

Learning Under Different Training and Testing Distributions

Dr Haider Raza

Postdoctoral Research Fellow
Institute for Analytics and Data Science, University of Essex,
Colchester, Essex, UK.
h.raza@essex.ac.uk



: @sagihaider



: sagihaider

“The Big Data & Analytics Summer School 2018”

Hashtag: [#IADSSummerSchool](#)

23th July 2018

Tweet us!

If you like, please Tweet



Use Hashtag **#IADSSummerSchool**

&

Use **@EssexIADS** to find, follow, and tag.



GitHub: Lab work and presentation



Follows the steps:

- ▶ Go to link
https://github.com/sagihaidar/IADS_SC_2018
- ▶ Look at right hand side in green color: **Clone or download**.
Click it and Download Zip
- ▶ When downloading is finished. Copy the Zip file and take to the location you want such any folder and paste it. Extract it.
- ▶ If you have **Anaconda3** installed. Go to terminal or command prompt and type “jupyter notebook”

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ $(X \rightarrow Y)$: Y is determined by values of X (e.g. credit card fraud detection) **Predictive models** (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ $(Y \rightarrow X)$: Y determines the values of X (e.g. medical diagnosis) **Generative models** (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ $(X \rightarrow Y)$: Y is determined by values of X (e.g. credit card fraud detection) Predictive models (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ $(Y \rightarrow X)$: Y determines the values of X (e.g. medical diagnosis) Generative models (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ $(X \rightarrow Y)$: Y is determined by values of X (e.g. credit card fraud detection) Predictive models (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ $(Y \rightarrow X)$: Y determines the values of X (e.g. medical diagnosis) Generative models (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ $(X \rightarrow Y)$: Y is determined by values of X (e.g. credit card fraud detection) **Predictive models** (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ $(Y \rightarrow X)$: Y determines the values of X (e.g. medical diagnosis) **Generative models** (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ ($X \rightarrow Y$): Y is determined by values of X (e.g. credit card fraud detection) **Predictive models** (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ ($Y \rightarrow X$): Y determines the values of X (e.g. medical diagnosis) **Generative models** (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ ($X \rightarrow Y$): Y is determined by values of X (e.g. credit card fraud detection) **Predictive models** (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ ($Y \rightarrow X$): Y determines the values of X (e.g. medical diagnosis) **Generative models** (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Notation

- ▶ A set of features or covariates X .
- ▶ A set of target or class variables Y .
- ▶ A joint distribution $P(Y,X)$ or $P(Y \cap X)$ (i.e. Probability of Y and X).
- ▶ ($X \rightarrow Y$): Y is determined by values of X (e.g. credit card fraud detection) **Predictive models** (e.g. Logistic Regression, SVM, and Neural Networks.)
- ▶ ($Y \rightarrow X$): Y determines the values of X (e.g. medical diagnosis) **Generative models** (e.g. GMM, HMM, and Naive Bayes).
- ▶ The joint distribution $P(Y,X)$ can be written as
 1. $P(Y|X)P(X)$ in $X \rightarrow Y$ problems.
 2. $P(X|Y)P(Y)$ in $Y \rightarrow X$ problems.
- ▶ P_{tr} : Data distribution in training
- ▶ P_{ts} : Data distribution in testing

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift

Learning Under Different Training and Testing Distributions

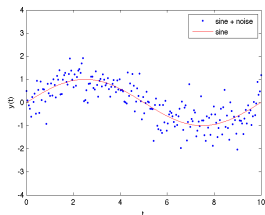
-
- Figure 1 consists of two side-by-side plots. The left plot is a histogram with a blue outline, showing the density function of a mixture of two distributions. The x-axis ranges from -5 to 10, and the y-axis (Density function) ranges from 0.00 to 0.15. The histogram has five bars: one at x ≈ -3.5 with height ≈ 0.08, one at x ≈ -1.5 with height ≈ 0.16, one at x ≈ 0.5 with height ≈ 0.08, one at x ≈ 4.5 with height ≈ 0.08, and one at x ≈ 6.5 with height ≈ 0.08. The right plot shows the same density function as a solid blue curve, which is bimodal with peaks at x ≈ -1.5 and x ≈ 6.5. Overlaid on this are three dashed red curves, each representing one of the three components of the mixture. These curves are unimodal and centered at x ≈ -3.5, x ≈ 0.5, and x ≈ 4.5, with heights of approximately 0.04, 0.04, and 0.08 respectively.

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Issues with Data in Machine Learning

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Issues with Data in Machine Learning

Dr Haider Raza

- ▶ Imbalanced dataset
- ▶ Overlapping dataset
- ▶ Density: Lack of data
- ▶ Noise in data
- ▶ Dataset Shift



- ▶ In learning theory **independent and identically distributed (i.i.d)** assumption (i.e. each random variable has the same probability distribution as the others and all are mutually independent).
- ▶ In practice **train** and **test** inputs have different distributions.
- ▶ The difference in distribution arises from operating in **non-stationary environments** in real-world application such as **finance**, **healthcare**, **brain signals**, much more...
- ▶ Learning in such non-stationary environment is difficult and we need an think before operating.

- ▶ In learning theory **independent and identically distributed (i.i.d)** assumption (i.e. each random variable has the same probability distribution as the others and all are mutually independent).
- ▶ In practice **train** and **test** inputs have different distributions.
- ▶ The difference in distribution arises from operating in **non-stationary environments** in real-world application such as **finance**, **healthcare**, **brain signals**, much more...
- ▶ Learning in such non-stationary environment is difficult and we need an think before operating.

- ▶ In learning theory **independent and identically distributed (i.i.d)** assumption (i.e. each random variable has the same probability distribution as the others and all are mutually independent).
- ▶ In practice **train** and **test** inputs have different distributions.
- ▶ The difference in distribution arises from operating in **non-stationary environments** in real-world application such as **finance**, **healthcare**, **brain signals**, much more...
- ▶ Learning in such non-stationary environment is difficult and we need an think before operating.

