

Domain Adaptation: Tackle Distribution Shift Without Access to Target Label

Théo Gnassounou

Reading Group Hi!Paris, 11-03-2025



What is Domain Adaptation (DA)?

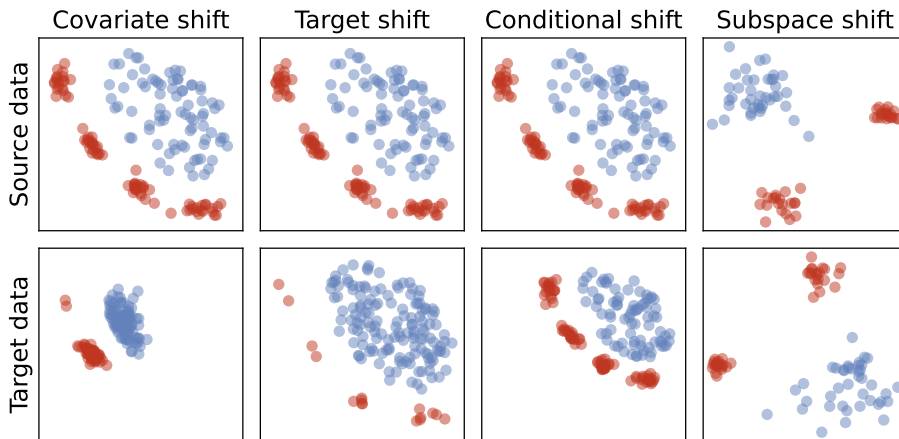
- Two type of domains: **Source** and **Target** .
- **Source** domains **with label** and **Target** domains **without label** .
- Assumption \rightarrow **shift** between the distribution of the domain's data

4 different types of shifts

- **Covariate Shift** (CS): $\mathcal{P}_x^s(x) \neq \mathcal{P}_x^t(x)$, $\mathcal{P}^s(y|X) = \mathcal{P}^t(y|X)$
- **Target Shift** (TS): $\mathcal{P}_y^s(x) \neq \mathcal{P}_y^t(x)$, $\mathcal{P}^s(X|y) = \mathcal{P}^t(X|y)$
- **Conditional Shift** (CondS): $\mathcal{P}^s(y|X) \neq \mathcal{P}^t(y|X)$ or $\mathcal{P}^s(X|y) \neq \mathcal{P}^t(X|y)$
- **Subspace Shift** (SS): $\mathcal{P}^s(X) \neq \mathcal{P}^t(X)$ but it exists a subspace projection W such that $\mathcal{P}^s(WX) = \mathcal{P}^t(WX)$

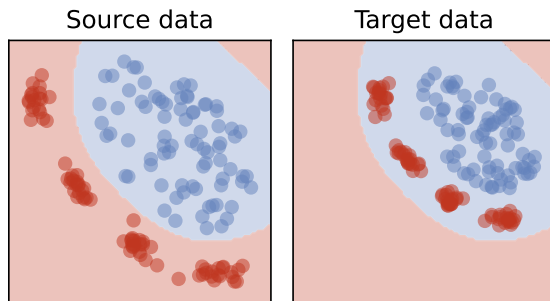
What is Domain Adaptation (DA)?

- Two type of domains: **Source** and **Target** .
- Source** domains **with label** and **Target** domains **without label** .
- Assumption → **shift** between the distribution of the domain's data



What is Domain Adaptation (DA)?

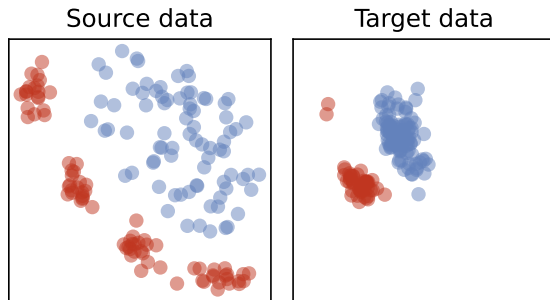
- Two type of domains: **Source** and **Target** .
- **Source** domains **with label** and **Target** domains **without label** .
- Assumption → **shift** between the distribution of the domain's data



→ **Problem: Drop in performance** when applying a model trained on the source to the target.

