

# Evaluación Honesta de Clasificadores en Clasificación Supervisada

Guzmán Santafé<sup>(1)</sup>, Iñaki Inza<sup>(2)</sup>

<sup>(1)</sup>Universidad Pública de Navarra

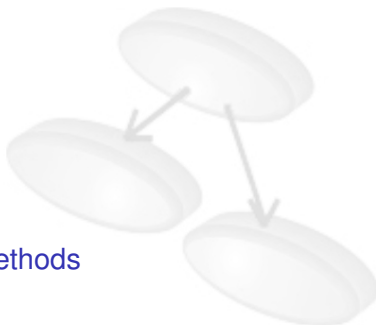
<sup>(2)</sup>Universidad del País Vasco

CAEPIA'11

7 de Noviembre, 2011

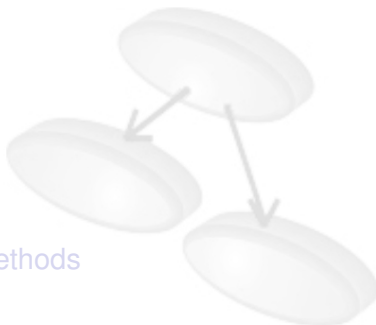
# Outline of the Tutorial

- 1 Introduction
- 2 Scores
- 3 Estimation Methods
- 4 Hypothesis Testing

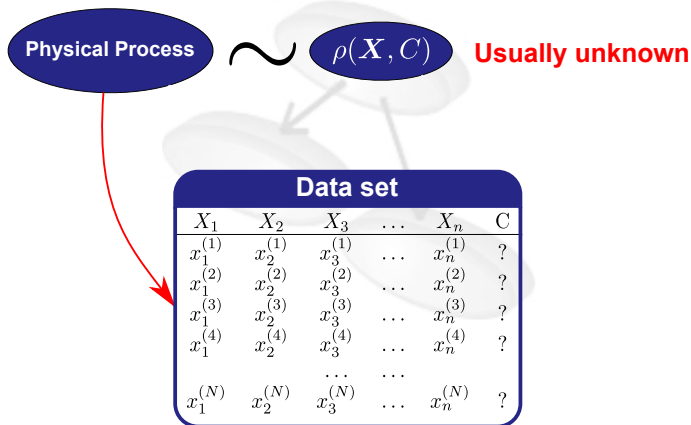


# Outline of the Tutorial

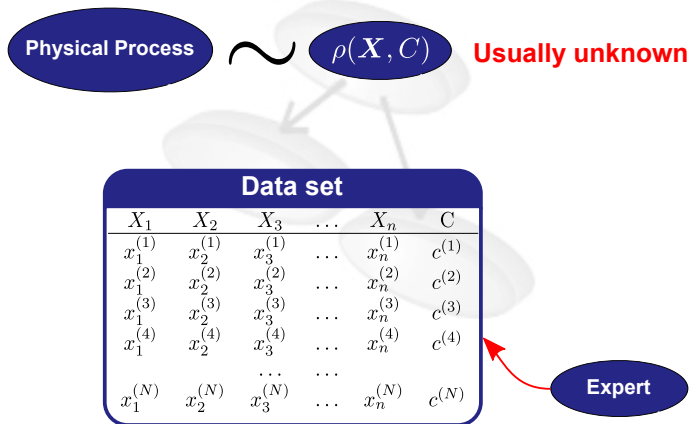
- 1 Introduction
- 2 Scores
- 3 Estimation Methods
- 4 Hypothesis Testing



# Classification Problem



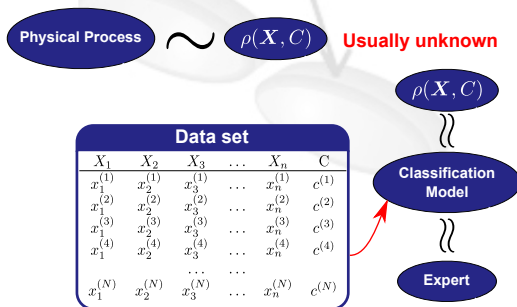
# Classification Problem



# Supervised Classification

## Learning from Experience

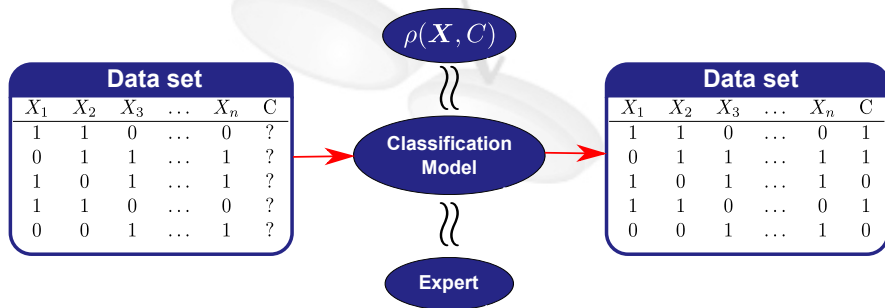
- “Automate the work of the expert”
- Tries to model  $\rho(\mathbf{X}, C)$



# Supervised Classification

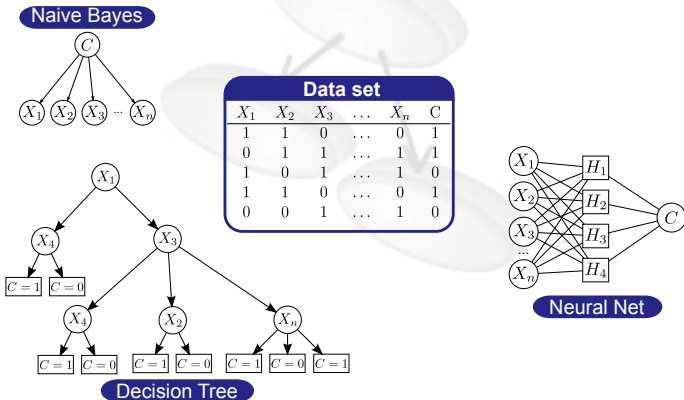
## Classification Model

- Classifier labels new data (unknown class value)



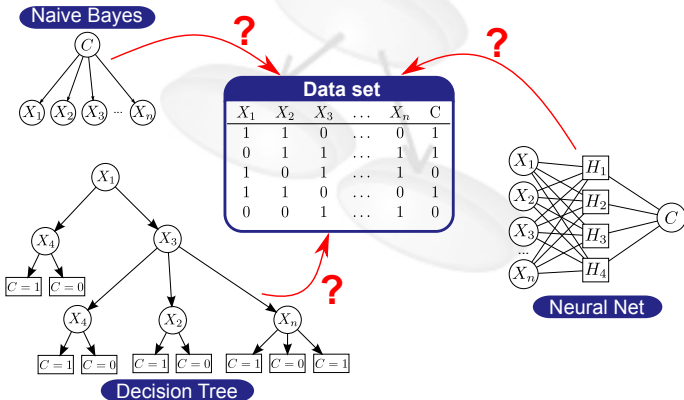
# Motivation for Honest Evaluation

- Many classification paradigms



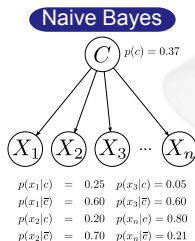
# Motivation for Honest Evaluation

- Which is the best paradigm for a classification problem?



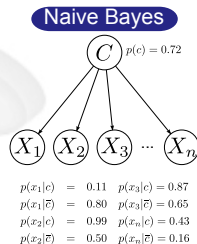
# Motivation for Honest Evaluation

- Many parameter configurations



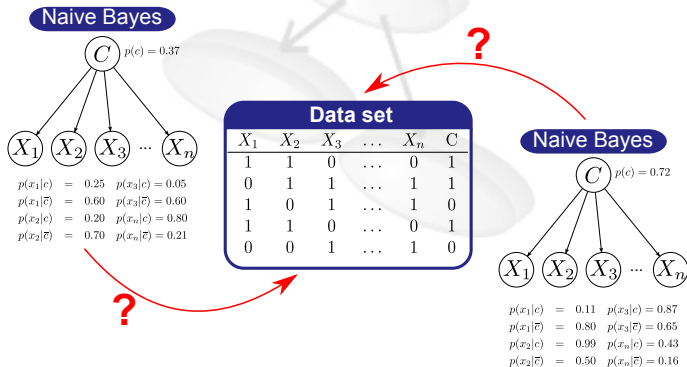
**Data set**

$X_1$	$X_2$	$X_3$	...	$X_n$	$C$
1	1	0	...	0	1
0	1	1	...	1	1
1	0	1	...	1	0
1	1	0	...	0	1
0	0	1	...	1	0



# Motivation for Honest Evaluation

- Which is the best parameter configuration for a classification problem?



# Motivation for Honest Evaluation

## Honest Evaluation

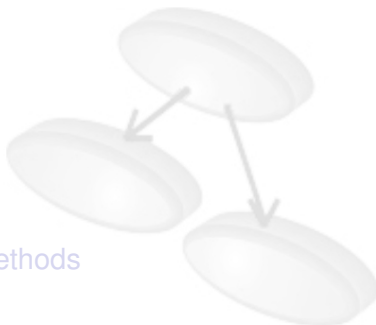
- Need to know the goodness of a classifier
- Methodology to compare classifiers
- Assess the validity of evaluation/comparison

## Steps for Honest Evaluation

- Scores: quality measures
- Estimation methods: estimate value of a score
- Statistical tests: comparison among different solutions

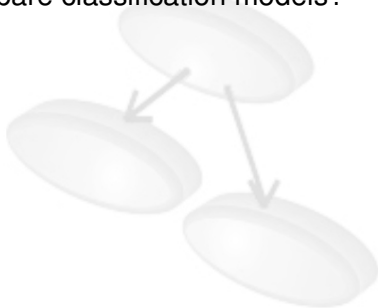
# Outline of the Tutorial

- 1 Introduction
- 2 Scores**
- 3 Estimation Methods
- 4 Hypothesis Testing



# Motivation

- How to compare classification models?



## Score

Function that provides a quality measure for a classifier when solving a classification problem

# Motivation

- How to compare classification models?

We need some way to measure the classification performance!!!

## Score

Function that provides a quality measure for a classifier when solving a classification problem

## Motivation

- How to compare classification models?

We need some way to measure the classification performance!!!

### Score

Function that provides a quality measure for a classifier when solving a classification problem

# Motivation

## What Does *Best Quality* Mean?

- What are we interested in?
- What do we want to optimize?
- Characteristics of the problem
- Characteristics of the data set

Different kind of scores

# Scores

## Based on Confusion Matrix

- Accuracy/Classification error
- Recall
- Specificity
- Precision
- F-Score

## Based on Receiver Operating Characteristics (ROC)

- Area under the ROC curve (AUC)

# Scores

## Based on Confusion Matrix

- Accuracy/Classification error → Classification
- Recall
- Specificity
- Precision
- F-Score

## Based on Receiver Operating Characteristics (ROC)

- Area under the ROC curve (AUC)

# Scores

## Based on Confusion Matrix

- Accuracy/Classification error → Classification
- Recall
- Specificity → Information Retrieval
- Precision
- F-Score

## Based on Receiver Operating Characteristics (ROC)

- Area under the ROC curve (AUC)

# Scores

## Based on Confusion Matrix

- Accuracy/Classification error → Classification
- Recall
- Specificity → Information Retrieval
- Precision
- F-Score

## Based on Receiver Operating Characteristics (ROC)

- Area under the ROC curve (AUC) → Medical Domains

# Confusion Matrix

## Two-Class Problem

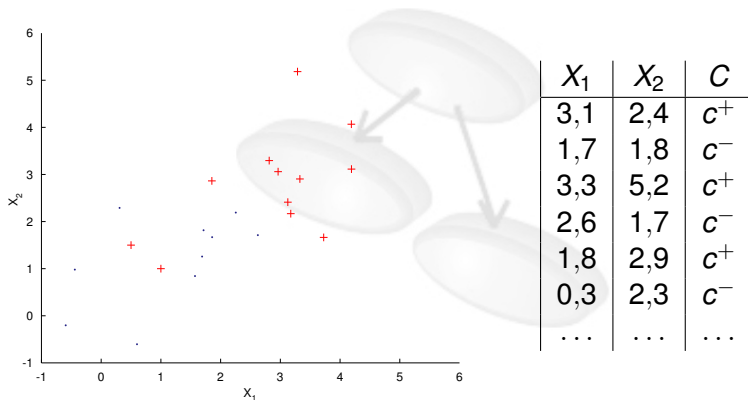
		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	$TP$	$FP$	$N^+$
	$c^-$	$FN$	$TN$	$N^-$
Total		$\hat{N}^+$	$\hat{N}^-$	$N$