- ▶ In learning theory **independent and identically distributed (i.i.d)** assumption (i.e. each random variable has the same probability distribution as the others and all are mutually independent).
- ▶ In practice **train** and **test** inputs have different distributions.
- ▶ The difference in distribution arises from operating in **non-stationary environments** in real-world application such as **finance**, **healthcare**, **brain signals**, much more...
- ▶ Learning in such non-stationary environment is difficult and we need an think before operating.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Shifts in Data

Outline

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹

1. “concept shift/drift” G. Widmer et al., 1996, 1998
2. “changes of classification” K. Wang et al., 2003
3. “changing environments” R. Alaiz-Rodriguez et al., 2008
4. “fracture point” N.V. Chawla et al., 2009
5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2006

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹

1. “concept shift/drift” G. Widmer et al., 1996, 1998
2. “changes of classification” K. Wang et al., 2003
3. “changing environments” R. Alaiz-Rodriguez et al., 2008
4. “fracture point” N.V. Chawla et al., 2009
5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2006

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹
 1. “concept shift/drift” G. Widmer et al., 1996, 1998
 2. “changes of classification” K. Wang et al., 2003
 3. “changing environments” R. Alaiz-Rodriguez et al., 2008
 4. “fracture point” N.V. Chawla et al., 2009
 5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2006

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹
 1. “concept shift/drift” G. Widmer et al., 1996, 1998
 2. “changes of classification” K. Wang et al., 2003
 3. “changing environments” R. Alaiz-Rodriguez et al., 2008
 4. “fracture point” N.V. Chawla et al., 2009
 5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2006

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹
 1. “concept shift/drift” G. Widmer et al., 1996, 1998
 2. “changes of classification” K. Wang et al., 2003
 3. “changing environments” R. Alaiz-Rodriguez et al., 2008
 4. “fracture point” N.V. Chawla et al., 2009
 5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2006

- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”¹
 1. “concept shift/drift” G. Widmer et al., 1996, 1998
 2. “changes of classification” K. Wang et al., 2003
 3. “changing environments” R. Alaiz-Rodriguez et al., 2008
 4. “fracture point” N.V. Chawla et al., 2009
 5. “fractures between data” J.G. Moreno-Torres et al., 2010

¹A. Storkey, *Dataset Shift in Machine Learning*, 2009

Dataset Shift: General Example

- ▶ Speech recognition system
- ▶ Training the speech recognition system
- ▶ Voice recognition systems fails

Dataset Shift: General Example

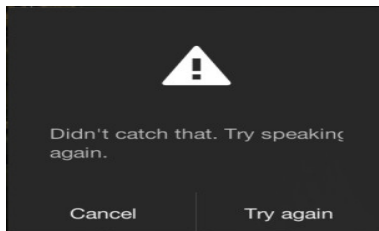
- ▶ Speech recognition system
- ▶ Training the speech recognition system
- ▶ Voice recognition systems fails

Dataset Shift: General Example

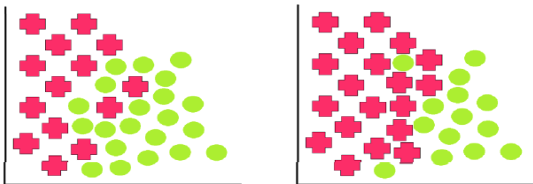
- ▶ Speech recognition system
- ▶ Training the speech recognition system
- ▶ Voice recognition systems fails

Dataset Shift: General Example

- ▶ Speech recognition system
- ▶ Training the speech recognition system
- ▶ Voice recognition systems fails

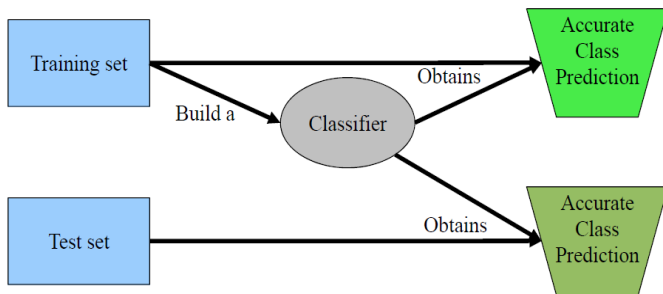


- $$P_{tr}(X, Y) \neq P_{ts}(X, Y)$$



Dataset Shift...cont

- Basic assumption for classification in operating under stationary environment



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

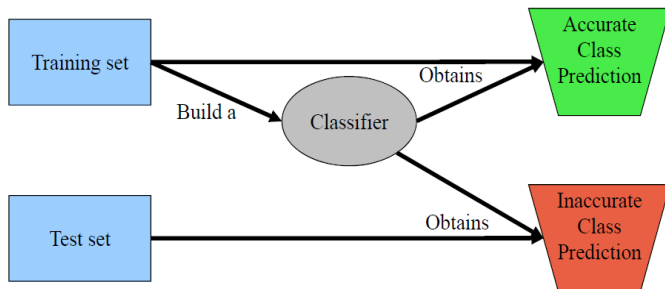
Learning in Dataset Shift

Summary

Dataset Shift...cont

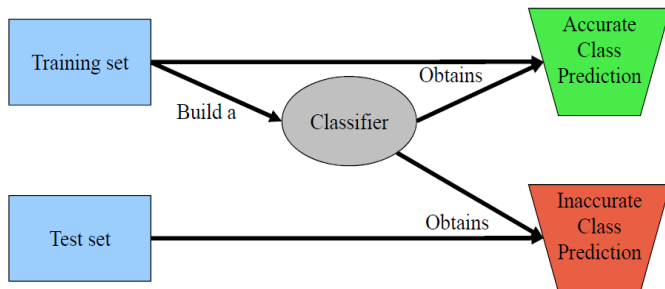
Learning Under Different Training and Testing Distributions

► But sometimes...



Dataset Shift...cont

- ▶ But sometimes...



- ▶ The classifier has overfitting problem

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between training and **test** sets: **then** possible actions ²
 - ▶ Train a new classifier for the **test** set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the **test** set

²Alippi et al., IEEE TNNLS, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between training and test sets: **then** possible actions ²
 - ▶ Train a new classifier for the test set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the test set

²Alippi et al., IEEE TNNLS, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between training and test sets: **then** possible actions ²
 - ▶ Train a new classifier for the test set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the test set

²Alippi et al., IEEE TNNLS, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between **training** and **test** sets: **then** possible actions ²
 - ▶ Train a new classifier for the **test** set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the **test** set

²Alippi et al., IEEE TNNLS, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between **training** and **test** sets: **then** possible actions ²
 - ▶ Train a new classifier for the **test** set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the **test** set

²Alippi et al., *IEEE TNNLS*, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between **training** and **test** sets: **then** possible actions ²
 - ▶ Train a new classifier for the **test** set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the **test** set

²Alippi et al., *IEEE TNNLS*, 2008.

The problem of Dataset Shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ **If** the classifier has an overfitting problem: **then** possible actions
 - ▶ Change the parameters of the algorithm
 - ▶ Use a more general learning method
- ▶ **If** there is a change in the data distribution between **training** and **test** sets: **then** possible actions ²
 - ▶ Train a new classifier for the **test** set
 - ▶ Adapt to classifier
 - ▶ Modify the data in the **test** set

²Alippi et al., *IEEE TNNLS*, 2008.

Dataset Shift: A literature review

- ▶ As an example, the following terms have been used in the literature to refer to Dataset Shift: “Concept shift”, “Changes of classification”, “Changing environments”, “Contract mining in classification learning”, “Fracture points”, and “Fractures between data”
- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”³

³Moreno-Torres et al., *Pattern Recognition*, 2011.

