

# Machine Learning

## Master of Data Science

Alexandre Aussem

LIRIS UMR 5205 CNRS  
Data Mining & Machine Learning Group (DM2L)  
University of Lyon 1  
Web: [perso.univ-lyon1.fr/alexandre.aussem](http://perso.univ-lyon1.fr/alexandre.aussem)

September 18, 2023

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material



# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material



## 1 Introduction

### ■ Machine Learning: An overview

- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Problem categories: Supervised learning

A learning problem tries to predict properties of unknown data. In **supervised learning**, the data comes with additional attributes that we want to predict.

We can separate supervised learning problems in a few large categories:

- **Classification:** Identifying to which category an object belongs to.  
Samples belong to two or more classes and we want to learn from already labeled data how to predict the class of unlabeled data. Variants:  
*multi-class, multi-label classification, label ranking, collaborative filtering.*
- **Regression:** Predicting a continuous-valued attribute associated with an object. Variants: *multi-output learning* when the desired output consists of more continuous variables.

# Problem categories: Supervised learning

- **Structured (output) learning:** techniques that involve predicting structured objects, rather than scalar discrete or real values, e.g. translating a natural language sentence into a syntactic representation such as a parse tree, learning social networks, gene networks etc.
- **Probabilistic graphical models:** large class of structured prediction models. Learn the statistical relationships between the variables. Training and prediction using a trained model are often computationally infeasible and approximate inference is used.
- **Time series prediction and classification:** Techniques that involve assigning time series patterns to specific categories or predicting future values based on past observations. The data is a *sequence* of discrete-time observations, not *i.i.d.* anymore.

# Problem categories: Semi-supervised learning

- **Semi-supervised learning** is a class of supervised learning tasks and techniques that also make use of unlabeled data for training – typically a small amount of labeled data with a large amount of unlabeled data.
- The problem arises whenever the acquisition of labeled data requires a skilled human agent (e.g. to transcribe an audio segment) or a physical experiment (e.g. determining the 3D structure of a protein).
- The **cost associated with the labeling process is prohibitive**, whereas acquisition of unlabeled data is relatively inexpensive.

# Problem categories: Unsupervised learning

In **unsupervised learning**, the target values are missing. The goal may be to:

- Discover groups of similar examples within the data, known as **clustering**,
- Determine the distribution of data within the input space, known as **density estimation**.
- Project the data into a lower-dimensional space for the purpose of **visualization** or **dimensionality reduction**.

# Problem categories: Reinforcement learning

In **Reinforcement learning**, the goal is to find suitable actions (seen as missing targets) to take in a given situation in order to maximize the future rewards.

- The optimal actions are not given but must instead be discovered by **interacting with the environment in a sequence of states and actions.**
- Actions not only affects the immediate reward but also has an impact on the reward at subsequent time steps.

# Machine learning tasks and applications

- **Preprocessing:** Feature extraction and normalization. Bias correction and handling of missing values. Transforming input data such as text/image/video for use with machine learning algorithms.
- **Dimensionality reduction:** Reducing the number of random variables to consider for parsimony and increased efficiency.
- **Clustering:** Customer segmentation, Grouping experiment outcomes, ...
- **Classification & Regression:** Spam detection, image recognition, medical diagnosis ...
- **Structured learning:** Learning gene regulatory networks, natural language processing, speech recognition, and computer vision.

# Machine learning models and algorithms

- **Collaborative filtering:** Recommender systems collect preferences or taste information from many users (collaborating) for making automatic predictions (filtering) about the interests of each user.
- **Model selection:** Comparing, validating and choosing parameters and models for improved accuracy via parameter tuning.
- **Classification & Regression:** Decision Trees, ensemble methods (bagging, Random Forests, Boosting etc.) support vector machines (SVM), logistic regression, linear regression, naive Bayes classification, ridge regression, Lasso,  $k$  nearest neighbors, (deep) neural networks ...

# Machine learning models and algorithms

- **Clustering:** k-Means, spectral clustering, mean-shift, ...
- **Structured learning:** Bayesian networks and random fields, structured SVMs, Markov logic networks etc. Approximate inference and learning methods are usually needed.
- **Dimensionality reduction:** PCA, SVD, feature extraction and selection, non-negative matrix factorization, latent Dirichlet allocation (LDA) ...
- **Model selection:** Grid search, cross validation ...

# Choosing the right estimator

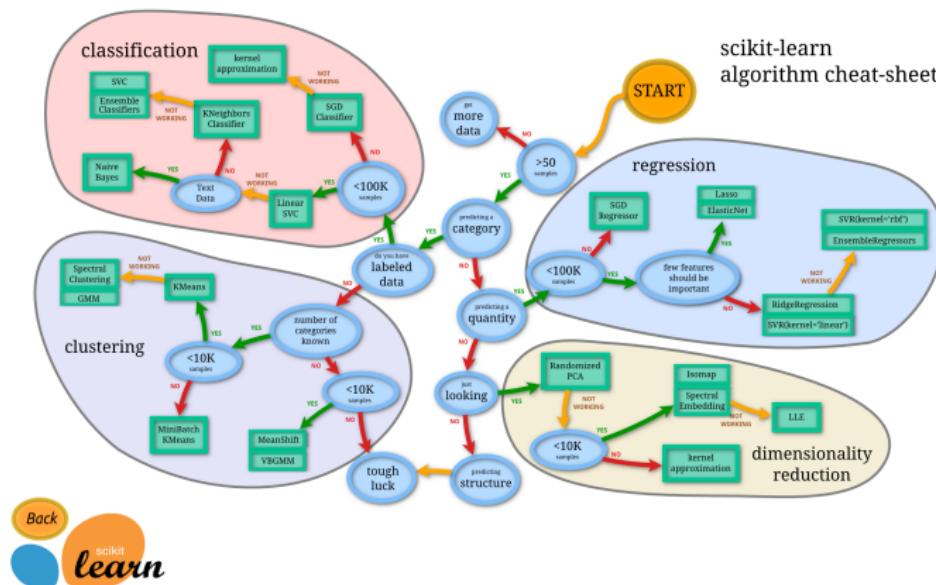


Figure: A road map on how to approach problems with regard to which estimators to try on your data.

# Model selection: illustration

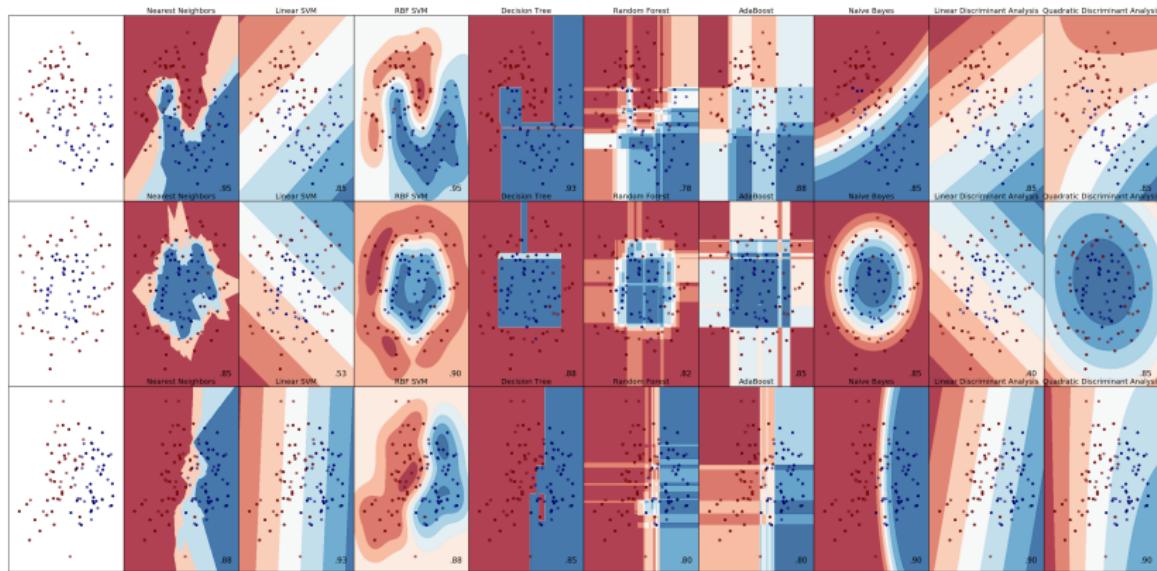


Figure: A comparison of several classifiers in scikit-learn on synthetic bivariate data.

# Model comparisons: illustration

APPROACH	ACC		AUC		RMS	
	$L$ IS TUNED	$L=200$	$L$ IS TUNED	$L=200$	$L$ IS TUNED	$L=200$
Ad	0.846±0.11	0.857±0.10	0.880±0.13	0.893±0.12	0.676±0.13	0.668±0.10
ADET	0.784±0.13	<b>0.862±0.09</b>	0.796±0.14	<b>0.898±0.12</b>	0.572±0.13	<b>0.667±0.09</b>
ADSt	0.847±0.12	0.833±0.11	0.898±0.12	0.874±0.13	0.603±0.11	0.598±0.08
ARCX4	0.866±0.09	0.852±0.09	<b>0.920±0.09</b>	0.892±0.11	<b>0.715±0.10</b>	0.686±0.10
ARCX4ET	0.866±0.09	0.868±0.08	<b>0.921±0.09</b>	0.901±0.10	0.715±0.10	0.693±0.09
BAG	<b>0.858±0.08</b>	0.823±0.10	<b>0.914±0.10</b>	0.875±0.12	<b>0.714±0.09</b>	0.660±0.10
BAGET	<b>0.865±0.10</b>	0.836±0.11	<b>0.916±0.11</b>	0.893±0.11	<b>0.717±0.10</b>	0.673±0.10
LOGB	0.848±0.10	0.845±0.10	0.880±0.12	0.884±0.13	<b>0.679±0.14</b>	0.635±0.09
RF	0.865±0.09	0.864±0.09	<b>0.915±0.11</b>	0.896±0.12	<b>0.716±0.10</b>	0.689±0.10
RADP	0.859±0.08	0.850±0.09	<b>0.915±0.09</b>	0.889±0.13	<b>0.714±0.09</b>	0.669±0.09
RADPET	0.864±0.10	0.861±0.09	0.915±0.11	0.908±0.10	<b>0.716±0.10</b>	0.680±0.09
ROT	0.864±0.09	0.865±0.08	0.916±0.10	0.903±0.11	<b>0.722±0.10</b>	0.700±0.10
ROTB	0.862±0.09	0.865±0.09	0.913±0.10	0.897±0.11	0.719±0.10	0.702±0.11
ROTBET	0.864±0.09	0.866±0.09	0.913±0.10	0.900±0.11	0.719±0.10	0.704±0.11
ROTE <sup>T</sup>	0.864±0.09	0.871±0.08	0.915±0.10	0.901±0.10	<b>0.721±0.10</b>	0.698±0.10
SWT	0.851±0.10	0.859±0.09	0.899±0.10	0.888±0.11	<b>0.692±0.08</b>	0.638±0.07
SWTET	0.864±0.10	0.866±0.08	0.913±0.11	0.890±0.11	<b>0.699±0.09</b>	0.649±0.08
VAD	0.812±0.10	<b>0.858±0.09</b>	0.817±0.13	<b>0.894±0.12</b>	0.628±0.14	<b>0.684±0.11</b>
VADET	0.791±0.13	<b>0.864±0.08</b>	0.792±0.15	<b>0.899±0.11</b>	0.601±0.15	<b>0.681±0.09</b>

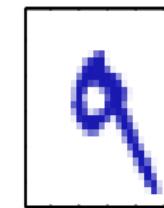
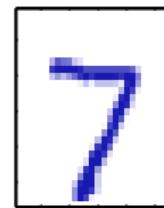
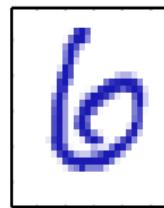
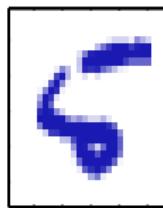
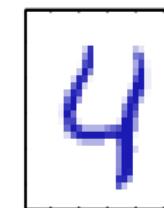
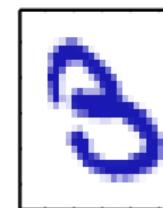
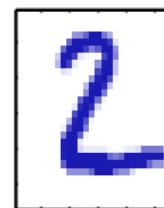
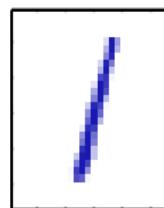
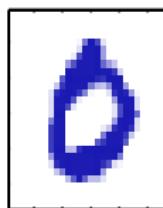
# Terminology

Common terms, synonyms or closely related topics:

- Machine Learning
- Statistical Learning
- Computational Learning Theory
- Knowledge Discovery in Databases
- Pattern Recognition
- Data Mining
- Artificial Intelligence

≠ Cognitive science ...

# Illustration: Handwritten digit recognition



## Illustration: Handwritten digit recognition

- Each digit corresponds to a  $28 \times 28$  pixel image, represented by a vector comprising 784 real numbers. The goal is to build a machine that will take such a vector  $x$  as input and that will produce the identity of the digit  $0, \dots, 9$  as the output.
- This is a nontrivial problem due to the **wide variability of handwriting**.
- The original input variables are typically preprocessed (feature extraction) to transform them into some new space of variables
- The precise form of the function  $f(x)$  is determined during the training phase,

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Image annotation/retrieval

Target 1: cloud yes/no

Target 2: sky yes/no

Target 3: tree yes/no

...

...

...



# Multi-label learning

- Training data:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ,  $y_i \in \mathcal{Y} = \{0, 1\}^m$
- **Predict** the vector  $y = (y_1, y_2, \dots, y_m)$  for a given  $x$ .

	$X_1$	$X_2$	$Y_1$	$Y_2$	...	$Y_m$
$x_1$	5.0	4.5	1	1		0
$x_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$x_n$	3.0	3.5	0	1		1
$x$	4.0	2.5	?	?		?

# Multivariate prediction

- Training data:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ,  $y_i \in \mathcal{Y} = \mathbb{R}^m$
- **Predict** the vector  $y = (y_1, y_2, \dots, y_m)$  for a given  $x$ .

	$X_1$	$X_2$	$Y_1$	$Y_2$	...	$Y_m$
$x_1$	5.0	4.5	14	0.3		9
$x_2$	2.0	2.5	15	1.1		4.5
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$x_n$	3.0	3.5	19	0.9		2
$x$	4.0	2.5	?	?		?

# Label ranking

- Training data:  $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ , where  $\mathbf{y}_i$  is a ranking (permutation) of a fixed number of labels/alternatives.
- **Predict** permutation  $(y_{\pi(1)}, y_{\pi(2)}, \dots, y_{\pi(m)})$  for a given  $\mathbf{x}$ .

---

					
	$X_1$	$X_2$	$Y_1$	$Y_2$	$Y_m$
$\mathbf{x}_1$	5.0	4.5	1	3	2
$\mathbf{x}_2$	2.0	2.5	2	1	3
$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	3	1	2
$\mathbf{x}$	4.0	2.5	?	?	?

---

# Multi-task learning

- Training data:  $\{(\mathbf{x}_{1j}, y_{1j}), (\mathbf{x}_{2j}, y_{2j}), \dots, (\mathbf{x}_{nj}, y_{nj})\}$ ,  $j = 1, \dots, m$ ,  
 $y_{ij} \in \mathcal{Y} = \mathbb{R}$ .
- **Predict**  $y_j$  for a given  $\mathbf{x}_j$ .

	$X_1$	$X_2$	$Y_1$	$Y_2$	...	$Y_m$
$x_1$	5.0	4.5	14			9
$x_2$	2.0	2.5		1.1		
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$x_n$	3.0	3.5				2
$x$	4.0	2.5				?

# Collaborative filtering

- Training data:  $\{(u_i, m_j, y_{ij})\}$ , for some  $i = 1, \dots, n$  and  $j = 1, \dots, m$ ,  $y_{ij} \in \mathcal{Y} = \mathbb{R}$ .
- **Predict**  $y_{ij}$  for a given  $u_i$  and  $m_j$ .

	$m_1$	$m_2$	$m_3$	$\cdots$	$m_m$
$u_1$	1			$\cdots$	4
$u_2$	3		1	$\cdots$	
$u_3$		2	5	$\cdots$	
$\cdots$				$\cdots$	
$u_n$		2		$\cdots$	1

# Dyadic prediction

							4	5	...	7		8	6
							10	14	...	9		21	12
instances			$y_1$	$y_2$	...	$y_m$		$y_{m+1}$	$y_{m+2}$				
1	1	$x_1$	10	?	...	1		?	?				
3	5	$x_2$		0.1	...	0					?		
7	0	$x_3$	?	?	...	1		?					
1	1	...		...	...	0					?		
3	1	$x_n$		0.9	...	1		?	?				
<hr/>													
2	3	$x_{n+1}$	?		...	?					?		
3	1	$x_{n+2}$		?	...	?		?	?				

# Further problems

- Structured output prediction
- Multi-task learning and transfer learning
- Matrix factorization
- Sequence learning, time series prediction and data stream mining
- Metric learning
- Topic Modelling
- Causal inference
- ...

# Softwares

Name	Advantages	Disadvantages	Open Source
R	Library support and Visualization	Steep learning curve	Yes
Matlab	Native matrix support, Visualization	Expensive, incomplete statistics support	No
Scientific Python	Ease and Simplicity	Heavy development	Yes
Excel	Easy, Visual, Flexible	Large datasets	No
SAS	Large Datasets	Expensive, outdated programming language	No
Stata	Easy Statistical Analysis		No
SPSS	Like Stata but more expensive and less flexible		

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Supervised learning

- In the setting of supervised learning, a function of  $f(., \mathbf{w}) : X \rightarrow Y$  is to be learned, that predicts well on instances that are drawn from a **joint probability distribution**  $p(\mathbf{x}, y)$  on  $X \times Y$ .
- The true distribution  $p(\mathbf{x}, y)$  is **unknown**. Instead, one has access to a training set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .
- In statistical learning models, the training sample  $(\mathbf{x}_i, y_i)$  are assumed to have been drawn i.i.d. from  $p(\mathbf{x}, y)$ .

# Generalization error

- Given a loss function  $L(f(\mathbf{x}, \mathbf{w}), y) : X \times X \rightarrow \mathbb{R}$  that measures the difference between the predicted value  $f(\mathbf{x}, \mathbf{w})$  and the true value  $y$ . The objective is to minimize the expected "risk":

**True loss, generalization error or expected risk :**

$$E(\mathbf{w}) = \mathbb{E}[L(f(\mathbf{X}, \mathbf{w}), Y)] = \int \int L(f(\mathbf{x}, \mathbf{w}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

- The True loss is expressed as an expectation over the unknown **joint probability distribution**  $p(\mathbf{x}, y)$ .