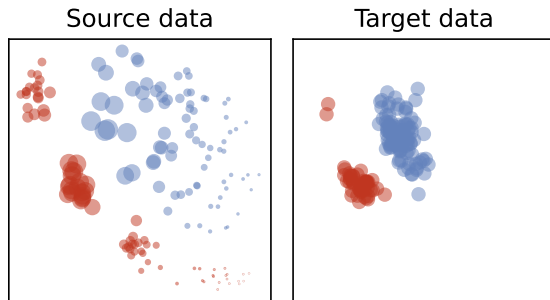
Traditional DA methods: Covariate Shift and target shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Reweighting**



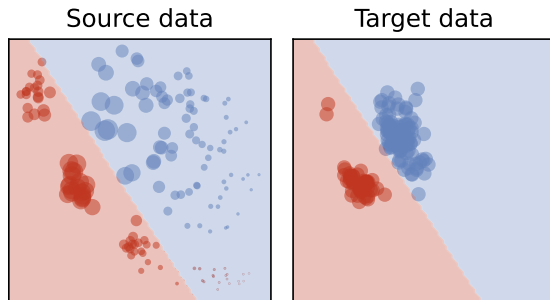
Traditional DA methods: Covariate Shift and target shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Reweighting**



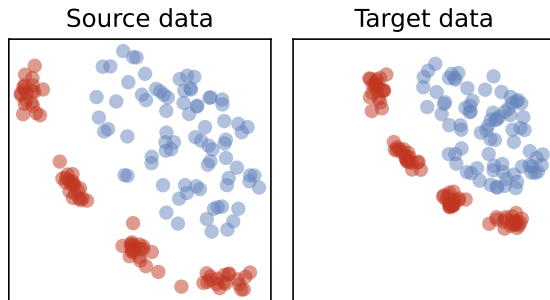
Traditional DA methods: Covariate Shift and target shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Reweighting**



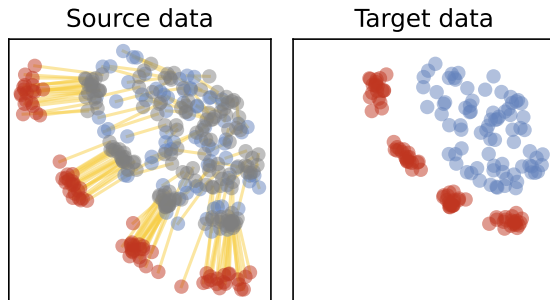
Traditional DA methods: Conditional Shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Mapping**



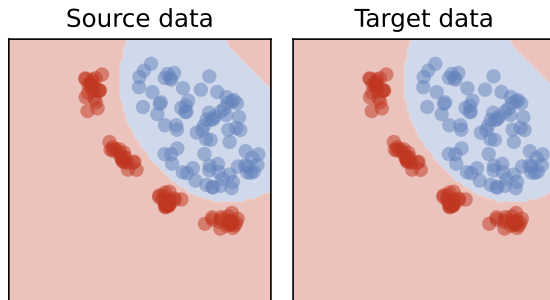
Traditional DA methods: Conditional Shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Mapping**



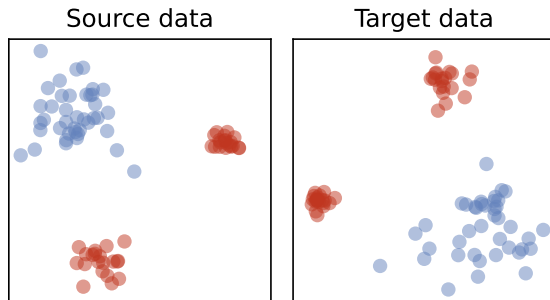
Traditional DA methods: Conditional Shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Mapping**



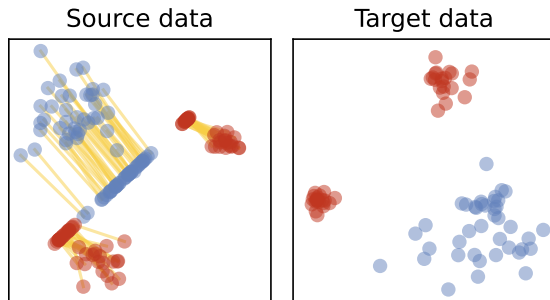
Traditional DA methods: Subspace shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Subspace**



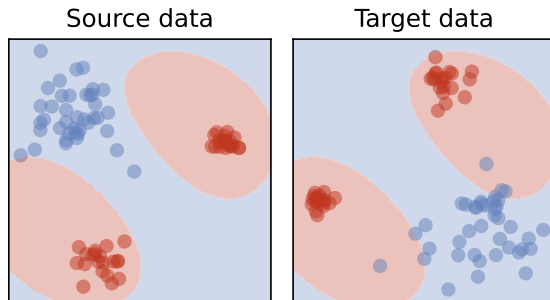
Traditional DA methods: Subspace shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Subspace**



