art method, while the Bin-DANN columns represent the domain adaptation method without the comparison of histograms.

| $\mathcal{S}$ | $\mathcal{T}$ | SAE | | | Bin-DANN | | |
|---|---|---|---|---|---|---|---|
| | | P | R | $F_1$ | P | R | $F_1$ |
| | Salzinnes | 96.8 | 94.9 | 95.9 | - | - | - |
| Salzinnes | Dibco | 60.1 | 96.2 | 74.0 | 74.8 | 93.3 | **83.0** |
| | Phi | 10.8 | 99.9 | 19.5 | 53.6 | 99.7 | **69.7** |
| | Dibco | 94.7 | 83.4 | 88.7 | - | - | - |
| Dibco | Salzinnes | 99.8 | 78.6 | **88.0** | 80.2 | 69.1 | 74.2 |
| | Phi | 13.3 | 99.8 | 23.5 | 64.0 | 98.4 | **77.6** |
| | Phi | 85.5 | 83.5 | 84.5 | - | - | - |
| Phi | Salzinnes | 99.9 | 40.4 | **57.5** | 82.8 | 37.3 | 51.4 |
| | Dibco | 99.2 | 30.8 | **47.0** | 88.7 | 30.4 | 45.3 |
| Average | | 63.9 | 74.3 | 51.6 | 74.0 | 71.3 | **66.9** |

target images, which, as we have indicated, is solved after the adaptation process.

After analyzing these preliminary results, we observe that $F_1$, the harmonic mean of P and R, is enough for the evaluation. Therefore, from here on we will report the results considering only this metric.

## 4.4. Impact of domain similarity hyper-parameterization

We shall now evaluate the proposed method with which to calculate the similarity between domains. As mentioned previously, it is necessary to establish a threshold $\rho_{\text{th}}$ for the correlation value $\rho_{\text{S,T}}$ in order to decide whether the model should be binarized with the supervised SAE or with the model based on DA. For the following experiments, we are going to consider the quotient $F_1^{\text{Bin-DANN}}/F_1^{\text{SAE}}$ as a measure of the relative improvement obtained by the DA process with respect to the conventional SAE.

As mentioned in Section 3.4, in addition to the Pearson's correlation, we considered a set of known metrics to compute the similarity or divergence between the two histograms $\mathcal{H}_{\text{S}}$ and $\mathcal{H}_{\text{T}}$. Specifically, we compared the Pearson's correlation [54], the Kullback–Leibler (KL) divergence [59], the Jensen-Shannon (JS) divergence [60] and the histogram intersection [61].
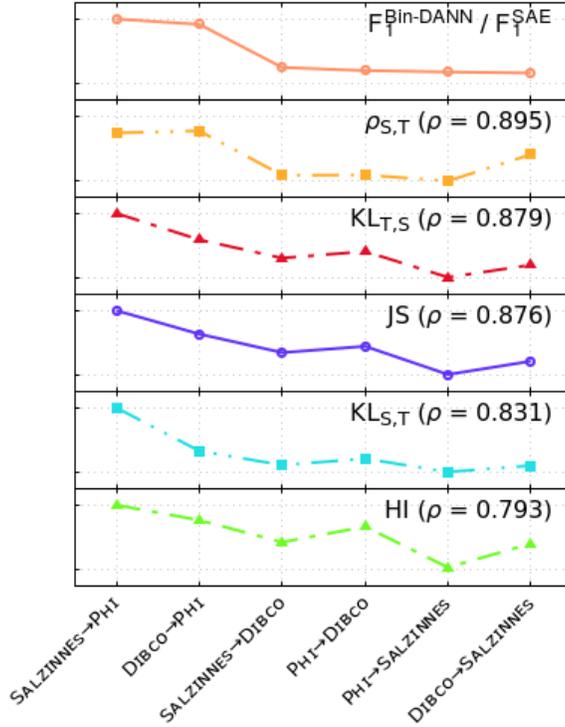
Figure 4: Comparison among different similarity/divergence metrics between $\mathcal{H}_S$ and $\mathcal{H}_T$. The horizontal axis indicates the different pairs of domains $\mathcal{S} \to \mathcal{T}$ studied. The first curve stands for the quotient $F_1^{\text{Bin-DANN}}/F_1^{\text{SAE}}$. The rest of graphs represent the different correlation/divergence metrics normalized between 0 and 1 for easy comparison. They are ordered according to the correlation ($\rho$) of each graph compared with the first curve, and whose values are included between brackets.

Figure 4 shows a comparison among all these metrics, where it is observed that all have a similar trend. We also compared this trend with the curve obtained by the relative improvement $F_1^{\text{Bin-DANN}}/F_1^{\text{SAE}}$ for each pair of domains $\mathcal{S} \to \mathcal{T}$, ordered by this quotient. For further analysis, the figure includes the Pearson's correlation between each curve compared with the degree of improvement of Bin-DANN with respect to SAE (i.e., the first curve).