# Supervised learning

- The ideal goal is to select a function  $f(., \mathbf{w}) \in \mathcal{H}$ , where  $\mathcal{H}$  is a space of functions called a **hypothesis space**, so that the true loss is minimised.
- As  $p(\mathbf{x}, y)$  is unknown, a common paradigm is to estimate a function  $\hat{f}$  through (regularized) **empirical risk minimization**.
- The loss function may be chosen for numerical reasons (e.g. convexity).
- Distinct loss functions may lead to functions  $\hat{f}$  that differ significantly in their predictions.

# Standard loss functions

- In **regression**,  $L(\cdot)$  is typically the **squared loss**:

$$\mathbb{E}[L(f(X, \mathbf{w}), Y)] = \int \int (f(\mathbf{x}, \mathbf{w}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

- In **classification**,  $L(\cdot)$  is typically the **zero-one loss**:

$$\begin{aligned}\mathbb{E}[L(f(X, \mathbf{w}), Y)] &= P_{(X, Y)}[f(X, \mathbf{w}) \neq Y] \\ &= \sum_y \int \mathbb{I}[f(\mathbf{x}, \mathbf{w}) \neq y] p(\mathbf{x}, y) d\mathbf{x}\end{aligned}$$

- In **probabilistic classification**,  $f((X, \mathbf{w}), y)$  is an estimate for  $P(y|\mathbf{x})$  and  $L(\cdot)$  is typically the **logarithmic loss**:

$$\mathbb{E}[L(f(X, \mathbf{w}), Y)] = - \int \int \log(f((\mathbf{x}, \mathbf{w}), y)) p(\mathbf{x}, y) d\mathbf{x} dy$$



# MAP estimate in classification

- In classification,  $L(f(\mathbf{x}), y) = \mathbb{I}(f(\mathbf{x}) \neq y)$  called the **0-1 loss**. ( $\mathbf{w}$  is omitted for conciseness). It is convenient to denote by  $C_k$  the class  $k$ . So

$$\arg \min_{f \in \mathcal{H}} \sum_y \int \mathbb{I}[f(\mathbf{x}, \mathbf{w}) \neq y] p(\mathbf{x}, y) dy$$

$$= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^K \int \mathbb{I}(f(\mathbf{x}) \neq C_k) p(\mathbf{x}, C_k) d\mathbf{x}$$

$$= \arg \min_{f \in \mathcal{H}} \int \left\{ \sum_{i=1}^K \mathbb{I}(f(\mathbf{x}) \neq C_k) p(C_k | \mathbf{x}) \right\} p(\mathbf{x}) d\mathbf{x}$$

- For all  $\mathbf{x}$ , the optimal solution  $f(\mathbf{x})$  satisfies

$$\arg \min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^K \mathbb{I}(f(\mathbf{x}) \neq C_k) p(C_k | \mathbf{x}) \right\}$$

# MAP estimate in classification

- When  $L(f(\mathbf{x}), y)$  is **0-1 loss**, we can prove that MAP estimate minimizes 0-1 loss,

$$\begin{aligned}\arg \min_{f \in \mathcal{H}} \sum_{i=1}^K L[f(\mathbf{x}_i), C_k] P(C_k | \mathbf{x}_i) &= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^K \mathbb{I}(f(\mathbf{x}_i) \neq C_k) P(C_k | \mathbf{x}_i) \\ &= \arg \min_{f \in \mathcal{H}} [1 - P(f(\mathbf{x}) = C_k | \mathbf{x})] \\ &= \arg \max_{f \in \mathcal{H}} P(f(\mathbf{x}) = C_k | \mathbf{x})\end{aligned}$$

- The optimal solution for the 0-1 loss is the most likely class (the mode of the distribution).

# Posterior mean minimizes quadratic loss

- For continuous outputs, a more appropriate loss function is **squared error** or **quadratic loss**.
- The posterior expected loss is given by

$$\begin{aligned}\int_y L[f(\mathbf{x}), y] p(y|\mathbf{x}) dy &= \int_y [y - f(\mathbf{x})]^2 p(y|\mathbf{x}) dy \\ &= \int_y \left[ y^2 - 2yf(\mathbf{x}) + f(\mathbf{x})^2 \right] p(y|\mathbf{x}) dy\end{aligned}$$

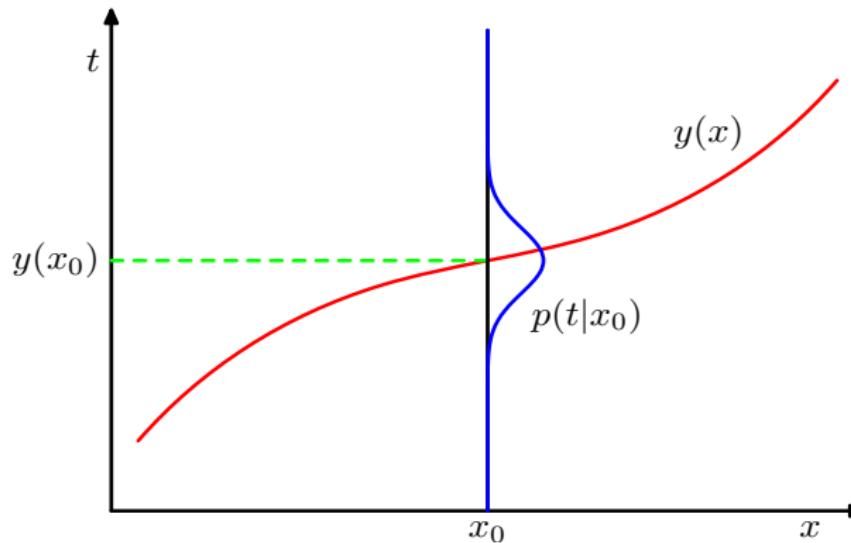
# Posterior mean minimizes quadratic loss

- Setting the derivative w.r.t.  $f$  equal to zero gives

$$\begin{aligned}\frac{\partial \rho}{\partial f} &= \int_y [-2y + 2f(\mathbf{x})] p(y|\mathbf{x}) dy = 0 \Rightarrow \\ \int_y f(\mathbf{x}) p(y|\mathbf{x}) dy &= \int_y y p(y|\mathbf{x}) dy \\ f(\mathbf{x}) \int_y p(y|\mathbf{x}) dy &= \mathbb{E}_{p(y|\mathbf{x})}[Y] = \mathbb{E}[Y|\mathbf{x}] \\ \Rightarrow f(\mathbf{x}) &= \mathbb{E}[Y|\mathbf{x}]\end{aligned}$$

- Hence the the **minimum mean squared error** solution or **MMSE** is the **posterior mean**.

# Squared loss



The regression function  $y(x)$ , which minimizes the expected squared loss, is given by the mean of the conditional distribution  $p(t|x)$ .

# Loss functions

- There are situations in which squared loss can lead to very poor results in regression and where more sophisticated approaches are needed. An important example concerns situations in which the conditional distribution  $p(y|x)$  is multimodal.
- A simple generalization of the squared loss, called the *Minkowski loss*, is given by

$$\mathbb{E}[L(f(\mathbf{X}, \mathbf{w}), Y)] = \int \int |f(\mathbf{x}, \mathbf{w}) - y|^q p(\mathbf{x}, y) d\mathbf{x} dy$$

- Is there a closed-form expression for the optimal solution  $f^*(\cdot)$ ?

# Loss functions

$$\mathbb{E}[L(f(\mathbf{X}, \mathbf{w}), Y)] = \int \int |f(\mathbf{x}, \mathbf{w}) - y|^q p(\mathbf{x}, y) d\mathbf{x} dy$$

One can show that the optimal solution,  $\hat{f}(\mathbf{x})$ , reduces to

- The conditional **mean**,  $\mathbb{E}[Y|\mathbf{x}]$ , for  $q = 2$  (i.e., quadratic loss function)
- The conditional **median** of  $p(y|\mathbf{x})$  for  $q = 1$  (i.e., absolute loss function)
- The conditional **mode** of  $p(y|\mathbf{x})$  for  $q \rightarrow 0$  (i.e., zero-one loss function)

Is there  $\mathbf{w}^*$  such that  $f(\cdot, \mathbf{w}^*) = \hat{f}(\cdot)$ ? Otherwise the model will be biased.

# Posterior median minimizes the absolute loss

- A more robust alternative to the quadratic loss is the absolute loss. The optimal estimate is the **posterior median**,

$$\begin{aligned}\rho &= \int_L[f(\mathbf{x}), y] p(y|\mathbf{x}) dy = \int_y |y - f(\mathbf{x})| p(y|\mathbf{x}) dy \\ &= \int_{y < f(\mathbf{x})} (f(\mathbf{x}) - y) p(y|\mathbf{x}) dy + \int_{y \geq f(\mathbf{x})} (y - f(\mathbf{x})) p(y|\mathbf{x}) dy\end{aligned}$$

$$\begin{aligned}\frac{\partial \rho}{\partial f} &= \int_{y < f(\mathbf{x})} p(y|\mathbf{x}) dy - \int_{y \geq f(\mathbf{x})} p(y|\mathbf{x}) dy = 0 \\ \Rightarrow P(Y < f(\mathbf{x})|\mathbf{x}) &= P(Y \geq f(\mathbf{x})|\mathbf{x}) = 0.5\end{aligned}$$

# Loss function and consistency

- The loss function should approximate the actual cost we are paying due to misclassification errors.
- There are many loss functions that are **non-convex and discontinuous**. It is difficult to optimize these losses directly (e.g. gradient-based optimization).
- The loss functions used for training classifiers (e.g. exponential loss, hinge loss) are usually approximations or **surrogates** of the zero-one-loss (misclassification rate).
- In practice, we consider a surrogate loss function which can be optimized by efficient algorithms.
- A learning algorithm is **consistent** if the expected risk of a learned function converges to the Bayes risk as the training sample size increases.

# More complex loss functions with multiple labels

- So far, each instance was associated with a single label  $Y$ . In real-world applications, one object is usually relevant to multiple labels  $\{Y_1, \dots, Y_L\}$ .
- Ranking these labels is a central part of many **information retrieval** problems, such as document retrieval, collaborative filtering, sentiment analysis, and online advertising.
- The **rank-loss** evaluates the label order induced by giving a numerical score for each item in some sense:

$$L_{rankloss}(f(\mathbf{x}), \mathbf{y}) = \frac{1}{n^+ \cdot n^-} \sum_{y_i \leq y_j} \mathbb{I}[f_i(\mathbf{x}) \geq f_j(\mathbf{x})]$$

where  $n^+ = \{y_i : y_i = +1\}$  and  $n^- = \{y_j : y_j = -1\}$ .

- It is difficult to find surrogate loss consistent with the ranking loss.

# Empirical loss

- The true loss  $E(\mathbf{w})$  can not be computed exactly as the distribution  $p(\mathbf{x}, y)$  is **unkown**. We may estimate the expected risk from  $\mathcal{D}$  using the approximation,

$$\mathbb{E}[f(\mathbf{X})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \simeq \frac{1}{n} \sum_{j=1}^n f(\mathbf{x}_j)$$

- So, the intent is to find a function  $f(., \mathbf{w})$  that minimizes the empirical loss,

## Empirical loss:

$$\hat{E}(\mathbf{w}) = \frac{1}{n} \sum_{j=1}^n L(f(\mathbf{x}_j, \mathbf{w}), y_j)$$

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- **The curse of dimensionality**
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

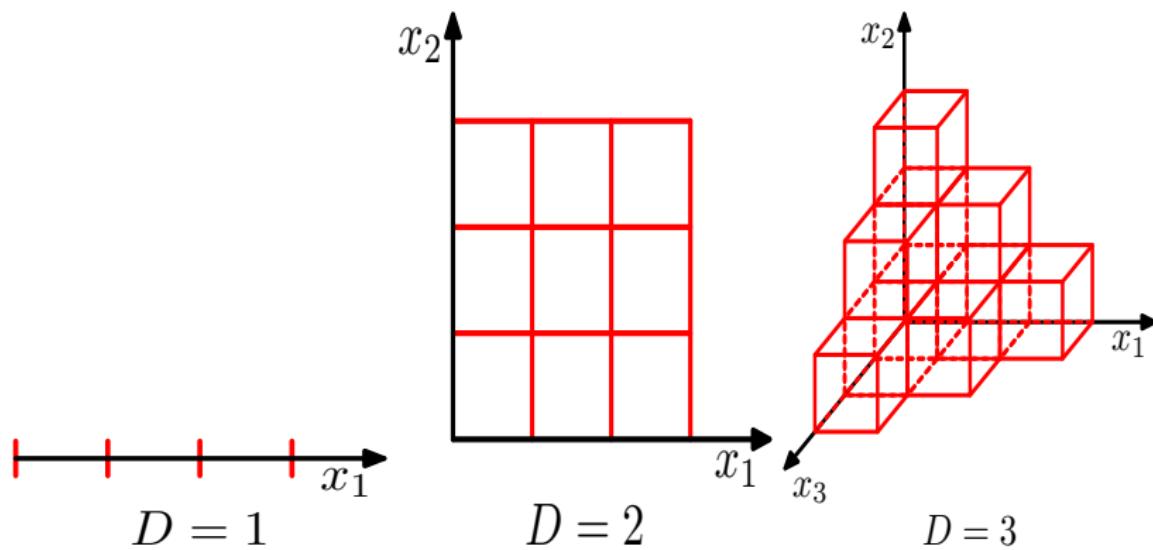
## 6 Neural networks

## 7 References & Further Material

# The curse of dimensionality

- In practical applications, we have to deal with spaces of high dimensionality comprising many input variables
- Suppose the input space is divided into cells and any new test point is assigned to the class that has a majority number of representatives in the same cell as the test point.

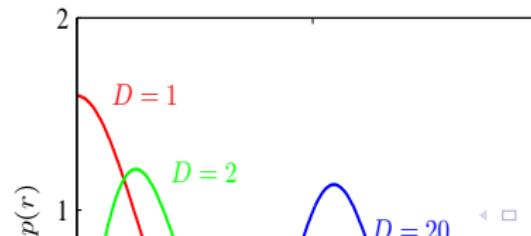
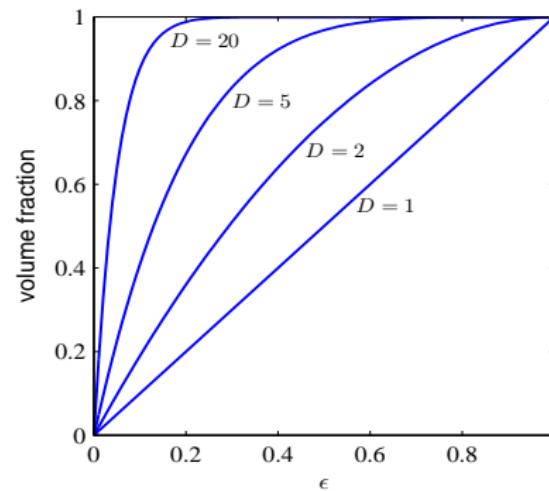
# The curse of dimensionality



# The curse of dimensionality

- As the number of such cells grows exponentially with  $D$ , we need an exponentially large quantity of training data in order to ensure that the cells are not empty
- For a polynomial of order  $M$ , the number of coefficients to estimate varies as  $O(D^M)$ . More **parsimonious** models are needed.
- Our geometrical intuitions fail badly when we consider spaces of higher dimensionality. Learning algorithms based on the euclidean metric suffer drastically from the curse of dimensionality.

# The curse of dimensionality



# The curse of dimensionality

- With high-dimensional data, the **Euclidean distance is not informative** because all vectors are almost equidistant to the search query vector.
- Local methods work best in low-dimensional embeddings.
- Example: The k-nearest neighbors algorithm (k-NN) is a non-parametric method used for classification and regression (majority vote of its  $k$  neighbors or average of their values) where the function is only approximated locally.
- Dimension reduction is usually performed prior to applying the k-NN algorithm to high-dimensional data ( $d > 10$ ).