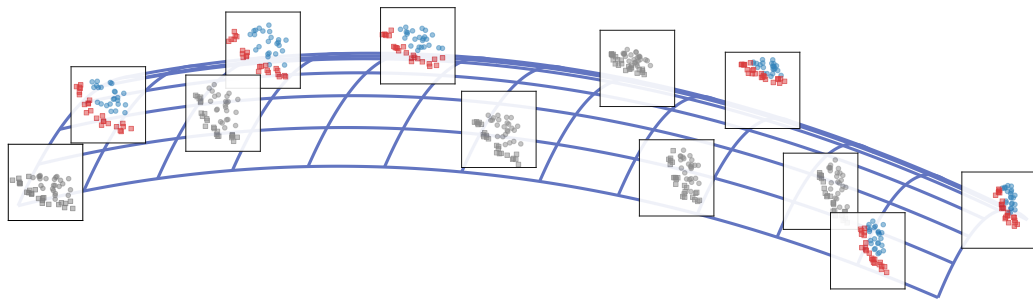
Traditional DA methods: Subspace shift

- **One** source (X_s, y_s) and **one** target (X_t, y_t)
- **Adapt** the source to the target via: **Subspace**



Multi-source multi-target Domain Adaptation

Domain manifold



Source-free Domain Adaptation (or Test-Time DA)

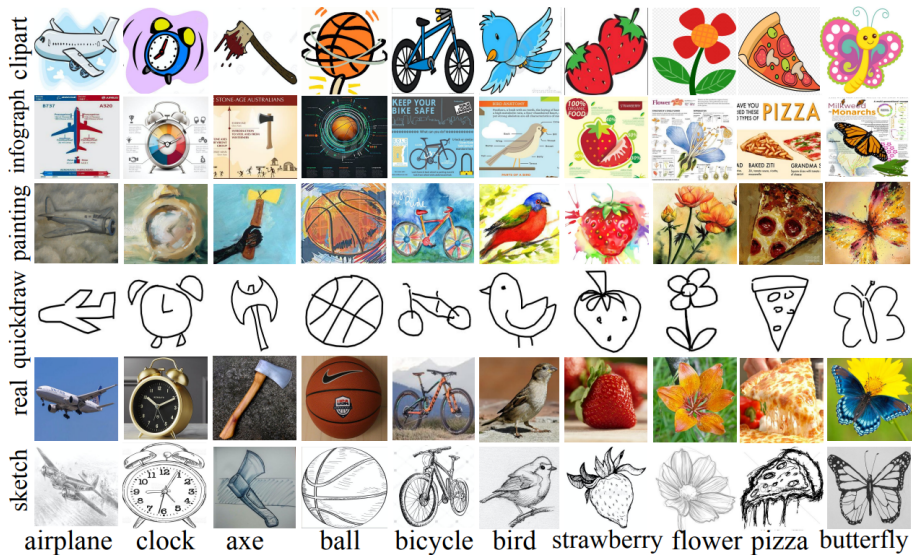
1. Train-time

- Access to **Source** domains with labels
- **No** access to **Target** domains
- **Train** a model on the source domains with labels

2. Test-time

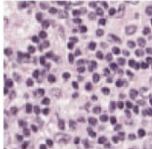
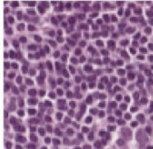
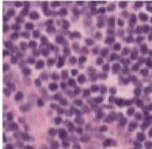
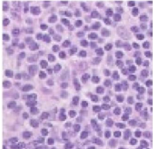
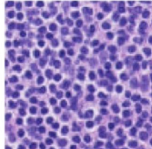
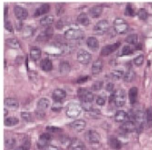
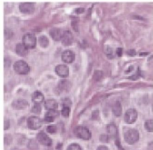
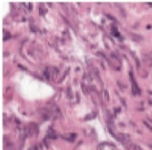
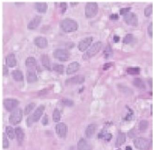
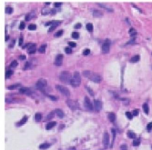
- **No** access to **Source** domains
- Access to **Target** domains **without** labels
- **Finetune** the model on the target domains without access to the target labels

Real-world applications: Computer Vision



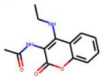
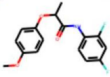
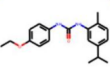
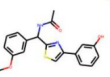
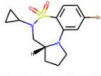
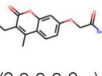
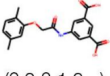
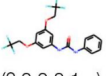
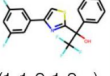
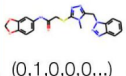
Source: Peng et. al., 2019

Real-world applications: Computer Vision

Train			Val (OOD)	Test (OOD)	
	d = Hospital 1	d = Hospital 2	d = Hospital 3	d = Hospital 4	d = Hospital 5
y = Normal					
y = Tumor					