Dataset Shift: A literature review

- ▶ As an example, the following terms have been used in the literature to refer to Dataset Shift: “Concept shift”, “Changes of classification”, “Changing environments”, “Contract mining in classification learning”, “Fracture points”, and “Fractures between data”
- ▶ “cases where the joint distribution of inputs and outputs differs between training and test stage”³

³Moreno-Torres et al., *Pattern Recognition*, 2011.

Outline

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Types of dataset shift

1. Covariate shift
2. Prior probability shift
3. Concept shift Concept

Types of dataset shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

1. Covariate shift
2. Prior probability shift
3. Concept shift Concept

Types of dataset shift

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

1. Covariate shift
2. Prior probability shift
3. Concept shift Concept

Covariate Shift

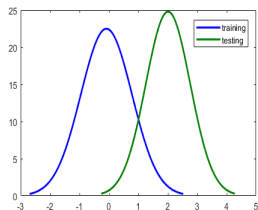
- Covariate shift appears only in $X \rightarrow Y$ problems⁴, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

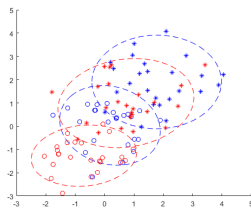
&

$$P_{tr}(X) \neq P_{ts}(X)$$

Uni-variate



Bi-variate



⁴Raza et al., Pattern Recognition, 2015.

- ▶ The term covariate shift was first defined 18 years ago by (Shimodaira, 2000⁵), where it refers to changes in the distribution of the input variables X .
- ▶ Covariate shift is probably the most studied type of shift, but there appears to be some confusion in the literature about the exact definition of the term. There are also some equivalent names, such as “population drift”, “changes in the data distributions”, “differing training and test distributions”.

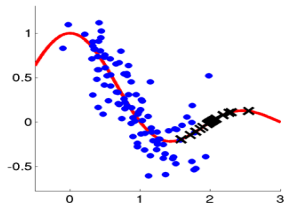
⁵Shimodaira et al., *Journal of Statistical Planning and Inference*, 2000

- ▶ The term covariate shift was first defined 18 years ago by (Shimodaira, 2000⁵), where it refers to changes in the distribution of the input variables X .
- ▶ Covariate shift is probably the most studied type of shift, but there appears to be some confusion in the literature about the exact definition of the term. There are also some equivalent names, such as “population drift”, “changes in the data distributions”, “differing training and test distributions”.

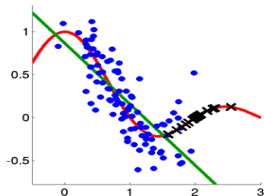
⁵ Shimodaira et al., *Journal of Statistical Planning and Inference*, 2000.

Covariate Shift: An Example

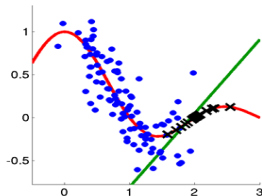
► A regression example ⁶



Training

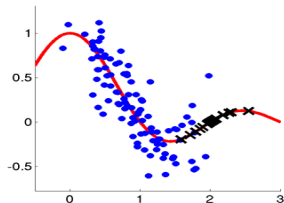


Testing

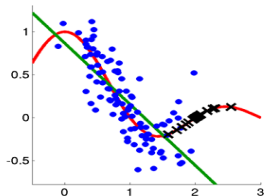


Covariate Shift: An Example

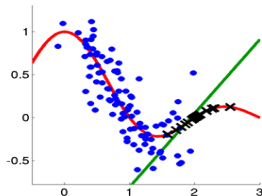
► A regression example ⁶



Training



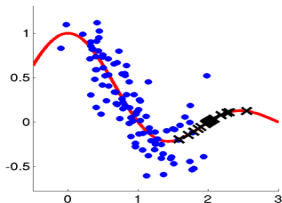
Testing



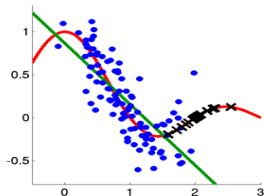
⁶Sugiyama et al., *Journal of Machine Learning Research*, 2007.

Covariate Shift: An Example

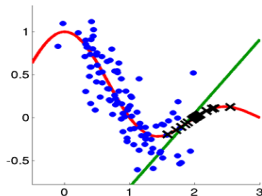
- A regression example ⁶



Training



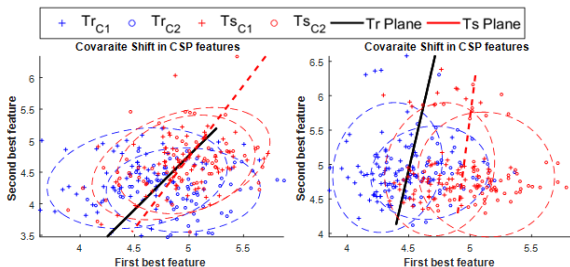
Testing



⁶Sugiyama et al., *Journal of Machine Learning Research*, 2007.

Covariate Shift: An Example

► A Classification example ⁷



- Covariate shift (CS) between the **training** and **test** distributions of the EEG signal from the healthy subject (a) illustrates the CS in the mu band [8-12] Hz and (b) shows the CS in the beta band [14-30] Hz.

Shifts in Data

Dataset Shift

Types of Dataset Shift

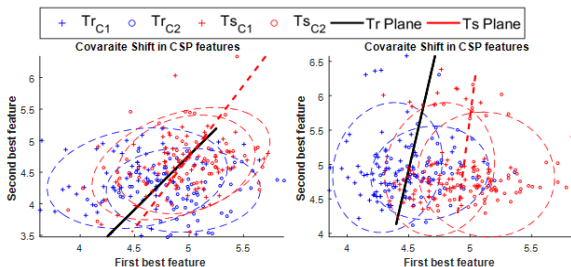
Causes of Dataset Shift

Learning in Dataset Shift

Summary

Covariate Shift: An Example

► A Classification example ⁷



- Covariate shift (CS) between the **training** and **test** distributions of the EEG signal from the healthy subject (a) illustrates the CS in the mu band [8-12] Hz and (b) shows the CS in the beta band [14-30] Hz.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

⁷Raza et al., *Soft Computing*.,2015 and *IEEE-IJCNN*., 2015.