We observe that Pearson has the highest correlation, with a figure of 0.895. The JS and KL divergences also obtain competitive figures between 0.879 and 0.831, whereas HI reports the worst correlation with 0.793. It should be noted that KL divergence is not symmetric, so it changes according to the order in which the operations are computed. Particularly, we

21

observe a correlation of 0.879 for $KL_{T,S}$, and a value of 0.831 for the opposite $KL_{S,T}$. For all this, we finally selected the Pearson's correlation for the rest of experiments, since in addition to reporting the best results, it also returns bounded values between -1 and 1, simplifying the search for the optimal threshold $\rho_{S,T}$. Anyway, since similar behaviors are observed for all cases, we can assume that our method is robust regardless of the particular metric used to compare the probability distributions.

In Figure 5, we plot the quotient $F_1^{Bin\text{-}DANN}/F_1^{SAE}$ with respect to $\rho_{S,T}$ in order to assess their relationship. This figure shows that when $\rho_{S,T}$ is near to the maximum value, the improvement ratio obtained with Bin-DANN is virtually non-existent, or even detrimental (below 1). This phenomenon would appear to be associated with those cases in which the SAE model provides similar histogram distributions for both the source and target domains. In these cases, and following our premise, the binarization can be performed by the SAE trained with $\mathcal{S}$.



Figure 5: Relationship between the correlation coefficient and the improvement obtained when applying domain adaptation. The horizontal dotted line marks the boundary of no improvement with respect to the SAE method. Higher figures are those cases in which Bin-DANN improves the results and the lower ones represent the cases in which it does not.

However, when the histogram distributions between $\mathcal{S}$ and $\mathcal{T}$ barely match (left-hand side of the horizontal axis), the improvement rate rises drastically, increasing to a maximum factor of over 3 times with respect to the SAE model. For example, in the case SALZINNES → PHI, we obtain an

SAE binarization performance of 19.5%, which is considerably improved by Bin-DANN to 69.7%, or in the case DIBCO → PHI, when the performance is increased from 23.5% to 77.6%.



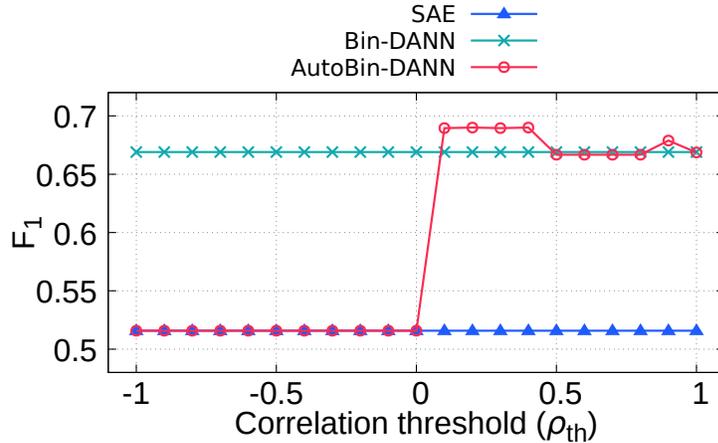Figure 6: Comparison between the averurlage performance obtained by SAE, Bin-DANN and AutoBin-DANN, depending on the threshold selected $\rho_{th}$. When the correlation $\rho_{S,T}$ is less than or equal to $\rho_{th}$, the method selects Bin-DANN; otherwise, it selects SAE.

We shall now evaluate the method that automatically selects the SAE or the Bin-DANN, henceforth denoted as AutoBin-DANN (Algorithm 1). We assess its performance with respect to $\rho_{th}$ in Figure 6. The curves reveal a clear increase in performance of AutoBin-DANN when $\rho_{th} \geq 0.1$. The best results are obtained in the range between 0.1 and 0.4, obtaining an $F_1$ value of 69% for AutoBin-DANN when compared with the SAE and the Bin-DANN, with only 51.6% or 66.9%, respectively.

The preliminary experiments, therefore, show that the best results are attained when $\rho_{S,T} \in [0.1, 0.4]$. For the final experiments, which are described in the next section, we consider the middle value within that range, i.e., $\rho_{th} = 0.25$.

To complement these experiments, Table 2 analyzes the values of the Pearson's correlation coefficient ($\rho_{S,S}$) when comparing samples within the same domain. We observe that, in all cases, the correlation is higher than 0.9 with very low variance. It must be kept in mind that, in these datasets, there are samples in which the color distribution is maintained but the ink density varies, as well as samples of the opposite. However, the values from Table 2

Table 2: Analysis of the similarity between images of the same domain $\rho_{S,S}$. We report the values in terms of Pearson's correlation (average $\pm$ std. deviation).

| Domain | Correlation $\rho_{S,S}$ |
|---|---|
| SALZINNES | $0.99 \pm 6 \times 10^{-8}$ |
| DIBCO | $0.95 \pm 7 \times 10^{-2}$ |
| PHI | $0.99 \pm 3 \times 10^{-4}$ |

indicate that the intra-domain correlation is consistent and, therefore, the proposed metric is suitable for the task at issue.