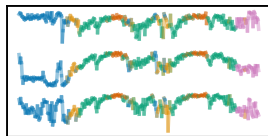
Source: Koh et. al., WILDS, 2020

Real-world applications: Biology

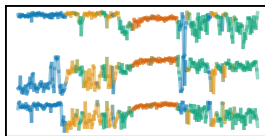
Train				Test
Scaffold 11	Scaffold 32	Scaffold 321	Scaffold 4413	Scaffold 54113
 (1,0,?,0,?,...)	 (?,0,0,0,?,...)	 (0,1,1,0,0,...)	 (?,0,0,0,?,...)	 (0,?,1,?,0,...)
 (?,0,0,0,?,...)	 (?,0,?,1,0,...)	 (?,0,0,0,1,...)	 (1,1,0,1,0,...)	Scaffold 65912
				 (0,1,0,0,0,...)

Source: Koh et. al., WILDS, 2020

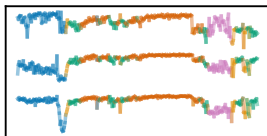
Real-world applications: Time series



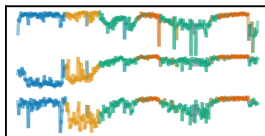
Domain 1
 (X_1, y_1)



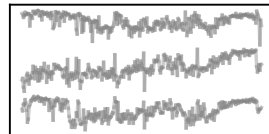
Domain 2
 (X_2, y_2)



Domain 3
 (X_3, y_3)



Domain 4
 (X_4, y_4)



Domain 5
 (X_5)

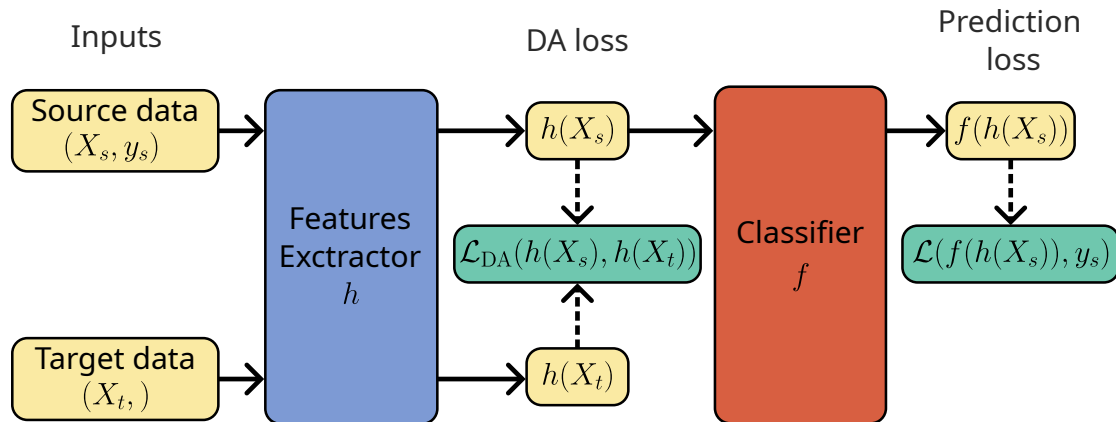
Example: Sleep stage classification from EEG signals¹

¹*Gnassounou et. al., 2023*

Deep Learning for DA

Why use Deep Learning?

- More **fine-grained feature** extraction
- End-to-end** learning: no need for adaptation step



Deep learning for DA

$$\mathcal{L}_{\text{tot}}(h, f) = \underbrace{\mathcal{L}(f(h(X_s)), y_s)}_{\text{Prediction Loss}} - \underbrace{\lambda}_{\text{Regularization}} \underbrace{\mathcal{L}_{\text{DA}}(f(h(X_s)), f(h(X_t)))}_{\text{DA loss}},$$

- **Loss** → **Cross-entropy** loss
- **Regularization** → **Threshold** between the loss and the DA loss
- **DA loss** → **Reduce divergence** between source and target features

$$(\hat{h}, \hat{f}) = \underset{h, f}{\operatorname{argmin}} \mathcal{L}_{\text{tot}}(h, f)$$

How to reduce the divergence between source and target features?

Domain-Adversarial Training of Neural Networks

Yaroslav Ganin
Evgeniya Ustinova

Skolkovo Institute of Science and Technology (Skoltech)
Skolkovo, Moscow Region, Russia

Hana Ajakan

Deep CORAL: Correlation Alignment for Deep Domain Adaptation

Hugo Larochelle

Baochen Sun* and Kate Saenko**

University of Massachusetts Lowell, Boston University

Abstract. Deep neural networks are able to learn powerful representations from large quantities of labeled input data, however they cannot always generalize well across changes in input distributions. Domain adaptation algorithms have been proposed to compensate for the degradation in performance due to domain shift. In this paper, we address the case when the target domain is unlabeled, requiring unsupervised adaptation. CORAL[1] is a “frustratingly easy” unsupervised domain adaptation method that aligns the second-order statistics of the source and target distributions with a linear transformation. Here, we extend CORAL to learn a nonlinear transformation that aligns correlations of layer activations in deep neural networks (Deep CORAL). Experiments on standard benchmark datasets show state-of-the-art performance.