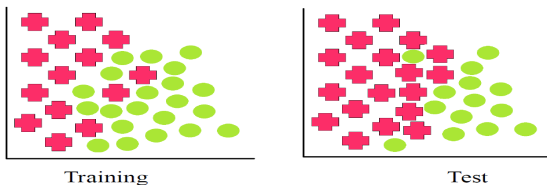
Prior probability shift

- Prior probability shift appears only in $Y \rightarrow X$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(Y) \neq P_{ts}(Y)$$



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

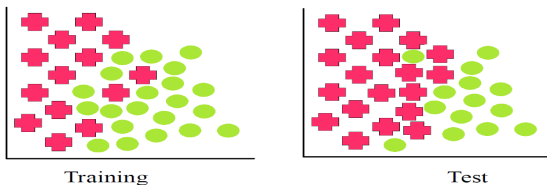
Prior probability shift

- Prior probability shift appears only in $Y \rightarrow X$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(Y) \neq P_{ts}(Y)$$



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

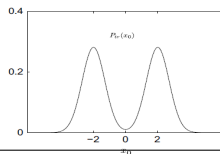
Prior probability shift: Example

Example⁸: $Y \rightarrow X$ problem with one covariate x_0 and a target y that may take the class value $y=0$ and $y=1$. In **training data**, $P_{tr}(y=0) = P_{tr}(y=1) = 0.5$ and $P_{tr}(x_0 | y)$ is defined as

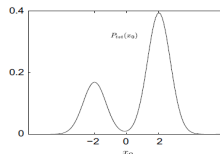
$$x_0 = \begin{cases} \mathcal{N}(2, 0.5), & \text{when } y = 1 \\ \mathcal{N}(-2, 0.5), & \text{otherwise} \end{cases} \quad (1)$$

Now consider that in the **test data** $P_{ts}(x_0 | y=0)$ and $P_{ts}(x_0 | y=1)$ remains unchanged, but the class prior probabilities vary, taking the values $P_{ts}(y=0) = 0.70$ and $P_{ts}(y=1) = 0.30$. This example is illustrated in the figure below

Training



Testing



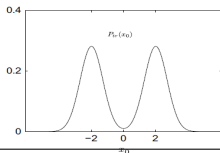
Prior probability shift: Example

Example⁸: $Y \rightarrow X$ problem with one covariate x_0 and a target y that may take the class value $y=0$ and $y=1$. In **training data**, $P_{tr}(y=0) = P_{tr}(y=1) = 0.5$ and $P_{tr}(x_0 | y)$ is defined as

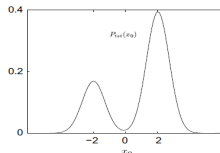
$$x_0 = \begin{cases} \mathcal{N}(2, 0.5), & \text{when } y = 1 \\ \mathcal{N}(-2, 0.5), & \text{otherwise} \end{cases} \quad (1)$$

Now consider that in the **test data** $P_{ts}(x_0 | y=0)$ and $P_{ts}(x_0 | y=1)$ remains unchanged, but the class prior probabilities vary, taking the values $P_{ts}(y=0) = 0.70$ and $P_{ts}(y=1) = 0.30$. This example is illustrated in the figure below

Training



Testing



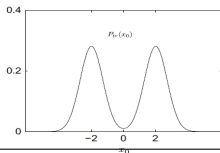
Prior probability shift: Example

Example⁸: $Y \rightarrow X$ problem with one covariate x_0 and a target y that may take the class value $y=0$ and $y=1$. In **training data**, $P_{tr}(y=0) = P_{tr}(y=1) = 0.5$ and $P_{tr}(x_0 | y)$ is defined as

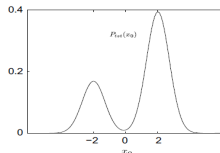
$$x_0 = \begin{cases} \mathcal{N}(2, 0.5), & \text{when } y = 1 \\ \mathcal{N}(-2, 0.5), & \text{otherwise} \end{cases} \quad (1)$$

Now consider that in the **test data** $P_{ts}(x_0 | y=0)$ and $P_{ts}(x_0 | y=1)$ remains unchanged, but the class prior probabilities vary, taking the values $P_{ts}(y=0) = 0.70$ and $P_{ts}(y=1) = 0.30$. This example is illustrated in the figure below

Training



Testing



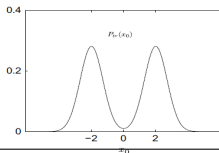
Prior probability shift: Example

Example⁸: $Y \rightarrow X$ problem with one covariate x_0 and a target y that may take the class value $y=0$ and $y=1$. In **training data**, $P_{tr}(y=0) = P_{tr}(y=1) = 0.5$ and $P_{tr}(x_0 | y)$ is defined as

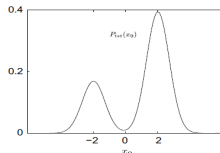
$$x_0 = \begin{cases} \mathcal{N}(2, 0.5), & \text{when } y = 1 \\ \mathcal{N}(-2, 0.5), & \text{otherwise} \end{cases} \quad (1)$$

Now consider that in the **test data** $P_{ts}(x_0 | y=0)$ and $P_{ts}(x_0 | y=1)$ remains unchanged, but the class prior probabilities vary, taking the values $P_{ts}(y=0) = 0.70$ and $P_{ts}(y=1) = 0.30$. This example is illustrated in the figure below

Training



Testing



⁸ *Moreno-Torres et al., Pattern Recognition, 2011.*

Concept Shift

- ▶ Concept shift is defined as:

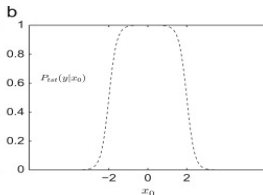
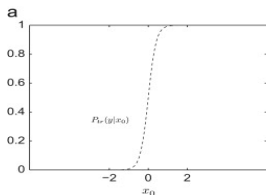
- ▶ $X \rightarrow Y$ problems

$$P_{tr}(Y | X) \neq P_{ts}(Y | X) \quad \text{and} \quad P_{tr}(X) = P_{ts}(X)$$

&

- ▶ $Y \rightarrow X$ problems

$$P_{tr}(X | Y) \neq P_{ts}(X | Y) \quad \text{and} \quad P_{tr}(Y) = P_{ts}(Y)$$



Concept Shift

- ▶ Concept shift is defined as:

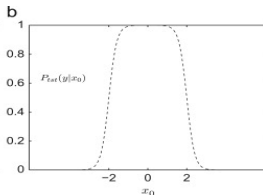
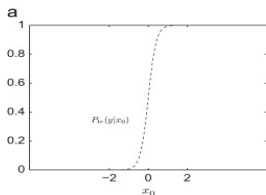
- ▶ $X \rightarrow Y$ problems

$$P_{tr}(Y | X) \neq P_{ts}(Y | X) \quad \text{and} \quad P_{tr}(X) = P_{ts}(X)$$

&

- ▶ $Y \rightarrow X$ problems

$$P_{tr}(X | Y) \neq P_{ts}(X | Y) \quad \text{and} \quad P_{tr}(Y) = P_{ts}(Y)$$