### 4.5. Final evaluation

We shall now employ the best configuration determined previously to study the results obtained for all the combinations of the datasets considered. The first columns of Table 3 show the $F_1$ obtained for the SAE, Bin-DANN, and AutoBin-DANN methods. We then include the correlation calculated for each combination of datasets, and finally, the improvement achieved with respect to the performance of SAE when applying Bin-DANN and AutoBin-DANN.

As expected, the domain adaptation method (Bin-DANN) is not always the best option, since there are several combinations of $\mathcal{S}$ and $\mathcal{T}$ that do not improve the performance of the SAE. If we focus on the correlation column $\rho_{S,T}$, we can observe that higher figures of this metric, i.e., those near to 1, are associated with both high binarization performance with SAE, a bare improvement or a decrease in the quality of the binarization in the adaptation model.

In general, the performance of the SAE is severely reduced when $\mathcal{S} \neq \mathcal{T}$. For example, if we focus on SALZINNES as a source, we observe a great difference in performance as regards the on-domain experiment (SALZINNES $\rightarrow$ SALZINNES)—with a result of 95.9%—and the cross-domain ones, with figures of between 12.1% and 74%. The same applies when PHI is used as a source: the on-domain experiment attains 84.5%, whereas the performance in the cross-domain experiments ranges from 21.8% to 62.2%.

When $\mathcal{S} \equiv$ DIBCO, the results of the music datasets (SALZINNES and EINSIEDELN) are also comparable with the on-domain case, with values near to 89%. This may be because SALZINNES and EINSIEDELN, in addition to other elements, also contain text. Indeed, in the opposite situation, i.e., when SALZINNES $\rightarrow$ DIBCO and EINSIEDELN $\rightarrow$ DIBCO, the performance decreases

Table 3: Final results for all corpora. The best results between the $F_1^{\text{SAE}}$, $F_1^{\text{Bin-DANN}}$ and $F_1^{\text{AutoBin-DANN}}$ columns are highlighted in bold type. The $\rho_{\text{S,T}}$ column is the correlation coefficient between the histograms obtained from the probabilistic map of $\mathcal{S}$ and $\mathcal{T}$ by means of the SAE model. Finally, the two last columns show the difference between our proposals (Bin-DANN and AutoBin-DANN) and the state-of-the-art model (SAE).

| $\mathcal{S}$ | $\mathcal{T}$ | $F_1^{\text{SAE}}$ (%) | $F_1^{\text{Bin-DANN}}$ (%) | $F_1^{\text{AutoBin-DANN}}$ (%) | $\rho_{\text{S,T}}$ | $F_1^{\text{Bin-DANN}}$ - $F_1^{\text{SAE}}$ | $F_1^{\text{AutoBin-DANN}}$ - $F_1^{\text{SAE}}$ |
|---|---|---|---|---|---|---|---|
| | SALZINNES | 95.9 | - | - | 1.00 | - | - |
| | EINSIEDELN | 67.3 | **92.1** | 67.3 | 0.60 | +24.8 | 0.0 |
| SALZINNES | DIBCO | 74.0 | **83.0** | 74.0 | 0.89 | +9.1 | 0.0 |
| | PHI | 19.5 | **69.7** | **69.7** | 0.08 | +50.2 | +50.2 |
| | PALM | 12.1 | **32.7** | **32.7** | 0.07 | +20.6 | +20.6 |
| | SALZINNES | **90.4** | 88.2 | **90.4** | 0.54 | −2.2 | 0.0 |
| | EINSIEDELN | 92.5 | - | - | 1.00 | - | - |
| EINSIEDELN | DIBCO | **86.1** | 84.2 | **86.1** | 0.99 | −1.9 | 0.0 |
| | PHI | 20.8 | **63.0** | **63.0** | −0.24 | +42.3 | +42.3 |
| | PALM | 16.0 | **30.6** | **30.6** | 0.25 | +14.6 | +14.6 |
| | SALZINNES | **88.0** | 74.2 | **88.0** | 0.49 | −13.7 | 0.0 |
| | EINSIEDELN | **89.6** | 83.2 | **89.6** | 0.99 | −6.5 | 0.0 |
| DIBCO | DIBCO | 88.7 | - | - | 1.00 | - | - |
| | PHI | 23.5 | **77.6** | **77.6** | 0.04 | +54.1 | +54.1 |
| | PALM | 15.8 | **30.4** | **30.4** | 0.13 | +14.6 | +14.6 |
| | SALZINNES | **57.5** | 51.4 | **57.5** | 0.99 | −6.2 | 0.0 |
| | EINSIEDELN | **62.2** | 60.2 | **62.2** | 0.77 | −2.0 | 0.0 |
| PHI | DIBCO | **47.0** | 45.3 | **47.0** | 0.88 | −1.7 | 0.0 |
| | PHI | 84.5 | - | - | 1.00 | - | - |
| | PALM | **21.8** | 20.1 | **21.8** | 0.92 | −1.7 | 0.0 |
| | SALZINNES | 70.8 | **76.1** | 70.8 | 0.32 | +5.2 | 0.0 |
| | EINSIEDELN | 73.8 | **80.7** | 73.8 | 0.43 | +6.9 | 0.0 |
| PALM | DIBCO | 58.2 | **68.7** | 58.2 | 0.66 | +10.5 | 0.0 |
| | PHI | **66.6** | 19.5 | **66.6** | 0.96 | −47.0 | 0.0 |
| | PALM | 39.4 | - | - | 1.00 | - | - |

slightly. We attribute this phenomenon to the wide variety of elements in music manuscripts, such as staff lines, ornaments and music notes, which do not appear in DIBCO.

