

Machine Learning for Time Series

(MLTS or MLTS-Deluxe Lectures)

Dr. Dario Zanca

Machine Learning and Data Analytics (MaD) Lab
Friedrich-Alexander-Universität Erlangen-Nürnberg
01.02.2024

- Time series fundamentals and definitions (2 lectures)
- Bayesian Inference (1 lecture)
- Gaussian processes (2 lectures)
- State space models (2 lectures)
- Autoregressive models (1 lecture)
- Data mining on time series (1 lecture)
- Deep learning on time series (4 lectures)
- Domain adaptation (1 lecture) ←

In this lecture...

- 1. Domain adaptation: overview**
- 2. Unsupervised domain adaptation**
- 3. Domain generalization (OOD generalization)**

References

Deep Learning Foundations by Soheil Feizi (University of Maryland):

<https://www.youtube.com/watch?v=E1760ZzsXN8>

https://www.youtube.com/watch?v=wwgt_ErD3vA



Domain adaptation

Domain adaptation: overview



The typical machine learning setup (so far)

The typical setup we have had so far included a training set

$$\{(x_i^{train}, y_i^{train})\}_{i=1}^m \sim Q_{X,Y}$$

training set.
→ distribution

Where $x_i \in X$, $y_i \in Y$, and where $Q_{X,Y}$ denotes the distribution from the training examples are sampled from.

Again, typically we want to learn an optimal mapping f_θ , for which we solve:

Learn
optimal
mapping
 f_θ .

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m L(f_\theta(x_i^{train}), y_i^{train}) \Rightarrow \theta^*$$

Loss is minimized.

We, then, evaluate our model on a hold out test set

$$\{(x_i^{test}, y_i^{test})\}_{i=1}^{m'} \sim Q_{X,Y}$$

by computing a test error

$$\epsilon_{test} = \frac{1}{m} \sum_{i=1}^{m'} L(f_{\theta^*}(x_i^{test}), y_i^{test})$$

(and we aim at a small ϵ_{test}).

small test error.

The typical machine learning setup (so far)

Summary:

1. $\{(x_i^{train}, y_i^{train})\}_{i=1}^m \sim Q_{X,Y}$
2. $\min_{\theta} \frac{1}{m} \sum_{i=1}^m L(f_{\theta}(x_i^{train}), y_i^{train}) \Rightarrow \theta^*$
3. $\{(x_i^{test}, y_i^{test})\}_{i=1}^{m'} \sim Q_{X,Y}$
4. $\epsilon_{test} = \frac{1}{m'} \sum_{i=1}^{m'} L(f_{\theta^*}(x_i^{test}), y_i^{test})$

$t_{train} \in t_{test}$ from
Same distribution.

Key assumption is that both the training and test set come from the same distribution.

Is it a realistic assumption?

Realistic?

In practice, the training distribution and the test distribution are often not the same.

→ We train an image classifier on a database of photos taken with a professional camera, and want our classifier to work on pictures taken with any smartphone camera.

→ Training distribution ≠ Test distribution

→ $Q_{X,Y} \neq P_{X,Y}$

Professional camera;
Smartphone camera.

Domain adaptation definitions

Source domain, target domain, and domain shift

$Q_{X,Y}$ (Source)

$(P_{X,Y} \text{ Target})$

We introduce some terminology from the Domain Adaptation domain:

