# Domain Invariant Representations with Deep Spectral Alignment

Christoph Raab[1], Peter Meier[1] and Frank-Michael Schleif[1] *

1- University of Applied Science Würzburg-Schweinfurt - Department of Computer Science, Sanderheinrichsleitenweg 20, Würzburg - Germany

**Abstract**.    Similar as traditional algorithms, deep learning networks struggle in generalizing across domain boundaries. A current solution is the simultaneous training of the classification model and the minimization of domain differences in the deep network. In this work, we propose a new *un*supervised deep domain adaptation architecture, which trains a classifier and minimizes the difference of spectral properties of the covariance matrix of the data. Evaluated against standard architectures and datasets, the approach shows an alignment with respect to the data variance between related domains.

## 1    Introduction

Recent developments in deep learning models have led to a significant increase in predictive performance in various applications, such as computer vision or speech recognition. Despite the outstanding results and multiple efforts to improve the generalization capabilities of deep learning networks, they still struggle to generalize across domain boundaries. This is caused due to differences in the distributions between the different domains [1].

A current solution to this problem is to over-generalize the networks by training on massive datasets and then specialize these networks by fine-tuning them to the desired target domain [1]. However, this solution assumes that there is a sufficient amount of labeled data in the target domain, but in many cases, this cannot be guaranteed [2]. So if fine-tuning is not feasible, unsupervised deep domain adaptation techniques are used to learn a network that achieves good classification results in the target domain [3]. Unsupervised deep domain adaptation makes the following assumptions: There exists a pre-trained over-generalized network, a small amount of labeled data in a source domain, and a small amount of unlabelled data in the target domain [4]. The goal is to learn a network that achieves good classification results and at the same time, reduces the differences between the distributions of the source and target domains, to apply it in the target domain.

In this work, we study the effect of the generalization properties of deep networks by adjusting the spectrum of the output distribution of higher layers of the network. We show that this is an implicit case of minimizing the difference of variance or second central moment of the distributions. Based on this, Section 3 introduces a new Deep Domain Adaptation architecture that simultaneously

learns a classifier and a domain invariant representation. We present the evaluation of our Spectral Alignment loss ($\mathcal{L}_{AS}$) in the Deep Spectral Network (DSN) against other established approaches and data sets in the field in Section 5. A summary and open problems are provided at the end of this paper.

## 2 Background and Related Work

For unsupervised deep domain adaptation [5, 6, 7, 8], we consider a labeled source dataset $D_s = \{\mathbf{X}_s, Y_s\} = \{\mathbf{x}_i, y_i\}_{i=1}^n \overset{i.i.d.}{\sim} p(\mathcal{S})$ in the source domain $\mathcal{S}$ and an unlabeled target dataset $D_t = \{\mathbf{X}_t, Y_t\} = \{\mathbf{x}_j, y_j\}_{j=1}^m \overset{i.i.d.}{\sim} p(\mathcal{T})$ in the target domain $\mathcal{T}$ with same label space $\forall i, j : y_i, y_j \in \mathcal{C}$ but different distributions $p(\mathcal{S}) \neq p(\mathcal{T})$. The overall goal is still to learn a classifier model, but additionally, it should generalize in a related target domain. Initially, we start with $D_s$ and learn parameters $\theta$ of a neural network $f_\theta : \mathcal{X} \times \mathcal{Y}$ minimizing the risk $R[L(f(\mathbf{x}; \theta), y)]$. The loss can be defined for example as the cross-entropy loss $L(f(\mathbf{x}; \theta), y) = -\sum_{c \in C} y_c log(f_\theta(\mathbf{x})_c)$. The loss is part of the empirical objective function

$$\arg \min_\theta E[L(f(\mathbf{X}_s; \theta), Y_s)]. \tag{1}$$

The network itself consists of multiple hidden layers and an output or classification layer. Consider $g(\mathbf{X}_s; \theta)_l = a(f(\mathbf{x}; \theta)_l)_l$ as the layer $l$ with an activation function $a(\cdot)_l$ and parameter layer $f(\cdot)_l$ given source data and for target data $g(\mathbf{X}_t; \theta)_l$ analogously.