With regard to the experiments with $\mathcal{S} \equiv$ PALM, the on-domain results are an exception. It is worth noting that this dataset contains poor-quality images, and this dataset does not, therefore, achieves good results, regardless of the $\mathcal{S}$ considered, with the best case being the on-domain situation with 39.4%, when compared to the cross-domain experiments with values of only 21.8% (PHI → PALM).

It should also be noted that Bin-DANN can greatly increase the performance of binarization in many cases, while the correlation-based selection allows the baseline performance to be maintained when the situation is not

suitable for adaptation. Furthermore, there are cases in which this decision is not the best. For example, in the case of SALZINNES → EINSIEDELN, the performance of Bin-DANN improves from 67.3% to 92.1%, but the correlation coefficient reaches a high value of 0.60. It is precisely for this reason that the substantial improvement made by the domain adaptation method is not selected when the correlation threshold is set to 0.25. This is not, of course, detrimental when compared to the baseline case, signifying that the binarization can be carried out with fair results. A similar situation can be found in the case of SALZINNES → DIBCO, with an improvement from 74.0% to 83.0%, and PALM → DIBCO, from 58.2% to 68.7%.

In spite of the above, the AutoBin-DANN is beneficial in a number of cases. For example, when EINSIEDELN → PHI or EINSIEDELN → PALM, the binarization results of Bin-DANN are improved from 20.8% to 63.0% and from 12.1% to 32.7%, respectively, and the correlation allows the eventual selection of the adapted method. There are similar cases for almost all the source scenarios, with the exception of PHI, since when it is used as a source, the adaptation is detrimental in all the target cases. However, the AutoBin-DANN manages to correct this issue by selecting the state-of-the-art method, and does not, therefore, compromise the binarization performance.

After the correlation thresholding, none of the 20 unsupervised experiments (without including those when the same domain is considered as source and target) obtained a loss of performance with respect to the SAE model. The low threshold value obtained in the preliminary experiments makes it possible to increasing the robustness of AutoBin-DANN, but, as mentioned previously, there are cases in which the DA approach leads to improvement, but which are not exploited owing to the correlation bias.

Moreover, note that the adaptation process is not suitable when the performance of the baseline is high. For example, in the case of EINSIEDELN → DIBCO, SAE attains 86.1%, while Bin-DANN obtains 84.2%, DIBCO → SALZINNES decreases from 88.0% for SAE to 74.2% with DA and DIBCO → EINSIEDELN worsens the binarization from 89.6% to 83.2%. We attribute this loss to the fact that the feature representation learned by SAE works properly with $\mathcal{T}$, but DA modifies and blurs it in the learning process. Furthermore, when $\mathcal{S}$ is PALM, the domain adaptation is slightly better for each target except for the combination PALM → PHI, in which it falls from 66.5% to 19.5%. However, the correlation bias makes it possible to avoid this drawback by selecting the SAE method. In addition, when PHI is used as a source, the domain adaptation is slightly detrimental for all the targets

considered. $\rho_{th}$, therefore, plays an important role as regards avoiding these situations in order to maximize the robustness of the combined approach.
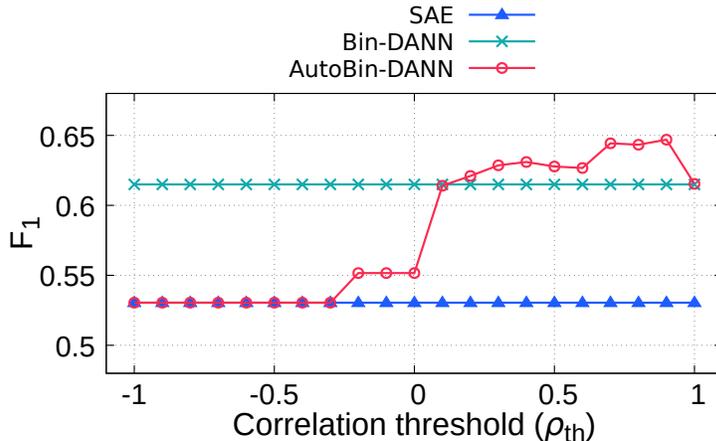


Figure 7: Comparison between SAE, Bin-DANN and AutoBin-DANN depending on the threshold $\rho_{th}$ selected.

All of the above will now be employed to analyze the importance of the correlation threshold in Figure 7. Similar to what occurred in the preliminary analysis, AutoBin-DANN improves the SAE model with $\rho_{th} \in [-0.2, 1.0]$, with this improvement being more relevant within the range $[0.1, 1.0]$. In these results, the same trend is obtained with the datasets used for the adjustment of the hyper-parameters, in which, for thresholds from 0.1, the model provided better results when applying the proposed AutoBin-DANN.