# Confusion Matrix

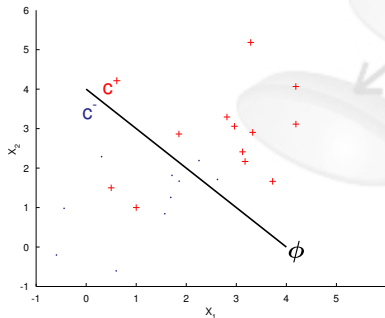
## Several-Class Problem

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$N_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$N_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$N_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$N_n$
Total		$\hat{N}_1$	$\hat{N}_2$	$\hat{N}_3$	...	$\hat{N}_n$	$N$

# Two-Class Problem - Example



## Two-Class Problem - Example

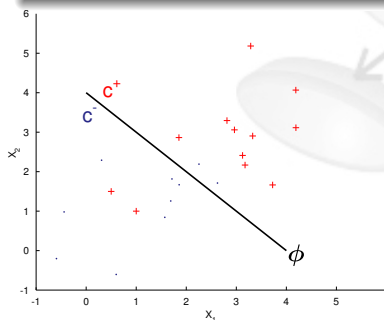


		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	10	2	12
	$c^-$	2	8	10
Total		12	10	<b>22</b>

# Accuracy/Classification Error

## Definition

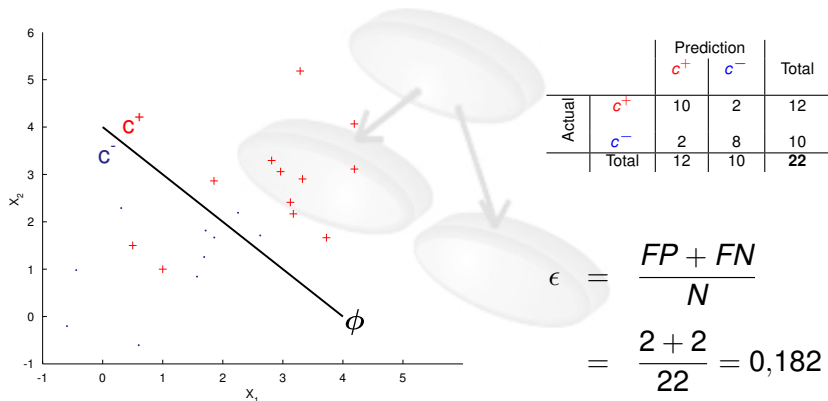
- Data samples classified correctly/incorrectly



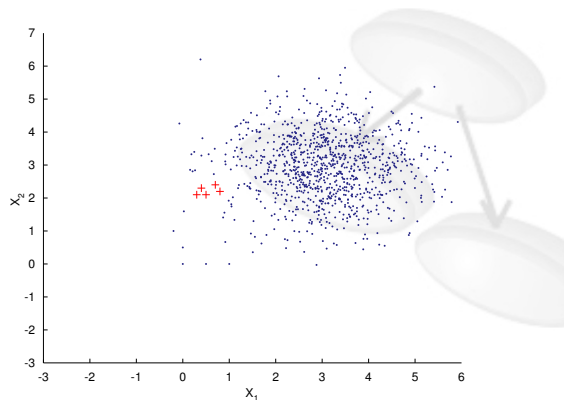
		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	10	2	12
	$c^-$	2	8	10
Total		12	10	<b>22</b>

$$\epsilon(\phi) = p(\phi(\mathbf{X}) \neq C) = E_{\rho(\mathbf{x}, c)}[1 - \delta(c, \phi(\mathbf{x}))]$$

# Accuracy/Classification Error

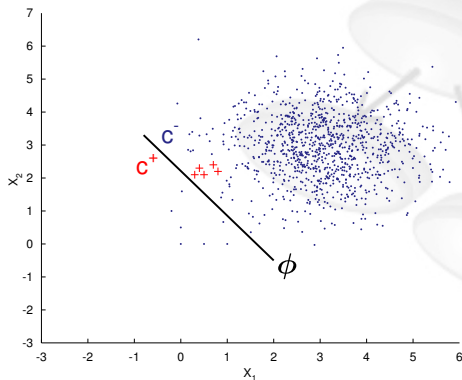


# Skew Data



$X_1$	$X_2$	$C$
0,8	2,2	$c^+$
0,47	2,3	$c^+$
0,5	2,1	$c^+$
2,4	2,9	$c^-$
3,1	1,2	$c^-$
2,5	3,1	$c^-$
...	...	...

# Skew Data - Classification Error

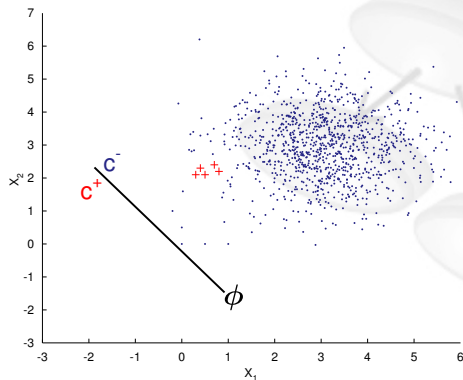


		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	7	993	1000
Total		7	998	<b>1005</b>

$$\epsilon = \frac{7 + 5}{1005} = 0,012$$

Very low  $\epsilon$ !!

# Skew Data - Classification Error

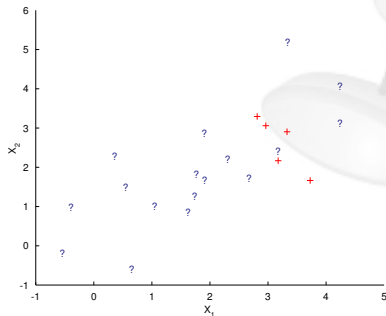


		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	0	1000	1000
Total		0	1005	<b>1005</b>

$$\epsilon = \frac{0 + 5}{1005} = 0,005$$

Better??

# Positive Unlabeled Learning



## Positive Labeled Data

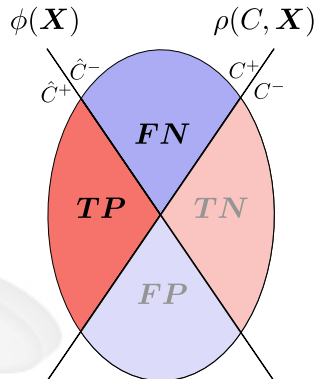
- Only positive samples labeled
- Many unlabeled samples:
  - Positive?
  - Negative?
- Classification error is useless

# Recall

## Definition

- Fraction of positive class samples correctly classified
- Other names  $\left\{ \begin{array}{l} \text{True positive rate} \\ \text{Sensitivity} \end{array} \right.$

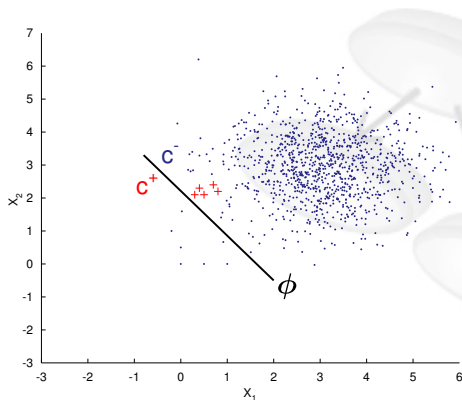
$$r(\phi) = \frac{TP}{TP + FN} = \frac{TP}{P}$$



## Definition Based on Probabilities

$$r(\phi) = p(\phi(\mathbf{x}) = c^+ | C = c^+) = E_{\rho(\mathbf{x}|C=c^+)}[\delta(\phi(\mathbf{x}), c^+)]$$

# Skew Data - Recall

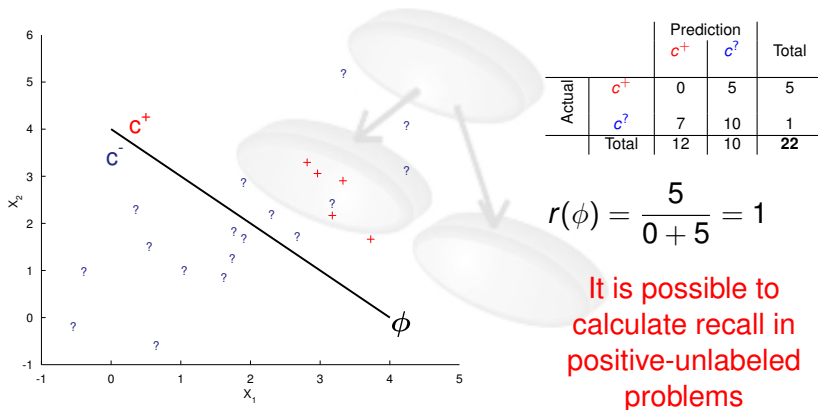


		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	7	993	1000
Total		7	998	<b>1005</b>

$$r(\phi) = \frac{0}{0 + 5} = 0$$

Very bad recall!!!

# Positive Unlabeled Learning - Recall

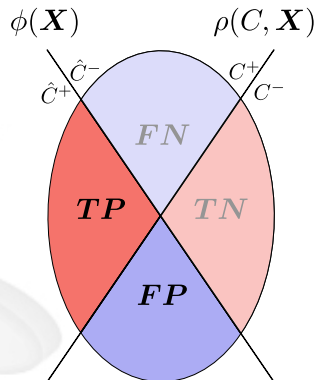


# Precision

## Definition

- Fraction of data samples classified as  $c^+$  which are actually  $c^+$

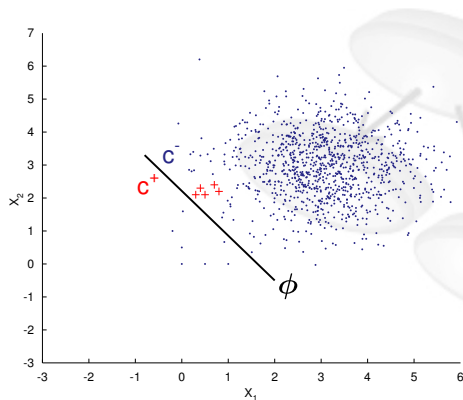
$$pr(\phi) = \frac{TP}{TP + FP} = \frac{TP}{\hat{p}}$$



## Definition Based on Probabilities

$$pr(\phi) = p(C = c^+ | \phi(\mathbf{x}) = c^+) = E_{\rho(\mathbf{x} | \phi(\mathbf{x}) = c^+)}[\delta(\phi(\mathbf{x}), c^+)]$$

# Skew Data - Precision

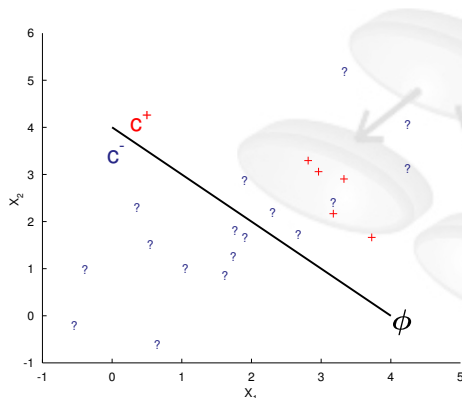


		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	7	993	1000
Total		7	998	<b>1005</b>

$$pr(\phi) = \frac{0}{0 + 7} = 0$$

Very bad precision!!

