A framework to train on labeled (simulated) and unlabeled (real) data

> Marcus Engsig - bgr911 Gor Nahapetyan - xzn779 Frederik Hansen - bcp695 Johan Sieborg - hnw224

### **Table of Contents**

- 1. Problem Statement
- 2. Dataset Structure
- 3. Domain Adversarial Neural Networks (DANN)
  - a. Architecture
  - b. Loss Function
  - c. Gradient Reversal Layer
- 4. Training Procedure
- 5. Results
  - a. MLP (Null-hypothesis)
  - b. DANN
  - c. Multi-DANN
- 6. Conclusion & Discussion

Frequently, **simulated data** is **easy** to obtain.



Frequently, **simulated data** is **easy** to obtain.

Conversely, **real data** is generally **difficult** and **expensive** to gather.





Frequently, **simulated data** is **easy** to obtain.

Conversely, **real data** is generally **difficult** and **expensive** to gather.

On top of that, **real data** tends to be **difficult** and sometimes impossible to **label**.





Frequently, **simulated data** is **easy** to obtain.

Conversely, **real data** is generally **difficult** and **expensive** to gather.

On top of that, **real data** tends to be **difficult** and sometimes impossible to **label**.

How can we **train** on **simulated data**, and **generalize** our model to **real data** without <u>overfitting</u> (to the simulated data).



Frequently, **simulated data** is **easy** to obtain.

Conversely, **real data** is generally **difficult** and **expensive** to gather.

On top of that, **real data** tends to be **difficult** and sometimes impossible to **label**.

How can we **train** on **simulated data**, and **generalize** our model to **real data** without <u>overfitting</u> (to the simulated data).





Our dataset consists of:

- 1. a set of labelled 'simulated' data.
- 2. a set of **unlabelled** 'real' data.

Our dataset consists of:

- 1. a set of **labelled** 'simulated' data.
- 2. a set of **unlabelled** 'real' data.



Our dataset consists of:

- 1. a set of **labelled** 'simulated' data.
- 2. a set of **unlabelled** 'real' data.



Our dataset consists of:

- 1. a set of **labelled** 'simulated' data.
- 2. a set of **unlabelled** 'real' data.

Even though the 'real' data is **unlabelled**, there is still some hidden and **mineable information**.



Our dataset consists of:

- 1. a set of **labelled** 'simulated' data.
- 2. a set of **unlabelled** 'real' data.

Even though the 'real' data is **unlabelled**, there is still some hidden and **mineable information**.

The label of being 'simulated' vs. 'real'.



How can we use the 'real' vs. 'simulated' **label** to **generalize** our model to the 'real' data?

How can we use the 'real' vs. 'simulated' **label** to **generalize** our model to the 'real' data?



Y. Ganin et. al., 2016.

The short answer... Adversarial Neural Nets.

How can we use the 'real' vs. 'simulated' **label** to **generalize** our model to the 'real' data?



The short answer... Adversarial Neural Nets.

1. The <u>feature extractor</u> tries to **mine relevant information** from the real or simulated data that is **indistinguishable**, despite the **origin** of the data (think auto-encoder).

How can we use the 'real' vs. 'simulated' **label** to **generalize** our model to the 'real' data?



#### The short answer... Adversarial Neural Nets.

- 1. The <u>feature extractor</u> tries to **mine relevant information** from the real or simulated data that is **indistinguishable**, despite the **origin** of the data (think auto-encoder).
- 2. The <u>discriminator</u> (domain classifier) tries to **classify** whether the data is **real** or **simulated**.

How can we use the 'real' vs. 'simulated' **label** to **generalize** our model to the 'real' data?



#### The short answer... Adversarial Neural Nets.

- 1. The <u>feature extractor</u> tries to **mine relevant information** from the real or simulated data that is **indistinguishable**, despite the **origin** of the data (think auto-encoder).
- 2. The <u>discriminator</u> (domain classifier) tries to **classify** whether the data is **real** or **simulated**.
- 3. The <u>label predictor</u> attempts to **classify** (main task) based on the **main-objective** labels from **only** the **simulated** data.

### Loss Function

$$E(\theta_f, \theta_y, \theta_d) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_y^i(\theta_f, \theta_y) - \lambda \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}_d^i(\theta_f, \theta_d) + \frac{1}{n'_{i=n+1}} \sum_{i=1}^N \mathcal{L}_d^i(\theta_f, \theta_d)\right)$$







Where the loss function is Binary Cross Entropy (BCE)

$$\mathcal{L}_i = -(y_i \log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i))$$



Where the loss function is Binary Cross Entropy (BCE)

$$\mathcal{L}_i = -(y_i \log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i))$$

How do we implement this?...