The relationship between the correlation $\rho_{S,T}$ and the improvement obtained with DA is shown in Figure 8. Note that the selected threshold of 0.25 splits the distribution in two parts: the first, with $\rho_{th} <= 0.25$, in which the improvements obtained range from 192% to 357% with respect to the baseline, and the second, with $\rho_{th} > 0.25$, which is composed of the lowest improvement cases and the scenarios in which the Bin-DANN is worse than the SAE. Note also that when $\rho_{th} \in [0.7, 0.9]$, the approach selects the Bin-DANN model for several combinations in which it is the best choice, but also others in which the SAE would be the best option.

Table 4 includes a summary of the average of the results shown in Table 3. While the state of the art obtains 53% of performance, Bin-DANN improves it with 61.5% and AutoBin-DANN outperforms both with 62.9%. To provide a reference of the best result that could be attained when the task is performed
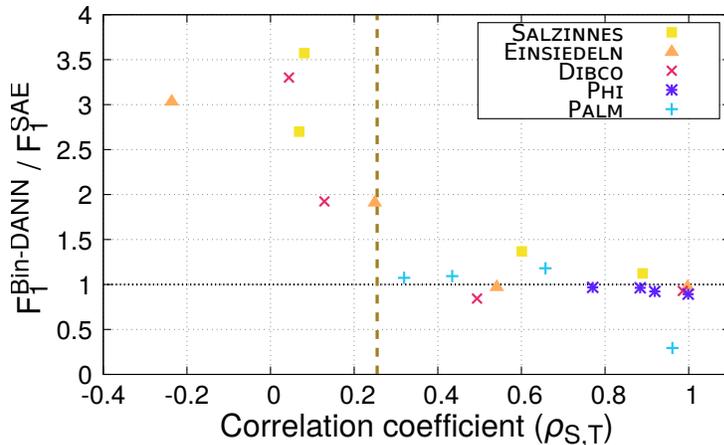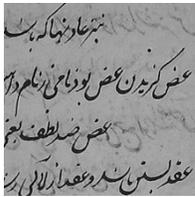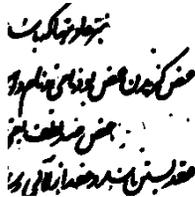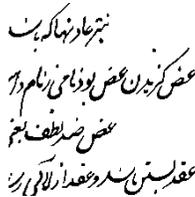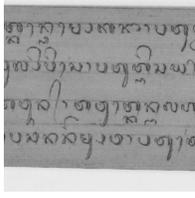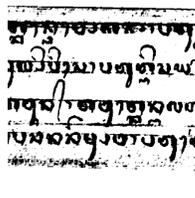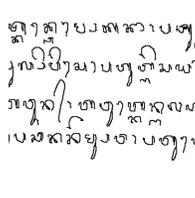
Figure 8: Relationship between the correlation coefficient and the improvement obtained when applying DA, using all datasets. The horizontal dotted line marks the boundary of no improvement with respect to the SAE method. Higher figures are the cases in which Bin-DANN improves the results, while the lower ones represent the cases in which it does not. The vertical dashed line symbolizes the correlation threshold applied by AutoBin-DANN.

with supervised learning, the table includes the performance achieved with SAE for all $\mathcal{T} \rightarrow \mathcal{T}$ (see last column), obtaining an average result of 76.5. This corresponds to the upper bound that our proposal attempts to attain. Both Bin-DANN and AutoBin-DANN are able to recover part of the loss caused by the SAE model when $\mathcal{S} \neq \mathcal{T}$. Specifically, Bin-DANN recovers 36% of the performance, while AutoBin-DANN achieves over 42%, thus getting closer to the ideal case by almost half.

Table 4: Average results for each source corpus.

| $\mathcal{S}$ | $F_1^{SAE}$ (%) | $F_1^{Bin\text{-}DANN}$ (%) | $F_1^{AutoBin\text{-}DANN}$ (%) | $F_1^{SAE}$ (ref. $\mathcal{T}$ %) |
|---|---|---|---|---|
| SALZINNES | 43.2 | **69.4** | 60.9 | 76.2 |
| EINSIEDELN | 53.3 | 66.5 | **67.5** | 77.1 |
| DIBCO | 54.2 | 66.3 | **71.4** | 78.1 |
| PHI | **47.1** | 44.2 | **47.1** | 79.1 |
| PALM | **67.4** | 61.2 | **67.4** | 72.3 |
| Average | 53.0 | 61.5 | **62.9** | 76.5 |

28

Table 5: Selected examples from DIBCO → PHI and DIBCO → PALM cases. We show the binarization obtained by the state-of-the-art model (SAE), that provided by our approach (Bin-DANN), and the ground truth. In both cases, AutoBin-DANN selects the DA approach.

| | Input | SAE | Bin-DANN | GT |
|---|---|---|---|---|
| DIBCO ↓ PHI | | | | |
| DIBCO ↓ PALM | | | | |

## 4.6. Qualitative evaluation

To conclude our experimental section, we present some representative examples of image binarization, comparing the different approaches. Table 5 shows the binarization results obtained using DIBCO as $\mathcal{S}$ and images from PHI and PALM as $\mathcal{T}$. In both cases, the correlation between $\mathcal{S}$ and $\mathcal{T}$ is low, with values of 0.13 and 0.04 for PHI and PALM, respectively. Therefore, the AutoBin-DANN selects the DA approach. As observed in Section 4.5, lower correlation coefficients are usually aligned to higher improvements in performance by Bin-DANN, with the SAE method providing poor-quality binarization.