# Binary classification

- Many ML problems can be cast as **binary classification problems**.
- The goal is to infer a function  $f : X \rightarrow \{-1, +1\}$  from  $\mathcal{D}$  such that the expected risk,  $P_{(X,Y)}[f(X, \mathbf{w}) \neq y]$ , is as low as possible.
- The model output value is often regarded as **score** of the positive class that is compared to some arbitrary **threshold**  $\theta$ .
- A **probabilistic** classifier output an estimate of  $P(Y = 1|x)$ , but **calibration** of the probabilities is often required when the loss is expressed as an expected cost.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- **Polynomial curve fitting**
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

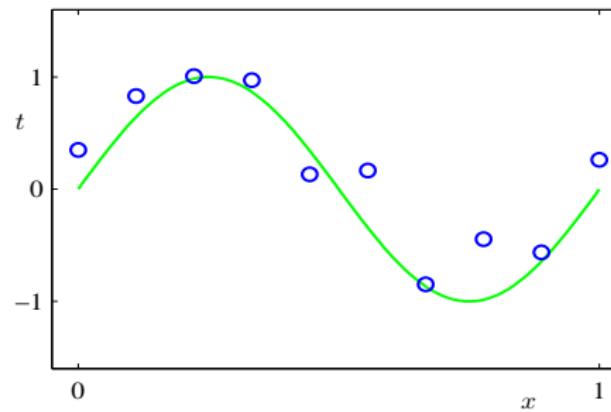
## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

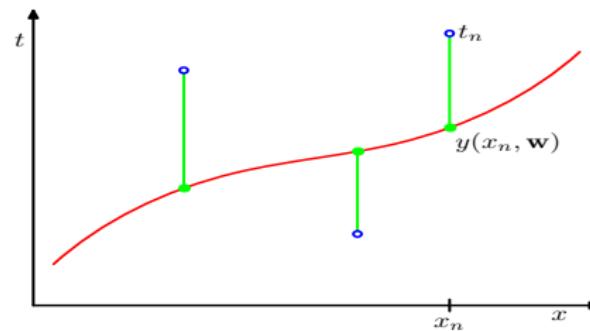
# Example: Polynomial curve fitting



$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_j x^j$$



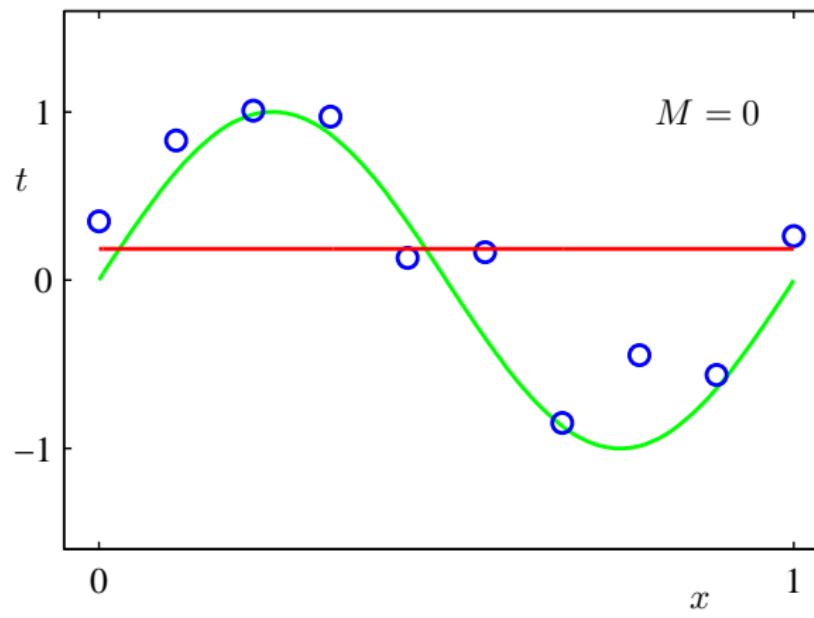
# Sum of the squares



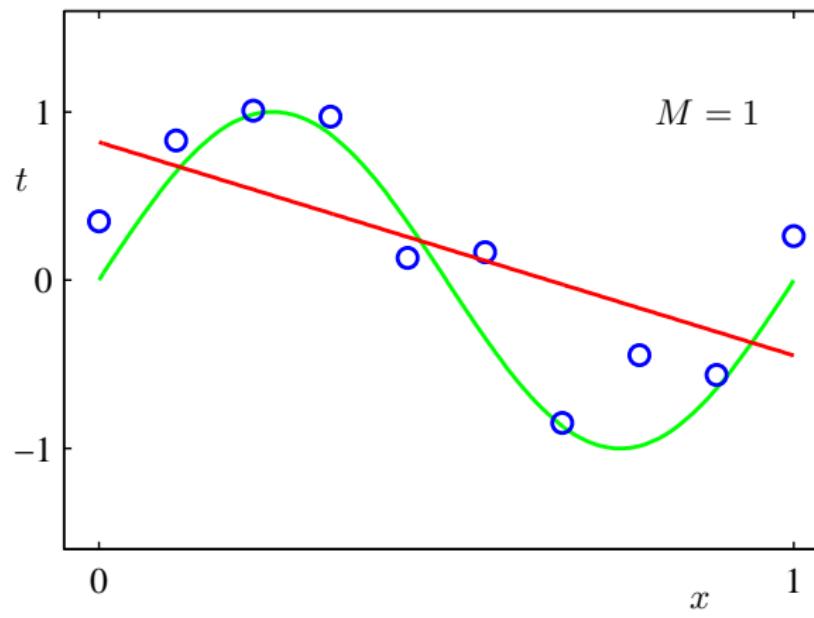
The sum of the squares of the errors measures the misfit between the function  $y(x, \mathbf{w})$ , for any given value of  $\mathbf{w}$ , and the training set data points,

$$\hat{E}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \{f(x_n, \mathbf{w}) - y_n\}^2$$

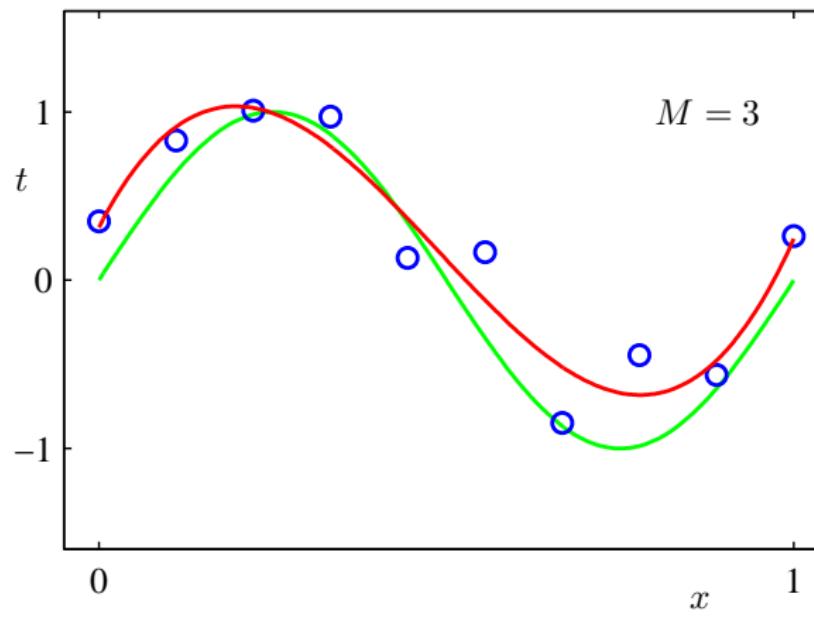
# Polynomial of order 0

 $M = 0$

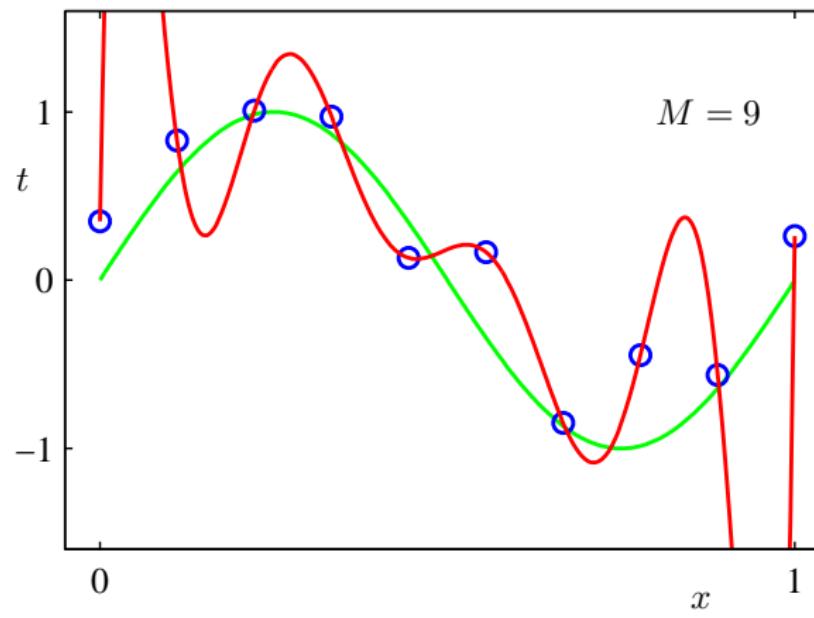
# Polynomial of order 1



# Polynomial of order 3



# Polynomial of order 9

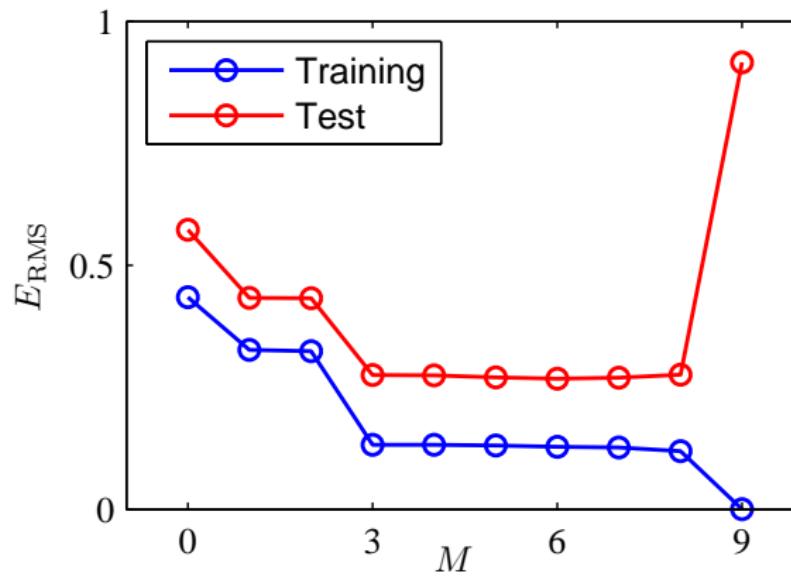


# Polynomial of order 9

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

We observe that, as  $M$  increases, the magnitude of the coefficients typically gets larger, i.e.  $|w_j| \rightarrow \infty$  when  $M \rightarrow \infty$ .

# Over-fitting



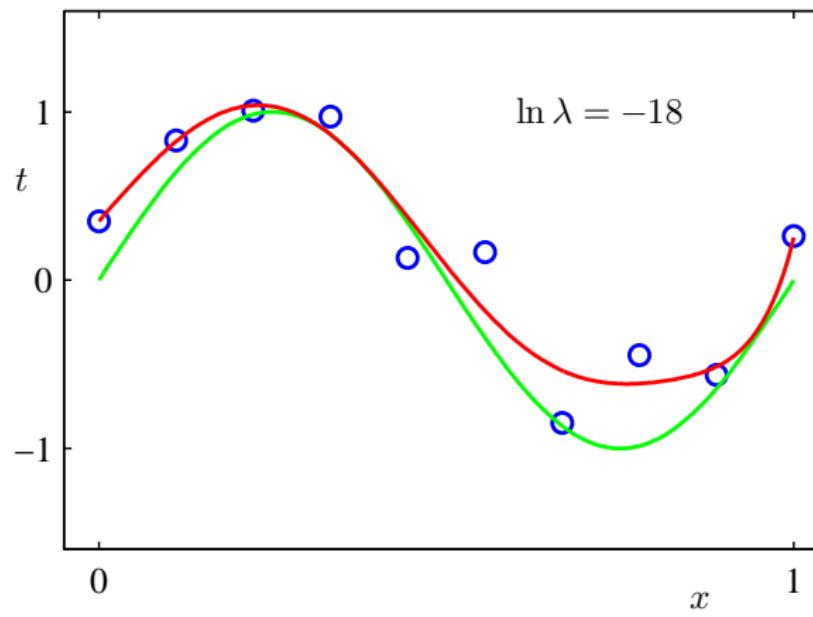
# Regularization

- Regularization is one technique often used to control the over-fitting phenomenon. It involves adding a penalty term to the error function to discourage the coefficients from reaching large values,

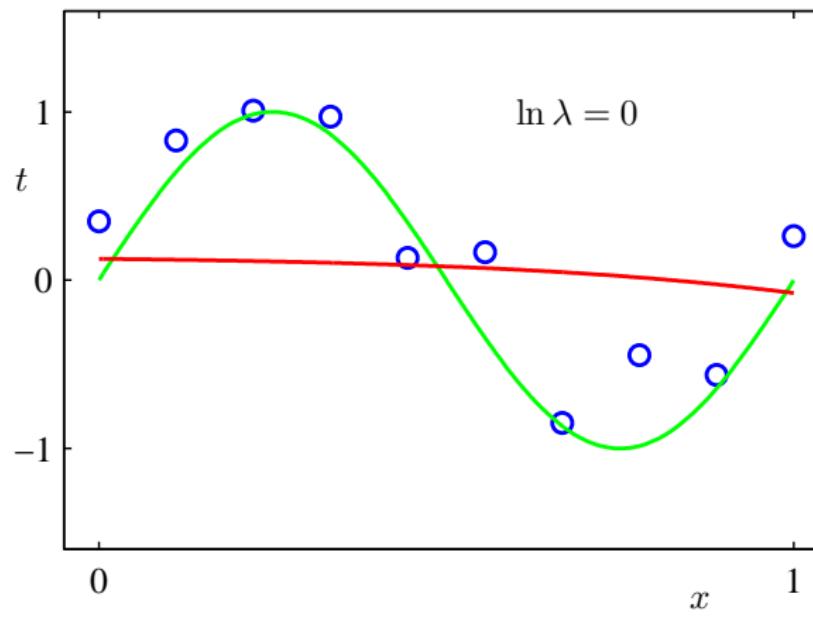
$$\hat{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^M \{f(x_n, \mathbf{w}) - y_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- $\lambda$  controls the effective complexity of the model and hence determines the degree of over-fitting.

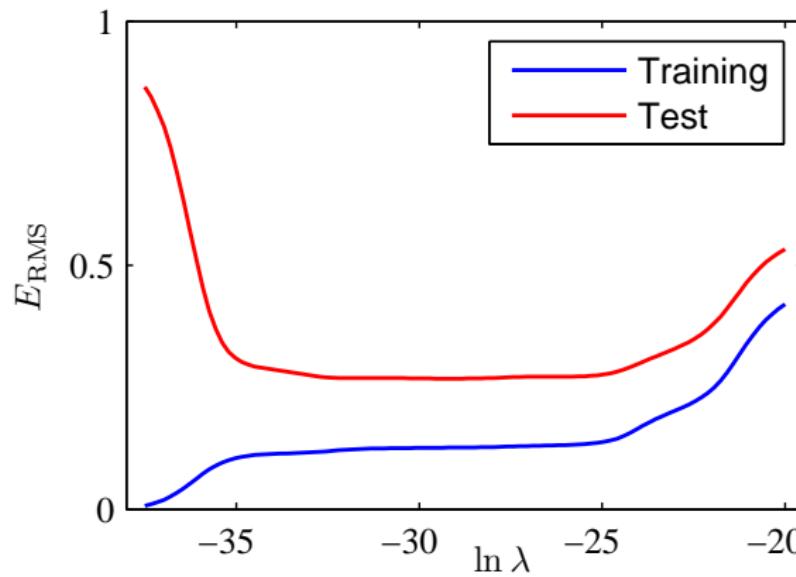
# Regularization $\ln \lambda = -18$



# Regularization $\ln \lambda = 0$



# Regularization : $E_{RMLS}$ vs. $\ln \lambda$



# Polynomial of order 9

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01

We see that, as  $\lambda$  increases, the magnitude of the coefficients typically gets larger.

# Bayes optimal prediction

- In regression, we typically assume that  $y = f(x) + \epsilon$  where  $\epsilon$  is a white noise with variance  $\sigma^2$ . the expected squared loss can be written in the form:

$$\begin{aligned}\mathbb{E}[\mathbf{w}] &= \int \int (f(x, \mathbf{w}) - y)^2 f_{X,Y}(x, y) dx dy \\ &= \int \int (f(x, \mathbf{w}) - \mathbb{E}[Y|x])^2 f_{X,Y}(x, y) dx dy \\ &+ \int \left\{ \int (\mathbb{E}[Y|x] - y)^2 f_{Y|X}(y) dy \right\} f_X(x) dx\end{aligned}$$

where the conditional expectation  $\mathbb{E}[Y|x] = \int yp(y|x))dy$ , is the **Bayes optimal prediction** for the squared loss function,

- The last term, equal to  $\sigma^2$ , is the **intrinsic noise** on the data and represents the minimum achievable value of the expected loss.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- **The Bias-Variance decomposition**

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# The Bias-Variance decomposition

- $f(x, \mathbf{w})$  is dependent of  $\mathcal{D}$ . Taking the expectation of this expression with respect to  $\mathcal{D}$ ,

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[(f(x, \mathbf{w}) - \mathbb{E}[Y|x])^2] &= (\mathbb{E}_{\mathcal{D}}[f(x, \mathbf{w})] - \mathbb{E}[Y|x])^2 \\ &\quad + \mathbb{E}_{\mathcal{D}}[(f(x, \mathbf{w}) - \mathbb{E}_{\mathcal{D}}[f(x, \mathbf{w})])^2]\end{aligned}$$

we get the following decomposition of the expected squared loss:

$$\text{Error} = (\text{bias})^2 + \text{variance} + \text{noise}$$

- The model with the optimal predictive capability is the one that leads to the best balance between bias and variance.