#### How do we implement this?...

### Loss Function

$$E(\theta_f, \theta_y, \theta_d) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_y^i(\theta_f, \theta_y) - \lambda \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}_d^i(\theta_f, \theta_d) + \frac{1}{n'_{i=n+1}} \sum_{i=1}^N \mathcal{L}_d^i(\theta_f, \theta_d)\right)$$

How do we implement this?...







1. To **reverse** the **gradient** from the <u>discriminator</u> (the minus sign).



- 1. To **reverse** the **gradient** from the <u>discriminator</u> (the minus sign).
- 2. To scale the gradient by a factor  $\lambda$



- 1. To **reverse** the **gradient** from the <u>domain classifier</u> (the minus sign).
- 2. To scale the gradient by a factor  $\lambda$
- 3. Train the <u>Feature Extractor</u> to 'fool' the domain classifier  $\rightarrow$  reduce simulated and real data into an **indistinguishable latent space**.



Y. Ganin et. al., 2016.

### An important anecdote

Y. Ganin et. al., showed that DANNs perform **best** where the is a **domain** shift between the 'real' and simulated data.

#### Domain-Adversarial Training of Neural Networks

Yaroslav Ganin Evgeniva Ustinova Skolkovo Institute of Science and Technology (Skoltech) Skolkovo, Moscow Region, Russia

GANIN@SKOLTECH.RU EVGENIYA.USTINOVA@SKOLTECH.RU

Hana Ajakan HANA.AJAKAN.1@ULAVAL.CA Pascal Germain PASCAL GERMAIN@IFT.ULAVAL CA Département d'informatique et de génie logiciel, Université Laval Québec, Canada, G1V 0A6

Hugo Larochelle Département d'informatique, Université de Sherbrooke Québec, Canada, J1K 2R1

François Laviolette

Québec, Canada, G1V 0A6

Mario Marchand

FRANCOIS.LAVIOLETTE@IFT.ULAVAL.CA MARIO, MARCHAND@IFT, ULAVAL, CA Département d'informatique et de génie logiciel, Université Laval

HUGO.LAROCHELLE@USHERBROOKE.CA

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

LEMPITSKY@SKOLTECH.RU

Editor: Urun Dogan, Marius Kloft, Francesco Orabona, and Tatiana Tommasi

#### Abstract

We introduce a new representation learning approach for domain adaptation, in which data at training and test time come from similar but different distributions. Our approach 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 cannot discriminate between the training (source) and test (target) domains.

The approach implements this idea in the context of neural network architectures that are trained on labeled data from the source domain and unlabeled data from the target domain (no labeled target-domain data is necessary). As the training progresses, the approach promotes the emergence of features that are (i) discriminative for the main learning task on the source domain and (ii) indiscriminate with respect to the shift between the domains. We show that this adaptation behaviour can be achieved in almost any feed-forward model by augmenting it with few standard layers and a new gradient reversal layer. The resulting augmented architecture can be trained using standard backpropagation and stochastic gradient descent, and can thus be implemented with little effort using any of the deep learning packages.

### An important anecdote

Y. Ganin et. al., showed that DANNs perform **best** where the is a **domain shift** between the 'real' and simulated data.



MLP

Adversarial



### An important anecdote

Y. Ganin et. al., showed that DANNs perform **best** where the is a **domain shift** between the 'real' and simulated data.

The data must be **fundamentally distinguishable**.



MLP

Random noise augmentation:



Adversarial

### **Our Data Augmentation**

We have Monte-Carlo ALEPH data.

- 1. 6 parameter input
- 2. 1 parameter classification output

### **Our Data Augmentation**

We have Monte-Carlo ALEPH data.

- 1. 6 parameter input
- 2. 1 parameter classification output

We augment our data X with a 'noise' parameter **t** with four functions  $f_i$ .



### **Our Data Augmentation**

We have Monte-Carlo ALEPH data.

- 1. 6 parameter input
- 2. 1 parameter classification output

We augment our data X with a 'noise' parameter **t** with four functions  $f_i$ .


# **Our Data Augmentation**

We have Monte-Carlo ALEPH data.

- 1. 6 parameter input
- 2. 1 parameter classification output

We augment our data X with a 'noise' parameter **t** with four functions  $f_i$ .

**Transformed data** has **no truth values** during training.



- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



We train an MLP: t = 8

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

We train it on 'simulated' data, and use it to predict the 'real' data.



We train an MLP: t = 8

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

We train it on 'simulated' data, and use it to predict the 'real' data.



We train an MLP: t = 8

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

We train it on 'simulated' data, and use it to predict the 'real' data.



We train an MLP: t = 8

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

We train it on 'simulated' data, and use it to predict the 'real' data.