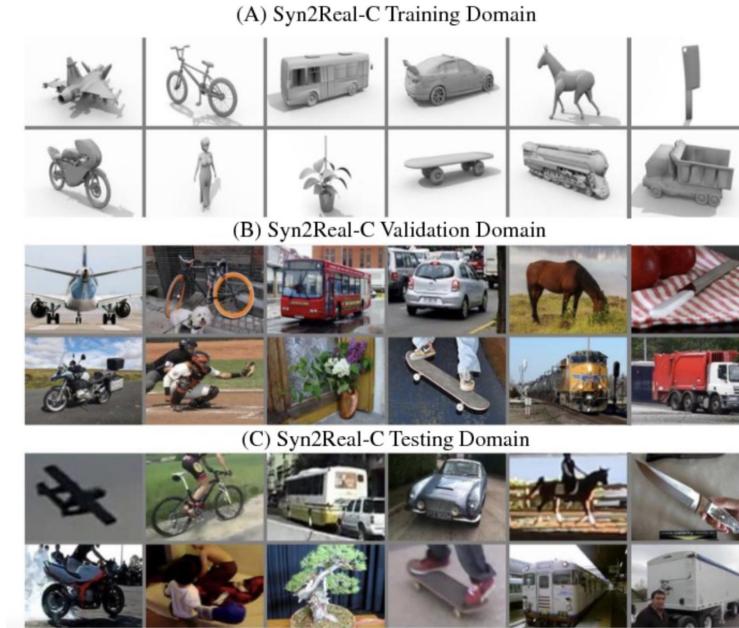
- **Source domain $Q_{X,Y}$.** The data distribution on which the model is trained using labeled examples.
→ photos taken with a professional camera.
different yet related.
- **Target domain $P_{X,Y}$.** A different, yet “related” distribution on which it is required to perform a similar task.
→ photos taken with a smartphone.
similar task $P_{X,Y}$
- **Domain shift.** It is the statistical difference between different domains.
→ statistical difference between $Q_{X,Y}$ and $P_{X,Y}$.
Domain shift
⇒ *statistical differences*.

We introduce some terminology from the Domain Adaptation domain:

- **Source domain $Q_{X,Y}$.** The data distribution on which the model is trained using labeled examples.
→ photos taken with a professional camera.
- **Target domain $P_{X,Y}$.** A different, yet “related” distribution on which it is required to perform a similar task.
→ photos taken with a smartphone.
- **Domain shift.** It is the statistical difference between different domains.
→ statistical difference between $Q_{X,Y}$ and $P_{X,Y}$.

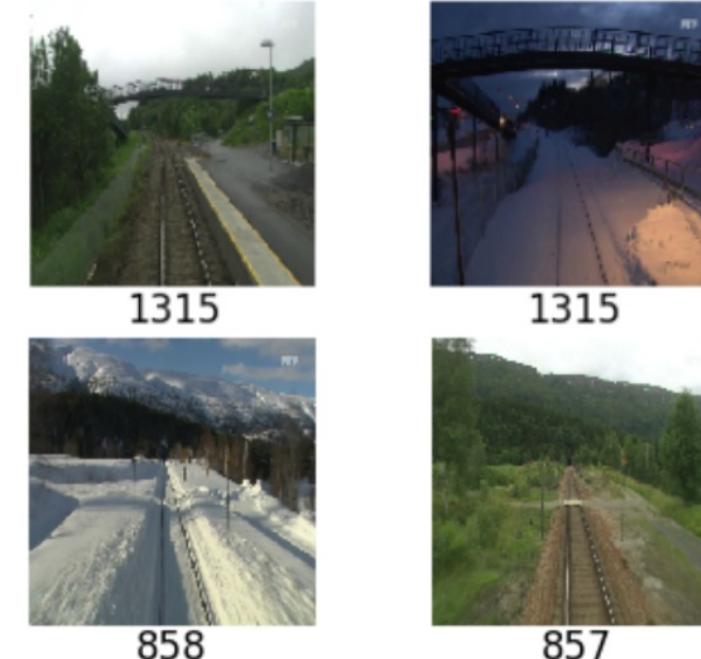
Domain adaptation examples

Same view different seasons.



Train on synthetic samples and use with
real samples.