Abstract

We introduce a new representation learning approach for domain adaptation. We show that data at training and test time come from similar but different distributions. This is directly inspired by the theory on domain adaptation suggesting that, for effective domain transfer to be achieved, predictions must be made based on features that represent

GANIN@SKOLTECH.RU

EVGENIYA.USTINOVA@SKOLTECH.RU

DeepJDOT: Deep Joint Distribution Optimal Transport for Unsupervised Domain Adaptation

¹ Bharath Bhushan Damodaran^{1*}, Benjamin Kellenberger^{2*}, Rémi Flamary³,
HUGO LAROCHELLE⁴, Devis Tuia², Nicolas Courty¹

¹ Université de Bretagne Sud, IRISA, UMR 6074, CNRS, France

² Wageningen University, the Netherlands

FRAN³ Université Côte d’Azur, OCA, UMR 7293, CNRS, Laboratoire Lagrange, France

MAIL {bharath-bhushan.damodaran@irisa.fr, benjamin.kellenberger@wur.nl}

Abstract. In computer vision, one is often confronted with problems of domain shifts, which occur when one applies a classifier trained on a source dataset to target data sharing similar characteristics (e.g. same classes), but also different latent data structures (e.g. different acquisition conditions). In such a situation, the model will perform poorly on the new data, since the classifier is specialized to recognize visual cues specific to the source domain. In this work we explore a solution, named DeepJDOT, to tackle this problem: through a measure of discrepancy on joint deep representations/labels based on optimal transport, we not only learn new data representations aligned between the source and target domain, but also simultaneously preserve the discriminative information used by the classifier. We applied DeepJDOT to a series of visual recognition tasks, where it compares favorably against state-of-the-art deep domain adaptation methods.

DeepCoral: Correlation Alignment¹

With d the dimension of the feature space, the **Coral loss** is defined as:

$$\mathcal{L}_{\text{DA}}(h) = \frac{1}{4d^2} \left\| \underbrace{C(h(X_s))}_{\text{Source Covariance}} - \underbrace{C(h(X_t))}_{\text{Target Covariance}} \right\|_F^2,$$

with $\|\cdot\|_F$ the Frobenius norm. The covariance matrices are defined as:

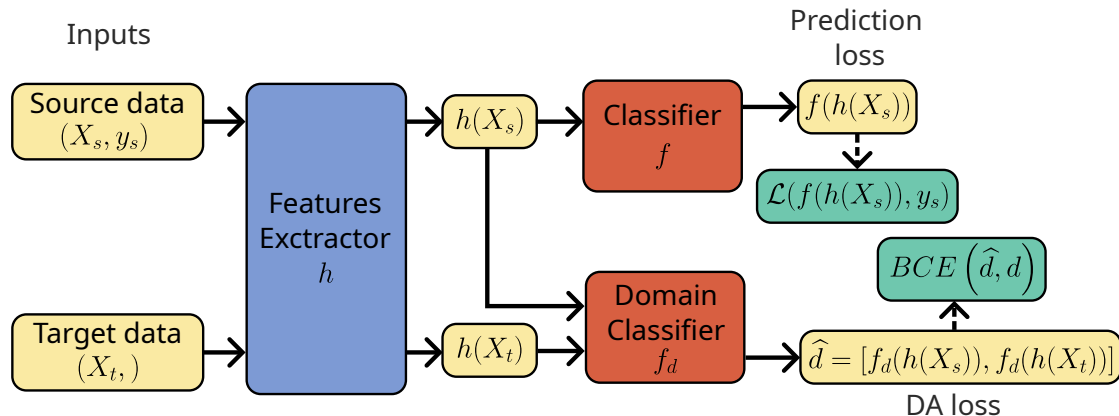
$$C(X) = \frac{1}{N-1} \left(X^\top X - \frac{1}{N} (\mathbf{1}^\top X)^\top (\mathbf{1}^\top X) \right),$$

with $\mathbf{1}$ a vector of ones and N the number of samples.