# Positive Unlabeled Learning - Precision



- Precision is not a good score for positive-unlabeled data samples
- **Not all the positive samples are labeled**

# Precision & Recall Application Domains

## Spam Filtering

- Decide if an email is spam or not
  - Precision: Proportion of real spam in the spam-box
  - Recall: Proportion of total spam messages identified by the system

## Sentiment Analysis

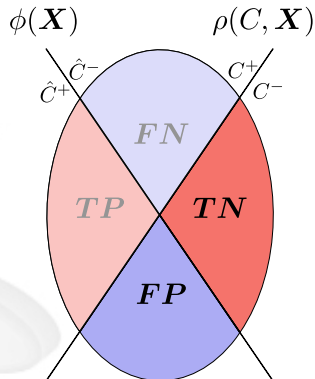
- Classify opinions about specific products given by users in blogs, webs, forums, etc.
  - Precision: Proportion of opinions classified as positive being actually positive
  - Recall: Proportion of positive opinions identified as positive

# Specificity

## Definition

- Fraction of negative class samples correctly identified
- *Specificity* =  $1 - \text{FalsePositiveRate}$

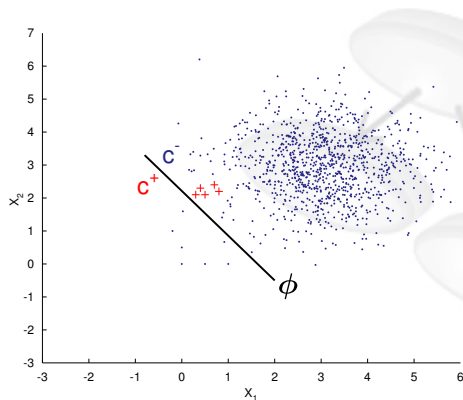
$$sp(\phi) = \frac{TN}{TN + FP} = \frac{TN}{N}$$



## Definition Based on Probabilities

$$sp(\phi) = p(\phi(\mathbf{x}) = c^- | C = c^-) = E_{p(\mathbf{x}|C=c^-)}[1 - \delta(\phi(\mathbf{x}), c^-)]$$

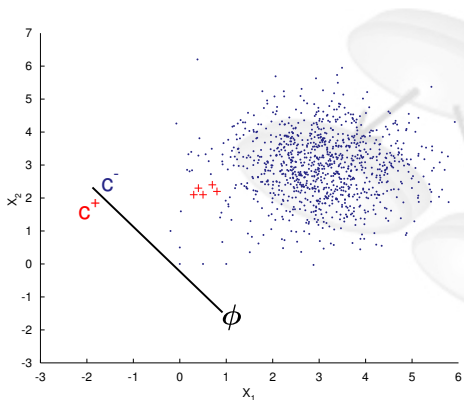
# Skew Data - Specificity



		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	7	993	1000
Total		7	998	<b>1005</b>

$$sp(\phi) = \frac{993}{993 + 7} = 0,99$$

# Skew Data - Specificity



		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	0	1000	1000
Total		0	1005	<b>1005</b>

$$sp(\phi) = \frac{1000}{1000 + 0} = 1,00$$

## Balanced Scores

- Balanced accuracy rate

$$Bal. acc = \frac{1}{2} \left( \frac{TP}{P} + \frac{TN}{N} \right) = \frac{recall + specificity}{2}$$

- Balanced error rate

$$Bal. \epsilon = \frac{1}{2} \left( \frac{FP}{P} + \frac{FN}{N} \right)$$

### Skew Data

		Prediction		Total
		$c^+$	$c^-$	
Actual	$c^+$	0	5	5
	$c^-$	7	993	1000
Total		7	998	<b>1005</b>

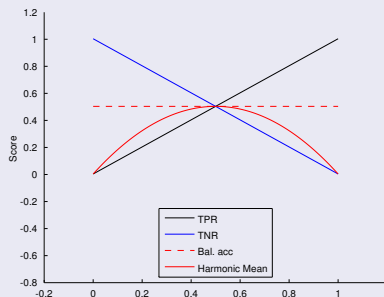
- $Bal. acc = \frac{1}{2} \left( \frac{0}{5} + \frac{993}{1000} \right) \approx 0,5$
- $Bal. \epsilon = \frac{1}{2} \left( \frac{7}{7} + \frac{5}{1000} \right) \approx 0,5$

## Balanced Scores

- $F - \text{Score} = \frac{(\beta^2 + 1) \text{Precision} \cdot \text{Recall}}{\beta^2 (\text{Precision} + \text{Recall})}$
- $F_1 - \text{Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \rightarrow \text{Harmonic Mean}$

### Harmonic Mean

- Maximized with balanced components
- Bal. acc  $\rightarrow$  arithmetic mean



## Classification Cost

- All misclassifications cannot be equally considered

### E.g. Medical Diagnosis Problem

It does not have the same cost diagnosing a healthy patient as ill rather than diagnosing an ill patient as healthy

### Classification Model

May be of interest to minimize the expected cost instead the classification error

# Dealing with Classification Cost

## Loss Function

Associate an economic/utility/etc. cost to each classification.

- Typical loss function in classification  $\rightarrow$  0/1 Loss
- We can use cost matrix to specify the associated cost:

		Prediction	
		$c^+$	$c^-$
Actual	$c^+$	0	1
	$c^-$	1	0

## Dealing with Classification Cost

### Loss Function

Associate an economic/utility/etc. cost to each classification.

- Typical loss function in classification  $\rightarrow$  0/1 Loss
- We can use cost matrix to specify the associated cost:

		Prediction	
		$c^+$	$c^-$
Actual	$c^+$	$Cost_{TP}$	$Cost_{FN}$
	$c^-$	$Cost_{FP}$	$Cost_{TN}$

## Dealing with Classification Cost

### Loss Function

Associate an economic/utility/etc. cost to each classification.

- Typical loss function in classification  $\rightarrow$  0/1 Loss
- We can use cost matrix to specify the associated cost:

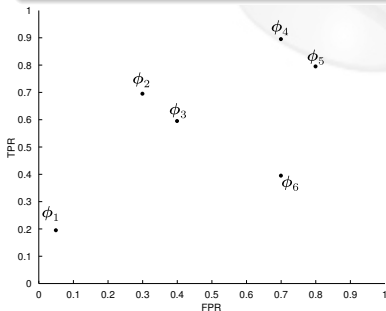
		Prediction	
		$c^+$	$c^-$
Actual	$c^+$	$Cost_{TP}$	$Cost_{FN}$
	$c^-$	$Cost_{FP}$	$Cost_{TN}$

Usually not easy to give an associated cost

# Receiver Operating Characteristics (ROC)

## ROC Space

Coordinate system used for visualizing classifiers performance where  $TPR$  is plotted on the  $Y$  axis and  $FPR$  is plotted on the  $X$  axis.

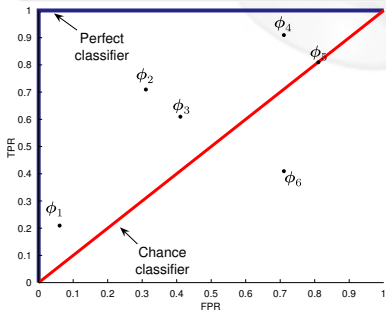


- $\phi_1$ : kNN
- $\phi_2$ : Neural network
- $\phi_3$ : Naive Bayes
- $\phi_4$ : SVM
- $\phi_5$ : Linear regression
- $\phi_6$ : Decision tree

# Receiver Operating Characteristics (ROC)

## ROC Space

Coordinate system used for visualizing classifiers performance where  $TPR$  is plotted on the  $Y$  axis and  $FPR$  is plotted on the  $X$  axis.

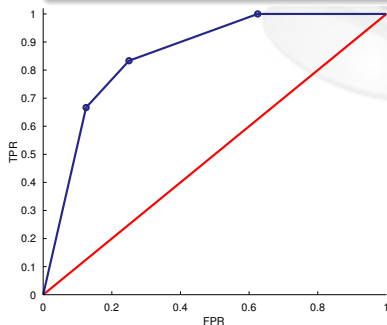


- $\phi_1$ : kNN
- $\phi_2$ : Neural network
- $\phi_3$ : Naive Bayes
- $\phi_4$ : SVM
- $\phi_5$ : Linear regression
- $\phi_6$ : Decision tree

# Receiver Operating Characteristics (ROC)

## ROC Curve

For a probabilistic/fuzzy classifier, a ROC curve is a plot of the TPR vs. FPR as its discrimination threshold is varied

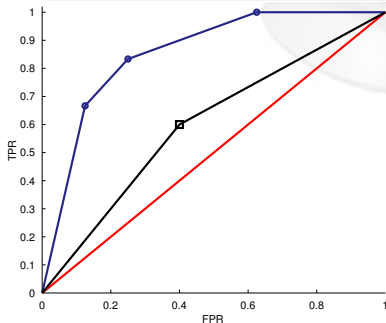


$p(c \mathbf{x})$	$T = 0,2$	$T = 0,5$	$T = 0,8$	$C$
0,99	$c^+$	$c^+$	$c^+$	$c^+$
0,90	$c^+$	$c^+$	$c^+$	$c^+$
0,85	$c^+$	$c^+$	$c^+$	$c^+$
0,80	$c^+$	$c^+$	$c^+$	$c^-$
0,78	$c^+$	$c^+$	$c^-$	$c^+$
0,70	$c^+$	$c^+$	$c^-$	$c^-$
0,60	$c^+$	$c^+$	$c^-$	$c^+$
0,45	$c^+$	$c^-$	$c^-$	$c^-$
0,40	$c^+$	$c^-$	$c^-$	$c^-$
0,30	$c^+$	$c^-$	$c^-$	$c^-$
0,20	$c^+$	$c^-$	$c^-$	$c^+$
0,15	$c^-$	$c^-$	$c^-$	$c^-$
0,10	$c^-$	$c^-$	$c^-$	$c^-$
0,05	$c^-$	$c^-$	$c^-$	$c^-$

# Receiver Operating Characteristics (ROC)

## ROC Curve

For a crisp classifier a ROC curve can be obtained by interpolation from a single point



$p(c \mathbf{x})$	$T = 0,2$	$T = 0,5$	$T = 0,8$	$C$
0,99	$c^+$	$c^+$	$c^+$	$c^+$
0,90	$c^+$	$c^+$	$c^+$	$c^+$
0,85	$c^+$	$c^+$	$c^+$	$c^+$
0,80	$c^+$	$c^+$	$c^+$	$c^-$
0,78	$c^+$	$c^+$	$c^-$	$c^+$
0,70	$c^+$	$c^+$	$c^-$	$c^-$
0,60	$c^+$	$c^+$	$c^-$	$c^+$
0,45	$c^+$	$c^-$	$c^-$	$c^-$
0,40	$c^+$	$c^-$	$c^-$	$c^-$
0,30	$c^+$	$c^-$	$c^-$	$c^-$
0,20	$c^+$	$c^-$	$c^-$	$c^+$
0,15	$c^-$	$c^-$	$c^-$	$c^-$
0,10	$c^-$	$c^-$	$c^-$	$c^-$
0,05	$c^-$	$c^-$	$c^-$	$c^-$

# Receiver Operating Characteristics (ROC)

## ROC Curve

- Insensitive to skew class distribution
- Insensitive to misclassification cost

## Dominance Relationship

A ROC curve  $A$  dominates another ROC curve  $B$  if  $A$  is always above and to the left of  $B$  in the plot

# Receiver Operating Characteristics (ROC)

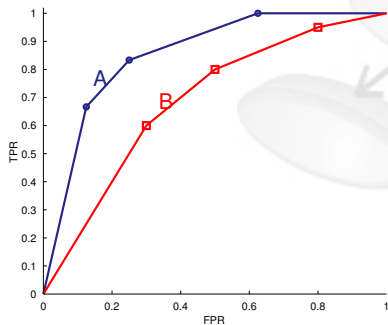
## ROC Curve

- Insensitive to skew class distribution
- Insensitive to misclassification cost

## Dominance Relationship

A ROC curve  $A$  dominates another ROC curve  $B$  if  $A$  is always above and to the left of  $B$  in the plot

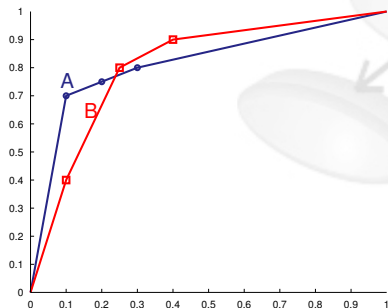
# Receiver Operating Characteristics (ROC)



## Dominance

- A dominates B throughout all the range of  $T$
- A has a better predictive performance over any condition of cost and class distribution

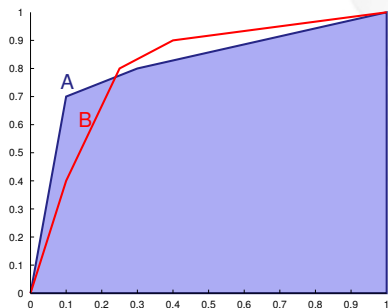
# Receiver Operating Characteristics (ROC)



## No-Dominance

- The dominance relationship may not be so clear
- No model is the best under all possible scenarios

# Receiver Operating Characteristics (ROC)



## Area Under ROC Curve

- Equivalent to Wilcoxon test
- If  $A$  dominates  $B$ :  
 $AUC(A) \geq AUC(B)$
- If  $A$  does not dominate  $B$   
 $AUC$  “cannot identify the best classifier”

# Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

## $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

# Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

## $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

# Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

## $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

# Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

## $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

# Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

## $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

## Generalization to Multilabel-Class

- Most of the presented scores are for binary classification
- Generalization to multilabel is possible
  - E.g. One-vs-All approach

		Prediction					Total
		$c_1$	$c_2$	$c_3$	...	$c_n$	
Actual	$c_1$	$TP_1$	$FN_{12}$	$FN_{13}$	...	$FN_{1n}$	$P_1$
	$c_2$	$FN_{21}$	$TP_2$	$FN_{23}$	...	$FN_{2n}$	$P_2$
	$c_3$	$FN_{31}$	$FN_{32}$	$TP_3$	...	$FN_{3n}$	$P_3$
	...	...	...	...	...	...	...
	$c_n$	$FN_{n1}$	$FN_{n2}$	$FN_{n3}$	...	$TP_n$	$P_n$
Total		$\hat{P}_1$	$\hat{P}_2$	$\hat{P}_3$	...	$\hat{P}_n$	

### $c_1$ vs. All ( $score_1$ )

- $TP$
- $TN$
- $FN$
- $FP$

$$score_{TOT} = \sum_{i=1}^n score_i \cdot p(c_i)$$

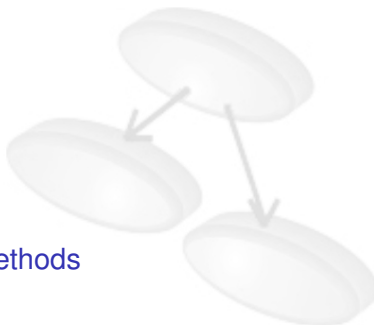
# Scores

## The Use of a Specific Score Depends on:

- Application domain
- Characteristics of the problem
- Characteristics of the data set
- Our interest when solving the problem
- etc.