Recent work in the field uses one or more higher layers, i.e. the fully connected layers of the network to adapt the output distributions of the (hidden) layers $g(\mathbf{X}_s; \theta)_l$ and $g(\mathbf{X}_t; \theta)_l$[5]. This leads to very individualized approaches. To measure the difference between the output distributions of the network, some type of dissimilarity measure $d : g(\mathbf{X}_s; \theta)_l \times g(\mathbf{X}_t; \theta)_l) \to \mathbb{R}^+$ is employed and added to the objective function:

$$\arg \min_\theta E[L(f(\mathbf{X}_s; \theta), Y_s)] + \lambda d(g(\mathbf{X}_s; \theta)_l, g(\mathbf{X}_t; \theta)_l). \tag{2}$$

The dissimilarity measure is used as a regularization. The parameter $\lambda \in [0, \infty]$ controls the trade-off between aligning the output distributions and minimizing the classification objective. By setting a high $\lambda$, the distributions are very closely aligned with each other, but the feature representations are degenerated, causing a high loss [2].

In the following, we will discuss prior related work, minimizing the statistical properties [1] of source and target distributions of the output of some layer. A commonly used dissimilarity measure is the Maximum Mean Discrepancy (MMD) [9], which is the difference in mean of two matrices in a reproducing kernel Hilbert space (RKHS). In [6], the gradient of the MMD with respect to the parameters is used additionally for the parameter updates. The work in [2] uses the MMD to directly propagate the update step with respect to data, causing the parameters to change accordingly. The original MMD considers

only the prior distribution, while in [8] (used in the JAN network), a joint-MMD was proposed minimizing the conditional distribution discrepancy. The minimization of MMD in the proposed networks can be seen as an alignment of statistical moments given a particular kernel, e.g. RBF-Kernel, of the two domains [10]. The authors of [5] proposed the Central Moment Discrepancy for domain adaptation, which strives for explicitly minimizing higher central moments. The CORAL Loss [7] is most related to us, minimizing the difference of the full covariance matrices between two domains. Our proposal is a particular case of CORAL, by aligning *only* the singular value spectrum of the domains. As a side effect, we are also aligning the covariances. However, our DSN minimizes a diagonal matrix, which is easier to compute. Further, we do not rely on a particular kernel matrix nor kernel function, but any positive semi definite (psd) kernel can be used. Due to this flexibility, it is also rather easy to integrate low-rank approximation techniques. An interpretation of our loss is the minimization of the second central moment between domains.

## 3  Deep Spectral Alignment

In this section, we present the spectral alignment and the modified network architecture. Let $\mathbf{X}_s^l$ and the $\mathbf{X}_t^l$ be the output of source and target from layer $k$, respectively. The Singular Value Decomposition (SVD) of these outputs is given with $\mathbf{X}_s^l = \mathbf{U\Sigma V}^T$ and $\mathbf{X}_t^l = \mathbf{LTR}^T$. Here $\mathbf{R}, \mathbf{V} \in \mathbb{R}^{d \times d}$, $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{L} \in \mathbb{R}^{m \times m}$ are matrices. Further, $\mathbf{U}, \mathbf{L}, \mathbf{V}$ and $\mathbf{R}$ are column-orthogonal. $\mathbf{\Sigma}$ is a $n \times d$ matrix wherein all entries $\sigma_{ij} = 0$ iff $i \neq j$ and $\mathbf{T}$ is a $m \times d$ matrix wherein all entries $t_{ij} = 0$ iff $i \neq j$. Furthermore, by $\sigma_k$ we denote the singular value in the k-th column of $\mathbf{\Sigma}$. For a linear covariance function, we can decompose the respective kernel with the eigenvalue decomposition and the SVD into

$$\mathbf{K} = \mathbf{CDC}^{-1} = \mathbf{X}^T\mathbf{X} = (\mathbf{V\Sigma U}^T)(\mathbf{U\Sigma V}^T) = \mathbf{C\Sigma}^2\mathbf{C}^{-1}, \qquad (3)$$