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Concept Shift

- ▶ Concept shift is defined as:

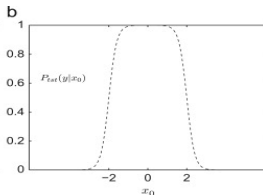
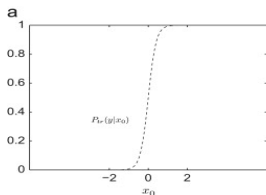
- ▶ $X \rightarrow Y$ problems

$$P_{tr}(Y | X) \neq P_{ts}(Y | X) \quad \text{and} \quad P_{tr}(X) = P_{ts}(X)$$

&

- ▶ $Y \rightarrow X$ problems

$$P_{tr}(X | Y) \neq P_{ts}(X | Y) \quad \text{and} \quad P_{tr}(Y) = P_{ts}(Y)$$



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Outline

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Causes of Dataset Shift

1. The main two causes of dataset are Sample **Selection Bias** and **Non-stationary environments**.
2. These concepts have created confusion at times, so it is important to remark that these terms are factors that can lead to the appearance of some of the shift explained, but they do not constitute Dataset Shift themselves.

Causes of Dataset Shift

1. The main two causes of dataset are Sample **Selection Bias** and **Non-stationary environments**.
2. These concepts have created confusion at times, so it is important to remark that these terms are factors that can lead to the appearance of some of the shift explained, but they do not constitute Dataset Shift themselves.

Causes of Dataset Shift

1. The main two causes of dataset are Sample **Selection Bias** and **Non-stationary environments**.
2. These concepts have created confusion at times, so it is important to remark that these terms are factors that can lead to the appearance of some of the shift explained, but they do not constitute Dataset Shift themselves.

Causes of Dataset Shift. . .

1. **Sample selection bias**: the discrepancy in distribution is due to the fact that the **training examples** have been **obtained through a biased method**, and thus do not represent reliably the operating environment where the classifier is to be deployed (In ML terms, would constitute the **test** set).
2. **Non-stationary environments**: It appears when the **training environment** is **different** from the **test** one, whether it is due to a temporal or a spatial change.

Causes of Dataset Shift. . .

1. **Sample selection bias**: the discrepancy in distribution is due to the fact that the **training examples** have been **obtained through a biased method**, and thus do not represent reliably the operating environment where the classifier is to be deployed (In ML terms, would constitute the **test** set).
2. **Non-stationary environments**: It appears when the **training environment** is **different** from the **test** one, whether it is due to a temporal or a spatial change.

Causes of Dataset Shift. . .

1. **Sample selection bias**: the discrepancy in distribution is due to the fact that the **training examples** have been **obtained through a biased method**, and thus do not represent reliably the operating environment where the classifier is to be deployed (In ML terms, would constitute the **test** set).
2. **Non-stationary environments**: It appears when the **training environment** is **different** from the **test** one, whether it is due to a temporal or a spatial change.

Sample selection bias

1. The term Sample selection bias refers to a systematic flaw in the process of data collection or labeling which causes **training** examples to be selected non-uniformly from the population to be modeled.
2. The term has been used as a synonym of covariate shift (which is not correct), but also on its own as a related problem to Dataset shift

Sample selection bias

1. The term Sample selection bias refers to a systematic flaw in the process of data collection or labeling which causes **training** examples to be selected non-uniformly from the population to be modeled.
2. The term has been used as a synonym of covariate shift (which is not correct), but also on its own as a related problem to Dataset shift

Sample selection bias: An example

1. **Example:** *Survivorship bias* is a common type of sample selection bias. When back-testing an investment strategy on a large group of stocks. Look for securities that have data for the entire sample period (i.e. 15 years).
2. Now, in testing strategy, we need 15 years of stock data.
3. However, eliminating a stock that stopped trading, or shortly left the market, would input a bias in our data sample. Since we are only including stocks that lasted the 15-year period, our final results would be flawed, as these performed well enough to survive the market.

Sample selection bias: An example

1. **Example:** *Survivorship bias* is a common type of sample selection bias. When back-testing an investment strategy on a large group of stocks. Look for securities that have data for the entire sample period (i.e. 15 years).
2. Now, in testing strategy, we need 15 years of stock data.
3. However, eliminating a stock that stopped trading, or shortly left the market, would input a bias in our data sample. Since we are only including stocks that lasted the 15-year period, our final results would be flawed, as these performed well enough to survive the market.

Sample selection bias: An example

1. **Example:** *Survivorship bias* is a common type of sample selection bias. When back-testing an investment strategy on a large group of stocks. Look for securities that have data for the entire sample period (i.e. 15 years).
2. Now, in **testing** strategy, we need 15 years of stock data.
3. However, eliminating a stock that stopped trading, or shortly left the market, would input a **bias** in our data sample. Since we are only including stocks that lasted the 15-year period, our final results would be **flawed**, as these performed well enough to survive the market.

Sample selection bias: An example

1. **Example:** *Survivorship bias* is a common type of sample selection bias. When back-testing an investment strategy on a large group of stocks. Look for securities that have data for the entire sample period (i.e. 15 years).
2. Now, in **testing** strategy, we need 15 years of stock data.
3. However, eliminating a stock that stopped trading, or shortly left the market, would input a **bias** in our data sample. Since we are only including stocks that lasted the 15-year period, our final results would be **flawed**, as these performed well enough to survive the market.

Learning Under Different Training and Testing Distributions

Causes of Dataset Shift

-
- Group 1 Group 2 ... Group k-1 Group k
- ● ● ● ● ● ● ● ● ● ● ● ●
- Training Validation
- $\hat{f}(x)$ $(\hat{f}(x_t) - y_t)^2$

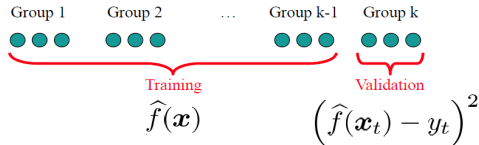
- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

Correcting: Dataset shift generated by Sample Selection Bias

Dr Haider Raza

1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.

Causes of Dataset Shift

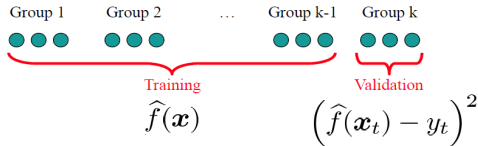


Correcting: Dataset shift generated by Sample Selection Bias

Dr Haider Raza

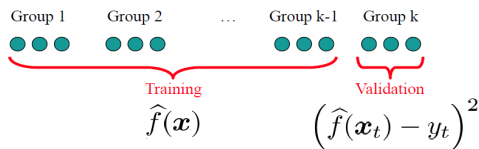
1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.