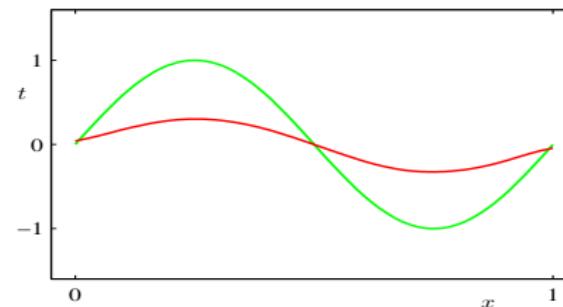
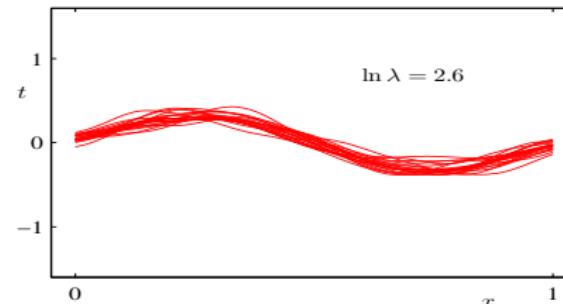
# Regularization

- Regularization is one technique often used to control the over-fitting phenomenon. It involves adding a penalty term to the error function to discourage the coefficients from reaching large values,

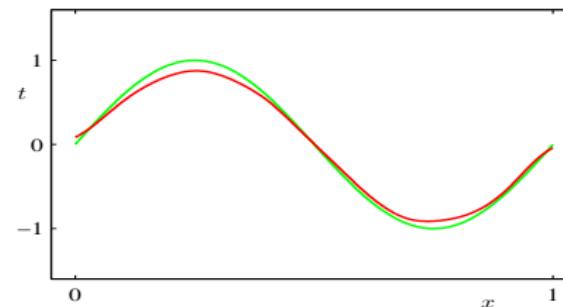
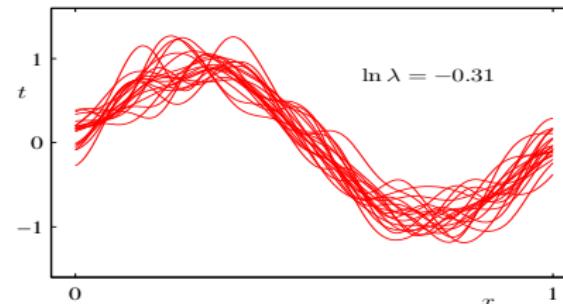
$$\hat{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^M \{f(x_n, \mathbf{w}) - y_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- **$\lambda$  controls the bias-variance trade-off.** Bias represents the extent to which the average prediction over all data sets differs from the desired regression function. Variance measures the extent to which the solutions for individual data sets vary around their average

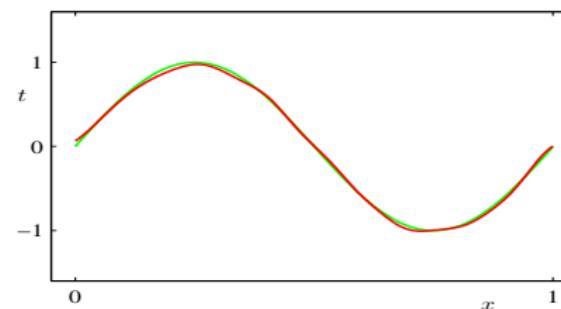
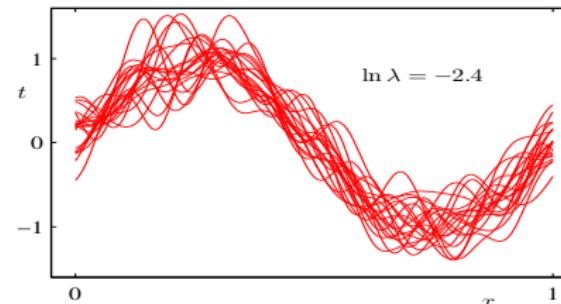
# The bias-variance decomposition as a fonction of $\ln \lambda$



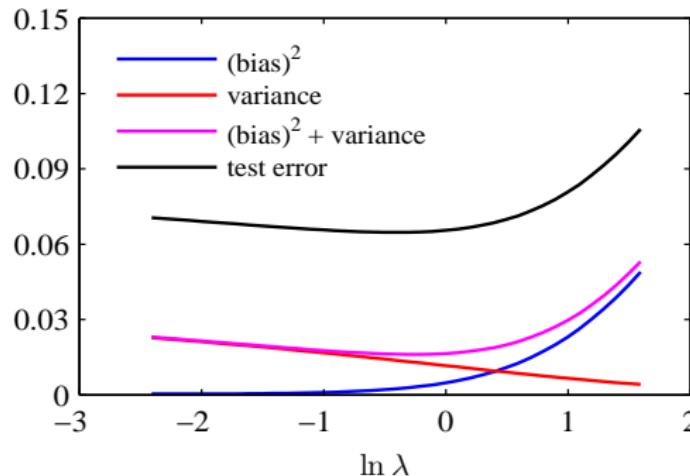
# The bias-variance decomposition as a fonction of $\ln \lambda$



# The bias-variance decomposition as a fonction of $\ln \lambda$



# The bias-variance decomposition



Very flexible models having low bias and high variance, and relatively rigid models having high bias and low variance.

Université Claude Bernard Lyon 1



# A bias-variance decomposition in classification

Several bias-variance decompositions in classification exist, one is given by

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[P(Y \neq f(X, \mathbf{w}))] &= 1 - P(j^*|X) \\ &+ \mathbb{E}_X[(P(j^*|X) - P(\hat{j}|X))P(\hat{j}|f, X)] \\ &+ \mathbb{E}_X[\sum_{j \neq \hat{j}} ((P(j^*|X) - P(j|X))P(j|f, X))]\end{aligned}$$

with notations

$$P(j|f, x)) = P_{\mathcal{D}}(f(x, \mathbf{w}) = j | X = x))$$

$$P(j|x)) = P(Y = j|x))$$

$$j^*(x) = \operatorname{argmax}_j P(j|x)$$

$$\hat{j}(x) = \operatorname{argmax}_j P(j|f, x)$$

# The bias-variance decomposition

We obtain the following decomposition of the expected squared loss

## Bias-variance decomposition in regression

$$\text{Error} = (\text{bias})^2 + \text{variance} + \text{noise}$$

## Bias-variance decomposition in classification

$$\text{Error} = \text{bias} + \text{spread} + \text{Bayes error}$$

The *spread* acts as variance en regression. Noise and Bayes error are irreducible.  
The goal is to find the best balance between bias and variance (or spread).



# Decision Theory

- Probability theory provides a consistent mathematical framework for quantifying and manipulating uncertainty.
- Determination of  $p(x, t)$  from a set of training data is typically a difficult inference problem in high dimensional spaces (curse of dimensionality).
- We also need a decision theory that, when combined with probability theory, allows us to make **optimal decisions in situations involving uncertainty**.
- We need a rule that assigns each value of  $x$  to one of the available classes given by  $p(x, t)$ . These decision should be optimal in some appropriate sense (**cost dependent**).



# Decision Theory

- We need a rule that divide the input space into regions  $\mathcal{R}_k$  called **decision regions**, one for each class, such that all points in  $\mathcal{R}_k$  are assigned to class  $\mathcal{C}_k$ .
- The boundaries between decision regions are called **decision boundaries** or decision surfaces.

# Decision Theory

- Suppose  $p(\mathbf{x}, y)$  is given with  $y \in \{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ , we are interested in  $p(\mathcal{C}_k|\mathbf{x})$ . Using Bayes' theorem,

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$$

- We can interpret  $p(\mathcal{C}_k)$  as the prior probability for the class  $\mathcal{C}_k$ , and  $p(\mathcal{C}_k|\mathbf{x})$  as the corresponding posterior probability.
- If our aim is to minimize the chance of assigning  $x$  to the wrong class, we choose the **class having the higher posterior probability**.
- This amounts to choose  $k$  s.t.  $p(\mathcal{C}_k|\mathbf{x}) > p(\mathcal{C}_j|\mathbf{x}), \forall j \neq k$ . In this case  $\mathcal{R}_k = \{\mathbf{x} | p(\mathcal{C}_k|\mathbf{x}) > p(\mathcal{C}_j|\mathbf{x}), \forall j \neq k\}$ .

# Minimizing the misclassification rate

- Note that

$$p(\mathcal{C}_k | \mathbf{x}) > p(\mathcal{C}_j | \mathbf{x})$$

is equivalent to

$$p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k) > p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)$$

- So, instead of learning  $p(\mathcal{C}_k | \mathbf{x})$ , we may learn  $K$  distinct models  $p(\mathbf{x} | \mathcal{C}_k)$  for  $k = 1, \dots, L$ . The  $p(\mathcal{C}_k)$  are easily estimated..

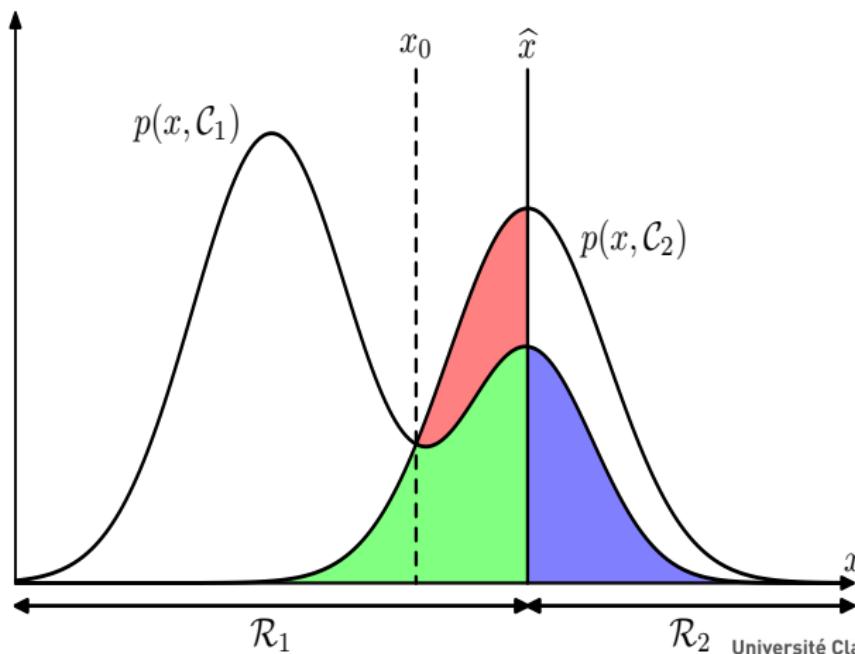
# Minimizing the misclassification rate

- In the 2-class problem, the probability an error occurring is given by

$$\begin{aligned} P(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) + p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) \\ &= \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} + \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} \end{aligned}$$

- Clearly to minimize  $P(\text{mistake})$  we should arrange that each  $\mathbf{x}$  is assigned to whichever class has the smaller value.

# Minimizing the misclassification rate



# Minimizing the expected loss

- The misclassification rate is not appropriate when the consequences of the mistakes can be dramatically different, we need an overall measure of loss (or cost function) incurred in taking any of the available decisions.
- By assigning a new value of  $x$ , with true class  $\mathcal{C}_k$ , to class  $\mathcal{C}_j$  (where  $j$  may not be equal to  $k$ ), we incur some level of loss that we denote by  $L_{kj}$ .
- The average loss is given by

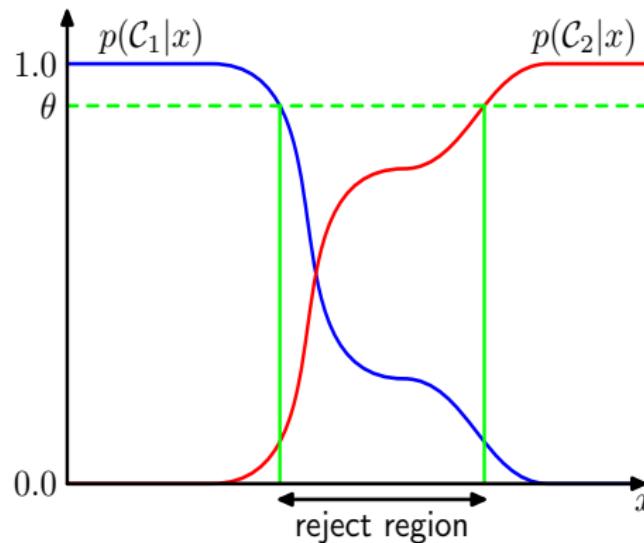
$$\mathbb{E}(L) = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

# Minimizing the expected loss

$$\mathbb{E}(L) = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

- The goal is to choose the regions  $\mathcal{R}_j$  in order to minimize  $\mathbb{E}(L)$ , which implies that for each  $\mathbf{x}$  we should minimize  $\sum_k L_{kj} p(\mathbf{x}, \mathcal{C}_k)$
- Thus the trivial decision rule that minimizes  $E(L)$  assigns each new  $\mathbf{x}$  to the class  $j$  for which the quantity  $\sum_k L_{kj} p(\mathbf{x}, \mathcal{C}_k)$  is minimum.

# The reject option



Inputs  $x$  such that  $\max_k p(\mathcal{C}_k|x) < \theta$ , are rejected.

# Inference and decision

There are three distinct approaches to solving decision problems, in decreasing order of complexity,

- 1 **Generative models:** Learn  $p(\mathbf{x}|\mathcal{C}_k)$  for each class  $\mathcal{C}_k$ ) individually. Also separately infer the prior class probabilities  $p(\mathcal{C}_k)$ . Then use Bayes' theorem to find  $p(\mathcal{C}_k|\mathbf{x})$ . One may also model the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  directly.
- 2 **Discriminative models:** Learn  $p(\mathcal{C}_k|\mathbf{x})$ , and then subsequently use decision theory to assign each new  $\mathbf{x}$  to one of the classes.
- 3 **Non probabilistic models:** Find a function  $f(\mathbf{x})$ , called a discriminant function, which maps each input  $\mathbf{x}$  directly onto a class label. Probabilities play no role here.

# Generative models

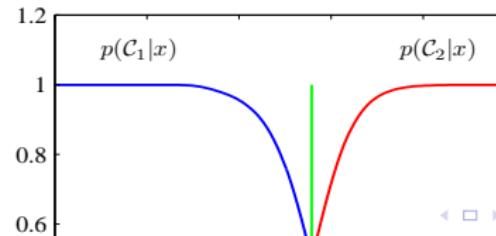
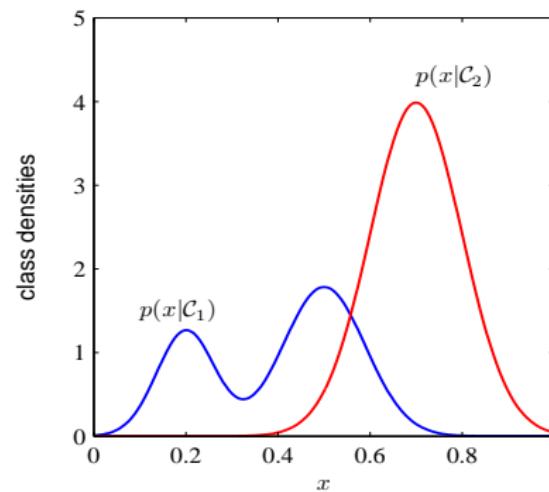
The relative merits of **Generative models**:

- If  $\mathbf{x}$  has high dimensionality, a large training set is needed to determine  $p(\mathcal{C}_k|\mathbf{x})$  from  $p(\mathbf{x}|\mathcal{C}_k)$ . to reasonable accuracy.
- Advantage: it allows the marginal density of data to be determined,

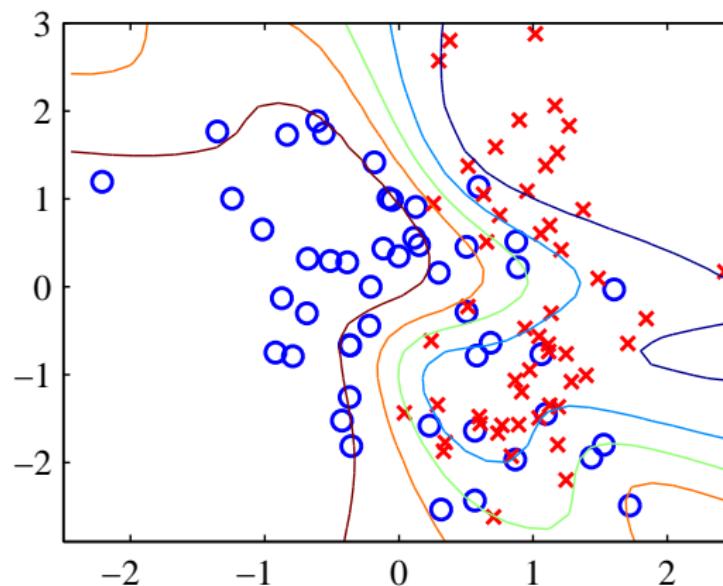
$$p(\mathbf{x}) = \sum_k p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

- Useful for detecting new data points that have low probability under the model known as (*outlier detection*) or (*novelty detection*).
- Disadvantage: wasteful of computational resources and excessively demanding of data if we only wish to make classification decisions.  $p(\mathbf{x}|\mathcal{C}_k)$  may contain a lot of structure that has little effect on the posterior probabilities

# Illustration



# Illustration: Decision boundaries



# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material



# Model Selection

- We have already seen that, in the maximum likelihood approach, the performance on the training set is not a good indicator of predictive performance on unseen data due to the problem of over-fitting
- If data is plentiful, the best approach is simply to decompose  $\mathcal{D}$  into 3 subsets:
  - A **training set**, some of the available data to train a range of models.
  - A **validation set** to estimate of the predictive performance of the models and compare them on independent data.
  - A **test set** on which the performance of the selected model is finally evaluated.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

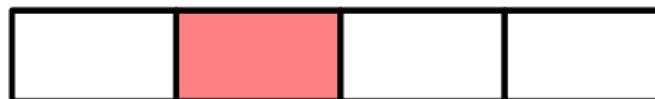
# Cross-validation

- In many applications, however, the supply of data for training and testing will be limited, and in order to build good models, we wish to use as much of the available data as possible for training.
- One solution is to use  $S$ -fold **cross-validation**. The technique involves taking the available data and partitioning it into  $S$  groups. Then  $S - 1$  of the groups are used to train a set of models that are then evaluated on the remaining group.
- The procedure is repeated for all  $S$  possible choices for the held-out group and the performance scores from the  $S$  runs are then averaged

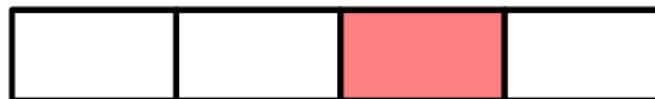
# Cross-validation



run 1



run 2



run 3



run 4

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Evaluation of binary classifiers

- There are many metrics that can be used to measure the performance of a classifier or predictor; different fields have different preferences for specific metrics due to different goals.
- To evaluate a classifier, one compares its output to a perfect classification and cross tabulates the data into a  $2 \times 2$  contingency table (confusion matrix),

	Positive	Negative
Positive prediction	TP	FP
Negative prediction	FN	TN

- **Prevalence** (i.e., the positive rate) has a significant impact on prediction values. A stupid classifier achieves an accuracy rate above 50% with (*imbalanced data sets*).
- The cost function may not be **symmetrical**, depending on the application



# Evaluation of binary classifiers

$$\text{Precision} = \frac{TP}{TP + FP}; \quad \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}; \quad \text{Sensibility} = \frac{TP}{TP + FN}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}; \quad \text{Balanced Acc.} = \frac{1}{2}(\text{Spe.} + \text{Sens.})$$

Each of the  $2 \times 2$  tables can be summarized as a pair of 2 numbers. In addition to the paired metrics, there are also single metrics that give a single score. An F-score is a combination of the precision and the recall:

$$\text{F-score} = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

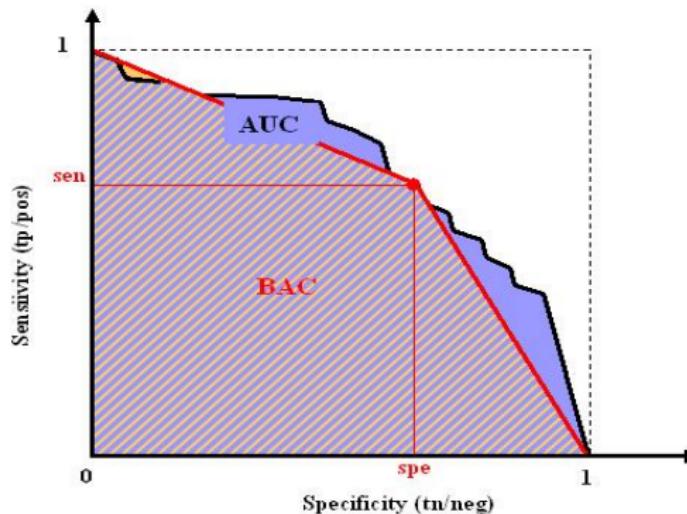
# Evaluation of binary classifiers

- **Prevalence** has a significant impact on prediction values.
- Suppose there is a test for a disease with 99% sensitivity and 99% specificity. If 2000 people are tested and the prevalence is 50%, 1000 of them are sick and 1000 healthy. Thus about 990 TP (and 990 TN) are likely, with 10 FP (and 10 FN). The P and N prediction values would be 99%.
- If the prevalence is only 5%, only 100 are sick. The likely result is 99 TP, 1 FN, 1881 TN and 19 FP. Of the 19+99 people tested positive, only 99 really have the disease, - only 84% chance that a patient has the disease given that his test result is positive, and only 1 chance in 1882, (0.05%) that the patient has the disease despite a negative test result.