where $\mathbf{C}$ are the eigenvectors (right singular values of $\mathbf{X}$) and $\mathbf{D}$ are the eigenvalues of $\mathbf{K}$. The singular values $\mathbf{\Sigma}$ of $\mathbf{X}$ are the square root eigenvalues of $\mathbf{K}$. Accordingly, the entries of the diagonal of $\mathbf{D}, \mathbf{\Sigma}^2$ give the variance of the columns of $\mathbf{K}$. Assuming that the expected values $\mu_s, \mu_t = 0$ of $\mathbf{X}_s^l$ and $\mathbf{X}_t^l$, minimizing the difference between $\mathbf{\Sigma}$ and $\mathbf{T}$ is the same as minimizing the variance of the covariance matrix. Due to the low number of parameters in comparison to other domain adaptation approaches [11, 7, 8], the spectral loss is easy to minimize.

For domain adaptation, we implement for the discrepancy measure $d(\cdot)$ in Eq. (2) the spectral loss

$$\mathcal{L}_{SA} = \frac{1}{n}||\mathbf{\Sigma} - \mathbf{T}||_F^2, \qquad (4)$$

to extend the classification loss. Where $\frac{1}{n}$ is a scaling constant and $||\cdot||_F^2$ denotes the squared Frobenius norm. Because we assume equal batch sizes for source and target data, we have $m = n$. Minimizing Eq. (4) will align the spectra of source and target. Accordingly, the variances of the outputs in layer $l$ are aligned to each other and a domain invariant representation for source and target is

learned. In the following, the derivative of Eq. (4) with respect to the data is given, which allows the optimization of the spectral parameters. Data-driven derivations are common in deep domain adaptation [5, 7, 8, 10]. Due to lack of space, we only give the derivative with respect to the source data and in the following $x_{ij} \in \mathbf{X}_s^l$. The derivative with respect to target data is straightforward. For computing $\frac{\partial \mathcal{L}_{SA}}{\partial x_{ij}}$, it is easy to see that the partial derivative of $||\mathbf{\Sigma} - \mathbf{T}||_F^2$ has the same form like $||\mathbf{\Sigma}||_F^2$, because $\mathbf{T}$ is constant and the problem reduces to compute the partial derivative of all $\sigma_k \in \mathbf{\Sigma}$ with respect to $x_{ij}$.

We follow the solution of [12] and analyze

$$\frac{\partial \mathbf{X}_s^l}{\partial x_{ij}} = \frac{\partial (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)}{\partial x_{ij}} = \frac{\partial \mathbf{U}}{\partial x_{ij}} \mathbf{\Sigma} \mathbf{V}^T + \mathbf{U} \frac{\partial \mathbf{\Sigma}}{\partial x_{ij}} \mathbf{V}^T + \mathbf{U} \mathbf{\Sigma} \frac{\partial \mathbf{V}^T}{\partial x_{ij}}. \quad (5)$$

Let $\mathbf{\Omega}_U^{ij} = \mathbf{U}^T \frac{\partial \mathbf{U}}{\partial x_{ij}}$ and let $\mathbf{\Omega}_V^{ij} = \frac{\partial \mathbf{V}^T}{\partial x_{ij}} \mathbf{V}$ depending on $x_{ij}$. By multiplying Eq. (5) with the orthogonal matrices $\mathbf{U}^T$ and $\mathbf{V}$ from left and right respectively, and inserting the definitions of $\mathbf{\Omega}_{(\cdot)}$ we obtain $\mathbf{U}^T \frac{\partial \mathbf{X}_s^l}{\partial x_{ij}} \mathbf{V} = \mathbf{\Omega}_U^{ij} \mathbf{\Sigma} + \frac{\partial \mathbf{\Sigma}}{\partial x_{ij}} + \mathbf{\Sigma} \mathbf{\Omega}_V^{ij}$. Because $\mathbf{\Omega}_U^{ij}$ and $\mathbf{\Omega}_V^{ij}$ are anti symmetric, zero on the diagonal and $\mathbf{\Sigma}$ is diagonal, the terms $\mathbf{\Omega}_U^{ij} \mathbf{\Sigma}$ and $\mathbf{\Omega}_V^{ij} \mathbf{\Sigma}$ are zero [12]. Also $\frac{\partial \mathbf{X}_s^l}{\partial x_{ij}} = 0$, for all entries $x_{kl}$ with $(k, l) \neq (i, j)$ and 1, otherwise. Putting all together we obtain $\frac{\partial \sigma_k}{\partial x_{ij}} = u_{ik} v_{jk}$.