Causes of Dataset Shift



Correcting: Dataset shift generated by Sample Selection Bias

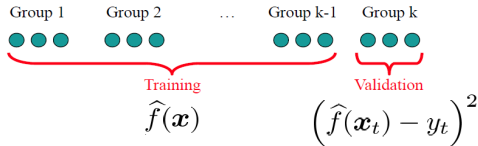
1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.



5. This method is **cross-validation (CV)** and is almost unbiased without covariate shift.
6. But, **CV is heavily biased under covariate shift.**

Correcting: Dataset shift generated by Sample Selection Bias

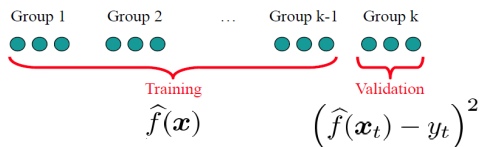
1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.



Correcting: Dataset shift generated by Sample Selection Bias

Dr Haider Raza

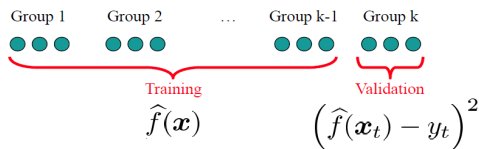
1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.



5. This method is **cross-validation (CV)** and is almost unbiased without covariate shift.
6. But, **CV is heavily biased under covariate shift.**

Correcting: Dataset shift generated by Sample Selection Bias

1. Divide the **training** samples into k groups.
2. **Train** a learning machine with $k - 1$ groups.
3. Validate the **trained** machine using the rest.
4. Repeat this for all the combination and output the mean validation error.



5. This method is **cross-validation (CV)** and is almost unbiased without covariate shift.
6. But, **CV is heavily biased under covariate shift.**

Non-stationary environments

1. In real-world applications, it is often the case that the data is not (time- or space-) stationary
2. One of the most relevant non-stationary scenarios involves adversarial classification problems, such as spam filtering, brain signal classification, and network intrusion detection.
3. This type of problem is receiving an increasing amount of attention in the machine learning field.

Non-stationary environments

1. In real-world applications, it is often the case that the data is not (time- or space-) stationary
2. One of the most relevant non-stationary scenarios involves adversarial classification problems, such as **spam filtering**, **brain signal classification**, and **network intrusion detection**.
3. This type of problem is receiving an increasing amount of attention in the machine learning field.

Non-stationary environments

1. In real-world applications, it is often the case that the data is not (time- or space-) stationary
2. One of the most relevant non-stationary scenarios involves adversarial classification problems, such as **spam filtering**, **brain signal classification**, and **network intrusion detection**.
3. This type of problem is receiving an increasing amount of attention in the machine learning field.

Non-stationary environments

1. In real-world applications, it is often the case that the data is not (time- or space-) stationary
2. One of the most relevant non-stationary scenarios involves adversarial classification problems, such as **spam filtering**, **brain signal classification**, and **network intrusion detection**.
3. This type of problem is receiving an increasing amount of attention in the machine learning field.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Learning in Dataset Shift

Learning in Non-stationary environments

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza



Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Ditzler, G et al., (2015). Learning in Nonstationary Environments : A Survey. *IEEE Computational Intelligence Magazine*, 10(4), 12–25.

Learning in Non-stationary environments

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

1. **Passive Approach:** continuously update the model over time (without requiring an explicit detection of the change)
2. **Active Approach:** rely on an explicit detection of the change in the data distribution to activate an adaptation mechanism

Learning in Non-stationary environments

1. **Passive Approach:** continuously update the model over time (without requiring an explicit detection of the change)
2. **Active Approach:** rely on an explicit detection of the change in the data distribution to activate an adaptation mechanism

1. Does not use a shift detection method to detect changes.
2. Perform a continuous adaptation of the model parameters every time new data arrive.
3. **Advantage:** Maintain an up-to-date model at all times.
4. **Advantage:** Avoiding the potential pitfall associated with the active approaches, that is, failing to detect a change or falsely detecting a non-existent change (false alarm).
5. **Disadvantage:** Update every time a new data arrives. Not suitable for real-time systems.

1. Does not use a shift detection method to detect changes.
2. Perform a continuous adaptation of the model parameters every time new data arrive.
3. **Advantage:** Maintain an up-to-date model at all times.
4. **Advantage:** Avoiding the potential pitfall associated with the active approaches, that is, failing to detect a change or falsely detecting a non-existent change (false alarm).
5. **Disadvantage:** Update every time a new data arrives. Not suitable for real-time systems.

1. Does not use a shift detection method to detect changes.
2. Perform a continuous adaptation of the model parameters every time new data arrive.
3. **Advantage:** Maintain an up-to-date model at all times.
4. **Advantage:** Avoiding the potential pitfall associated with the active approaches, that is, failing to detect a change or falsely detecting a non-existent change (false alarm).
5. **Disadvantage:** Update every time a new data arrives. Not suitable for real-time systems.

1. Does not use a shift detection method to detect changes.
2. Perform a continuous adaptation of the model parameters every time new data arrive.
3. **Advantage:** Maintain an up-to-date model at all times.
4. **Advantage:** Avoiding the potential pitfall associated with the active approaches, that is, failing to detect a change or falsely detecting a non-existent change (false alarm).
5. **Disadvantage:** Update every time a new data arrives. Not suitable for real-time systems.

1. Does not use a shift detection method to detect changes.
2. Perform a continuous adaptation of the model parameters every time new data arrive.
3. **Advantage:** Maintain an up-to-date model at all times.
4. **Advantage:** Avoiding the potential pitfall associated with the active approaches, that is, failing to detect a change or falsely detecting a non-existent change (false alarm).
5. **Disadvantage:** Update every time a new data arrives. Not suitable for real-time systems.

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error.**
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble.**
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers.**

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

1. Single Classifier:

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. Ensemble Classifier:

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Passive Approaches Methods

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error.**
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble.**
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers.**

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Passive Approaches Methods

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

1. **Single Classifier:**

- ▶ Provide lower computational cost.
- ▶ **Decision trees** are the most common classifier for data stream mining.
- ▶ **Very Fast Decision Tree (VFDT)** and **Online Information Network (ONI)** are very popular approaches based on sliding window method.
- ▶ Recently, **Extreme Learning Machine (ELM)** based on neural networks gaining popularity for learning non-stationary data.

2. **Ensemble Classifier:**

- ▶ More accurate than single classifier due to **reduction in the variance of the error**.
- ▶ Flexible to incorporate new data, simply by **adding new members to ensemble**.
- ▶ Provide mechanism to **forget irrelevant knowledge**, simply by **removing old classifiers**.