# Outline of the Tutorial

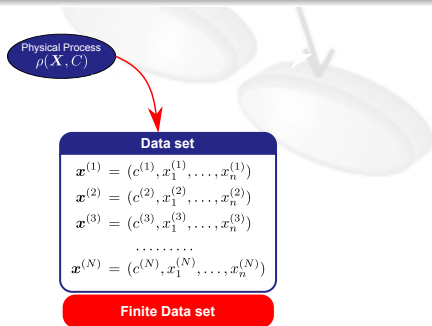
- 1 Introduction
- 2 Scores
- 3 Estimation Methods**
- 4 Hypothesis Testing



# Introduction

## Estimation

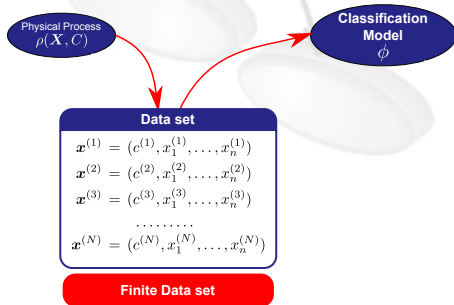
- Select a score to measure the quality
- Calculate the true value of the score
- Limited information is available



# Introduction

## Estimation

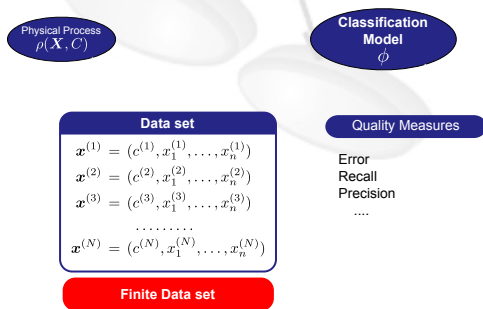
- Select a score to measure the quality
- Calculate the true value of the score
- Limited information is available



# Introduction

## Estimation

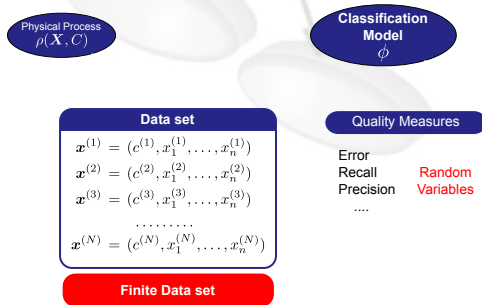
- Select a score to measure the quality
- Calculate the true value of the score
- Limited information is available



# Introduction

## Estimation

- Select a score to measure the quality
- Calculate the true value of the score
- Limited information is available



# Introduction

True Value -  $\epsilon_N$

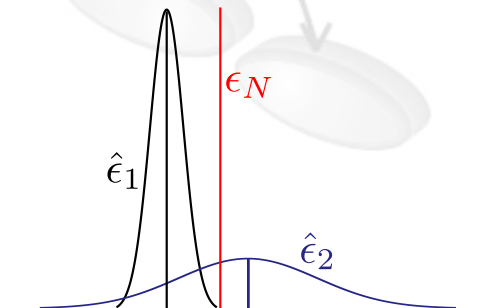
Expected value of the score for a set of  $N$  data samples  
sampled from  $\rho(\mathbf{X}, C)$

# Introduction

True Value -  $\epsilon_N$

Expected value of the score for a set of  $N$  data samples sampled from  $\rho(\mathbf{X}, C)$

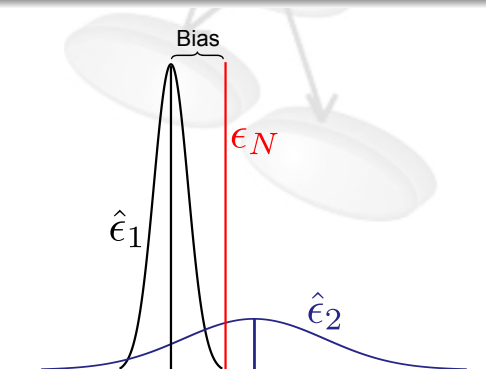
$\rho(\mathbf{X}, C)$  unknown  $\rightarrow$  Point estimation of the score ( $\hat{\epsilon}$ )



# Introduction

## Bias

Difference between the estimation of the score and its true value:  $E_{\rho}[\hat{\epsilon}] - \epsilon_N$

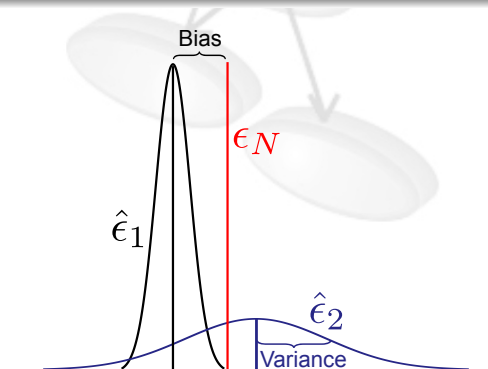


# Introduction

## Variance

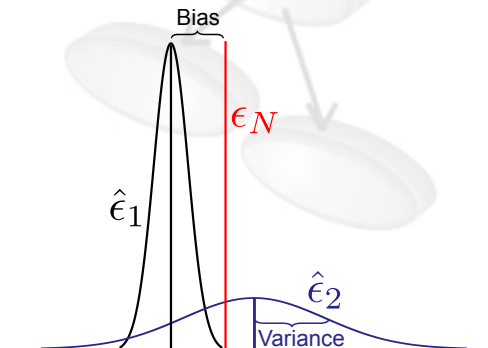
Deviation of the estimated value from its expected value:

$$\text{var}(\hat{\epsilon}) = E[\hat{\epsilon} - E_{\rho}[\hat{\epsilon}]]$$



# Introduction

- Bias and variance depend on the estimation method
- Trade-off between bias and variance needed



# Introduction

## Data set

$$\mathbf{x}^{(1)} = (c^{(1)}, x_1^{(1)}, \dots, x_n^{(1)})$$

$$\mathbf{x}^{(2)} = (c^{(2)}, x_1^{(2)}, \dots, x_n^{(2)})$$

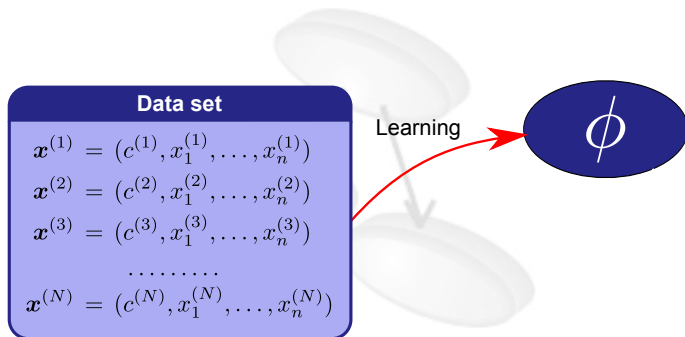
$$\mathbf{x}^{(3)} = (c^{(3)}, x_1^{(3)}, \dots, x_n^{(3)})$$

.....

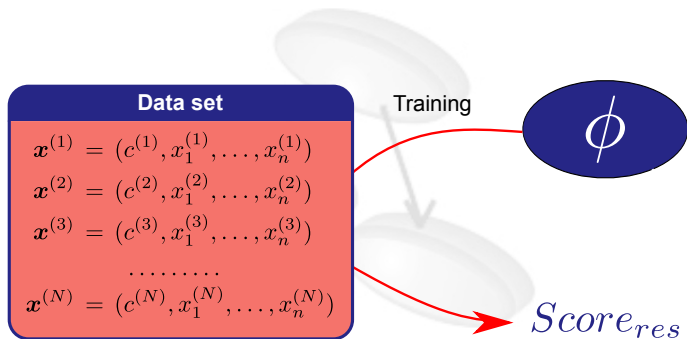
$$\mathbf{x}^{(N)} = (c^{(N)}, x_1^{(N)}, \dots, x_n^{(N)})$$

- Finite data set to estimate the score
- Several choices depending on how this data set is dealt with

# Resubstitution



# Resubstitution

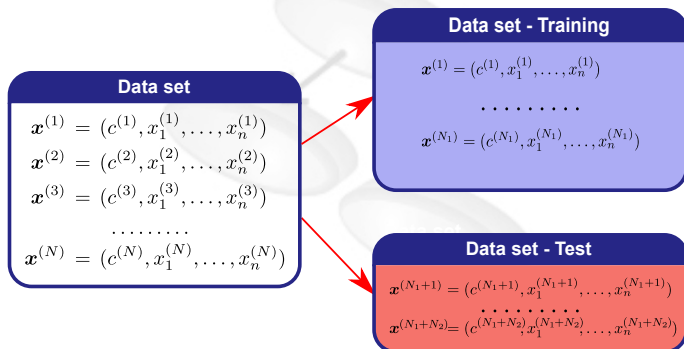


# Resubstitution

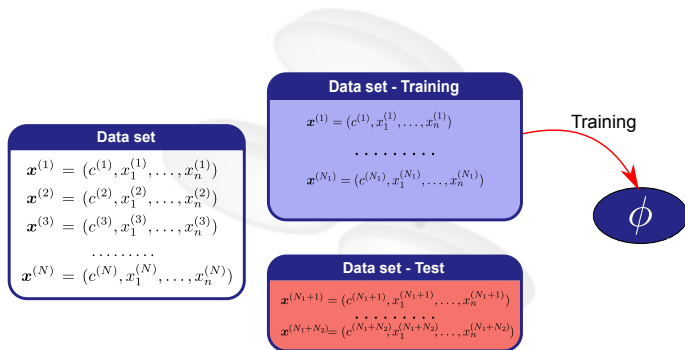
## Classification Error Estimation

- The simplest estimation method
- Biased estimation  $\epsilon_N$
- Smaller variance
- Too optimistic (overfitting problem)
- Bad estimator of the true classification error