If we examine the input image of the first example, i.e. DIBCO → PHI, we see that the ink bleeds through the paper from the reverse side but, given that this is noise, these pixels are labeled as background in the ground truth. However, the state-of-the-art approach is not able to differentiate this from the actual text of the page. Bin-DANN, by adapting to the new domain, deals better with the situation, obtaining a binarized image that is much closer to the ground truth. Something similar happens in the second example. In this case, it is observed that PALM documents have a very low contrast between the ink and the background, which does not happen in DIBCO. This confuses the state-of-the-art method, which produces many

false positives, failing to differentiate the background from the ink. Again, the proposed method manages to better deal with the issue thanks to the adaptation process.

If we visually analyze the domains, we may observe graphic similarities and differences between them, such as in DIBCO and PHI, which present similar intensity in ink pixels but have differences in background color (see Figure 3). However, the performance of binarization depends on the features learned by the network from the source domain, which do not have to coincide with our visual appreciation. For example, according to Table 3, PHI → DIBCO obtains a correlation of 0.88, which means that both domains are quite similar (according to our measure); however, if we analyze the inverse case (DIBCO → PHI), the correlation drops to 0.04. Therefore, the features on which the network is focusing on, can be, and certainly are, different to the features that the human eye is perceiving. This reinforces the use of the proposed AutoBin-DANN method, since it is based on the result obtained by the network when using the learned features, analyzing whether these features are also suitable to process a certain target domain.

## 5. Conclusions

In this paper, we propose an unsupervised neural network approach with which to binarize images by means of adversarial training from a domain whose ground truth is not available. The approach employs a state-of-the-art method based on Selectional Auto-Encoder (SAE) as its basis and makes use of a Domain Adaptation (DA) artifact denominated as Gradient Reversal Layer (GRL). This makes it possible to learn a common feature representation in order to binarize both the labeled and unlabeled domains by penalizing those features that differentiate the domain of the input image. The model is then able to learn to binarize images of different domains with respect to that used for training.

The results suggest that this DA approach can be employed to address the binarization issue, and that in most cases it obtains a clear improvement. However, it has been also shown that the adaptation is not always suitable, depending on the pair of corpora considered. In order to solve this problem, we propose a method that makes it possible to compare the similarity between the domains and determine whether applying the DA process is appropriate. This process is performed by means of a comparison between the histogram

obtained from the probability maps provided by the SAE (trained with $\mathcal{S}$) for both domains $\mathcal{S}$ and $\mathcal{T}$.

The experiments were carried out with five different domains and, therefore, 20 possible combinations of pairs of domains. The results reveal that the decision to use DA or SAE is essential if a robust model for binarization is to be obtained in unsupervised scenarios, since there was a substantial increase in the average results from 53% to 62.9% of performance for all the study cases considered, approaching the upper bound by over 42% when compared with the state of the art.

After analyzing the results, we realized that our proposal has room for improvement. Although the decision algorithm proposed in this paper can make robust decisions, it is not, in some cases, perfect. Our future work will, therefore, involve studying other decision strategies based on machine learning.

## Acknowledgment

## References

[1] D. Doermann, K. Tombre, et al., Handbook of document image processing and recognition, Springer, 2014.

[2] V. B. Campos, A. H. Toselli, E. Vidal, Natural language inspired approach for handwritten text line detection in legacy documents, in: Proceedings of the 6th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities, LaTeCH '12, Association for Computational Linguistics, USA, 2012, p. 107–111.

[3] D. Bainbridge, T. Bell, The challenge of optical music recognition, Computers and the Humanities 35 (2) (2001) 95–121.

[4] M. Krallinger, A. Valencia, Text-mining and information-retrieval services for molecular biology, Genome biology 6 (7) (2005) 224.

[5] K. Wang, B. Babenko, S. Belongie, End-to-end scene text recognition, in: 2011 International Conference on Computer Vision, IEEE, 2011, pp. 1457–1464.

[6] A. Hankinson, J. A. Burgoyne, G. Vigliensoni, I. Fujinaga, Creating a large-scale searchable digital collection from printed music materials, in: Proceedings of the 21st International Conference on World Wide Web, ACM, 2012, pp. 903–908.

[7] G. Louloudis, B. Gatos, I. Pratikakis, C. Halatsis, Text line detection in handwritten documents, Pattern Recognition 41 (12) (2008) 3758–3772.

[8] S. He, M. Wiering, L. Schomaker, Junction detection in handwritten documents and its application to writer identification, Pattern Recognition 48 (12) (2015) 4036–4048.