Looking again at Eq. (4), the derivative of $\mathcal{L}_{\mathcal{SA}}$ is given by the derivative of $\sigma_k$ w.r.t. $x_{ij}$. Hence, the derivative of the loss function is

$$\frac{\partial \mathcal{L}_{SA}}{\partial x_{ij}} = \frac{1}{n} \frac{\partial \sigma_k}{\partial x_{ij}} \left[ \sqrt{\Sigma_{k=1}^n (\sigma_k - t_k)^2} \right]^2 = \frac{2}{n} \Sigma_{k=1}^n (\sigma_k - t_k) \frac{\partial \sigma_k}{\partial x_{ij}} \quad (6)$$

$$= \frac{2}{n} \Sigma_{k=1}^n (\sigma_k - t_k) \cdot u_{ik} v_{jk}. \quad (7)$$

The loss can be integrated into any layer as regularization term or simultaneously used in multiple layers, as suggested by [8]. Here we use the proposed approach in the last layer.

## 4 Experiments

We evaluate our network against other recent approaches on the standard image dataset Office-31. Source code and datasets are available at `github.com/ChristophRaab/DSN`. The study follows the standard protocol for evaluating unsupervised deep domain adaptation [8] and utilizes all available source data for learning and all target data for knowledge transfer and evaluation. Other datasets or sampling strategies [13] are omitted due to space issues. The **Office-31** dataset contains images from three separated domains, namely Amazon, webcam, and digital single-lens reflex camera (DSLR). Each domain has 31 classes with objects frequently located in the office. Since all three domains are acquired with different settings and photo cameras, the adaptation problem is to train on one domain and to test on another. The results are a summary from five

| Dataset | AlexNet [15] | DDC [15] | DAN [11] | JAN [8] | CORAL [7] | DSN |
|---|---|---|---|---|---|---|
| A → W | 48.0 (2.4) | 49.3 (0.8) | 62.8 (0.8) | **63.0** (2.5) | 59.9 (1.0) | 59.6 (2.8) |
| A → D | 55.6 (2.2) | 59.0 (2.0) | **63.6** (0.8) | 63.1 (1.4) | 59.6 (1.0) | 62.3 (0.8) |
| W → A | 54.9 (0.5) | 55.2 (1.6) | 55.6 (0.8) | 55.8 (1.3) | 53.4 (1.3) | **56.2** (0.5) |
| W → D | 98.4 (0.3) | 98.5 (0.5) | 98.8 (0.2) | 98.5 (0.4) | **99.0** (0.3) | 98.4 (0.2) |
| D → A | 53.4 (0.7) | 53.7 (0.8) | 55.2 (1.3) | 56.1 (0.5) | 52.3 (1.1) | **56.5** (1.6) |
| D → W | 94.3 (0.7) | 94.7 (0.7) | **95.8** (0.2) | 95.7 (0.5) | 94.8 (0.3) | 95.0 (0.7) |
| Mean | 67.4 (1.1) | 68.4 (1.1) | 72.0 (0.7) | 72.0 (1.1) | 69.8 (0.8) | 71.3 (1.1) |

Table 1: Mean prediction accuracy with standard deviation in brackets on the Office-31 dataset over five runs.

test runs on the dataset combinations A→D (Amazon to DSLR), A→W, D→A, D→W, W→A, and W→D and presented as mean accuracy with standard deviation. The parameters are optimized as shown in [14]. The baseline network is Alexnet [15], which is pre-trained on the Imagenet dataset. For the Alexnet, we follow [7] and modify the classifier layer to have 4096 input dimensions and 31 output dimensions. For a fair comparison, we implement our loss and all competitive losses in the classification layer of the network. Hence, the loss is propagated through all layers of the network. All approaches are trained with a learning rate of $1e^{-3}$ and momentum stochastic gradient descent with a decay of 0.9. As suggested in [7], the classifier layer has ten times the learning rate as the remaining network, i.e. 0.01. We used 50 epochs to retrain the network and $\lambda = 0.5$ for all networks.