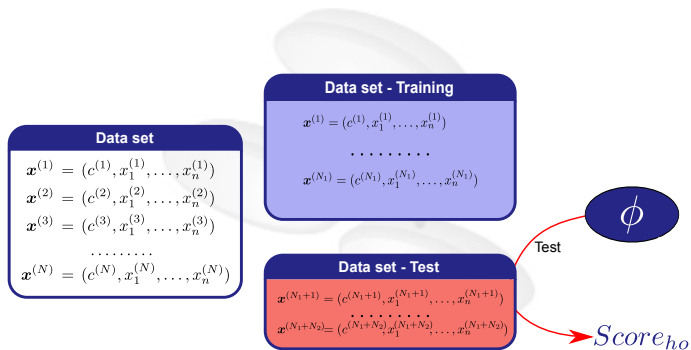
# Hold-Out



# Hold-Out



# Hold-Out

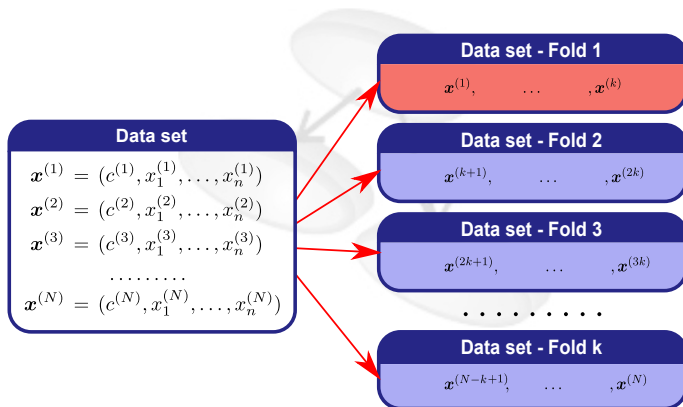


# Hold-Out

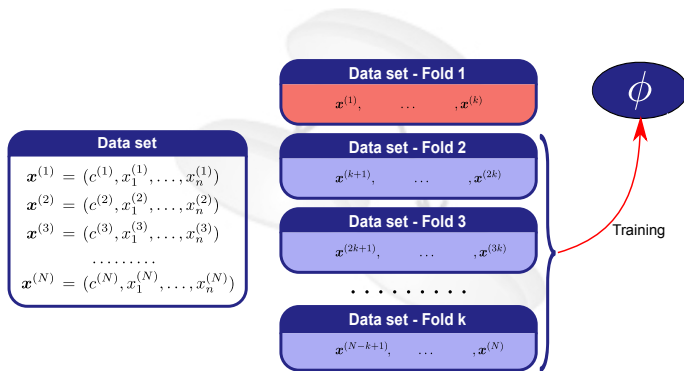
## Classification Error Estimation

- Unbiased estimator of  $\epsilon_{N_1}$
- Biased estimator of  $\epsilon_N$
- Large bias (pessimistic estimation of the true classification error)
- Bias related to  $N_1$  and  $N_2$

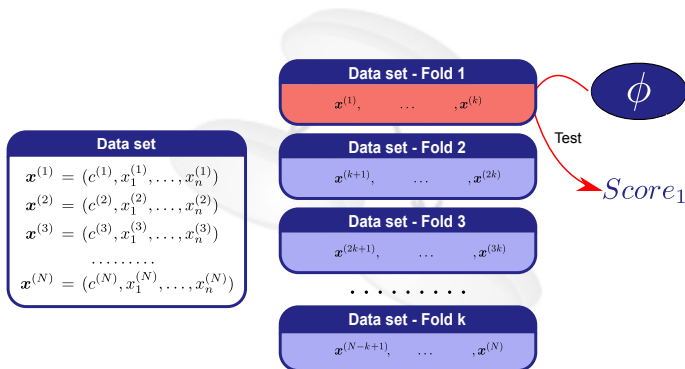
# k-Fold Cross-Validation



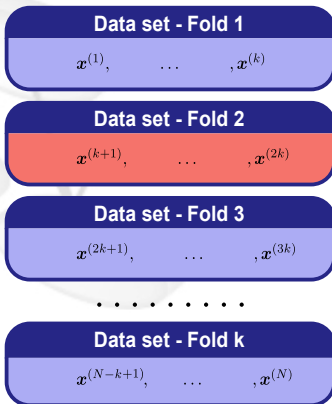
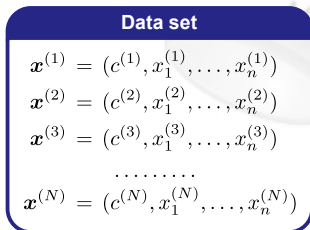
# k-Fold Cross-Validation



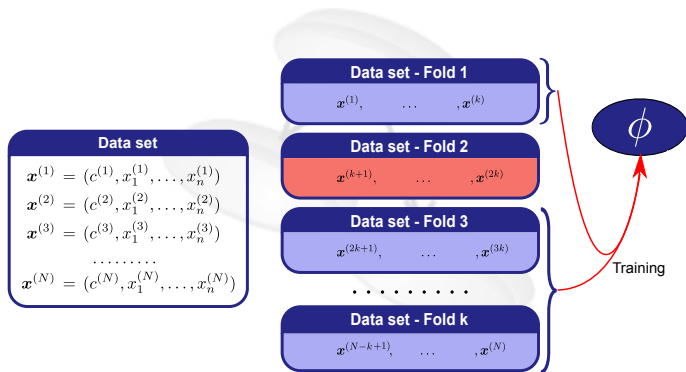
# k-Fold Cross-Validation



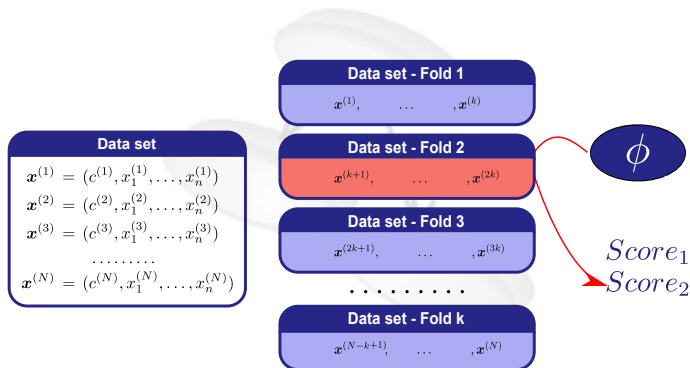
# k-Fold Cross-Validation



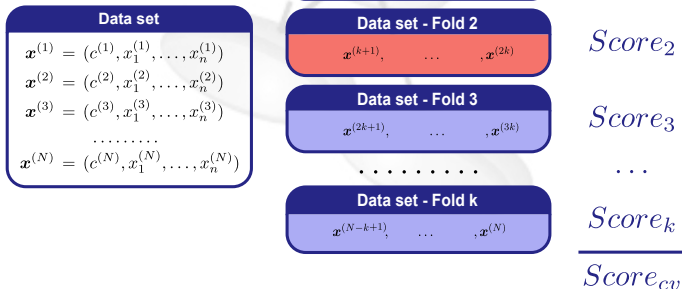
# k-Fold Cross-Validation



# k-Fold Cross-Validation



# k-Fold Cross-Validation



# $k$ -Fold Cross-Validation

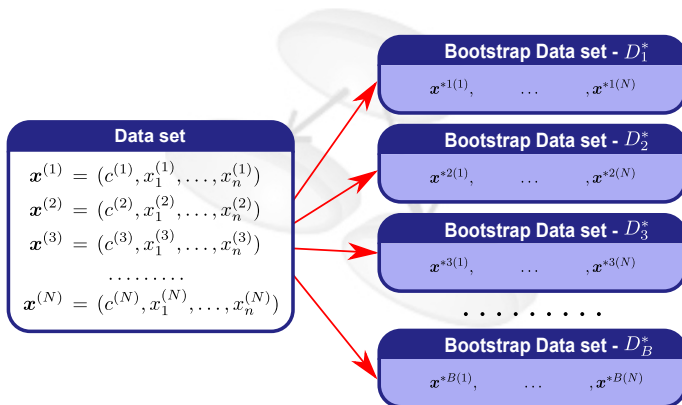
## Classification Error Estimation

- Unbiased estimator of  $\epsilon_{N-\frac{N}{k}}$
- Biased estimation of  $\epsilon_N$
- Smaller bias than Hold-Out

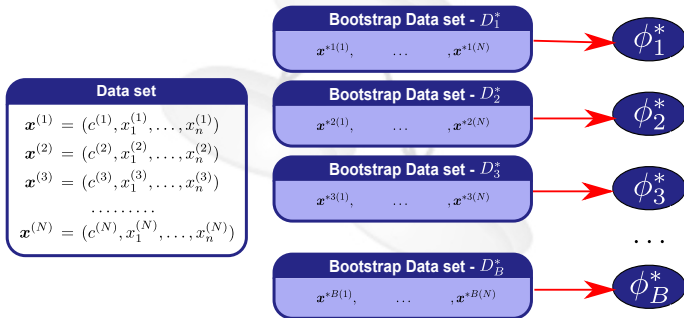
## Leaving-One-Out

- Special case of  $k$ -fold Cross-Validation ( $k = N$ )
- Quasi unbiased estimation for  $N$
- Improves the bias with respect to CV
- Increases the variance  $\rightarrow$  more unstable
- Higher computational cost

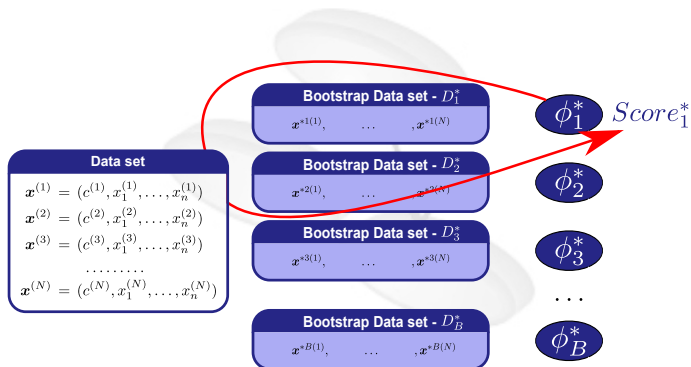
# Bootstrap



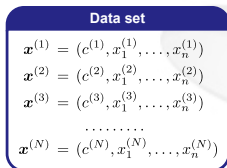
# Bootstrap



# Bootstrap



# Bootstrap



$$\phi_1^* \text{ Score}_1^*$$



$$\phi_2^* \text{ Score}_2^*$$



$$\phi_3^* \text{ Score}_3^*$$



$$\phi_B^* \text{ Score}_B^*$$

---


$$\text{Score}_{boot}$$

# Bootstrap

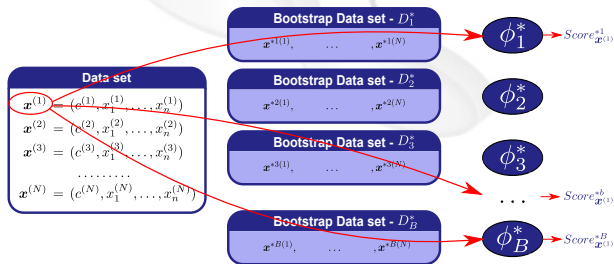
## Classification Error Estimation

- Biased estimation of the classification error
- Variance improved because of resampling
- Uses for testing part of the data used for learning
- “Similar to resubstitution”
- Problem of overfitting

# Leaving-One-Out Bootstrap

- Mimics Cross-Validation

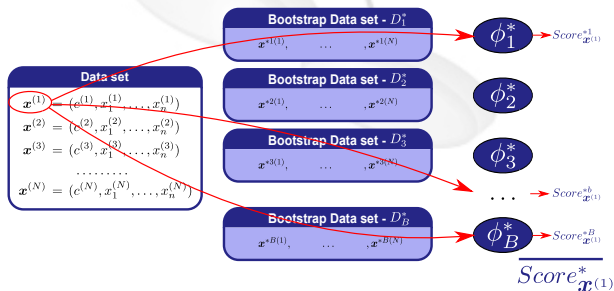
- Each  $\mathbf{x}^{(i)}$  is only evaluated by  $\phi_j^* \begin{cases} j = 1, \dots, N \\ \mathbf{x}^{(i)} \notin D_j^* \end{cases}$



# Leaving-One-Out Bootstrap

- Mimics Cross-Validation

- Each  $\mathbf{x}^{(i)}$  is only evaluated by  $\phi_j^*$   $\begin{cases} j = 1, \dots, N \\ \mathbf{x}^{(i)} \notin D_j^* \end{cases}$



# Leaving-One-Out Bootstrap

- Mimics Cross-Validation

- Each  $\mathbf{x}^{(i)}$  is only evaluated by  $\phi_j^*$   $\begin{cases} j = 1, \dots, N \\ \mathbf{x}^{(i)} \notin D_j^* \end{cases}$

Data set
$\mathbf{x}^{(1)} = (c^{(1)}, x_1^{(1)}, \dots, x_n^{(1)})$
$\mathbf{x}^{(2)} = (c^{(2)}, x_1^{(2)}, \dots, x_n^{(2)})$
$\mathbf{x}^{(3)} = (c^{(3)}, x_1^{(3)}, \dots, x_n^{(3)})$
.....
$\mathbf{x}^{(N)} = (c^{(N)}, x_1^{(N)}, \dots, x_n^{(N)})$

Bootstrap Data set -  $D_1^*$

$\mathbf{x}^{*1(1)}, \dots, \mathbf{x}^{*1(N)}$

$Score_{\mathbf{x}^{(1)}}^*$

Bootstrap Data set -  $D_2^*$

$\mathbf{x}^{*2(1)}, \dots, \mathbf{x}^{*2(N)}$

$Score_{\mathbf{x}^{(2)}}^*$

Bootstrap Data set -  $D_3^*$

$\mathbf{x}^{*3(1)}, \dots, \mathbf{x}^{*3(N)}$

...