[9] A. P. Giotis, G. Sfikas, B. Gatos, C. Nikou, A survey of document image word spotting techniques, Pattern Recognition 68 (2017) 310 – 332.

[10] J. Sauvola, M. Pietikäinen, Adaptive document image binarization, Pattern recognition 33 (2) (2000) 225–236.

[11] S. He, L. Schomaker, Deepotsu: Document enhancement and binarization using iterative deep learning, Pattern recognition 91 (2019) 379–390.

[12] K. Kita, T. Wakahara, Binarization of color characters in scene images using k-means clustering and support vector machines, in: 2010 20th International Conference on Pattern Recognition, IEEE, 2010, pp. 3183–3186.

[13] J. Pastor-Pellicer, S. E. Boquera, F. Zamora-Martínez, M. Z. Afzal, M. J. C. Bleda, Insights on the use of convolutional neural networks for document image binarization, in: Advances in Computational Intelligence - 13th International Work-Conference on Artificial Neural Networks, IWANN 2015, Palma de Mallorca, Spain, June 10-12, 2015. Proceedings, Part II, 2015, pp. 115–126.

[14] J. Calvo-Zaragoza, A.-J. Gallego, A selectional auto-encoder approach for document image binarization, Pattern Recognition 86 (2019) 37 – 47.

[15] M. Wang, W. Deng, Deep visual domain adaptation: A survey, Neurocomputing 312 (2018) 135–153.

[16] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, V. Lempitsky, Domain-adversarial training of neural networks, The Journal of Machine Learning Research 17 (1) (2016) 2096–2030.

[17] N. Otsu, A threshold selection method from gray-level histograms, IEEE transactions on systems, man, and cybernetics 9 (1) (1979) 62–66.

[18] A. Perez, R. C. Gonzalez, An iterative thresholding algorithm for image segmentation, IEEE transactions on pattern analysis and machine intelligence PAMI-9 (6) (1987) 742–751.

[19] R. F. Moghaddam, M. Cheriet, Adotsu: An adaptive and parameterless generalization of otsu's method for document image binarization, Pattern Recognition 45 (6) (2012) 2419–2431.

[20] W. Niblack, An introduction to digital image processing, Strandberg Publishing Company, 1985.

[21] C. Wolf, J. M. Jolion, F. Chassaing, Text localization, enhancement and binarization in multimedia documents, in: Proceedings of the International Conference on Pattern Recognition, Vol. 2, 2002, pp. 1037–1040.

[22] A. Mishra, K. Alahari, C. Jawahar, An mrf model for binarization of natural scene text, in: 2011 International Conference on Document Analysis and Recognition, IEEE, 2011, pp. 11–16.

[23] F. Jia, C. Shi, K. He, C. Wang, B. Xiao, Degraded document image binarization using structural symmetry of strokes, Pattern Recognition 74 (2018) 225–240.

[24] S. S. Lokhande, N. A. Dawande, A survey on document image binarization techniques, in: 2015 International Conference on Computing Communication Control and Automation, 2015, pp. 742–746. `doi: 10.1109/ICCUBEA.2015.148`.

[25] A. Sulaiman, K. Omar, M. Nasrudin, Degraded historical document binarization: A review on issues, challenges, techniques, and future directions, Journal of Imaging 5 (4) (2019) 48.

[26] C. Tensmeyer, T. R. Martinez, Historical document image binarization: A review, SN Comput. Sci. 1 (3) (2020) 173.

[27] M. A. Hearst, S. T. Dumais, E. Osuna, J. Platt, B. Scholkopf, Support vector machines, IEEE Intelligent Systems and their applications 13 (4) (1998) 18–28.

[28] C.-H. Chou, W.-H. Lin, F. Chang, A binarization method with learning-built rules for document images produced by cameras, Pattern Recognition 43 (4) (2010) 1518–1530.

[29] W. Xiong, J. Xu, Z. Xiong, J. Wang, M. Liu, Degraded historical document image binarization using local features and support vector machine (svm), Optik 164 (2018) 218–223.

[30] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, nature 521 (7553) (2015) 436–444.

[31] M. Z. Afzal, J. Pastor-Pellicer, F. Shafait, T. M. Breuel, A. Dengel, M. Liwicki, Document image binarization using lstm: A sequence learning approach, in: Proceedings of the 3rd International Workshop on Historical Document Imaging and Processing, 2015, pp. 79–84.

[32] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, J. Schmidhuber, Lstm: A search space odyssey, IEEE transactions on neural networks and learning systems 28 (10) (2016) 2222–2232.

[33] F. Milletari, N. Navab, S.-A. Ahmadi, V-net: Fully convolutional neural networks for volumetric medical image segmentation, in: 2016 Fourth International Conference on 3D Vision (3DV), IEEE, 2016, pp. 565–571.

[34] K. R. Ayyalasomayajula, F. Malmberg, A. Brun, Pdnet: Semantic segmentation integrated with a primal-dual network for document binarization, Pattern Recognition Letters 121 (2019) 52–60.

[35] H. Yan, Y. Ding, P. Li, Q. Wang, Y. Xu, W. Zuo, Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 2272–2281.