0.7

0.6

 $\cos^{0.5}$ 

0.4

0.3

0

1000

2000

Epoch No.

3000

4000

5000

For small transformations MLPs do great!



0.2

0.0

0.0

0.2

0.4

Ealse Positive Rate

ROC curve (area = 0.92)

0.8

1.0

0.6

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





0.7

0.6

sso 0.5

0.4

0.3

- 1. Architecture: 6x10x10x10x5x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



#### A quick reminder on DANNs



- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.









- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.

- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.





We train an DANN: **t** = 3, **λ** = 150

- 1. Architecture: Main task 6x10x10x10x5x1, Domain Discriminator 10x20x10x1.
- 2. LR: 0.001 with scheduler.
- 3. Validation loss early-stopper.



MLP: Real Accuracy = 0.693

DANN: Real Accuracy = 0.894

Hybrid-DANN: Real Accuracy = 0.899



 DANNs provide a methodology to generalize simulated data to real data.



- DANNs provide a methodology to generalize simulated data to real data.
- Relative increase in performance increases as quality of simulation decreases.



- DANNs provide a methodology to generalize simulated data to real data.
- Relative increase in performance increases as quality of simulation decreases.
- Multiple noise profiles can be generalized via a Hybrid-DANN approach, taking advantage of pairwise learning.



- DANNs provide a methodology to generalize simulated data to real data.
- Relative increase in performance increases as quality of simulation decreases.
- Multiple noise profiles can be generalized via a Hybrid-DANN approach, taking advantage of pairwise learning.



We provided theoretical results with theoretical domain shifts.

- The target Aleph Bjet dataset had no domain shift.
- Indistinguishable from MC data.

Consequently, DANNs and Hybrid-DANNs are ineffective, as the fundamental adversarial cannot discriminate.



# Thank you for your time (:

## Appendix:

Ganin, Yaroslav, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2015. "Domain-Adversarial Training of Neural Networks." http://arxiv.org/abs/1505.07818

## GitHub:

https://github.com/mengsig/APML\_DANN

## Real Data - Aleph data results

• Unfortunately, due to the lack of a domain shift between the Monte Carlo and real data, DANNs were not applicable. In summary we found for t = 3:

ML-Architecture	Naive MLP	DANN	Hybrid-DANN
Real Data Accuracy	0.853	0.732	0.841

- Showcasing that the Domain Adversarial Neural Networks could not discriminate the real and simulated data. Moreover, the Hybrid-DANN only does better than the DANN due to the pairwise learning that exists in this architecture.
- That however, does not mean that DANNs are not applicable to real data! It just means that they are applicable when there is a real domain shift in the data!

# Different DANN backpropagation

- In our initial implementation (see Github \*\*\_novel\_learning.py files), we treated the entire model as a single entity, and did backpropagation in the following manner.
  - Forward pass of simulated data with full gradient checkpointing.
  - Forward pass of real data without gradient checkpointing in the Label Classifier.
  - Gradient reversal layer between Feature Extractor and Domain Classifier.
  - Full back-propagation based on the sum off

total\_loss = label\_loss + lambda \* domain\_loss

• This resulted in similar results to the real DANN, but for different reasons, as the Domain Classifier had substantially higher accuracy. Occasionally, this model outperformed the normal DANN.

• Since we worked with MC data, we had labels for finding accuracy.

- Since we worked with MC data, we had labels for finding accuracy.
- Real data have no labels, so we need some other way to measure performance.

- Since we worked with MC data, we had labels for finding accuracy.
- Real data have no labels, so we need some other way to measure performance.
- We use a metric called **consistency**:

Consistency = 
$$1 - \frac{1}{N} \sum_{i=1}^{N} \left| \hat{y}_i - \frac{1}{k} \sum_{j \in k - \text{NN}(x_i)} \hat{y}_j \right|$$

- Since we worked with MC data, we had labels for finding accuracy.
- Real data have no labels, so we need some other way to measure performance.
- We use a metric called **consistency**:

Consistency = 
$$1 - \frac{1}{N} \sum_{i=1}^{N} \left| \hat{y}_i - \frac{1}{k} \sum_{j \in k - \text{NN}(x_i)} \hat{y}_j \right|$$

Prediction of i-th observation

- Since we worked with MC data, we had labels for finding accuracy.
- Real data have no labels, so we need some other way to measure performance.
- We use a metric called **consistency**:

Consistency = 
$$1 - \frac{1}{N} \sum_{i=1}^{N} \left| \hat{y}_i - \frac{1}{k} \sum_{j \in k - \text{NN}(x_i)} \hat{y}_j \right|$$
  
Predictions of its k-nearest neighbors
## Consistency metric - sinusoidal transformation



## **Consistency metric - Gaussian transformation**