Image from : Peng X. et al., "Syn2Real: A New Benchmark for Synthetic-to-Real Visual Domain Adaptation"



Same view from different seasons

Image from : Olid D. et al., "Single-View Place Recognition under Seasonal Changes"

Unsupervised domain adaptation

- Labeled samples for the source domain
- $Q_{X,Y} \sim \{(x_i^S, y_i^S)\}_{i=1}^{m_S} := (X^S, Y^S)$
- Only unlabeled samples available for the target domain

$$P_X \sim \{x_i^T\}_{i=1}^{m_T} := X^T$$

only unlabeled for target domain

Semi-supervised

Semi-supervised domain adaptation

- Labeled samples for the source domain
- $Q_{X,Y} \sim \{(x_i^S, y_i^S)\}_{i=1}^{m_S} := (X^S, Y^S)$
- Unlabeled target samples + “Few” labeled target samples

Unlabeled +
“Few” labeled target.

(multiple source domains)

Domain generalization

- Labeled samples for the multiple source domains
 $Q_{X,Y}^1 \sim \{(x_i^{S_1}, y_i^{S_1})\}_{i=1}^{m_{S_1}} := (X^{S_1}, Y^{S_1})$
 $Q_{X,Y}^2 \sim \dots$
 \dots
- **No** samples from the target domain available during training
- This problem is also called “out-of-distribution generalization”

out of distribution generalization

Notice:

- We use both the samples from the source domain and from the target domain during training
- The target domain is different than what we use to call test set
- We need labelled samples from the target domain for testing, in all three scenarios

Need labelled sample from target
domain for
testing in all
scenarios.



Domain adaptation

Unsupervised domain adaptation



Let's assume, for simplicity and without loss of generalization, that $m_S = m_T = m$, i.e.,

- Source domain: $(X^S, Y^S) = \{(x_i^S, y_i^S)\}_{i=1}^m \sim Q_{X,Y}$
- Target domain: $X^T = \{x_i^T\}_{i=1}^m \sim P_X$

The goal in unsupervised domain adaptation is that of, given a hypothesis class H , to pick a function $h \in H$ such that

$$\epsilon_T(h) = \mathbb{E}[L(h(x), y)]$$

with $(x, y) \sim P_{X,Y}$

*Unsupervised domain
adaptation*

Let's assume, for simplicity and without loss of generalization, that $m_S = m_T = m$, i.e.,

- Source domain: $(X^S, Y^S) = \{(x_i^S, y_i^S)\}_{i=1}^m \sim Q_{X,Y}$
- Target domain: $X^T = \{x_i^T\}_{i=1}^m \sim P_X$

Hypothesis class H ;
pick a function $h \in H$

The goal in **unsupervised domain adaptation** is that of, given a hypothesis class H , to pick a function $h \in H$ such that

$$\epsilon_T(h) = \mathbb{E}[L(h(x), y)]$$

with $(x, y) \sim P_{X,Y}$
Unlabelled
sample from target.

Unsupervised domain adaptation

Assumptions

if conditional label distribution
does not change source & target

1. **Covariate shifts.** P and Q satisfy the covariate shift assumption if the conditional label distribution does not change between source and target distribution.

$$\forall x \in X, y \in \{0, 1\} \Rightarrow P(y | x) = Q(y | x)$$

2. **Similarity of distributions.** Source and target (marginal) distribution should be similar.

$$Q_X \dots < = > \dots P_X$$

3. **Small joint error.** If I “had” labeled samples, the joint error should be small.

$$\epsilon_{joint} = \min \left[\frac{1}{m} \sum_{i=1}^m L(h(x_i^S), y^S) + \frac{1}{m} \sum_{i=1}^m L(h(x_i^T), y^T) \right] \approx 0$$

small joint error

H-divergence is defined as:

$$2 \sup_{h \in H} |p_{x \in Q_X}(h(x) = 1) - p_{x \in P_X}(h(x) = 1)| \triangleq d_H(Q_X, P_X)$$

Lemma. The H-divergence $d_H(Q_X, P_X)$ can be estimated by $n_S = m_T = m$ samples from source and target domains, $VC(H) = d$, with probability $1 - \delta$,

$$d_H(Q_X, P_X) \leq d_H(Q_X^{(m)}, P_X^{(m)}) + 4 \sqrt{\frac{d \log(2m) - \log(\frac{2}{5})}{m}}$$

Estimate H-divergence

Methods for estimating the H-divergence

→ classifier → separate
source domain from
target domain

The H-divergence can be computed by finding a classifier to separate source domain from target domain.