# ROC Curve

- The overall performance of the classifier can be visualized and studied using the Receiver Operating Characteristic curve of **ROC curve**.
- A ROC curve, is a graphical plot that illustrates the performance of a binary classifier system as its discrimination **threshold** is varied. The curve is created by plotting the sensitivity against the specificity at various threshold settings.
- The best possible prediction method would yield a point in the upper left corner or coordinate  $(1, 1)$  of the ROC space,
- The diagonal divides the ROC space. Points above the diagonal represent good classification results (better than random), points below the line represent poor results (worse than random)

# Courbe ROC



The area under the curve (AUC) is equal to the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

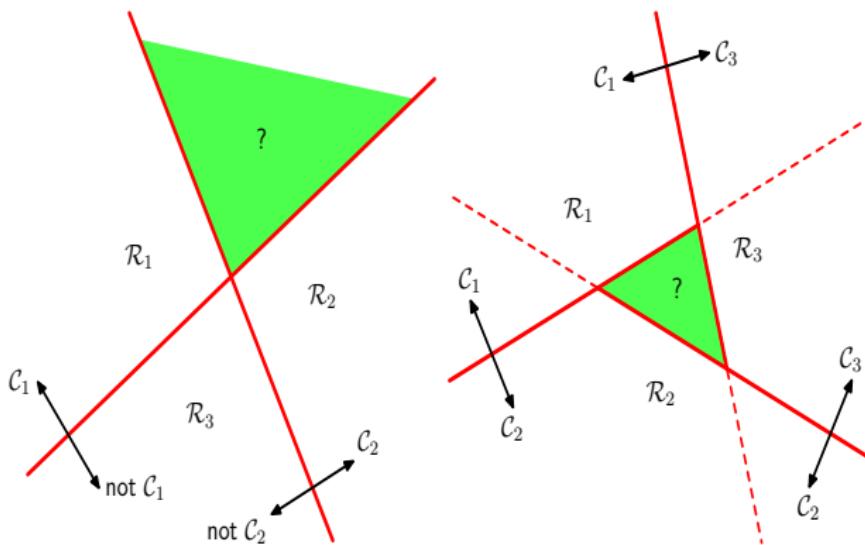
## 7 References & Further Material



# Multiclass classification

- Multiclass or multinomial classification is the problem of classifying instances into one of  $K$  classes ( $K > 2$ ).
- Many binary classification algorithms can be turned into multiclass classifiers by a variety of strategies.
  - 1 *One-against-all*: involves training a single classifier per class, with the samples of that class as positive samples and all other samples as negatives. Base classifiers produce real-valued confidence scores rather than just a class label.
  - 2 *One-against-one*:  $K(K - 1)/2$  binary classifiers are trained, each with the samples of the pair of classes. At prediction time, all classifiers are applied; the class with the highest number of votes is usually predicted.
  - 3 *Error-Correcting Output Codes* (ECOC) is an ensemble method designed for multi-class classification problem where each class is assigned a unique binary string of length  $n$ , called a codeword.

# Illustration



## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# ECOC - Illustration

Classe	vl	hl	dl	cc	ol	or
0	0	0	0	1	0	0
1	1	0	0	0	0	0
2	0	1	1	0	1	0
3	0	0	0	0	1	0
4	1	1	0	0	0	0
5	1	1	0	0	1	0
6	0	0	1	1	0	1
7	0	0	1	0	0	0
8	0	0	1	1	0	0
9	0	0	0	1	1	0

Table: A 6 bit error-correcting output code for a ten-class problem

# ECOC - Illustration

- In digit recognition task, we need to map each hand written digit to one of  $K = 10$  classes. Every column is given an interpretation

Column	notation	meaning
1	vl	contains a vertical line
2	hl	contains a horizontal line
3	dl	contains a diagonal line
4	cc	contains a closed curve
5	ol	contains an open curve on the right side
6	or	contains an open curve on the left side

# ECOC - Illustration

- During training, one binary classifier is learned for each column.
- To classify a new data point  $x$ , all  $n$  binary classifiers are evaluated to obtain a 6-bit string. Finally, we choose the class whose *codeword* is closest to  $x$ 's output string as the predicted label.
- Example : 110001 is closest to 110000, thus the output class is 4.
- The rows have more bits than is necessary ( $\log_2(10)$ ). Using some redundant "error-correcting" bits, we can tolerate some error.
- If  $d$  is the minimum Hamming distance between any pair of code words, then the code can correct at least  $[(d - 1)/2]$  single bit errors.
- In the previous code,  $d = 2$  thus no error is allowed.

# ECOC

- There are many ways to design the error-correcting output code.
- When  $K$  is small, one can use exhaustive codes. Each code has length  $2^{K-1} - 1$ . Row 1 contains only ones. Row 2 consists of  $2^{K-2}$  zeros followed by  $2^{K-2} - 1$  ones, and so on. This code has the largest inter-row Hamming distance.
- When  $K$  is large, random codes can be used. Random code works as well as optimally constructed code.
- The major benefit of error-corrective coding is variance reduction via model averaging.

# One-against-one

- One-against-one suffers from ambiguities (just as one-against-all) in that some regions of its input space may receive the same number of votes.
- Instead of taking the class with the highest number of votes, it is possible to combine the  $K(K - 1)/2$  outputs values of the *One-against-one* strategy to obtain a posterior class probabilities.
- Let  $\mathcal{C}_{ij}$  denotes " $x$  is in class  $\mathcal{C}_i$  or  $\mathcal{C}_j$ " and let

$$P_{ij} = P(\mathcal{C}_i | \mathcal{C}_{ij}, X = x)$$

- Can  $P(\mathcal{C}_i | X = x)$  be written as a function of  $P_{ij}$  ?

# One-against-one

- The idea:

$$\begin{aligned} P\left(\bigcup_{j=1}^K \mathcal{C}_j | X = x\right) &= 1 \\ &= P\left(\bigcup_{j=1, i \neq j}^K \mathcal{C}_{ij} | X = x\right) \\ &= \sum_{j=1, i \neq j}^K P(\mathcal{C}_{ij} | X = x) - (K-2) \cdot P(\mathcal{C}_i | X = x) \end{aligned}$$

with

$$P_{ij} = P(\mathcal{C}_i | \mathcal{C}_{ij}, X = x) = \frac{P(\mathcal{C}_i | X = x)}{P(\mathcal{C}_{ij} | X = x)}$$

# One-against-one

- One obtains the  $K$  a posteriori probabilities given the  $K(K - 1)/2$  probabilities  $P_{ij}$  :

## Recombination

$$P(\mathcal{C}_i | X = x) = \frac{1}{\sum_{j=1, i \neq j}^K \frac{1}{P_{ij}} - (K - 2)}$$

- With  $P_{ij} = 1, \forall i$ , we get  $P(\mathcal{C}_i | X = x) = 1$  as expected.
- No clear advantage between *one-against-one* and *one-against-all* techniques.

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Multi-label classification

- The multi-label classification (MLC) problem differs from conventional classification, as multiple labels can be assigned simultaneously to the same instance.
- This situation is encountered in many recent real-world problems, including image indexing and annotation, facial expression analysis, text categorization, sentiment analysis, fault-control, drug side effects prediction, genome-wide protein function assignment, early detection of chronic diseases to cite just a few.
- For example, a document can be classified as both "spam" and "politics". Or, a patient can be diagnosed with both "cancer" and "diabetes".
- Multilabel learning is a challenging problem because it is difficult to learn the relationships between multiple labels.
- The MLC problem has received increasing attention from the ML community.



# Learning setting

- Learning in MLC amounts to finding a mapping from a space of features  $\mathcal{X} = \mathbb{R}^d$  to a space of labels  $\mathcal{Y} = \{0, 1\}^m$ , given a set of training samples in  $\mathcal{X} \times \mathcal{Y}$  and a loss function  $L$ .
- Formally, the training set  $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^s$  consists in i.i.d. samples drawn from the joint distribution  $p(\mathbf{x}, \mathbf{y})$  ( $\mathbf{x} \in \mathcal{X}$ ,  $\mathbf{y} \in \mathcal{Y}$ ), and the loss function  $L : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$  is a distance on the label space.
- The goal of learning is then to build a function  $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$  which maps any new input  $\mathbf{x}$  to its proper set of labels, as closely as possible given the loss function  $L$ .
- There are many different cost functions that can be used for multilabel learning.

# Learning setting

- From a probabilistic perspective, solving this problem amounts to modeling the conditional joint distribution  $p(\mathbf{y}|\mathbf{x})$ , and inferring Bayes-optimal predictions by minimizing the expected loss,  
$$\mathbf{h}^*(\mathbf{x}) = \operatorname{argmin}_{\hat{\mathbf{y}}} \mathbb{E}_{\mathbf{y}|\mathbf{x}} [L(\hat{\mathbf{y}}, \mathbf{y})].$$
- Multi-label classification raises a number of computational and statistical challenges, mainly due the size of the output space which grows exponentially with the number of labels. Typically, modeling naively  $p(\mathbf{y}|\mathbf{x})$  requires estimating  $O(2^m)$  parameters, while inferring a Bayes-optimal prediction for an arbitrary loss function requires  $O(2^{2m})$  evaluations of  $L$ .
- Label dependencies have to be incorporated into the learning process in order to improve the classification performance.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# MLC cost functions

There are many different cost functions that can be used for multilabel learning. Some of the most common cost functions include:

- **Hamming loss:** This is the most common cost function for multilabel learning. It measures the number of labels that are misclassified.
- **Zero-one loss:** This is a stricter cost function than Hamming loss. It measures the number of labels that are misclassified, as well as the number of labels that are not classified at all.
- **F-measure:** This cost function balances the precision and recall of the predictions.

The choice of cost function depends on the specific application. For example, if the application requires a high precision, then the zero-one loss function may be a good choice. If the application requires a high recall, then the F-measure may be a good choice.



## Further MLC cost functions

- The **Jaccard index**, is defined as the number of correctly predicted labels divided by the union of predicted and true labels,  $\frac{|T \cap P|}{|T \cup P|}$  where P and T are sets of predicted labels and true labels respectively.
- **Precision** is  $\frac{|T \cap P|}{|P|}$ ,
- **Recall** is  $\frac{|T \cap P|}{|T|}$ ,
- $F_1$  is their harmonic mean

# Hamming loss

- Optimizing decomposable loss functions such as the popular Hamming loss requires only to model the marginal distribution of each label.

$$L_H(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y_i = h_i(\mathbf{x})],$$

where  $\mathbb{I}[\cdot]$  is the standard  $\{False, True\} \rightarrow \{0, 1\}$  mapping.

- Binary relevance** is a simple method for multilabel learning. It treats each label as a separate binary classification problem. This means that the model is trained to predict each label independently of the other labels.
- Binary relevance can be inaccurate if the labels are not independent of each other.

# The Zero-one loss

- Contrary to decomposable loss functions, optimizing complex performance metrics such as the subset 0/1 loss, the F-measure or the Jaccard index requires to model the joint distribution of the labels (at least to some extent),
- The *subset* 0/1 loss, which generalizes the well-known 0/1 loss from the conventional to the multi-label setting, i.e.,

$$L_S(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \mathbb{I}[\mathbf{y} = \mathbf{h}(\mathbf{x})].$$

- The risk-minimizing prediction for subset 0/1 loss is given by the mode of the conditional distribution  $p(\mathbf{y}|\mathbf{x})$ , a.k.a. the maximum a posterior probability estimate (MAP), or the most probable explanation (MPE),

$$\mathbf{h}^*(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}).$$

- Therefore, one iteration through the  $2^m$  label combinations suffices to compute  $\mathbf{h}^*(\mathbf{x})$  for the subset 0/1 loss.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

### ■ MLC loss functions

### ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Learning Methods

- Probabilistic classifier chains (PCC) explicitly learn a maximum-likelihood estimate of  $p(\mathbf{y}|\mathbf{x})$  with logistic classifiers for MAP inference (computationally expensive).
- Label ranking is a method that ranks the labels for each instance. The top-ranked labels are then considered to be the predicted labels (difficult to train).
- Ensemble learning is a method for combining the predictions from multiple models. This can be done by averaging the predictions, or by using a voting system (computationally expensive).
- Weighted majority voting combines the predictions from multiple binary classifiers. The weights for the classifiers are determined by their accuracy (computationally expensive).

# Challenges of multilabel learning

- There are many challenges associated with multilabel learning. Some of the most common challenges include:
- The curse of dimensionality: This is the problem of having too many features relative to the number of instances. This can make it difficult to learn the relationships between the features and the labels.
- Label imbalance: This is the problem of having some labels that are much more common than others. This can make it difficult to learn the minority labels.
- Class overlap: This is the problem of having labels that are similar to each other. This can make it difficult to distinguish between the labels.
- These challenges can be addressed using a variety of techniques, such as feature selection,

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Problem

- When dealing with two random variables  $X$  and  $Y$ , one can use
  - 1 a *generative* model, i.e. which models the joint distribution  $p(X, Y)$ ,
  - 2 a *conditional* model which models the conditional probability of the output, given the input  $p(Y|X)$ .
  - 3 a *discriminative* which outputs the class directly.
- *Linear regression* and *logistic regression* are conditional models.
- The *SVM* is a discriminative model.
- The *Naive Bayes* is a generative model.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## MLC loss functions

## MLC learning methods

## 5 Classification and Regression Models

### Linear regression

- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Problem

- Suppose that  $Y \in \mathbb{R}$  depends linearly on  $X \in \mathbb{R}^p$ . Let  $w \in \mathbb{R}^p$  be a weighting vector and  $\sigma^2 > 0$ .
- We make the following assumption:

$$Y | X \sim \mathcal{N}(w^\top X, \sigma^2),$$

which can be rewritten as

$$Y = w^\top X + \epsilon,$$

with  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ .

- If an offset  $w_0 \in \mathbb{R}^p$  is added, i.e.  $Y = w^\top X + w_0 + \epsilon$ , the weighting vector is augmented  $\tilde{w} \in \mathbb{R}^{p+1}$  such that

$$Y = \tilde{w}^\top \begin{pmatrix} X \\ 1 \end{pmatrix} + \epsilon.$$

# Linear regression

- Let  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  be a training set of i.i.d. random variables. Each  $y_i$  is a *label* (a decision) on observation  $\mathbf{x}_i$ .
- The *conditional* distribution of all outputs given all inputs:

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n; \mathbf{w}, \sigma^2) = \prod_{i=1}^n p(y_i | \mathbf{x}_i; \mathbf{w}, \sigma^2).$$

- The associated log-likelihood has the following expression:

$$-\ell(\mathbf{w}, \sigma^2) = -\sum_{i=1}^n \log p(y_i | \mathbf{x}_i) = \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{\sigma^2}.$$

- The minimization problem with respect to  $\mathbf{w}$  can be reformulated as:

$$\text{find } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2.$$

# Linear regression

- Define the so-called *design matrix*  $\mathbf{X}$  as

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times p}$$

- Denote by  $\mathbf{y}$  the vector of coordinates  $(y_1, \dots, y_n)$ . The minimization problem over  $w$  can be rewritten in a more compact way as:

$$\text{find } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2.$$

# Linear regression

- Define

$$f : \mathbf{w} \mapsto \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = \frac{1}{2n} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w})$$

- $f$  is strictly convex iff its Hessian matrix is invertible. This is never the case when  $n < p$  (underdetermined problems).
- Usually, the Hessian matrix is invertible when  $n \geq p$ .
- Otherwise, Tikhonov regularization may be performed. It consists of adding a penalization term of the  $\ell_2$ -norm of  $\mathbf{w}$  by minimizing  $f(\mathbf{w}) + \lambda \|\mathbf{w}\|^2$  with some hyperparameter  $\lambda > 0$ .

# Linear regression

- The gradient of  $f$  is

$$\nabla f(\mathbf{w}) = \frac{1}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0 \iff \mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

- The equation

$$\boxed{\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}}$$

is known as the **normal equation**.

# Linear regression

- If  $\mathbf{X}^\top \mathbf{X}$  is invertible, then the optimal weighting vector is

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^\dagger \mathbf{y}$$

- $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is the *Moore-Penrose pseudo-inverse* of  $\mathbf{X}$ .
- The optimal weighting vector from  $\mathbf{X}$  and  $\mathbf{y}$  can be computed in  $O(p^3)$  (use a Cholesky decomposition of matrix  $\mathbf{X}^\top \mathbf{X}$  and solve two triangular systems).
- If  $\mathbf{X}^\top \mathbf{X}$  is not invertible, the solution is not unique anymore, and for any  $\mathbf{h} \in \ker(\mathbf{X})$ ,  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y} + \mathbf{h}$  is an admissible solution. In this case, regularization is necessary.