Bootstrap Data set -  $D_B^*$

$\mathbf{x}^{*B(1)}, \dots, \mathbf{x}^{*B(N)}$

$Score_{\mathbf{x}^{(N)}}^*$

---

$Score_{boot}$

## Leaving-One-Out Bootstrap

- Mimics Cross-Validation
- Each  $\mathbf{x}^{(i)}$  is only evaluated by  $\phi_j^*$   $\left\{ \begin{array}{l} j = 1, \dots, N \\ \mathbf{x}^{(i)} \notin D_j^* \end{array} \right.$

### Tries to Avoid the Overfitting Problem

- Expected number of distinct samples on bootstrap data set  $\approx 0,632N$
- Similar to repeated Hold-Out
- Biased upwards:
  - Tends to be a pessimistic estimation of the score

## Improving the Estimation - Bias

- Bias correction terms can be used for error estimation

### Hold-Out/Cross-Validation

- Several proposals
- Improves bias estimation
- Surprisingly not very extended

### Bootstrap

- Improves bias estimation
- Well established methods

## Improving the Estimation - Bias

Corrected Hold-Out ( $\hat{\epsilon}_{ho}^+$ ) - (*Burman, 1989*)

$$\hat{\epsilon}_{ho}^+ = \hat{\epsilon}_{ho} + \hat{\epsilon}_{res} - \hat{\epsilon}_{ho-N}$$

Where

- $\hat{\epsilon}_{ho}$  = standard Hold-Out estimator
- $\hat{\epsilon}_{res}$  = resubstitution error
- $\hat{\epsilon}_{ho-N} = \phi$  learned on Hold-Out learning set but tested on  $D$ .

# Improving the Estimation - Bias

## Corrected Hold-Out ( $\hat{\epsilon}_{ho}^+$ ) - (Burman, 1989)

$$\hat{\epsilon}_{ho}^+ = \hat{\epsilon}_{ho} + \hat{\epsilon}_{res} - \hat{\epsilon}_{ho-N}$$

## Improvement

- $Bias_{\hat{\epsilon}_{ho}} \approx Cons_0 \frac{N_2}{N_1 \cdot N}$
- $Bias_{\hat{\epsilon}_{ho}^+} \approx Cons_1 \frac{N_2}{N_1 \cdot N^2}$

# Improving the Estimation - Bias

## Corrected Cross-Validation ( $\hat{\epsilon}_{cv}^+$ ) - (Burman, 1989)

$$\hat{\epsilon}_{cv}^+ = \hat{\epsilon}_{cv} + \hat{\epsilon}_{res} - \hat{\epsilon}_{cv-N}$$

## Improvement

- $Bias_{\hat{\epsilon}_{cv}} \approx Cons_0 \frac{1}{(k-1) \cdot N}$
- $Bias_{\hat{\epsilon}_{cv}^+} \approx Cons_1 \frac{1}{(k-1) \cdot N^2}$

## Improving the Estimation - Bias

### 0.632 Bootstrap ( $\hat{\epsilon}_{boot}^{.632}$ )

$$\hat{\epsilon}_{boot}^{.632} = 0.368\hat{\epsilon}_{res} + 0.632\hat{\epsilon}_0$$

### Improvement

- $\hat{\epsilon}_0$  is similar to  $\hat{\epsilon}_{100-boot}$  estimator
- Tries to balance optimism (resubstitution) and pessimism ( $\hat{\epsilon}_0$ )
- Works well with “light-fitting” classifiers
- With overfitting classifiers  $\hat{\epsilon}_{boot}^{.632}$  is still too optimistic

## Improving the Estimation - Bias

### 0.632+ Bootstrap ( $\hat{\epsilon}_{boot}^{.632+}$ ) - (Efron & Tibshirani, 1997)

- Correct bias when there is great amount of overfitting
- Based on the non-information error rate ( $\gamma$ ):

$$\hat{\gamma} = \sum_{i=1}^N \sum_{j=1}^N \delta(\mathbf{c}_i, \phi_{\mathbf{x}}(\mathbf{x}_j)) / N^2$$

- Uses the relative overfitting to correct the bias:

$$\hat{R} = \frac{\hat{\epsilon}_0 - \hat{\epsilon}_{res}}{\hat{\gamma} - \hat{\epsilon}_{res}}$$

# Improving the Estimation - Bias

0.632+ Bootstrap ( $\hat{\epsilon}_{boot}^{.632+}$ ) - (Efron & Tibshirani, 1997)

$$\hat{\epsilon}_{boot}^{.632} = (1 - \hat{w})\hat{\epsilon}_{res} + \hat{w}\hat{\epsilon}_0$$

- $\hat{w} = \frac{0.632}{1 - 0.638\hat{R}}$
- $\hat{\gamma} = \sum_{i=1}^N \sum_{j=1}^N \delta(\mathbf{c}_i, \phi_{\mathbf{x}}(\mathbf{x}_j)) / N^2$
- $\hat{R} = \frac{\hat{\epsilon}_0 - \hat{\epsilon}_{res}}{\hat{\gamma} - \hat{\epsilon}_{res}}$

## Improving the Estimation - Variance

### Stratification

- Keeps the proportion of each class in the train/test data
  - Hold-Out: Stratified splitting
  - Cross-Validation: Stratified splitting
  - Bootstrap: Stratified sampling

May improve the variance of the estimation

## Improving the Estimation - Variance

### Repeated Methods

- Applicable to Hold-Out and Cross-Validation
- Bootstrap already includes sampling

### Repeated Hold-Out/Cross-Validation

- Repeat estimation process  $t$ -times
- Simple average over results

### Classification Error Estimation

- Same bias as standard estimation methods
- Reduces the variance with respect  
Hold-Out/Cross-Validation

# Estimation Methods

- Which estimation method is better?

## May Depend on Many Aspects

- The size of the data set
- The classification paradigm used
- The stability of the learning algorithm
- The characteristics of the classification problem
- The bias/variance/computational cost trade-off
- ...

# Estimation Methods

- Which estimation method is better?

## Large Data Sets

- Hold-out may be a good choice
  - Computationally not so expensive
  - Larger bias but depends on the data set size

## Smaller Data Sets

- Repeated Cross-Validation
- Bootstrap 0.632

# Estimation Methods

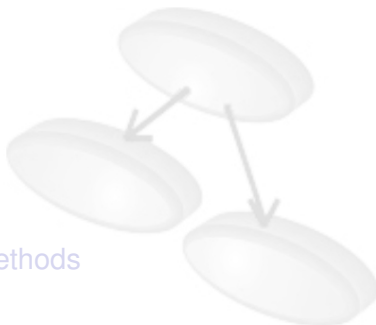
- Which estimation method is better?

## Small Data Sets

- Bootstrap and repeated Cross-Validation may not be informative
- Permutation test (*Ojala & Garriga, 2010*):
  - Can be used to ensure the validity of the estimation
- Confidence intervals (*Isaksson et al., 2008*):
  - May provide more reliable information about the estimation

# Outline of the Tutorial

- 1 Introduction
- 2 Scores
- 3 Estimation Methods
- 4 Hypothesis Testing**



# Motivation

## Basic Concepts

- Hypothesis testing form the basis of scientific reasoning in experimental sciences
- They are used to set scientific statements
- A hypothesis  $H_0$  called **null hypothesis** is tested against another hypothesis  $H_1$  called alternative
- The two hypotheses are not at the same level: reject  $H_0$  does not mean acceptance of  $H_1$
- The objective is to know when **the differences in  $H_0$  are due to randomness or not**

# Hypothesis Testing

## Possible Outcomes of a Test

- Given a sample, a decision is taken about the null hypothesis ( $H_0$ )
- The decision is taken under uncertainty

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

# Hypothesis Testing: An Example

## A Simple Hypothesis Test

- A process is given in nature that follows a Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$
- We have a sample of this process  $\{x_1, \dots, x_N\}$  and a decision must be taken about the following hypotheses:

$$\begin{cases} H_0 : \mu = 60 \\ H_1 : \mu = 50 \end{cases}$$

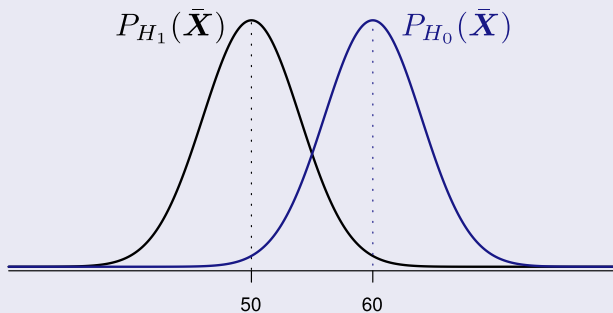
- A **statistic** (function) of the sample is used to take the decision. In our example  $\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$
- The probability distribution of the statistic is known:

$$\bar{X} \rightsquigarrow \mathcal{N}(\mu, \sigma^2/N)$$

# Hypothesis Testing: An Example

## Accept and Reject Regions

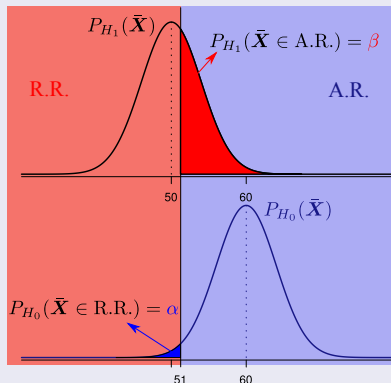
- The sample statistic has a different probability distribution under  $H_0$  and  $H_1$



# Hypothesis Testing: An Example

## Accept and Reject Regions

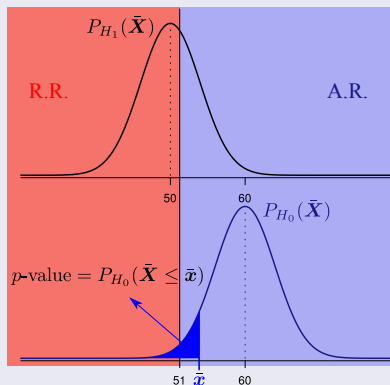
- By controlling  $\alpha$  we set the A.R. and R.R.



# Hypothesis Testing: An Example

## Accept and Reject Regions

- Given a sample and the specific value of the test statistic,  $\bar{x}$ :  $p\text{-value} = P_{H_0}(\bar{X} \leq \bar{x})$



## Hypothesis Testing: Remarks

### Power: $(1 - \beta)$

- Depending on the hypotheses the type II error ( $\beta$ ) can not be calculated:

$$\begin{cases} H_0 : \mu = 60 \\ H_1 : \mu \neq 60 \end{cases}$$

- In this case we do not know the value of  $\mu$  for  $H_1$  so we can not calculate the power  $(1 - \beta)$
- A good hypothesis test: given an  $\alpha$  the test maximises the power  $(1 - \beta)$

### Parametric test vs non-parametric test

# Hypothesis Testing in Supervised Classification

## Scenarios

- Two classifiers (algorithms) vs More than two
- One dataset vs More than one dataset
- Score
- Score estimation method known vs unknown
- The classifiers are trained and tested in the same datasets
- .....