- Label all source samples as +1 → all source + 1
- Label all target samples as 0 → all target 0
- Train a classifier to minimize the classification error:

$$\epsilon_{class} = \min_{h \in H} \left[\frac{1}{m} \sum_{i=1}^m 1(h(x_i^S) = 0) + \frac{1}{m} \sum_{i=1}^m 1(h(x_i^T) = 1) \right]$$

The classification loss is inversely proportional to the H-divergence,

$$\frac{1}{2} d_H(Q_X^{(m)}, P_X^{(m)}) = 1 - \epsilon_{class}$$

Classification loss
 $\propto \frac{1}{H\text{-divergence}}$

Symmetric difference hypothesis space

Definition

Definition. For the hypothesis class H , the symmetric difference hypothesis space $H\Delta H$ is the set of disagreements between any two hypothesis in H .

Hypothesis space
 $H\Delta H$

$$H\Delta H = \{g(x) = h(x) \oplus h'(x) | h, h' \in H\}$$

Set of disagreements b/w any two hypothesis in H .

Main result.

The following “main result” has inspired many practical methods in domain adaptation.

Main result. H is a hypothesis class with $\underline{VC}(H) = d$. We are given unlabeled samples from the target $P_X^{(m)}$ and labeled samples from the sources $Q_{X,Y}^{(m)}$. With probability $1 - \delta$, for any $h \in H$,

$$\epsilon_T(h) \leq \epsilon_S(h) + \frac{1}{2} d_{H\Delta H} \left(Q_X^{(m)}, P_X^{(m)} \right) + \epsilon_{joint}$$

(Target error \leq source error + divergence + joint error)

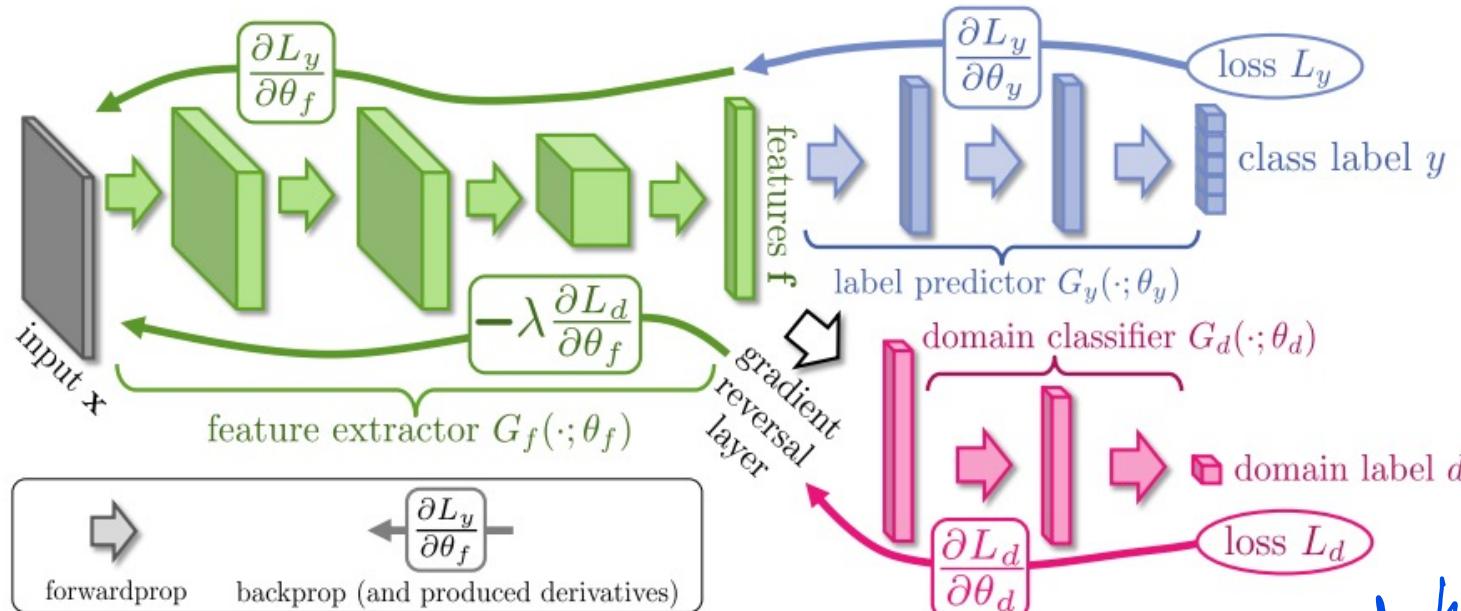
The main result resulted in many practical methods (approximation methods) in order to use the concept of divergence in the training itself.