# Linear regression

- Now, let's differentiate  $\ell(\mathbf{w}, \sigma^2)$  with respect to  $\sigma^2$ : we have

$$\nabla_{\sigma^2} \ell(\mathbf{w}, \sigma^2) = \frac{n}{2\sigma^2} - \frac{n}{2\sigma^4} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2.$$

- Setting  $\nabla_{\sigma^2} \ell(\mathbf{w}, \sigma^2)$  to zero yields,

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2.$$

- In practice,  $\mathbf{X}$  may be badly conditioned. To avoid numerical issues, columns are normalized.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## MLC loss functions

## MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- **Logistic regression**
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Logistic regression

- Let  $X \in \mathbb{R}^p$ ,  $Y \in \{0, 1\}$ . We assume that  $Y$  follows a Bernoulli distribution with parameter  $\theta$ . The problem is to find  $\theta$ .
- Let's define the *sigmoid* function  $\sigma$  defined on the real axis and taking values in  $[0, 1]$ , such that

$$\forall z \in \mathbb{R}, \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

- One can easily prove that  $\forall z \in \mathbb{R}$ ,

$$\begin{aligned}\sigma(-z) &= 1 - \sigma(z), \\ \sigma'(z) &= \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z).\end{aligned}$$

# Logistic regression

- Assume that  $Y|X = \mathbf{x}$  follows a Bernoulli distribution with parameter  $\theta = \sigma(\mathbf{w}^\top \mathbf{x})$ , where  $\mathbf{w}$  is a weighting vector. An offset  $\mathbf{w}^\top \mathbf{x} + w_0$  is added.
- The conditional distribution is given by

$$p(Y = y|X = \mathbf{x}) = \theta^y(1 - \theta)^{1-y} = \sigma(\mathbf{w}^\top \mathbf{x})^y \sigma(-\mathbf{w}^\top \mathbf{x})^{1-y}.$$

- Given a training set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  of *i.i.d.* random variables, we can compute the log-likelihood

$$l(\mathbf{w}) = \sum_{i=1}^n y_i \log \sigma(\mathbf{w}^\top \mathbf{x}_i) + (1 - y_i) \log \sigma(-\mathbf{w}^\top \mathbf{x}_i).$$

# Logistic regression

- In order to minimize the log-likelihood, since  $z \mapsto \log(1 + e^{-z})$  is a convex function and  $\mathbf{w} \mapsto \mathbf{w}^\top \mathbf{x}_i$  is linear,
- Let  $\eta_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$ , the gradient is:

$$\begin{aligned}\nabla_{\mathbf{w}} l(\mathbf{w}) &= \sum_{i=1}^n y_i \mathbf{x}_i \frac{\sigma(\mathbf{w}^\top \mathbf{x}_i) \sigma(-\mathbf{w}^\top \mathbf{x}_i)}{\sigma(\mathbf{w}^\top \mathbf{x}_i)} \\ &\quad - (1 - y_i) \mathbf{x}_i \frac{\sigma(\mathbf{w}^\top \mathbf{x}_i) \sigma(-\mathbf{w}^\top \mathbf{x}_i)}{\sigma(-\mathbf{w}^\top \mathbf{x}_i)} \\ &= \sum_{i=1}^n \mathbf{x}_i (y_i - \eta_i)\end{aligned}$$

# Logistic regression

- The stationary solution is,

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) = 0 \iff \sum_{i=1}^n \mathbf{x}_i (y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)) = 0$$

- The Hessian matrix  $H\ell(\mathbf{w})$  is given by:

$$\begin{aligned} H\ell(\mathbf{w}) &= \sum_{i=1}^n \mathbf{x}_i (0 - \sigma'(\mathbf{w}^\top \mathbf{x}_i) \sigma'(-\mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i^\top) \\ &= \sum_{i=1}^n (-\eta_i(1-\eta_i)) \mathbf{x}_i \mathbf{x}_i^\top \\ &= -\mathbf{X}^\top \text{Diag}(\eta_i(1-\eta_i)) \mathbf{X} \end{aligned}$$

## Brief review: First-order methods

- Let  $\ell : \mathbb{R}^p \rightarrow \mathbb{R}$  be the convex  $C^1$  function that we want to minimize. A *descent direction* at point  $\mathbf{x}$  is a vector  $d$  such that  $\langle \mathbf{d}, \nabla_{\mathbf{w}} \ell(\mathbf{w}) \rangle < 0$ .
- The minimization of  $f$  may be achieved by *descent methods*, which are adjust the current solution iteratively,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \varepsilon^{(k)} \mathbf{d}^{(k)},$$

where  $\varepsilon^{(k)}$  is the *stepsize*.

- We write the first-order Taylor-expansion of  $\ell$ :

$$\ell(\mathbf{w}) = \ell(\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^\top \nabla \ell(\mathbf{w}^t) + o(\|\mathbf{w} - \mathbf{w}^t\|)$$

- Minimization of the first terms leads to,

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \nabla_{\mathbf{w}} \ell(\mathbf{w}).$$

- The is called (pure) gradient descent..

## Brief review: First-order methods

There are several choices for  $\varepsilon^{(k)}$ :

- 1 Constant step:  $\varepsilon^{(k)} = \varepsilon$ . But the scheme does not necessarily converge.
- 2 Decreasing step size:  $\varepsilon^{(k)} \propto \frac{1}{k}$  with  $\sum_k \varepsilon^{(k)} = \infty$  and  $\sum_k (\varepsilon^{(k)})^2 < \infty$ .  
The scheme is guaranteed to converge.
- 3  $\varepsilon^{(k)}$  may be determined by *line search* which seeks  $\min_{\varepsilon} f(\mathbf{x}^{(k)} + \varepsilon \mathbf{d}^{(k)})$ :
  - either exactly but this is computationally expensive and often useless.
  - or approximately (e.g. Armijo line search) which usually performs better.

## Brief review: Newton's method

- Now, let  $\ell : \mathbb{R}^p \rightarrow \mathbb{R}$  be the  $C^2$  function that we want to minimize. We write the second-order Taylor-expansion of  $f$ :

$$\ell(\mathbf{w}) = \ell(\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^\top \nabla \ell(\mathbf{w}^t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^t)^\top H\ell(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + o(\|\mathbf{w} - \mathbf{w}^t\|)^2$$

- A local minimum  $\mathbf{x}^*$  is then reached whenever,

$$\begin{cases} \nabla_{\mathbf{w}} \ell(\mathbf{w}^*) = 0 \\ H\ell(\mathbf{w}^*) \succeq 0 \end{cases}$$

- Let  $\mathbf{h} = \mathbf{w} - \mathbf{w}^t$ , the minimization problem becomes:

$$\min_{\mathbf{h}} \mathbf{h}^\top \nabla_{\mathbf{w}} \ell(\mathbf{w}) + \frac{1}{2} \mathbf{h}^\top H\ell(\mathbf{w}) \mathbf{h}.$$

- This leads to set

$$\mathbf{w}^{t+1} = \mathbf{w}^t + H\ell(\mathbf{w}^t)^{-1} \nabla_{\mathbf{w}} \ell(\mathbf{w}).$$

## Brief review: Newton's method

- If the Hessian matrix is not invertible, we can regularize the problem and minimize  $\ell(\mathbf{w}) + \lambda \|\mathbf{w} - \mathbf{w}^t\|^2$  instead.
- The *Pure Newton step* does not converge always. It is better to use the so-called *Damped Newton method*:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \varepsilon^t \cdot H\ell(\mathbf{w}^t)^{-1} \nabla_{\mathbf{w}} \ell(\mathbf{w}).$$

where  $\varepsilon^t$  is set with the Armijo *Line Search*

# Logistic regression

- Newton's algorithm for logistic regression takes the form:

$$\begin{aligned}\ell(\mathbf{w}) &= \sum_{i=1}^n y_i \log \sigma(\mathbf{w}^\top \mathbf{x}_i) + (1 - y_i) \log \sigma(-\mathbf{w}^\top \mathbf{x}_i) \\ \nabla_{\mathbf{w}} \ell(\mathbf{w}) &= \sum_{i=1}^n \mathbf{x}_i (y_i - \eta_i) = \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\eta}) \\ H\ell(\mathbf{w}) &= -\mathbf{X}^\top \text{Diag}(\eta_i(1 - \eta_i)) \mathbf{X}\end{aligned}$$

# Logistic regression

- The minimization problem can also be seen as some *weighted* linear regression over  $h$  of some function of the form

$$\sum_i \frac{(\tilde{y}_i - \tilde{\mathbf{x}}_i^\top \mathbf{h})^2}{\sigma_i^2}$$

where  $\tilde{y}_i = y_i - \eta_i$  and  $\sigma_i^2 = [\eta_i(1 - \eta_i)]^{-1}$ .

- This method is often referred to as the *iterative reweighted least squares* algorithm (IRLS).

# Multiclass logistic regression

- In the logistic regression model for binary classification, we have that  $Y|X = \mathbf{x} \sim \text{Ber}(\mu(\mathbf{x}))$  and linear combination of parameters:

$$\mu(\mathbf{x}) = \sigma(\eta(\mathbf{x})) = 1 + \exp(-\eta(\mathbf{x}))^{-1}$$

$$\eta(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + \mathbf{b}.$$

- It is then natural to consider the generalization to a multiclass classification setting.
- In that case,  $Y|X = \mathbf{x}$  is multinomial distribution with natural parameters  $(\eta_1(\mathbf{x}), \dots, \eta_K(\mathbf{x}))$ .

# Multiclass logistic regression

- We may consider a classification problem with more than two classes :  
 $Y \in \{1, \dots, K\}$  with  $Y \sim \mathcal{M}(1, \pi_1(\mathbf{x}), \dots, \pi_K(\mathbf{x}))$  where

$$\pi_k(\mathbf{x}) = p(Y = k \mid \mathbf{x}) = \frac{e^{\mathbf{w}_k^\top \mathbf{x}}}{\sum_{j=1}^K e^{\mathbf{w}_j^\top \mathbf{x}}}$$

- We need to define a rule over the classifiers (softmax function, one-versus-all, etc.) in order to make a decision.

# Multiclass logistic regression

- Let  $w_k \in \mathbb{R}^p$  and  $b_k \in \mathbb{R}$ , for all  $1 \leq k \leq K$  and set  $\eta_k(x) = w_k^\top x + b_k$ .  
We have

$$\mathbb{P}(Y_k = 1 | X = x) = \frac{e^{\eta_k(x)}}{\sum_{k'=1}^K e^{\eta_{k'}(x)}} = \frac{e^{w_k^\top x + b_k}}{\sum_{k'=1}^K e^{w_{k'}^\top x + b_{k'}}},$$

and thus

$$\log \mathbb{P}(Y_k = y | X = x) = \sum_{k=1}^K y_k (w_k^\top x + b_k) - \log \left[ \sum_{k'=1}^K e^{w_{k'}^\top x + b_{k'}} \right].$$

- Like for binary logistic regression, the maximum likelihood principle can be used to learn  $(w_k, b_k)_{1 \leq k \leq K}$  using numerical optimization methods such as the IRLS algorithm.



# Multiclass logistic regression

- Note that the form of the parameterization obtained is the same as for the Naive Bayes model.
- However, the Naive Bayes model is learnt as a generative model, while the logistic regression is learnt as conditional model.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## MLC loss functions

## MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- **Shrinkage methods**
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Shrinkage methods

- We have seen that regularization control the over-fitting phenomenon by adding a penalty term to the error function,
- This technique, also called ridge regression, is an example of shrinkage method applied to least squares regression, to improve a least-squares estimator. Least-squares solutions having a large number of coefficients different from zero are penalized.
- Only the variables whose impact on the empirical risk is considerable have a coefficient bigger than zero and consequently appear in the fitted linear model.
- Doing shrinkage is therefore an implicit embedded manner of doing feature selection since only a subset of variables contributes to the final predictor.

# Ridge regression

- An equivalent way to write the ridge problem is

$$\hat{w}_r = \arg \min_w \sum_{i=1}^N (y_i - x_i^T w)^2,$$

subject to  $\sum_{j=1}^p w_j^2 \leq L$

- The ridge regression solution is

$$\hat{w}_r = (X^T X + \lambda I)^{-1} X^T Y$$

where  $I$  is the identity matrix of size  $p$

# Lasso

- Another well known shrinkage method is lasso where the estimate of the linear parameters is returned by

$$\hat{w}_r = \arg \min_w \sum_{i=1}^N (y_i - x_i^T w)^2,$$

subject to  $\sum_{j=1}^p |w_j| \leq L$

- The 1-norm penalty allows a stronger constraint on the coefficients, however it makes the solution nonlinear and requires a quadratic programming algorithm.

- If  $L > \sum_{j=1}^p \hat{w}_j$  the lasso returns the common least-squares solution. The penalty factor  $L$  is typically set by cross-validation.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- **Naive Bayes classifiers**
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Naive Bayes classifier

- The **Naive Bayes** (NB) classifier is a baseline model just as decision trees.
- Consider a multi-class learning problem.  $Y$  that takes values in the set  $\{c_1, \dots, c_K\}$ . The Bayes optimal classifier should return

$$c^*(x) = \arg \max_{j=1, \dots, K} \text{Prob}\{Y = c_j | x\}$$

- Using the Bayes theorem,

$$\begin{aligned} c^*(x) &= \arg \max_{j=1, \dots, K} \frac{\text{Prob}\{x|Y = c_j\} \text{Prob}\{Y = c_j\}}{\text{Prob}\{x\}} \\ &= \arg \max_{j=1, \dots, K} \text{Prob}\{x|Y = c_j\} \text{Prob}\{Y = c_j\} \end{aligned}$$

# Naive Bayes classifier

- It is easy to estimate each of the a priori probabilities  $\text{Prob}\{Y = c_j\}$ . The estimation of  $\text{Prob}\{\mathbf{x}|Y = c_j\}$  is much harder.
- The NB classifier is based on the **simplifying assumption** that the input values are **conditionally independent** given the target value:

$$\text{Prob}\{\mathbf{x}|Y = c_j\} = \text{Prob}\{x_1, \dots, x_n|Y = c_j\} = \prod_{h=1}^n \text{Prob}\{x_h|Y = c_j\}$$

- The NB classification is then:

$$c_{NB}(\mathbf{x}) = \arg \max_{j=1, \dots, K} \text{Prob}\{Y = c_j\} \prod_{h=1}^n \text{Prob}\{x_h|Y = c_j\}$$

# Naive Bayes classifier

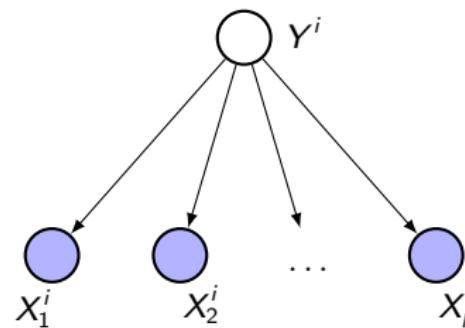
- The naive Bayes classifier is relevant when modeling the joint distribution of  $p(x|y)$  is difficult.
- Example : classification of documents based on a bag-of-word representation; i.e. each word of a reference dictionary is present in the document or not.
- The document  $i$  is represented by a vector of binary random variables.  
 $X^i : \Omega \mapsto \{0, 1\}^P$ , with  $x_j^i = 1$  iff word  $j$  of the dictionary is present in the  $i$ th document.

# Naive Bayes classifier

- It is possible to approach the problem using directly a *conditional model* of  $p(y | x)$  or using a *generative model* of the joint distribution modeling separately  $p(y)$  and  $p(x|y)$  and computing  $p(y|x)$  using Bayes rule.
- The naive Bayes model is a generative model.
- $Y^i$  is naturally modeled as a multinomial distribution with
$$p(y^i) = \prod_{k=1}^K \pi_k^{y_k^i}.$$
 However  $p(x^i|y^i) = p(x_1^i, \dots, x_p^i|y^i)$  has a priori  $2^p - 1$  parameters. The key assumption made is that  $X_1^i, \dots, X_p^i$  are all independent conditionally on  $Y^i.$
- While ignoring the correlations between words it works well in practice.

# Naive Bayes classifier

- These conditional independence assumptions correspond to the following graphical model:



# Naive Bayes classifier

- The distribution of  $Y^i$  is a multinomial distribution which we parameterize with  $(\pi_1, \dots, \pi_K)$ , and we write

$$\mu_{jk} = P(X_j^{(i)} = 1 | Y_k^{(i)} = 1)$$

- We then have

$$\begin{aligned} p(X^i = x^i, Y^i = y^i) &= p(x_i, y_i) \\ &= p(x^i | y^i) p(y^i) \\ &= \prod_{j=1}^p p(x_j^i | y^i) p(y^i) \end{aligned}$$

# Naive Bayes classifier

- The last expression leads to

$$p(x^i, y^i) = \left[ \prod_{j=1}^p \prod_{k=1}^K \mu_{jk}^{x_j^i y_k^i} (1 - \mu_{jk})^{(1-x_j^i)y_k^i} \right] \prod_{k=1}^K \pi_k^{y_k^i}$$

$$\begin{aligned} \log p(x^i, y^i) &= \sum_{k=1}^K \left( \sum_{j=1}^p (x_j^i y_k^i \log \mu_{jk} + (1 - x_j^i)y_k^i \log(1 - \mu_{jk})) \right. \\ &\quad \left. + y_k^i \log(\pi_k) \right) \end{aligned}$$

- Note that, in spite of the name the naive Bayes classifier is not a Bayesian approach to classification.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- **Support vector machines**
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Support vector machines (linear case)

- Consider a binary classification task. The problem of separating input data by a linear boundary is ill-posed as there are infinitely many possible separating hyperplanes characterized by the equation

$$\beta_0 + x^T \beta = 0$$

- The vector normal to the hyperplane is given by  $\beta^* = \frac{\beta}{|\beta|}$
- The signed distance of a point  $x$  to the hyperplane is called the geometric margin and is given by

$$\beta^{*T} (x - x_0) = \frac{x^T \beta - \beta x_0^T}{|\beta|} = \frac{1}{|\beta|} (x^T \beta + \beta_0)$$

# Support vector machines (linear case)

- This technique relies on an optimization approach to compute separating hyperplanes and was shown to lead to good classification performance on real data.
- The **SVM approach computes the (unique) maximal margin hyperplane** for a training set, i.e. the optimal separating hyperplane which separates the two classes by maximizing the distance to the closest point from either class.
- The problem is modeled as the optimization problem

$$\max_{\beta, \beta_0} C \text{ subject to } \frac{1}{|\beta|} y_i (x_i^T \beta + \beta_0) \geq C \quad \text{for } i = 1, \dots, N$$

So that all the points are at least a distance  $C$  from the decision boundary.

# Margin

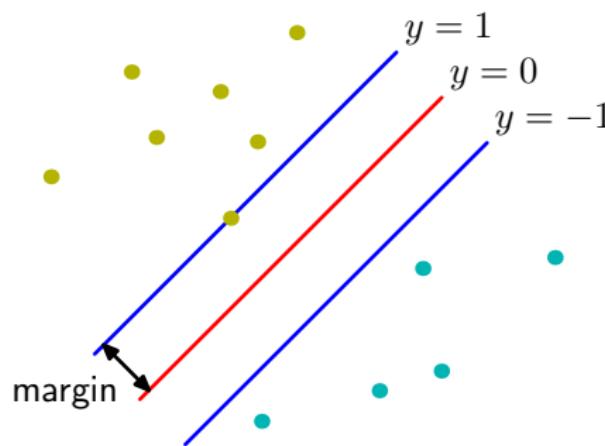


Figure: The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points.