# Testing Two Algorithms in a Dataset

## The General Approach

$$\left\{ \begin{array}{l} H_0 : \text{classifier } \phi \text{ has the same score value as} \\ \quad \text{classifier } \phi' \text{ in } \rho(\mathbf{x}, c) \\ \\ H_1 : \text{they have different values} \end{array} \right.$$

# Testing Two Algorithms in a Dataset

## The General Approach

$$\left\{ \begin{array}{l} H_0 : \text{classifier } \phi \text{ has the same score value as} \\ \quad \text{classifier } \phi' \text{ in } \rho(\mathbf{x}, c) \\ \\ H_1 : \text{they have different values} \end{array} \right.$$

$$\left\{ \begin{array}{l} H_0 : \text{algorithm } \phi \text{ has the same average score value as} \\ \quad \text{algorithm } \phi' \text{ in } \rho(\mathbf{x}, c) \\ \\ H_1 : \text{they have different values} \end{array} \right.$$

# Testing Two Algorithms in a Dataset

## An Ideal Context: We Can Sample $\rho(\mathbf{x}, c)$

- 1 Sample i.i.d.  $2n$  datasets from  $\rho(\mathbf{x}, c)$
- 2 Learn  $2n$  classifiers  $\phi_i^1, \phi_i^2$  for  $i = 1, \dots, n$
- 3 For each classifier obtain enough i.i.d. samples  $\{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_N, c_N)\}$  from  $\rho(\mathbf{x}, c)$
- 4 For each data set calculate the error of each algorithm in the test set

$$\epsilon_i^1 = \frac{1}{N} \sum_{j=1}^N \text{error}_i^1(\mathbf{x}_j) \quad \epsilon_i^2 = \frac{1}{N} \sum_{j=1}^N \text{error}_i^2(\mathbf{x}_j)$$

- 5 Calculate the average values over the  $n$  training datasets:

$$\bar{\epsilon}^1 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^1 \quad \bar{\epsilon}^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$$

## Testing Two Algorithms in a Dataset

### An Ideal Context: We Can Sample $\rho(\mathbf{x}, c)$

- Our test rejects the null hypothesis if  $|\bar{\epsilon}^1 - \bar{\epsilon}^2|$  (the statistic) is big
- Fortunately, by the central limit theorem:

$$\bar{\epsilon}^i \rightsquigarrow \mathcal{N}(\text{score}(\phi^i), \frac{\sigma_i^2}{N}) \quad i = 1, 2$$

- Therefore, under the null hypothesis (known  $\sigma_i^2$ ):

$$\hat{Z} = \frac{\bar{\epsilon}^1 - \bar{\epsilon}^2}{\sqrt{\frac{\sigma_1^2 + \sigma_2^2}{n}}} \rightsquigarrow \mathcal{N}(0, 1)$$

- ... and finally we reject  $H_0$  when  $|\hat{Z}| > z_{1-\alpha/2}$

## Testing Two Algorithms in a Dataset

### Properties of Our Ideal Framework

- Training datasets are **independent**
- Testing datasets are **independent**

### The Sad Reality

- We can not get i.i.d. **training** samples from  $\rho(\mathbf{x}, c)$
- We can not get i.i.d. **testing** samples from  $\rho(\mathbf{x}, c)$
- We have only one sample from  $\rho(\mathbf{x}, c)$

# Testing Two Algorithms in a Dataset

## McNemar Test (non-parametric)

- Compare two **classifiers** in a dataset after a Hold-Out process
- It is a paired non-parametric test

	$\phi^2$ error	$\phi^2$ ok
$\phi^1$ error	$n_{00}$	$n_{01}$
$\phi^1$ ok	$n_{10}$	$n_{11}$

- Under  $H_0$  we have  $n_{10} \approx n_{01}$  and the statistic

$$\frac{(|n_{01} - n_{10}| - 1)^2}{n_{01} + n_{10}}$$

follows a  $\chi^2$  distribution with 1 degree of freedom

- When  $n_{01} + n_{10}$  is small (<25) the binomial dist. can be used

# Testing Two Algorithms in a Dataset

## Tests Based on Resampling: Resampled t-test (parametric)

- The dataset is randomly divided  $n$  times in training and test
- Let  $\hat{\epsilon}_i$  be the difference between the performance of both algorithms in run  $i$  and  $\bar{\epsilon}$  the average. When it is assumed that  $\hat{\epsilon}_i$  are Gaussian and independent, under the null

$$t = \frac{\bar{\epsilon}}{\sqrt{\frac{1}{n} \frac{\sum_{i=1}^n (\hat{\epsilon}_i - \bar{\epsilon})^2}{n-1}}}$$

follows a  $t$  student distribution with  $n - 1$  degree of freedom

- **Caution:**
  - $\hat{\epsilon}_i$  are not Gaussian as  $\hat{\epsilon}_i^1$  and  $\hat{\epsilon}_i^2$  are not independent
  - $\hat{\epsilon}_i$  are not independent (overlap in training and testing)

## Testing Two Algorithms in a Dataset

### Resampled t-test Improved (Nadeau & Bengio, 2003)

- The variance in this case is too optimistic
- Two alternatives
  - Corrected resampled  $t$ :

$$t = \frac{\bar{\epsilon}}{\sqrt{\left(\frac{1}{n} + \frac{n_2}{n_1}\right) \frac{\sum_{i=1}^n (\hat{\epsilon}_i - \bar{\epsilon})^2}{n-1}}}$$

- Conservative  $Z$  (overestimation of the variance)

## Testing Two Algorithms in a Dataset

### t-test for k-fold Cross-validation

- It is similar to  $t$ -test for resampling
- In this case the testing datasets are independent
- The training datasets are still dependent

## Testing Two Algorithms in a Dataset

### 5x2 fold Cross-Validation (Dietterich 1998, Alpaydin 1999)

- Each cross-validation process has independent training and testing datasets
- The following statistic:

$$\frac{\sum_{i=1}^5 \sum_{j=1}^2 (\epsilon_i^{(j)})^2}{2 \sum_{i=1}^5 S_{\epsilon_i}^2}$$

follows a  $F$  distribution with 10 and 5 degrees of freedom under the null hypothesis

# Testing Two Algorithms in Several Datasets

## Initial Approaches

- Averaging Over Datasets
- Paired t-test

- $\epsilon^i = \epsilon_1^i - \epsilon_2^i$  and  $\bar{\epsilon} = \frac{1}{N} \sum_{i=1}^N \epsilon^i$

$$\frac{\bar{\epsilon}}{S_{\bar{\epsilon}}/\sqrt{N}} \rightsquigarrow t_{N-1}$$

## Problems

- Commensurability
- Outlier susceptibility
- (t-test) Gaussian assumption

# Testing Two Algorithms in Several Datasets

## Wilcoxon Signed-Ranks Test

- It is a non-parametric test that works as follows:
  - 1 Rank the module of the performance differences between both algorithms
  - 2 Calculate the sum of the ranks  $R^+$  and  $R^-$  where the first (resp. the second) algorithm outperforms the other
  - 3 Calculate  $T = \min(R^+, R^-)$
- For  $N \leq 25$  there are tables with critical values
- For  $N > 25$

$$z = \frac{T - \frac{1}{4}N(N+1)}{\sqrt{\frac{1}{24}N(N+1)(2N+1)}} \rightsquigarrow \mathcal{N}(0, 1)$$

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598		
Dataset2	0.599	0.591		
Dataset3	0.954	0.971		
Dataset4	0.628	0.661		
Dataset5	0.882	0.888		
Dataset6	0.936	0.931		
Dataset7	0.661	0.668		
Dataset8	0.583	0.583		
Dataset9	0.775	0.838		
Dataset10	1.000	1.000		

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591		
Dataset3	0.954	0.971		
Dataset4	0.628	0.661		
Dataset5	0.882	0.888		
Dataset6	0.936	0.931		
Dataset7	0.661	0.668		
Dataset8	0.583	0.583		
Dataset9	0.775	0.838		
Dataset10	1.000	1.000		

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971		
Dataset4	0.628	0.661		
Dataset5	0.882	0.888		
Dataset6	0.936	0.931		
Dataset7	0.661	0.668		
Dataset8	0.583	0.583		
Dataset9	0.775	0.838		
Dataset10	1.000	1.000		

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971	+0.017	
Dataset4	0.628	0.661	+0.033	
Dataset5	0.882	0.888	+0.006	
Dataset6	0.936	0.931	-0.005	
Dataset7	0.661	0.668	+0.007	
Dataset8	0.583	0.583	0.000	
Dataset9	0.775	0.838	+0.063	
Dataset10	1.000	1.000	0.000	

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971	+0.017	
Dataset4	0.628	0.661	+0.033	
Dataset5	0.882	0.888	+0.006	
Dataset6	0.936	0.931	-0.005	
Dataset7	0.661	0.668	+0.007	
Dataset8	0.583	0.583	0.000	
Dataset9	0.775	0.838	+0.063	
Dataset10	1.000	1.000	0.000	

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971	+0.017	
Dataset4	0.628	0.661	+0.033	
Dataset5	0.882	0.888	+0.006	
Dataset6	0.936	0.931	-0.005	
Dataset7	0.661	0.668	+0.007	
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	
Dataset10	1.000	1.000	0.000	1.5

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971	+0.017	
Dataset4	0.628	0.661	+0.033	
Dataset5	0.882	0.888	+0.006	
Dataset6	0.936	0.931	<b>-0.005</b>	
Dataset7	0.661	0.668	+0.007	
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	
Dataset10	1.000	1.000	0.000	1.5

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	
Dataset2	0.599	0.591	-0.008	
Dataset3	0.954	0.971	+0.017	
Dataset4	0.628	0.661	+0.033	
Dataset5	0.882	0.888	+0.006	
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	
Dataset10	1.000	1.000	0.000	1.5

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ =$$

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 7 + 8 + 4 + 5 + 9 + 1/2(1,5 + 1,5)$$

## Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 34.5$$

## Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 34.5 \quad R^- = 10 + 6 + 3 + 1/2(1,5 + 1,5)$$

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 34.5 \quad R^- = 20.5$$

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 34.5$$

$$R^- = 20.5$$

$$T = \min(R^+, R^-)$$

# Wilcoxon Signed-Ranks Test: Example

	$\phi^1$	$\phi^2$	diff	rank
Dataset1	0.763	0.598	-0.165	10
Dataset2	0.599	0.591	-0.008	6
Dataset3	0.954	0.971	+0.017	7
Dataset4	0.628	0.661	+0.033	8
Dataset5	0.882	0.888	+0.006	4
Dataset6	0.936	0.931	-0.005	3
Dataset7	0.661	0.668	+0.007	5
Dataset8	0.583	0.583	0.000	1.5
Dataset9	0.775	0.838	+0.063	9
Dataset10	1.000	1.000	0.000	1.5

$$R^+ = 34.5 \quad R^- = 20.5 \quad T = \min(R^+, R^-) = 20.5$$

# Testing Two Algorithms in Several Datasets

## Wilcoxon Signed-Ranks Test

- It also suffers from commensurability but only qualitatively
- When the assumptions of the  $t$  test are met, Wilcoxon is less powerful than  $t$  test

## Testing Two Algorithms in Several Datasets

### Signed Test

- It is a non-parametric test that counts the number of losses, ties and wins
- Under the null the number of wins follows a binomial distribution  $B(1/2, N)$
- For large values of  $N$  the number of wins follows  $\mathcal{N}(N/2, \sqrt{N}/2)$  under the null
- This test does not make any assumptions
- It is weaker than Wilcoxon

# Testing Several Algorithms in Several Datasets

Dataset (Demšar, 2006)

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79	0.84	0.89	0.72
$D_2$	0.57	0.88	0.88	0.79
$D_3$	0.71	0.87	0.88	0.62
$D_4$	0.65	0.81	0.69	0.72
$D_5$	0.89	0.89	0.91	0.67
$D_6$	0.65	0.63	0.98	0.55

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Testing all possible pairs of hypotheses  $\epsilon_{\phi^i} = \epsilon_{\phi^j} \quad \forall i, j$ .  
Multiple hypothesis testing
- Testing the hypothesis  $\epsilon_{\phi^1} = \epsilon_{\phi^2} = \dots = \epsilon_{\phi^k}$

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Testing all possible pairs of hypotheses  $\epsilon_{\phi^i} = \epsilon_{\phi^j} \quad \forall \quad i, j$ .  
Multiple hypothesis testing
- Testing the hypothesis  $\epsilon_{\phi^1} = \epsilon_{\phi^2} = \dots = \epsilon_{\phi^k}$

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Testing all possible pairs of hypotheses  $\epsilon_{\phi^i} = \epsilon_{\phi^j} \quad \forall \quad i, j$ .  
Multiple hypothesis testing
- Testing the hypothesis  $\epsilon_{\phi^1} = \epsilon_{\phi^2} = \dots = \epsilon_{\phi^k}$

## ANOVA vs Friedman

- *Repeated measures* ANOVA: Assumes Gaussianity and sphericity
- Friedman: Non-parametric test

# Testing Several Algorithms in Several Datasets

## Freidman Test

- 1 Rank the algorithms for each dataset separately (1-best). In case of ties assigned average ranks
- 2 Calculate the average rank  $R_j$  of each algorithm  $\phi^j$
- 3 The following statistic:

$$\chi_F^2 = \frac{12N}{k(k+1)} \left[ \sum_j R_j^2 - \frac{k(k+1)^2}{4} \right]$$

follows a  $\chi^2$  with  $k - 1$  degrees of freedom ( $N > 10$ ,  $k > 5$ )