[36] B. Sun, J. Feng, K. Saenko, Return of frustratingly easy domain adaptation, arXiv preprint arXiv:1511.05547 (2015).

[37] J. Shen, Y. Qu, W. Zhang, Y. Yu, Wasserstein distance guided representation learning for domain adaptation, arXiv preprint arXiv:1707.01217 (2017).

[38] A. Rozantsev, M. Salzmann, P. Fua, Beyond sharing weights for deep domain adaptation, IEEE transactions on pattern analysis and machine intelligence 41 (4) (2018) 801–814.

[39] B. Sun, K. Saenko, Deep coral: Correlation alignment for deep domain adaptation, in: European conference on computer vision, Springer, 2016, pp. 443–450.

[40] B. Bhushan Damodaran, B. Kellenberger, R. Flamary, D. Tuia, N. Courty, Deepjdot: Deep joint distribution optimal transport for unsupervised domain adaptation, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 447–463.

[41] C. Villani, Optimal transport: old and new, Vol. 338, Springer Science & Business Media, 2008.

[42] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, W. Li, Deep reconstruction-classification networks for unsupervised domain adaptation, in: European Conference on Computer Vision, Springer, 2016, pp. 597–613.

[43] P. Isola, J.-Y. Zhu, T. Zhou, A. A. Efros, Image-to-image translation with conditional adversarial networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1125–1134.

[44] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in neural information processing systems, 2014, pp. 2672–2680.

[45] A.-J. Gallego, J. Calvo-Zaragoza, R. B. Fisher, Incremental unsupervised domain-adversarial training of neural networks, IEEE Transactions on Neural Networks and Learning Systems (2020).

[46] D. Zou, Q. Zhu, P. Yan, Unsupervised domain adaptation with dual-scheme fusion network for medical image segmentation, in: C. Bessiere (Ed.), Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20, International Joint Conferences on Artificial Intelligence Organization, 2020, pp. 3291–3298.

[47] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2223–2232.

[48] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, CyCADA: Cycle-consistent adversarial domain adaptation, in: J. Dy, A. Krause (Eds.), Proceedings of the 35th International Conference on Machine Learning, Vol. 80 of Proceedings of Machine Learning Research, PMLR, 2018, pp. 1989–1998.

[49] Y. Li, L. Yuan, N. Vasconcelos, Bidirectional learning for domain adaptation of semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 6936–6945.

[50] M. M. Haq, J. Huang, Adversarial domain adaptation for cell segmentation, in: Medical Imaging with Deep Learning, PMLR, 2020, pp. 277–287.

[51] X. Peng, Y. Li, K. Saenko, Domain2vec: Domain embedding for unsupervised domain adaptation, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, Springer International Publishing, Cham, 2020, pp. 756–774.

[52] K. Osumi, T. Yamashita, H. Fujiyoshi, Domain adaptation using a gradient reversal layer with instance weighting, in: 2019 16th International Conference on Machine Vision Applications (MVA), 2019, pp. 1–5. `doi:10.23919/MVA.2019.8757975`.

[53] T. Salimans, I. J. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, Improved techniques for training gans, CoRR abs/1606.03498 (2016). `arXiv:1606.03498`.
URL `http://arxiv.org/abs/1606.03498`

[54] J. Benesty, J. Chen, Y. Huang, I. Cohen, Pearson correlation coefficient, in: Noise reduction in speech processing, Springer, 2009, pp. 1–4.

[55] F. Chollet, et al., Keras, `https://github.com/fchollet/keras` (2015).

[56] B. Gatos, K. Ntirogiannis, I. Pratikakis, Icdar 2009 document image binarization contest (dibco 2009), in: 2009 10th International conference on document analysis and recognition, IEEE, 2009, pp. 1375–1382.

[57] S. M. Ayatollahi, H. Z. Nafchi, Persian heritage image binarization competition (phibc 2012), in: 2013 First Iranian Conference on Pattern Recognition and Image Analysis (PRIA), IEEE, 2013, pp. 1–4.

[58] J.-C. Burie, M. Coustaty, S. Hadi, M. W. A. Kesiman, J.-M. Ogier, E. Paulus, K. Sok, I. M. G. Sunarya, D. Valy, Icfhr2016 competition on the analysis of handwritten text in images of balinese palm leaf manuscripts, in: 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR), IEEE, 2016, pp. 596–601.

[59] J. Shlens, Notes on kullback-leibler divergence and likelihood, arXiv preprint arXiv:1404.2000 (2014).

[60] J. Briët, P. Harremoës, Properties of classical and quantum jensen-shannon divergence, Physical review A 79 (5) (2009) 052311.

[61] S. Lee, J. H. Xin, S. Westland, Evaluation of image similarity by histogram intersection, Color Research & Application: Endorsed by Inter-Society Color Council, The Colour Group (Great Britain), Canadian Society for Color, Color Science Association of Japan, Dutch Society for the Study of Color, The Swedish Colour Centre Foundation, Colour Society of Australia, Centre Français de la Couleur 30 (4) (2005) 265–274.