(approximation methods)

- Classical domain adaptation methods
 - Metric learning
 - Sample re-weighting
 - Subspace alignment
 - ...
- Deep Learning-based methods
 - Nowadays an hot topic of research

Domain adaptation: Ganin & Lempitsky method

Ganin & Lempitsky.



Learn Mapping \rightarrow

task performance \rightarrow maximised
penalised domain classifier.

Image from: Ganin & Lempitsky, "Unsupervised Domain Adaptation by Backpropagation"



Domain adaptation
Domain generalization (OOD generalization)



Domain generalization (also called, Out-of-distribution (OOD) generalization)

out-of distribution.

The problem of domain generalization (also called, out-of-distribution (OOD) generalization) can be formalized as follows:

➤ Training: $K = |E|$ training domains

➤ $P^{(e)} \sim \{(x_i^e, y_i^e)\}_{i=1}^{m_e}$

➤ $1 \leq e \leq |E|$

➤ Goal: find $h \in H$ that performs well in an unseen domain $|E| + 1$

➤ $P^{(K+1)} \sim \{(x_i^{(K+1)}, y_i^{(K+1)})\}_{i=1}^{m_{(K+1)}}$

➤ Minimize the risk in the new environment

➤ $R^{(K+1)}(h) = \mathbb{E}_{(x,y) \sim P^{(K+1)}} [L(h(x), y)]$

Risk

"related"

find $h \in H$; perform well in unseen domain $|E| + 1$

minimize risk in new environment.

Note: also in this setup, different environments need to be "related" to each other.



- <http://ai.bu.edu/M3SDA/>
- 345 classes
- Domains: clipart, real, sketch, infograph, paintings, drawings
- <https://paperswithcode.com/dataset/pacs>
- 7 categories
- Domains: photo, paintings, cartoon, sketch

Baseline method

Method 1: Baseline method.

We call “Baseline” method the approach that consists simply on minimizing the error on the available domains.

- **Training:** $\min_f \frac{1}{K} \sum_{j=1}^k \mathbb{E}_{(x,y) \sim P^{(j)}} [L(f(x), y)]$
- **Test:** $\mathbb{E}_{(x,y) \sim P^{(K+1)}} [L(f(x), y)]$
- “Do nothing” method

(Do nothing method)

minimize error on domains.

Method 2: Invariant representation.

Learn a representation that is invariant across different domains

- Use domain adversarial neural networks (DANN)

Learn representations that are invariant.
Domain adversarial NN

- ϕ (feature extraction)

- $\omega \circ \phi$ (label classification)

- $c \circ \phi$ (domain classification)

$$loss = \frac{1}{K} \sum_{j=1}^K L(\omega \circ \phi(x), y) - \lambda \frac{1}{K} \sum_{j=1}^K L(c \circ \phi(x), y)$$

Minimize f.e loss & label classifier loss

- $\min_{\phi, \omega} loss$ & $\max_c loss$

Maximize domain classification loss.



Lecture title

Recap



- **Domain adaptation**
 - **Unsupervised domain adaptation**
 - Main result
 - Practical methods
 - **Semi-supervised domain adaptation**
 - **Domain generalization**
 - Baseline method
 - Invariant representations method