¹*Sun et. al., 2016*

DANN: Domain Adversarial Neural Network¹

- **Adversarial** training with a **domain classifier**
- **Binary** classification: **d = 0** for **source** and **d = 1** for **target**



¹Ganin et. al., 2016

DANN: Domain Adversarial Neural Network¹

Adversarial loss with binary cross entropy (BCE) loss: Binary cross entropy

$$\mathcal{L}_{\text{DA}}(g, f, f^d) = - \text{BCE}([f_d(h(X_s)), f_d(h(X_t))], [\mathbf{0}, \mathbf{1}]) .$$

Reverse gradient

$$(\hat{h}, \hat{f}) = \underset{h, f}{\operatorname{argmin}} \mathcal{L}_{\text{tot}}(g, f, f^d) ,$$

$$\hat{f}_d = \underset{f_d}{\operatorname{argmax}} \mathcal{L}_{\text{DA}}(g, f, f^d) .$$

In practice:

- **Joint** optimization of the feature extractor and the domain classifier
- **Reverse Gradient Layer**
- $f_d \rightarrow 3$ layers of fully connected layers

¹Ganin et. al., 2016

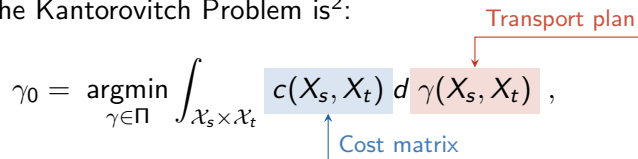
DeepJDOT: Joint Distribution Optimal Transport¹

→ **Optimal Transport** to align the source and target distributions

Brief reminder on Optimal Transport

The relaxed version of the Kantorovitch Problem is²:

$$\gamma_0 = \operatorname{argmin}_{\gamma \in \Pi} \int_{\mathcal{X}_s \times \mathcal{X}_t} c(X_s, X_t) d\gamma(X_s, X_t),$$



where Π is the set of all the **couplings** between marginal distributions μ_s and μ_t .

For discrete OT, introducing the cost matrix $(C)_{i,j} = c(X_s^i, X_t^j)$ the Kantorovitch Problem becomes:

$$\gamma_0 = \operatorname{argmin}_{\gamma \in \mathcal{B}} \langle \gamma, \mathbf{C} \rangle_F, \quad (1)$$

¹Damodaran et. al., 2018

²Peyré et. al., 2019

DeepJDOT: Joint Distribution Optimal Transport¹

The **Joint Distribution Optimal Transport** loss is defined as:

$$\mathcal{L}_{\text{DA}}(h, f) = \langle \gamma, \mathbf{C}_{h,f} \rangle_F ,$$

$C_{h,f}$ the cost matrix defines with the **feature extractor** and the **classifier** :

$$C_{h,f}(X_s^i, X_t^j) = \underbrace{\alpha}_{\text{Regularization}} \underbrace{\|h(X_s^i) - h(X_t^j)\|^2}_{\text{Distance}} + \underbrace{\beta}_{\text{Regularization}} \underbrace{\mathcal{L}(f(h(X_t^i)), y_s^i)}_{\text{Pseudo-labeling loss}} ,$$

- Cost matrix \rightarrow reduce **Distance** between source and target features
- Cost matrix \rightarrow Map target with **same** predicted **label** as source
- **Regularization** between distance and pseudo-labeling loss

¹*Damodaran et. al., 2018*

DeepJDOT: Joint Distribution Optimal Transport¹

Two steps optimization:

1. Compute the **optimal transport plan** γ between the source and target **batches**

$$\gamma = \underset{\gamma}{\operatorname{argmin}} \langle \gamma, \mathbf{C}_{\mathbf{h},\mathbf{f}} \rangle_F ,$$

Can be done using **POT**² library.

2. Update the **feature extractor** and the **classifier** by minimizing the **total loss**

$$(\hat{h}, \hat{f}) = \underset{h, f}{\operatorname{argmin}} \mathcal{L}_{\text{tot}}(h, f)$$

¹*Damodaran et. al., 2018*

²*Flamary et. al., 2017*

In practice: How to choose the best regularization?

No labels in the target domain \rightarrow **No** way to **tune** the **hyperparameters** .

- DeepCoral \rightarrow 1 hyperparameter
- DANN \rightarrow 1 hyperparameter
- DeepJDOT \rightarrow 2 hyperparameters

Solutions in paper:

- DeepCoral $\rightarrow \lambda$ fixed: $\mathcal{L} \sim \mathcal{L}_{DA}$ at the end of the training
- DANN \rightarrow Reversed cross-validation¹
- DeepJDOT \rightarrow "Fixed experimentally" ...