Active Approaches Methods

Shifts in Data

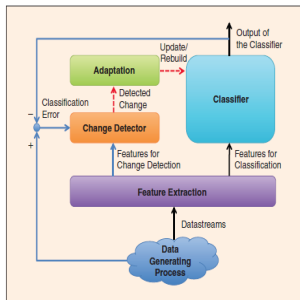
Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary



It is based on **change detection mechanism** that triggers, whenever advisable, an adaptation mechanism aiming at reacting to the detected change by updating or building new classifier.

1. **Change/Shift Detection:**

- ▶ Hypothesis Test, Change-point methods, Sequential hypothesis test, and Change-detection test.
- ▶ Popular methods: EWMA, CUSUM, JIT, ICI, DDM and many more.

2. **Adaptation:**

- ▶ Supervised adaptation, unsupervised adaptation, semi-supervised adaptation, and transduction.
- ▶ Popular methods: Learn⁺⁺.NSE, COMPOSE, JIT adaptive classifier, MOA, and many more.

1. **Change/Shift Detection:**

- ▶ Hypothesis Test, Change-point methods, Sequential hypothesis test, and Change-detection test.
- ▶ Popular methods: EWMA, CUSUM, JIT, ICI, DDM and many more.

2. **Adaptation:**

- ▶ Supervised adaptation, unsupervised adaptation, semi-supervised adaptation, and transduction.
- ▶ Popular methods: Learn⁺⁺, NSE, COMPOSE, JIT adaptive classifier, MOA, and many more.

1. **Change/Shift Detection:**

- ▶ Hypothesis Test, Change-point methods, Sequential hypothesis test, and Change-detection test.
- ▶ Popular methods: EWMA, CUSUM, JIT, ICI, DDM and many more.

2. **Adaptation:**

- ▶ Supervised adaptation, unsupervised adaptation, semi-supervised adaptation, and transduction.
- ▶ Popular methods: Learn⁺⁺.NSE, COMPOSE, JIT adaptive classifier, MOA, and many more.

1. **Change/Shift Detection:**

- ▶ Hypothesis Test, Change-point methods, Sequential hypothesis test, and Change-detection test.
- ▶ Popular methods: EWMA, CUSUM, JIT, ICI, DDM and many more.

2. **Adaptation:**

- ▶ Supervised adaptation, unsupervised adaptation, semi-supervised adaptation, and transduction.
- ▶ Popular methods: Learn⁺⁺.NSE, COMPOSE, JIT adaptive classifier, MOA, and many more.

1. **Change/Shift Detection:**

- ▶ Hypothesis Test, Change-point methods, Sequential hypothesis test, and Change-detection test.
- ▶ Popular methods: EWMA, CUSUM, JIT, ICI, DDM and many more.

2. **Adaptation:**

- ▶ Supervised adaptation, unsupervised adaptation, semi-supervised adaptation, and transduction.
- ▶ Popular methods: Learn⁺⁺.NSE, COMPOSE, JIT adaptive classifier, MOA, and many more.

Other Learning Approaches

Learning Under
Different Training
and Testing
Distributions

Dr Haider Raza

Shifts in Data

Dataset Shift

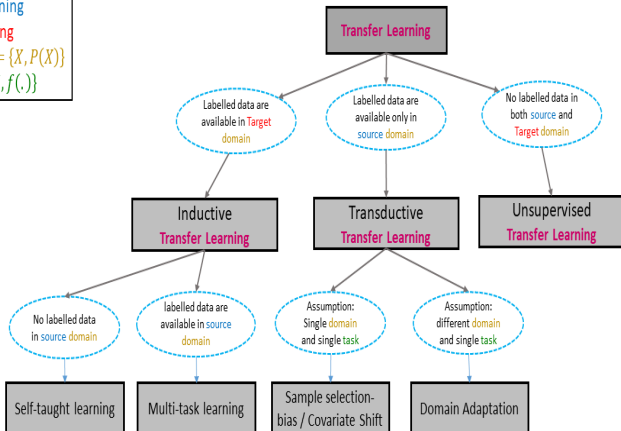
Types of Dataset Shift

Causes of Dataset Shift

Learning in
Dataset Shift

Summary

Source: Training
Target: Testing
Domain: $D = \{X, P(X)\}$
Task: $T = \{Y, f(\cdot)\}$



Domain Adaptation

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

- ▶ Learning from a **source (training)** data distribution a well performing model on a different (but related) **target (testing)** data distribution.
- ▶ Example, one of the tasks of the **common spam filtering problem** consists in adapting a model from one user (the **source** distribution) to a new one who receives significantly different emails (the **target** distribution).
- ▶ Note that, when more than one **source** distribution is available the problem is referred to as **multi-source domain adaptation**.
- ▶ Iterative Domain Adaptation Algorithm
 1. a model h is learned from the labeled examples;
 2. h automatically labels some target examples;
 3. a new model is learned from the new labeled examples.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. Spam-filtering, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish spam emails from legitimate ones, for example an Persian speaker may find that all emails in French are spam, not so for French speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to money transfer.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. Spam-filtering, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish spam emails from legitimate ones, for example an Persian speaker may find that all emails in French are spam, not so for French speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to money transfer.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. Spam-filtering, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish spam emails from legitimate ones, for example an Persian speaker may find that all emails in French are spam, not so for French speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to money transfer.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. **Spam-filtering**, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish **spam emails** from **legitimate ones**, for example an **Persian** speaker may find that all emails in **French** are spam, not so for **French** speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to **money transfer**.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. **Spam-filtering**, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish **spam emails** from **legitimate ones**, for example an **Persian** speaker may find that all emails in **French** are spam, not so for **French** speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to **money transfer**.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. **Spam-filtering**, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish **spam emails** from **legitimate ones**, for example an **Persian** speaker may find that all emails in **French** are spam, not so for **French** speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to **money transfer**.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Multi-task Learning (MLT)

- ▶ MTL is sub-field of learning in which multiple learning tasks are solved at the same time, while utilizing commonalities and differences across tasks.
- ▶ It aims to improve the performance of multiple classification tasks by learning them jointly
- ▶ Example:
 1. **Spam-filtering**, which can be treated as distinct but related classification tasks across different users.
 2. To make this more concrete, consider that different people have different distributions of features which distinguish **spam emails** from **legitimate ones**, for example an **Persian** speaker may find that all emails in **French** are spam, not so for **French** speakers.
 3. Yet there is a definite commonality in this classification task across users, for example one common feature might be text related to **money transfer**.
 4. Solving each user's spam classification problem jointly via MTL can let the solutions inform each other and improve performance.