# Maximizing the margin

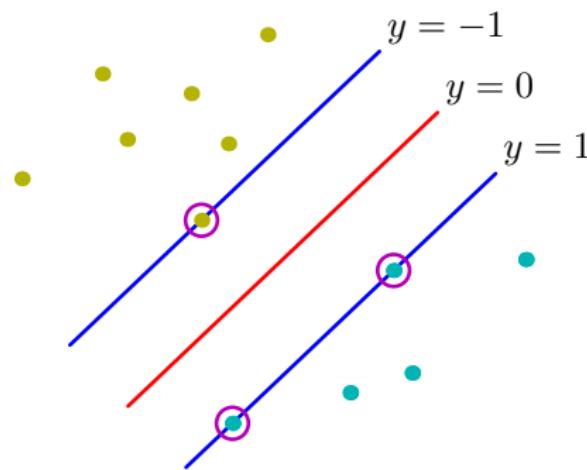


Figure: The location of this boundary is determined by a subset of the data points, known as support vectors.

# Support vector machines

- We can set  $|\beta| = 1/C$ . The maximization problem can be reformulated in a minimization form

$$\min_{\beta, \beta_0} \frac{1}{2} |\beta|^2 \text{ subject to } y_i(x_i^T \beta + \beta_0) \geq 1 \quad \text{for } i = 1, \dots, N$$

- The constraints impose a margin around the linear decision of thickness  $1/|\beta|$ .
- This optimization problem is a **convex optimization problem** (quadratic criterion with linear inequality constraints) where the primal Lagrangian is

$$L_P(\beta, \beta_0) = \frac{1}{2} |\beta|^2 - \sum_{i=1}^N \alpha_i [y_i(x_i^T \beta + \beta_0) - 1]$$

$\alpha_i \geq 0$  are the Lagrangian multipliers.

# Support vector machines

- Setting the derivatives  $\beta$  and  $\beta_0$  to zero we obtain:

$$\beta = \sum_{i=1}^N \alpha_i y_i x_i, \quad 0 = \sum_{i=1}^N \alpha_i y_i$$

- Substituting these in the primal form, we obtain the dual

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k$$

subject to  $\alpha_i \geq 0$ .

# Support vector machines

- The dual optimization problem is now

$$\max_{\alpha} L_D = \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k \langle x_i, x_k \rangle$$

subject to

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad \text{for } \alpha_i \geq 0, \quad i = 1, \dots, N$$

- $\langle x_i, x_k \rangle$  is the inner product of  $x_i$  and  $x_k$ .

# Support vector machines

- It can be shown that the optimal solution must satisfy the Karush-Kuhn-Tucker (KKT) condition

$$\alpha_i[y_i(x_i^T \beta + \beta_0) - 1] = 0, \quad \forall i$$

- This means that either of these two situations holds
  - 1  $y_i(x_i^T \beta + \beta_0) = 1$ , i.e. the point is on the boundary of the margin, then  $\alpha_i > 0$
  - 2  $y_i(x_i^T \beta + \beta_0) > 1$ , i.e. the point is not on the boundary of the margin, then  $\alpha_i = 0$
- The training points having an index  $i$  such that  $\alpha_i > 0$  are called the **support vectors**.

# Support vector machines

- The **decision function** can be written as

$$h(x, \beta, \beta_0) = \text{sign}[x^T \beta + \beta_0] = \text{sign}\left[\sum_{\text{SV}} y_i \alpha_i \langle x_i, x \rangle + \beta_0\right]$$

- Attractive property of SVM: it is expressed as a function of a limited number training points (support vectors) which are on the boundaries. Points far from the class boundary do not play a major role, unlike linear discriminant models.
- In the separable case

$$C = \frac{1}{|\beta|} = \frac{1}{\sqrt{\sum_{i=1}^N \alpha_i}} \quad (1)$$



# Support vector machines (nonlinear case)

- The extension of the SVM to nonlinear classification relies on the transformation of the input variables input space by using **basis functions**.
- Consider a regression problem where  $x \in \mathcal{X} \subset \mathbb{R}^n$  and define  $m$  new transformed variables  $z_j = z_j(x), j = 1, \dots, m$ , where  $z_j(\cdot)$  is a pre-defined nonlinear transformation (e.g.  $z_j(x_1, x_2) = \log x_1 + \log x_2$ ). This is equivalent to mapping the input space  $\mathcal{X}$  into a new space, known as the **feature space**,  $\mathcal{Z} = \{z = z(x) | x \in \mathcal{X}\}$ ,
- If  $m < n$ , this boils down to a dimensionality reduction.
- We may fit a linear model  $y = \sum_{j=1}^m \beta_m z_m$  to the training data in the new input space  $z \in \mathbb{R}^m$ . By doing this, we carry out a **nonlinear fitting** of data using conventional linear techniques.

# Kernels

- A **dot-product kernel** is a function  $K$ , such that for all  $x, x' \in \mathcal{X}$

$$K(x, x') = \langle z(x), z(x') \rangle$$

where  $\langle z_1, z_2 \rangle = z_1^T z_2$  stands for the inner product and  $z(\cdot)$  is the mapping from the original to the feature space  $\mathcal{Z}$ .

- Consider a simple kernel function given by  $K(x, x') = \langle x, x' \rangle^2 = \langle z(x), z(x') \rangle$ . The feature mapping takes the form  $z(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$  and therefore comprises second order terms.
- Stationary kernels  $K(x, x') = k(x - x')$  are invariant to translations in input space, e.g. radial basis functions or Gaussian kernels.
- Extensions of kernels to handle symbolic objects greatly expand the range of problems that can be addressed.

# Support vector machines

- The parametric identification step for binary classification by SVM in the separable case requires the solution of a quadratic programming problem in the space  $\mathcal{Z}$ :

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k \langle z_i, z_k \rangle$$

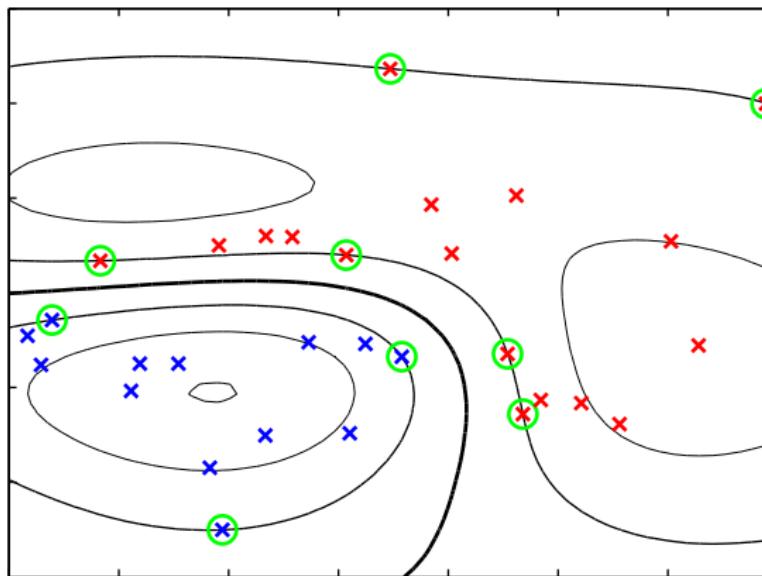
$$= \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k K(x_i, x_k)$$

$$\text{subject to } \sum_{i=1}^N \alpha_i y_i = 0, \text{ and } \alpha_i \geq 0, \quad i = 1, \dots, N$$

# Support vector machines

- The resolution of this problem differs from the linear one by the replacement of the quantities  $\langle x_i, x_k \rangle$  with  $K(x_i, x_k)$
- Whatever the dimensionality  $m$ , the SVM computation requires only the availability of the kernel matrix  $K$ .
- We don't need to know the underlying feature transformation function  $z(x)$ .
- The use of a kernel function is an attractive computational short-cut. In practice, the approach consists in defining a kernel function directly, hence implicitly defining the feature space.

# Support vector machines



**Figure:** Binary classification in two dimensions showing contours of constant  $f(x)$  obtained from a SVM with Gaussian kernel. The support vectors are in circles.

# Support vector machines

- Although the data set is not linearly separable in the two-dimensional data space  $\mathcal{X}$ , it is linearly separable in the nonlinear feature space  $\mathcal{Z}$  defined implicitly by the nonlinear kernel function  $K$ .
- Thus the training data points are perfectly separated in the original data space.

# Overlapping class distributions

- So far, we assumed that the training data points are **linearly separable** in the feature space  $\mathcal{Z}$ .
- In practice, the class-conditional distributions may overlap, in which case exact separation of the training data can lead to poor generalization.
- We need to reformulate the maximization problem so as to allow some of the training points to be misclassified. The goal is now to maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary,

$$\min_{\beta, \beta_0} \frac{1}{2} |\beta|^2 + C \sum |\xi_i| \text{ subject to } y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i, \forall i$$

- The hyperparameter  $C > 0$  controls the trade-off between the slack variable penalty and the margin.

# Overlapping class distributions

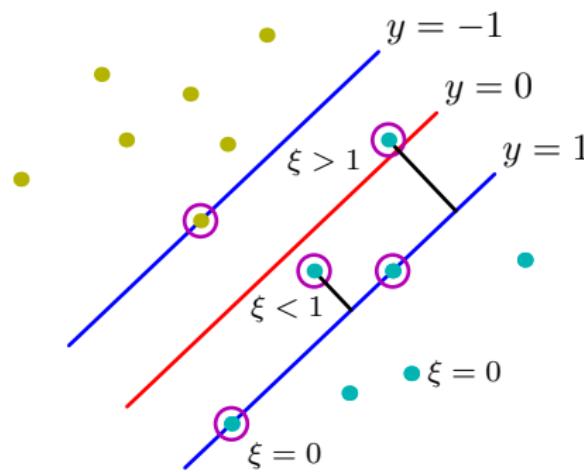


Figure: Illustration of the slack variables. Relaxing the hard margin constraint gives a soft margin and allows some of the training set data points to be misclassified.

## Overlapping class distributions

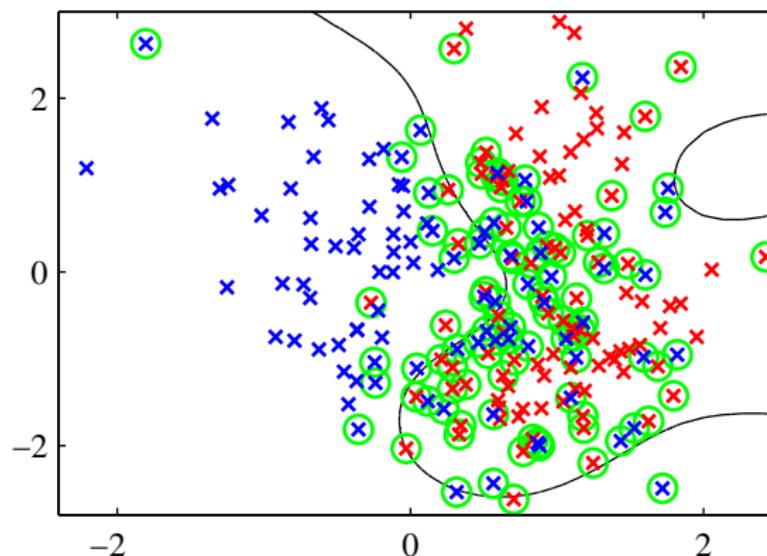


Figure: Illustration of the SVM with Gaussian kernels applied to a non separable data set in two dimensions.

Université Claude Bernard Lyon 1



## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- **Decision and regression trees**
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Decision and Regression Trees

- **Tree induction** from samples have been an active topic in the machine learning, the most representative methods of decision-tree induction are the ID3, C4 and the CART (Classification and Regression Trees) algorithm.
- A decision tree partitions the input space into **mutually exclusive regions**, each of which is assigned a procedure to characterize its data points
- An internal node is a decision-making unit that evaluates a decision function to determine which child node to visit next. A terminal node or leaf has no child nodes and is associated with one of the partitions of the input space.
- In classification trees each terminal node contains a label that indicates the class for the associated input region.

# Classification and Regression Trees

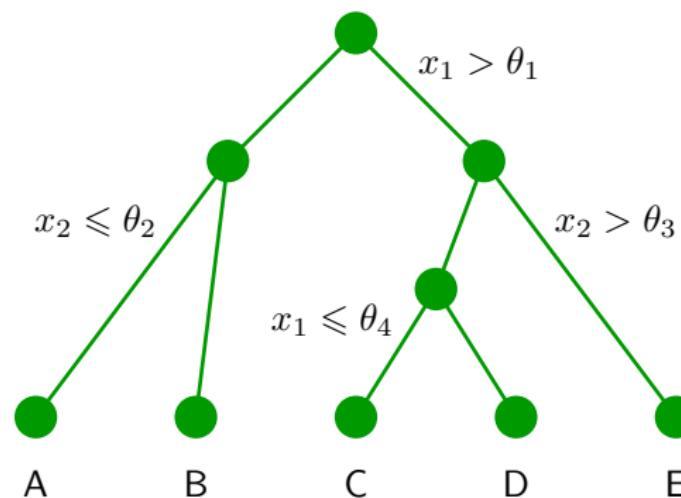


Figure: A binary decision tree.

# Decision and Regression Trees

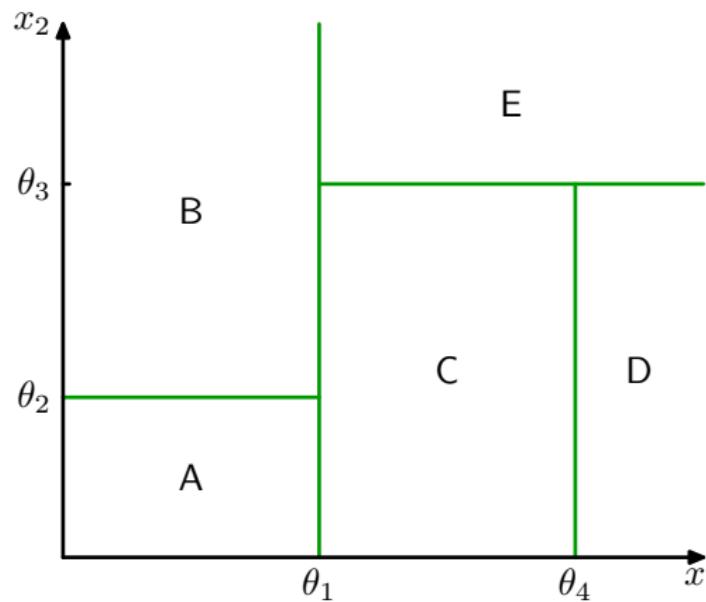


Figure: Input space partitioning induced on the input space by the binary tree.

# Decision and Regression Trees

- The structural identification in binary regression trees addresses the problem of choosing the optimal partitioning of the input space.
- Two steps of the procedure to construct an appropriate decision tree: CART first grows the tree on the basis of the training set, and then prunes the tree back based on a minimum cost-complexity principle.
- Tree growing : CART makes a succession of splits that partition the training data into disjoint subsets. Starting from the root node that contains the whole dataset, an exhaustive search is performed to find the split that best reduces a certain cost function.

# Regression Trees

- Consider a node  $t$  and let  $D(t)$  be the corresponding subset of the original  $D_N$ . The empirical error of the local model fitting the  $N(t)$  data contained in the node  $t$  is:

$$R_{emp}(t) = \min_{\alpha_t} \sum_{i=1}^{N(t)} L(y_i, h_t(x_i, \alpha_t))$$

- For any possible split  $s$  of node  $t$  into the two children  $t_r$  and  $t_l$  (and  $N(t_r) + N(t_l) = N(t)$ ), define

$$\Delta E(s, t) = R_{emp}(t) - (R_{emp}(t_l) + R_{emp}(t_r))$$

- The best split is  $s^* = \arg \max_s \Delta E(s, t)$

# Decision Trees

- Let  $p(\mathcal{C}_k)$  the probability of class  $\mathcal{C}_k$  in some node.
- Heterogeneity in the candidate node can measured by several criteria:
  - **Entropy** (*ID3, C4.5*) :  $H = - \sum_k p(\mathcal{C}_k) \log_2(p(\mathcal{C}_k))$ , bounded by  $\log_2(nb\text{classes})$  with equally balanced classes.
  - **Gini index** (*CART*) :  $Gini = 1 - \sum_k p^2(\mathcal{C}_k)$
  - **Error rate** :  $Error = 1 - \max_k(p(\mathcal{C}_k))$
- For example, with the Gini index at node  $t$ ,  $Gini(t)$ , the **homogeneity gain** of a split  $s$  at node  $t$  is given by,

$$\Delta E(s, t) = Gini(t) - \sum_{j=1}^m p(t_j) Gini(t_j)$$

where  $p(t_j)$  is the probability of reaching node  $t_j$  from node  $t$ .

# Decision and Regression Trees

- Once the best split is attained, the dataset is partitioned into the two disjoint subsets of length  $N(t_r)$  and  $N(t_l)$ , respectively. The same method is **recursively applied** to all the leaves.
- The procedure terminates either when the error measure associated with a node falls below a certain tolerance level, or when the error reduction  $\Delta E$  resulting from further splitting does not exceed a threshold value.
- The tree that the growing procedure yields is typically too large and presents a serious risk of overfitting the dataset. For that reason a **pruning procedure** is often adopted.

# Tree pruning

- Consider a fully expanded tree  $T_{max}$  characterized by  $L$  terminal nodes. Let us introduce a complexity based measure of the tree performance

$$R_\lambda(T) = R_{emp}(T) + \lambda|T| \quad (2)$$

where  $\lambda$  is a hyper-parameter that accounts for the tree's complexity and  $|T|$  is the number of terminal nodes of the tree  $T$ . For a fixed  $\lambda$  we define with  $T(\lambda)$  the tree structure which minimizes  $R_\lambda(T)$

- $\lambda$  is gradually increased in order to generate a sequence of tree configurations with decreasing complexity

$$T_L = T_{max} \supset T_{L-1} \supset \cdots \supset T_2 \supset T_1 \quad (3)$$

where  $T_i$  has  $i$  terminal nodes.

# Tree pruning

- In practice, this requires a **sequence of shrinking steps** where for each step we select the value of  $\lambda$  leading from a tree to a tree of inferior complexity.
- At the end of the shrinking process we have a sequence of candidate trees which have to be properly assessed in order to perform the structural selection.
- Cross-validation or independent testing can be used for model selection.