¹*Zhong et. al., 2010*

DA scorer: Metric to set the best hyperparameters

- **Import Weighted**¹ → Score as a **reweighted** accuracy on labeled sources data
- **Depp Embedded Validation**² → IW strategy in the **latent space** with variance reduction strategy
- **Prediction Entropy**³ → Reduce **Entropy** of the **prediction** on the target domain to reduce uncertainty
- **Circular Validation**⁴ → Adapt Source to Target then Target to Source with predicted labels

¹*Sugiyama et. al., 2007*

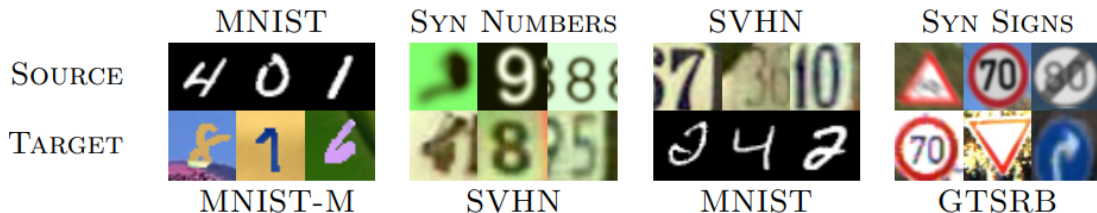
²*You et. al., 2019*

³*Morerio et. al., 2017*

⁴*Bruzzone et. al., 2010*

Experimental results: Digits dataset

Digits dataset



Source: Ganin et. al., 2016

- Classification of **10** classes over **5** domains
- Shift between the domains: **Font** , **Color** , **Style**

Experimental results: Digits dataset

Method	Adaptation:source→target			
	MNIST → USPS	USPS → MNIST	SVHN → MNIST	MNIST → MNIST-M
Source only	94.8	59.6	60.7	60.8
DeepCORAL [6]	89.33	91.5	59.6	66.5
MMD [14]	88.5	73.5	64.8	72.5
DANN [8]	95.7	90.0	70.8	75.4
ADDA [21]	92.4	93.8	76.0 ⁵	78.8
AssocDA [16]	-	-	95.7	89.5
Self-ensemble ⁴ [42]	88.14	92.35	93.33	-
DRCN [40]	91.8	73.6	81.9	-
DSN [41]	91.3	-	82.7	83.2
CoGAN [9]	91.2	89.1	-	-
UNIT [18]	95.9	93.5	90.5	-
GenToAdapt [19]	95.3	90.8	92.4	-
I2I Adapt [20]	92.1	87.2	80.3	-
StochJDOT	93.6	90.5	67.6	66.7
DeepJDOT (ours)	95.7	96.4	96.7	92.4
target only	95.8	98.7	98.7	96.8

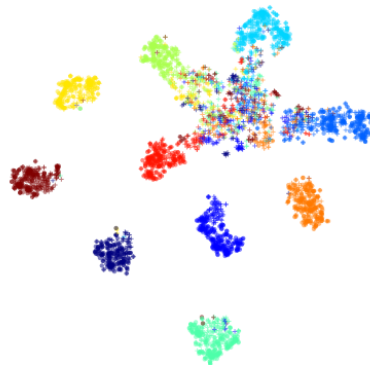
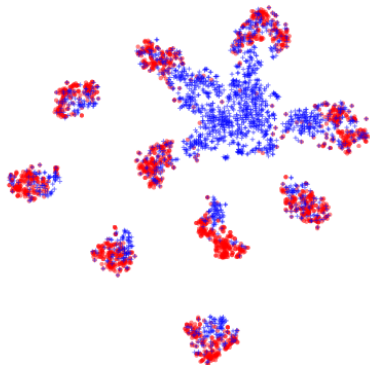
Source: Damodaran et. al., 2018

Experimental results: TSNE visualization with Source Only

Source (red) VS target (blue)