Covariate Shift Adaptation

- Covariate shift appears only in $X \rightarrow Y$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(X) \neq P_{ts}(X)$$

- Under covariate shift, the ratio $\frac{P_{tr}(X, Y)}{P_{ts}(X, Y)}$ can be re-written as follows:

$$\frac{P_{ts}(X, Y)}{P_{tr}(X, Y)} = \frac{P_{ts}(X) P_{ts}(Y | X)}{P_{tr}(X) P_{tr}(Y | X)} = \frac{P_{ts}(X)}{P_{tr}(X)}$$

- We want to weight each training instance with $\frac{P_{ts}(X)}{P_{tr}(X)}$.
- A major challenge is how to estimate the ratio $\frac{P_{ts}(X)}{P_{tr}(X)}$ for each x in the training set. In some work, a principled method of using non-parametric kernel density estimation is explored by Sugiyama in his work (i.e. density ratio estimation without direct estimation of P_{tr} and P_{ts}) Sugiyama et al., 2013, Neural Computation.

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Covariate Shift Adaptation

- Covariate shift appears only in $X \rightarrow Y$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(X) \neq P_{ts}(X)$$

- Under covariate shift, the ratio $\frac{P_{tr}(X,Y)}{P_{ts}(X,Y)}$ can be re-written as follows:

$$\frac{P_{ts}(X, Y)}{P_{tr}(X, Y)} = \frac{P_{ts}(X)}{P_{tr}(X)} \frac{P_{ts}(Y | X)}{P_{tr}(Y | X)} = \frac{P_{ts}(X)}{P_{tr}(X)}$$

- We want to weight each training instance with $\frac{P_{ts}(X)}{P_{tr}(X)}$.
- A major challenge is how to estimate the ratio $\frac{P_{ts}(X)}{P_{tr}(X)}$ for each x in the training set. In some work, a principled method of using non-parametric kernel density estimation is explored by Sugiyama in his work (i.e. density ratio estimation without direct estimation of P_{tr} and P_{tr}) Sugiyama et al., 2013, Neural Computation.

Covariate Shift Adaptation

- Covariate shift appears only in $X \rightarrow Y$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(X) \neq P_{ts}(X)$$

- Under covariate shift, the ratio $\frac{P_{tr}(X, Y)}{P_{ts}(X, Y)}$ can be re-written as follows:

$$\frac{P_{ts}(X, Y)}{P_{tr}(X, Y)} = \frac{P_{ts}(X) P_{ts}(Y | X)}{P_{tr}(X) P_{tr}(Y | X)} = \frac{P_{ts}(X)}{P_{tr}(X)}$$

- We want to weight each training instance with $\frac{P_{ts}(X)}{P_{tr}(X)}$.
- A major challenge is how to estimate the ratio $\frac{P_{ts}(X)}{P_{tr}(X)}$ for each x in the training set. In some work, a principled method of using non-parametric kernel density estimation is explored by Sugiyama in his work (i.e. [density ratio estimation](#) without direct estimation of P_{tr} and P_{ts}) [Sugiyama et al., 2013](#), [Neural Computation](#).

Covariate Shift Adaptation

- Covariate shift appears only in $X \rightarrow Y$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(X) \neq P_{ts}(X)$$

- Under covariate shift, the ratio $\frac{P_{tr}(X, Y)}{P_{ts}(X, Y)}$ can be re-written as follows:

$$\frac{P_{ts}(X, Y)}{P_{tr}(X, Y)} = \frac{P_{ts}(X) P_{ts}(Y | X)}{P_{tr}(X) P_{tr}(Y | X)} = \frac{P_{ts}(X)}{P_{tr}(X)}$$

- We want to weight each training instance with $\frac{P_{ts}(X)}{P_{tr}(X)}$.
- A major challenge is how to estimate the ratio $\frac{P_{ts}(X)}{P_{tr}(X)}$ for each x in the training set. In some work, a principled method of using non-parametric kernel density estimation is explored by Sugiyama in his work (i.e. [density ratio estimation](#) without direct estimation of P_{tr} and P_{ts}) [Sugiyama et al., 2013](#), [Neural Computation](#).

Covariate Shift Adaptation

- Covariate shift appears only in $X \rightarrow Y$ problems, and is defined as the case where

$$P_{tr}(Y | X) = P_{ts}(Y | X)$$

&

$$P_{tr}(X) \neq P_{ts}(X)$$

- Under covariate shift, the ratio $\frac{P_{tr}(X, Y)}{P_{ts}(X, Y)}$ can be re-written as follows:

$$\frac{P_{ts}(X, Y)}{P_{tr}(X, Y)} = \frac{P_{ts}(X) P_{ts}(Y | X)}{P_{tr}(X) P_{tr}(Y | X)} = \frac{P_{ts}(X)}{P_{tr}(X)}$$

- We want to weight each training instance with $\frac{P_{ts}(X)}{P_{tr}(X)}$.
- A major challenge is how to estimate the ratio $\frac{P_{ts}(X)}{P_{tr}(X)}$ for each x in the training set. In some work, a principled method of using non-parametric kernel density estimation is explored by Sugiyama in his work (i.e. **density ratio estimation** without direct estimation of P_{tr} and P_{tr}) [Sugiyama et al., 2013, Neural Computation.](#)

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Why is it difficult to learn from Data.
- ▶ Dataset shift and types of dataset shift.
- ▶ Causes of dataset shift.
- ▶ How to handle sample selection bias.
- ▶ Approaches to non-stationary learning (i.e. Passive and Active Approaches).

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Why is it difficult to learn from Data.
- ▶ Dataset shift and types of dataset shift.
- ▶ Causes of dataset shift.
- ▶ How to handle sample selection bias.
- ▶ Approaches to non-stationary learning (i.e. Passive and Active Approaches).

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Why is it difficult to learn from Data.
- ▶ Dataset shift and types of dataset shift.
- ▶ Causes of dataset shift.
- ▶ How to handle sample selection bias.
- ▶ Approaches to non-stationary learning (i.e. Passive and Active Approaches).

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Why is it difficult to learn from Data.
- ▶ Dataset shift and types of dataset shift.
- ▶ Causes of dataset shift.
- ▶ How to handle sample selection bias.
- ▶ Approaches to non-stationary learning (i.e. Passive and Active Approaches).

Summary

Shifts in Data

Dataset Shift

Types of Dataset Shift

Causes of Dataset Shift

Learning in Dataset Shift

Summary

- ▶ Why is it difficult to learn from Data.
- ▶ Dataset shift and types of dataset shift.
- ▶ Causes of dataset shift.
- ▶ How to handle sample selection bias.
- ▶ Approaches to non-stationary learning (i.e. Passive and Actives Approaches).

Contact me:

h.raza@essex.ac.uk

sagihaid@gmail.com

Webpage:

<https://sites.google.com/site/whoishraza/home>

<https://www.essex.ac.uk/people/razah72409>

Follow me:



: @sagihaider



LinkedIn : sagihaider

THANK YOU!