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79 (3)	0.84 (2)	0.89 (1)	0.72 (4)
$D_2$	0.57 (4)	0.88 (1.5)	0.88 (1.5)	0.79 (3)
$D_3$	0.71 (3)	0.87 (2)	0.88 (1)	0.62 (4)
$D_4$	0.65 (4)	0.81 (1)	0.69 (3)	0.72 (2)
$D_5$	0.89 (2.5)	0.89 (2.5)	0.91 (1)	0.67 (4)
$D_6$	0.65 (2)	0.63 (3)	0.98 (1)	0.55 (4)
avr. rank	3.08	2	1.42	3.5

## Testing Several Algorithms in Several Datasets

## Friedman Test: Example

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79 (3)	0.84 (2)	0.89 (1)	0.72 (4)
$D_2$	0.57 (4)	0.88 (1.5)	0.88 (1.5)	0.79 (3)
$D_3$	0.71 (3)	0.87 (2)	0.88 (1)	0.62 (4)
$D_4$	0.65 (4)	0.81 (1)	0.69 (3)	0.72 (2)
$D_5$	0.89 (2.5)	0.89 (2.5)	0.91 (1)	0.67 (4)
$D_6$	0.65 (2)	0.63 (3)	0.98 (1)	0.55 (4)
avr. rank	3.08	2	1.42	3.5

$$\chi_F^2 = \frac{12N}{k(k+1)} \left[ \sum_j R_j^2 - \frac{k(k+1)^2}{4} \right] =$$

## Testing Several Algorithms in Several Datasets

## Friedman Test: Example

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79 (3)	0.84 (2)	0.89 (1)	0.72 (4)
$D_2$	0.57 (4)	0.88 (1.5)	0.88 (1.5)	0.79 (3)
$D_3$	0.71 (3)	0.87 (2)	0.88 (1)	0.62 (4)
$D_4$	0.65 (4)	0.81 (1)	0.69 (3)	0.72 (2)
$D_5$	0.89 (2.5)	0.89 (2.5)	0.91 (1)	0.67 (4)
$D_6$	0.65 (2)	0.63 (3)	0.98 (1)	0.55 (4)
avr. rank	3.08	2	1.42	3.5

$$\chi_F^2 = \frac{12N}{k(k+1)} \left[ \sum_j R_j^2 - \frac{k(k+1)^2}{4} \right] = 9,8$$

## Testing Several Algorithms in Several Datasets

### Iman & Davenport, 1980

- An improvement of Friedman test:

$$F_F = \frac{(N - 1)\chi_F^2}{N(k - 1) - \chi_F^2}$$

follows a F-distribution with  $k - 1$  and  $(k - 1)(N - 1)$  degrees of freedom

# Testing Several Algorithms in Several Datasets

## Post-hoc Tests

- Decision on the null hypothesis
- In case of rejection use of **post-hoc** tests to:
  - 1 Compare all pairs
  - 2 Compare all classifiers with a control

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Several related hypothesis simultaneously  $H_1, \dots, H_n$

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Several related hypothesis simultaneously  $H_1, \dots, H_n$

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Several related hypothesis simultaneously  $H_1, \dots, H_n$

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

- Family-wise error: Probability of rejecting at least one hypothesis assuming that ALL ARE TRUE

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Several related hypothesis simultaneously  $H_1, \dots, H_n$

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

- Family-wise error: Probability of rejecting at least one hypothesis assuming that ALL ARE TRUE
- False discovery rate

# Testing Several Algorithms in Several Datasets

## Multiple Hypothesis Testing

- Several related hypothesis simultaneously  $H_1, \dots, H_n$

	$H_0$ TRUE	$H_0$ FALSE
Decision: ACCEPT	✓	Type II error ( $\beta$ )
Decision: REJECT	Type I error ( $\alpha$ )	✓

- **Family-wise error: Probability of rejecting at least one hypothesis assuming that ALL ARE TRUE**
- False discovery rate

# Testing Several Algorithms in Several Datasets

## Designing Multiple Hypothesis Test

- Controlling family-wise error
- If each test  $H_i$  has a type I error  $\alpha$  then the family-wise error (FWE) in  $n$  tests is:

$$\begin{aligned} & P(\text{accept } H_1 \cap \text{accept } H_2 \cap \dots \cap \text{accept } H_n) \\ &= P(\text{accept } H_1) \times P(\text{accept } H_2) \times \dots \times P(\text{accept } H_n) \\ &= (1 - \alpha)^n \end{aligned}$$

and therefore

$$\text{FWE} = 1 - (1 - \alpha)^n \approx 1 - (1 - \alpha n) = \alpha n$$

- In order to have FWE  $\alpha$  we need to modify the threshold at each test

# Testing Several Algorithms in Several Datasets

## Comparing with a Control

- The statistic for comparing  $\phi^i$  and  $\phi^j$  is:

$$z = \frac{(R_i - R_j)}{\sqrt{\frac{k(k+1)}{6N}}} \rightsquigarrow \mathcal{N}(0, 1)$$

## Bonferroni-Dunn Test

- It is a one-step method
- Modify  $\alpha$  by taking into account the number of comparisons:

$$\frac{\alpha}{k - 1}$$

# Testing Several Algorithms in Several Datasets

## Comparing with a Control

- Methods based on ordered  $p$ -values
- The  $p$ -values are ordered  $p_1 \leq p_2 \leq \dots \leq p_{k-1}$

## Holm Method

- It is a step-down procedure
- Starting from  $p_1$  check the first  $i = 1, \dots, k - 1$  such that  $p_i > \alpha / (k - i)$
- The hypothesis  $H_1, \dots, H_{i-1}$  are rejected. The rest of hypotheses are kept

## Testing Several Algorithms in Several Datasets

Friedman Test: Example ( $\alpha = 0.015$ )

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79 (3)	0.84 (2)	0.89 (1)	0.72 (4)
$D_2$	0.57 (4)	0.88 (1.5)	0.88 (1.5)	0.79 (3)
$D_3$	0.71 (3)	0.87 (2)	0.88 (1)	0.62 (4)
$D_4$	0.65 (4)	0.81 (1)	0.69 (3)	0.72 (2)
$D_5$	0.89 (2.5)	0.89 (2.5)	0.91 (1)	0.67 (4)
$D_6$	0.65 (2)	0.63 (3)	0.98 (1)	0.55 (4)
avr. rank	3.08	2	1.42	3.5

## Testing Several Algorithms in Several Datasets

Friedman Test: Example ( $\alpha = 0.015$ )

	$\phi^1$	$\phi^2$	$\phi^3$	$\phi^4$
$D_1$	0.79 (3)	0.84 (2)	0.89 (1)	0.72 (4)
$D_2$	0.57 (4)	0.88 (1.5)	0.88 (1.5)	0.79 (3)
$D_3$	0.71 (3)	0.87 (2)	0.88 (1)	0.62 (4)
$D_4$	0.65 (4)	0.81 (1)	0.69 (3)	0.72 (2)
$D_5$	0.89 (2.5)	0.89 (2.5)	0.91 (1)	0.67 (4)
$D_6$	0.65 (2)	0.63 (3)	0.98 (1)	0.55 (4)
avr. rank	3.08	2	1.42	3.5

$$z = \frac{(R_i - R_j)}{\sqrt{\frac{k(k+1)}{6N}}}$$

# Testing Several Algorithms in Several Datasets

Friedman Test: Example ( $\alpha = 0.015$ )

$$z = \frac{(R_i - R_j)}{\sqrt{\frac{k(k+1)}{6N}}}$$

	$z$
$z_{14}$	-0.76
$z_{24}$	-2.7
$z_{34}$	-3.74

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value
$z_{14}$	-0.76	0.447
$z_{24}$	-2.7	0.007
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )
$z_{14}$	-0.76	0.447	0.005
$z_{24}$	-2.7	0.007	0.005
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )
$z_{14}$	-0.76	0.447	0.005
$z_{24}$	-2.7	0.007	0.005
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )	Holm ( $\alpha/(4 - i)$ )
$z_{14}$	-0.76	0.447	0.005	
$z_{24}$	-2.7	0.007	0.005	
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005	

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )	Holm ( $\alpha/(4 - i)$ )
$z_{14}$	-0.76	0.447	0.005	
$z_{24}$	-2.7	0.007	0.005	
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005	0.005

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )	Holm ( $\alpha/(4 - i)$ )
$z_{14}$	-0.76	0.447	0.005	
$z_{24}$	-2.7	0.007	0.005	0.0075
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005	0.005

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )	Holm ( $\alpha/(4 - i)$ )
$z_{14}$	-0.76	0.447	0.005	0.015
$z_{24}$	-2.7	0.007	0.005	0.0075
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005	0.005

# Testing Several Algorithms in Several Datasets

## Friedman Test: Example ( $\alpha = 0.015$ )

	$z$	$p$ -value	Bonferroni ( $\alpha/3$ )	Holm ( $\alpha/(4 - i)$ )
$z_{14}$	-0.76	0.447	0.005	0.015
$z_{24}$	-2.7	0.007	0.005	0.0075
$z_{34}$	-3.74	$1,8 \cdot 10^{-4}$	0.005	0.005

# Testing Several Algorithms in Several Datasets

## Hochberg Method

- It is a step-up procedure
- Starting with  $p_{k-1}$  check the first  $i = k - 1, \dots, 1$  such that  $p_i < \alpha / (k - i)$
- The hypothesis  $H_1, \dots, H_{i-1}$  are rejected. The rest of hypotheses are kept

## Hommel Method

- Find the largest  $j$  such that  $p_{n-j+k} > k\alpha/j$  for all  $k = 1, \dots, j$
- Reject all hypotheses  $i$  such that  $p_i \leq \alpha/j$

# Testing Several Algorithms in Several Datasets

## Comments on the Tests

- Holm, Hochberg and Hommel tests are more powerful than Bonferroni
- Hochberg and Hommel are based on Simes conjecture and can have a higher than  $\alpha$  FWE
- In practice Holm obtains very similar results to the other

# Testing Several Algorithms in Several Datasets

## All Pairwise Comparisons

- Differences with Comparing with a Control
- The all pairwise hypotheses are logically related: not all combinations of true and false hypotheses are possible

$\phi_1$  better than  $\phi_2$     and     $\phi_2$  better than  $\phi_3$

and     $\phi_1$  equal to  $\phi_3$

# Testing Several Algorithms in Several Datasets

## Shaffer Static Procedure

- It is a modification of Homl's procedure
- Starting from  $p_1$  check the first  $i = 1, \dots, k(k-1)/2$  such that  $p_i > \alpha/t_i$
- The hypothesis  $H_1, \dots, H_{i-1}$  are rejected. The rest of hypotheses are kept
- $t_i$  is the maximum number of hypotheses that can be true given that  $(i-1)$  are false
- It is a static procedure:  $t_i$  is determined given the hypotheses independently of the  $p$ -values

## Testing Several Algorithms in Several Datasets

### Shaffer Dynamic Procedure

- It is similar to the previous procedure but  $t_i$  is changed by  $t_i^*$
- $t_i^*$  considers the maximum number of hypotheses that can be true given that the previous  $(i - 1)$  hypotheses are false
- It is a dynamic procedure as  $t_i^*$  depends on the hypotheses already rejected
- It is more powerful than the Shaffer Static Procedure

# Testing Several Algorithms in Several Datasets

## Bregmann & Hommel

- More powerful alternative than Shaffer Dynamic Procedure
- Difficult implementation

## Remarks

- Adjusted p-values

## Conclusions

### Two Classifiers in a Dataset

- The complexity of the estimation of the scores makes it difficult to carry out good statistical testing

### Two Classifiers in Several Datasets

- Wilcoxon Signed-Ranks Test is a good choice
- In case of many datasets and to avoid the commensurability problem the Signed test could be used

# Conclusions

## Several Classifiers in Several Datasets

- Friedman or Iman & Davenport are required
- Post-hoc test more powerful than Bonferroni:
  - Comparison with a control: Holm method
  - All-to-all comparison: Shaffer Static method

## An Idea for Future Work

- To consider the variability of the score in each classifier and dataset

# Evaluación Honesta de Clasificadores en Clasificación Supervisada

Guzmán Santafé<sup>(1)</sup>, Iñaki Inza<sup>(2)</sup>

<sup>(1)</sup>Universidad Pública de Navarra

<sup>(2)</sup>Universidad del País Vasco

CAEPIA'11

7 de Noviembre, 2011