Class discrimination

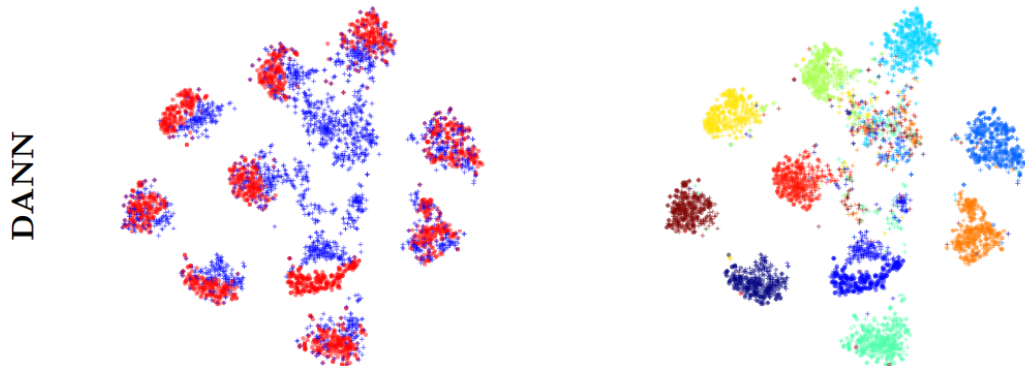
Source Only



■ **Target** domain samples are **not clustered**

Source: Damodaran et. al., 2018

Experimental results: TSNE visualization with DANN

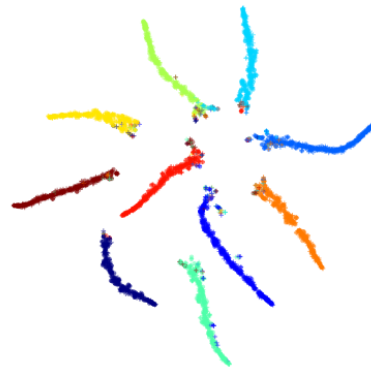
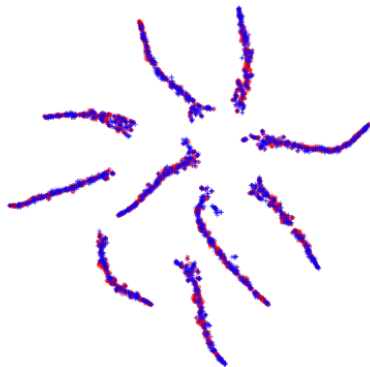


- **Target** domain samples are **more clustered** but few are misclassified

Source: Damodaran et. al., 2018

Experimental results: TSNE visualization with DeepJDOT

DeepJDOT



■ **Target** domain samples are **perfectly clustered** !

Source: Damodaran et. al., 2018

How to use DA?

Skada¹ is a **Python** library to **easily** use DA methods.

- **Homogeneous API** for all DA methods (Shallow and Deep learning).
- **Sklearn-like API** with estimator class (`.fit`, `.predict`, ...), pipeline, grid search ...
- **DA scorer** to validate hyper-parameters without using target label.



¹*Gnassounou et. al., 2024*

Data format in Skada

- $X \rightarrow$ 2D array of shape $(n_samples, n_features)$
- $y \rightarrow$ 1D array of shape $(n_samples,)$
- $sample_domain \rightarrow$ 1D array of shape $(n_samples,)$ giving the **domain** of each **sample**

```
1 from skada.datasets import make_shifted_datasets
2
3 X, y, sample_domain = make_shifted_datasets(
4     20, 20, shift='covariate_shift', random_state=42
5 )
```

- All shift are available in `make_shifted_datasets` function

Shallow DA in Skada

- Initialize the estimator
- Fit the model
- Don't forget to give the **sample domain**

```
1 from skada import LinOT
2
3 estimator = LinOT()
4 estimator.fit(X, y, sample_domain=sample_domain)
```

- ~ 20 shallow methods available in Skada

Pipeline DA in Skada

- Can be used with **Pipeline**

```
1 from skada import make_da_pipeline
2 from skada import LinOTAdapter, GaussianReweightAdapter
3 from sklearn.linear_model import LogisticRegression
4
5 pipeline = Pipeline(
6     LinOTAdapter(),
7     LogisticRegression()
8 )
9 pipeline.fit(X, y, sample_domain=sample_domain)
```

- Possibility to mixed DA adapters

```
1 pipeline = Pipeline(
2     LinOTAdapter(),
3     GaussianReweightAdapter(),
4     LogisticRegression()
5 )
```

DA scorer in Skada

- Possibility to use `cross_val_score` with **DA scorers**
- DA scorers are used to **validate** the **hyperparameters** without using the target labels

```
1 from skada.scorers import ImportanceWeightedScorer
2
3 scorer = ImportanceWeightedScorer()
4 score = cross_val_score(pipeline, X, y, sample_domain=sample_domain,
    ↪ scoring=scorer)
```

- 6 DA scorers available in Skada

Deep DA method in Skada

- Use **Skorch** → **Pytorch** wrapper for **Sklearn**
- Give an **architecture** and **hyperparameters**

```
1 from skada.deep import DeepCoral
2 from skada.deep.modules import ToyCNN
3
4 model = DeepCoral(
5     ToyCNN(),
6     batch_size=32,
7     max_epochs=5,
8     lr=1e-3,
9     reg=1,
10    layer_name="feature_extractor",
11)
12 model.fit(X, y, sample_domain=sample_domain)
```

- ~ 10 Deep DA methods available in Skada

Conclusion

- Distribution shift is a **challenging** problem
- Deep learning methods reduce the shift in the **feature space**
- Modern DA methods are more focus on **Test-Time DA**
- Try **Skada** to easily use DA methods
- Don't hesitate to contribute to the library!