Results are reported in Tab. 1. Our DSN is compared against the unsupervised deep domain adaptation algorithms DDC [15], DAN [11], JAN [8], and CORAL [7] and vanilla AlexNet as baseline. The best performing network is DAN, but the performance gap in accuracy from DAN to us is, on average, only 0.7 %. Because CORAL is worse than DSN in the performance experiments, we conclude that minimizing the spectral loss contributes more to the adaptation process of the network as minimizing the difference of a covariance matrix. The minimization of non-linear kernel measures as in DDC or DAN does not necessarily lead to better adaptation capacities. Finally, the tuning of the DSN is not expensive because it has no tunable parameters besides the regularization parameter, which makes it easier to apply compared to other approaches. The convergence behavior is not shown due to space issues.

## 5   Conclusion

We studied the effect of the spectral alignment in deep networks between two domains. The proposed Deep Spectral Network has competitive performance to state of the art unsupervised deep domain adaptation approaches. The proposed spectral loss is easy to optimize and the loss has no free parameters. This makes the approach more favorable compared to other networks, as the DSN is more comfortable to apply while maintaining competitiveness. In the future, non-linear kernels and multi-layer losses should be studied.

# References

[1] Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 312:135–153, 2018.

[2] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep Domain Confusion: Maximizing for Domain Invariance. *CoRR*, abs/1412.3, 2014.

[3] Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang, Chao Yang, and Chunfang Liu. A survey on deep transfer learning. In Věra Krurková, Yannis Manolopoulos, Barbara Hammer, Lazaros Iliadis, and Ilias Maglogiannis, editors, *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 11141 LNCS, pages 270–279, Cham, 2018. Springer International Publishing.

[4] Mingsheng Long, Jianmin Wang, Yue Cao, Jiaguang Sun, and Philip S. Yu. Deep learning of transferable representation for scalable domain adaptation. *IEEE Transactions on Knowledge and Data Engineering*, 28(8):2027–2040, 2016.

[5] Werner Zellinger, Thomas Grubinger, Edwin Lughofer, Thomas Natschläger, and Susanne Saminger-Platz. Central Moment Discrepancy (CMD) for Domain-Invariant Representation Learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.

[6] Muhammad Ghifary, W Bastiaan Kleijn, and Mengjie Zhang. Domain Adaptive Neural Networks for Object Recognition. In Duc-Nghia Pham and Seong-Bae Park, editors, *PRICAI 2014: Trends in Artificial Intelligence*, pages 898–904, Cham, 2014. Springer International Publishing.

[7] Baochen Sun and Kate Saenko. Deep CORAL: Correlation Alignment for Deep Domain Adaptation. In Gang Hua and Hervé Jégou, editors, *Computer Vision – ECCV 2016 Workshops*, pages 443–450, Cham, 2016. Springer International Publishing.

[8] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep Transfer Learning with Joint Adaptation Networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pages 2208–2217. JMLR.org, 2017.

[9] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A Kernel Two-sample Test. *J. Mach. Learn. Res.*, 13:723–773, 2012.

[10] Yujia Li, Kevin Swersky, and Rich Zemel. Generative Moment Matching Networks. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1718–1727, Lille, France, 2015. PMLR.

[11] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning Transferable Features with Deep Adaptation Networks. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 97–105, Lille, France, 2015. PMLR.

[12] Théodore Papadopoulo and Manolis I A Lourakis. Estimating the Jacobian of the Singular Value Decomposition: Theory and Applications. In *Computer Vision - ECCV 2000*, pages 554–570, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.

[13] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2066–2073, 2012.

[14] Erheng Zhong, Wei Fan, Qiang Yang, Olivier Verscheure, and Jiangtao Ren. Cross validation framework to choose amongst models and datasets for transfer learning. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 6323 LNAI(PART 3):547–562, 2010.

[15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In F Pereira, C J C Burges, L Bottou, and K Q Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.