# Regression Trees

- Regression trees are a very easy-to-interpret representation of a nonlinear input/output mapping.
- However, these methods are characterized by rough **discontinuity at the decision boundaries** which might bring undesired effects to the overall generalization.
- Dividing the data by partitioning the input space shows typically small estimator bias but at the cost of an increased variance. This is particularly problematic in high-dimensional spaces where data become sparse.
- One response to the problem is the adoption of simple local models (e.g. constant or linear to minimize variance at the cost of an increased bias) or to use of soft splits, allowing data to lie simultaneously in multiple regions.

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

- MLC loss functions

- MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

# Combination of two estimators

- Consider two unbiased estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$  of the same parameter  $\theta$ ,

$$E[\hat{\theta}_1] = \theta \quad E[\hat{\theta}_2] = \theta$$

having equal and non zero variance

$$\text{Var}[\hat{\theta}_1] = \text{Var}[\hat{\theta}_2] = v$$

and being uncorrelated, i.e.  $\text{Cov}[\hat{\theta}_1, \hat{\theta}_2] = 0$ .

- Let  $\hat{\theta}_{\text{cm}}$  be the combined estimator

$$\hat{\theta}_{\text{cm}} = \frac{\hat{\theta}_1 + \hat{\theta}_2}{2}$$



## Combination of two estimators

- This estimator has the nice properties of being unbiased

$$E[\hat{\theta}_{\text{cm}}] = \frac{E[\hat{\theta}_1] + E[\hat{\theta}_2]}{2} = \theta$$

and with a smaller variance than the original estimators

$$\text{Var}[\hat{\theta}_{\text{cm}}] = \frac{1}{4}\text{Var}[\hat{\theta}_1 + \hat{\theta}_2] = \frac{\text{Var}[\hat{\theta}_1] + \text{Var}[\hat{\theta}_2]}{4} = \frac{\nu}{2}$$

- This trivial computation shows that the simple average of two unbiased estimators with a non zero variance returns a combined estimator with reduced variance.

# Combination of $m$ estimators

- Now, we want to estimate the unknown parameter  $\theta$  by combining a set of  $m$  estimators  $\hat{\theta}_j, j = 1, \dots, m$ . having expected values and variances given by,

$$E[\hat{\theta}_j] = \mu_j \quad \text{Var} [\hat{\theta}_j] = v_j$$

- Suppose we are interested in estimating  $\theta$  by forming a linear combination

$$\hat{\theta}_{\text{cm}} = \sum_{j=1}^m w_j \hat{\theta}_j = w^T \hat{\theta}$$

- $\hat{\theta} = [\hat{\theta}_1, \dots, \hat{\theta}_m]^T$  is the vector of estimators and  $w = [w_1, \dots, w_m]^T$  is the weighting vector.

# Combination of $m$ estimators

- The mean-squared error of the combined system is

$$\begin{aligned}\text{MSE} &= E[(\hat{\theta}_{\text{cm}} - \theta)^2] \\ &= E[(w^T \hat{\theta} - E[w^T \hat{\theta}])^2] + (E[w^T \hat{\theta}] - \theta)^2 \\ &= E[(w^T (\hat{\theta} - E[\hat{\theta}]))^2] + (w^T \mu - \theta)^2 \\ &= w^T \Omega w + (w^T \mu - \theta)^2\end{aligned}$$

where  $\Omega$  is the covariance matrix whose  $ij^{\text{th}}$  term is

$\Omega_{ij} = E[(\hat{\theta}_i - \mu_i)(\hat{\theta}_j - \mu_j)]$  and  $\mu = (\mu_1, \dots, \mu_m)^T$  is the vector of expected values.

- The MSE error has a variance and a bias terms dependent the covariance and a bias of the single estimators.



# Linear constrained combination

- A commonly used constraint is

$$\sum_{j=1}^m w_j = 1, \quad w_j \geq 0, \quad j = 1, \dots, m$$

- This means that the combined estimator is unbiased if the individual estimators are unbiased.
- The constraint can be enforced in minimizing the MSE by writing  $w$  as

$$w = (u^T g)^{-1} g$$

where  $u = (1, \dots, 1)^T$  is an  $m$ -dimensional vector of ones,  
 $g = (g_1, \dots, g_m)^T$  and  $g_j > 0, \forall j = 1, \dots, m$ .

# Linear constrained combination

- The Lagrangian function writes

$$L = w^T \Omega w + (w^T \mu - \theta)^2 + \lambda(w^T u - 1)$$

with  $\lambda$  Lagrange multiplier.

- The optimum is achieved if we set

$$g^* = [\Omega + (\mu - \theta u)(\mu - \theta u)^T]^{-1} u$$

# Linear constrained combination

- With unbiased estimators ( $\mu = \theta$ ) we obtain

$$g^* = \Omega^{-1} u$$

- With unbiased and uncorrelated estimators

$$g_j^* = \frac{1}{v_j} \quad j = 1, \dots, m$$

This means that the optimal term  $g_j^*$  of each estimator is inversely proportional to its own variance.

# Model averaging approaches

- In model selection the winner-takes-all approach is intuitively the approach which should work the best.
- However, the final model can be improved not by choosing the one that performs apparently best but by creating a model whose output is the combination of the output of the models.
- The reason is that every hypothesis  $f(\cdot)$  is only an estimate of the real target and, like any estimate, is affected by a bias and a variance term.
- A variance reduction can be obtained by simply combining uncorrelated estimators. This simple idea underlies some of the most effective techniques recently proposed in machine learning.

# Stacked regression

- Suppose we have  $m$  distinct predictors  $f_j(\cdot)$ , for  $j = 1 \dots, m$  obtained from a given training set  $D$ .
- The idea of averaging models is to design an average estimator  $\sum_{j=1}^m \beta_j f_j(\cdot)$
- Once computed the least-squares solution  $\hat{\beta}$ , the combined estimator is

$$\hat{f}(x) = \sum_{j=1}^m \hat{\beta}_j f_j(x)$$

# Stacked regression

- The  $f_j$ s are correlated because all of them are estimated on the same training set  $D$ .
- Stacked regression is an idea for combining estimators without suffering of the correlation problem. It consists in estimating the  $\hat{\beta}_j$ s by solving the following optimization task,

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N \left( y_i - \sum_{j=1}^m \beta_j f_j^{(-i)}(x_i) \right)^2 \quad \text{s.t. } \beta_j \geq 0, \forall j.$$

- $f_j^{(-i)}(x_i)$  is the leave-one-out estimate of the  $j$ th model, i.e. the predicted outcome in  $x_i$  of the  $j$ th model trained on  $D$  with the  $i$ th sample  $(x_i, y_i)$  set aside.
- Avoids giving unfairly high weight to models with higher complexity.

# Bagging

- A learning algorithm is informally called **unstable** if small changes in the training data lead to significantly different models and relatively large changes in accuracy.
- Unstable learners can have low bias but have typically high variance. Unstable methods can have their accuracy improved by **perturbing** (i.e. generating multiple versions of the predictor by perturbing the training set or learning method) and **combining**.
- **Combining multiple estimator is a variance reduction technique.**  
Bagging aims to improve accuracy for unstable learners by averaging over such discontinuities.

# Bagging

- The idea of bagging or **bootstrap aggregating** is to imitate the stochastic process underlying the realization of  $D$
- A set of  $m$  repeated bootstrap samples are taken from  $D$ . A model is built for each bootstrap sample and a final predictor is built by aggregating the  $m$  models
- In the regression case the bagging predictor is the average, in classification it is the majority vote.

# Random Forest

- **Random forests** differ from bagging in only one way: RF use a modified tree learning algorithm that selects, at each candidate split in the learning process, a random subset of the features. This process is sometimes called "feature bagging".
- The reason is to reduce the correlation of the trees in ordinary bootstrap samples.
- Typically, for a classification problem with  $p$  features,  $\sqrt{p}$  (rounded down) features are used in each split.
- The **out-of-bag** (oob) error for each data point is recorded and averaged over the forest.

# Feature importance and dissimilarity measure

- Random forests can be used to rank the **importance of variables** in a regression or classification problem in a natural way. To measure the importance of the  $j$ -th feature after training, the values of the  $j$ -th feature are permuted among the training data and the out-of-bag error is again computed on this perturbed data set.
- Random forests naturally lead to a **dissimilarity measure** between the (labeled) observations and also between unlabeled data: the idea is to construct a random forest predictor that distinguishes the “observed” data from (suitably) generated synthetic data.
- For missing values, novelty detection, variable interaction etc. Consult [www.stat.berkeley.edu/~breiman/RandomForests/cc\\_home.htm](http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm)
- Variant: Extremely randomized trees. The top-down splitting in the tree learner is randomized further. Instead of computing the locally optimal feature/split combination, a random feature + value is selected for the split.



# Boosting

- **Boosting** is one of the most powerful learning ideas introduced in the last ten years.
- Boosting is a general method which attempts to boost the accuracy of any given learning algorithm. It was originally designed for classification problems but it can profitably be extended to regression as well.
- Boosting encompasses a family of methods. The focus of boosting methods is to produce a series of “weak” learners in order to produce a powerful combination.
- A **weak learner** is a learner that has accuracy only slightly better than chance.

# Boosting

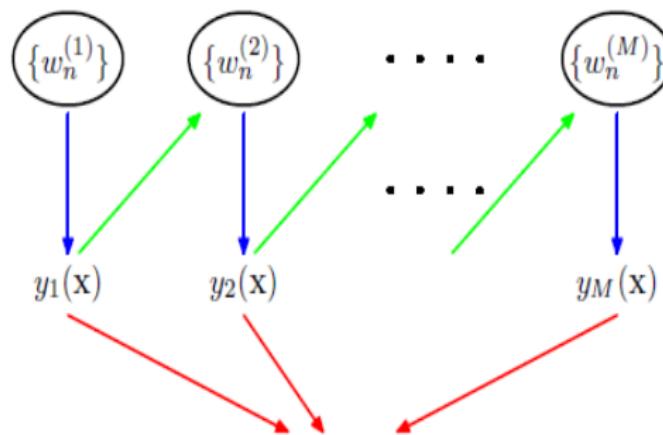
- Examples that are incorrectly predicted by previous classifiers in the series are chosen more often than examples that were correctly predicted.
- Thus Boosting attempts to produce new classifiers that are better able to predict examples for which the current ensemble's performance is poor.
- Unlike Bagging, the resampling of the training set is dependent on the performance of the earlier classifiers.
- The two most important types of boosting algorithms are the Ada Boost (Adaptive Boosting) algorithm (Freund, Schapire, 1997) and the Arcing algorithm (Breiman, 1996).

# Ada Boost

- Consider a binary classification problem where the output take values in  $\{-1, 1\}$ .
- A weak classifier is one whose misclassification error rate is only slightly better than random guessing.
- The purpose of boosting is to sequentially apply the weak classification algorithm to repeatedly modified versions of the data, thereby producing a sequence of classifiers  $f_j(\cdot)$ .
- The predictions of the  $m$  weak classifiers are then combined through a weighted majority vote to produce the final prediction

$$\hat{f}(x) = \text{sign} \left( \sum_{j=1}^m \alpha_j f_j(x) \right)$$

# Ada Boost



$$Y_M(x) = \text{sign} \left( \sum_m^M \alpha_m y_m(x) \right)$$

# Ada Boost

- The weights  $\alpha_j$  of the different classifiers are computed by the algorithm. The idea is to give higher influence to the more accurate classifiers in the sequence.
- At each step, the boosting algorithm samples  $N$  times from a distribution that puts a weight  $w_i$  on each sample  $(x_i, y_i)$  of  $D$ .
- Initially the weights are all set to  $w_i = 1/N$ . Then, the weights are individually modified and the classification algorithm is re-applied to the resampled training set.

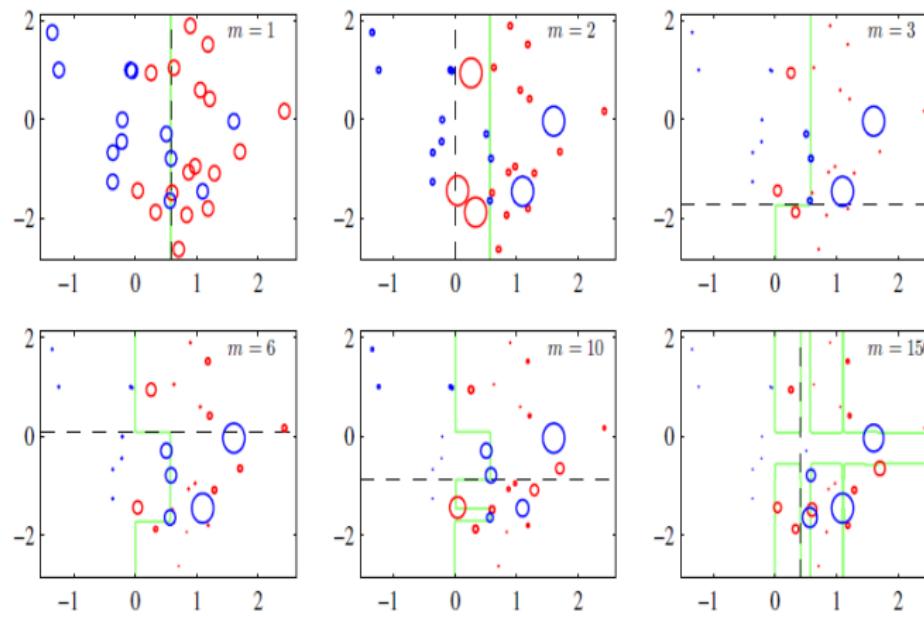
# Ada Boost

- Initialize the weights to  $w_i = 1/N$
- For  $j = 1$  to  $m$ ,
  - Fit a classifier  $f_j(\cdot)$  on  $D$  using weights  $w_i$  and compute the misclassification error on the training set

$$\hat{\epsilon}_j = \frac{\sum_{i=1}^N w_i \mathbb{I}[y_i \neq f_j(x_i)]}{\sum_{i=1}^N w_i}$$

- Let  $\alpha_j = \log((1 - \hat{\epsilon}_j)/\hat{\epsilon}_j)$
- Set  $w_i \leftarrow w_i \begin{cases} \exp[-\alpha_j] & \text{if correctly classified} \\ \exp[\alpha_j] & \text{if incorrectly classified} \end{cases}$
- Normalize the weights.
- Output  $\hat{f}(x) = \text{sign} \left( \sum_{j=1}^m \alpha_j f_j(x) \right).$

# Ada Boost



# Ada Boost

- Boosting has its roots in a theoretical framework for studying machine learning called the PAC learning model.
- The empirical error of the final hypothesis  $f_{\text{boo}}$  is at most

$$\prod_{j=1}^m \left[ 2\sqrt{\hat{\epsilon}_j \cdot (1 - \hat{\epsilon}_j)} \right]$$

- Boosting is simple and easy to program. Moreover, it has few parameters (e.g. max number of classifiers) to tune.
- Instead of trying to design a learning algorithm which should be accurate over the entire space, one can instead focus on finding weak algorithms that only need to be better than random.
- A nice property of Ada Boost is its ability to identify outliers (hard samples).

# Gradient boosting

- Like other boosting methods, **gradient boosting** combines weak "learners" into a single strong learner in an iterative fashion
- The gradient boosting method assumes a real-valued  $y$  and seeks an approximation  $\hat{f}(x)$  in the form of a weighted sum of functions

$$\hat{f}(x) = \sum_{j=1}^m \alpha_j f_j(x)$$

- The idea is to apply a steepest descent step to this minimization problem with the following equations

$$f_m(x) = f_{m-1}(x) - \gamma_m \sum_{i=1}^n \nabla_{f_{m-1}} L(y, f_{m-1}(x))$$

where the derivatives are taken with respect to the functions  $f_i$

# Gradient boosting

Given a training set  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable loss function  $L(y, f(x))$ , the generic gradient boosting method is:

- 1 Initialize model with a constant value  $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$ .
- 2 For  $m = 1$  to  $M$ :

- 1 Compute so-called pseudo-residuals:

$$r_{im} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x)=f_{m-1}(x)} \quad \forall i = 1, \dots, n.$$

- 2 Fit a base learner  $h_m(x)$  to pseudo-residuals  $\{(x_i, r_{im})\}_{i=1}^n$ .
- 3 Compute  $\gamma_m$  by solving the 1D optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, f_{m-1}(x_i) + \gamma h_m(x_i)).$$

- 4 Update the model  $f_m(x) = f_{m-1}(x) + \gamma_m h_m(x)$ .
- 3 Output:  $f_M(x)$ .

# Gradient boosting

- Gradient boosting is typically used with decision trees (especially CART trees) of a fixed size as base learners.
- $J$ , the number of terminal nodes in trees, can be adjusted for a data set at hand. It controls the level of interaction between variables in the model. With  $J = 2$  (decision stumps), no interaction between variables is allowed. Typically  $4 \leq J \leq 8$  work well.
- *XGBoost* is an open-source software library which provides the gradient boosting framework. It now has integrations with scikit-learn for Python users and supports the distributed processing frameworks Apache Hadoop, Apache Spark, and Apache Flink.
- *XGBoost* became popular among the *Kaggle* community where it has been used for a large number of competitions.

# Arcing

- **Adaptive Resampling and Combining** (ARCing) was proposed as a modification of the original Ada Boost algorithms by Breiman.
- It is based on the idea that the success of boosting is related to the adaptive resampling property where increasing weight is placed on those samples more frequently misclassified.
- The complex updating equations of Ada Boost are replaced by much simpler formulations. The final classifier is obtained by unweighted voting

# Arcing

- Initialize the weights to  $w_i = 1/N$
- For  $j = 1$  to  $m$ ,
  - Fit a classifier  $f_j(\cdot)$  to the training data obtained by resampling using weights  $w_i$
  - Compute the  $e_i$ : the number of misclassifications of the  $i$ th sample by the  $j$  classifiers  $f_1, \dots, f_j$
  - The updated weights are defined by

$$w_i = \frac{1 + e_i^4}{\sum_{i=1}^N (1 + e_i^4)}$$

- The output is obtained by unweighted voting of the  $m$  classifiers  $h_j$ ,

$$f_{\text{arc}}(x) = \text{sign} \left( \sum_{j=1}^m f_j(x) \right)$$
.Université Claude Bernard Lyon 1

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## ■ MLC loss functions

## ■ MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material

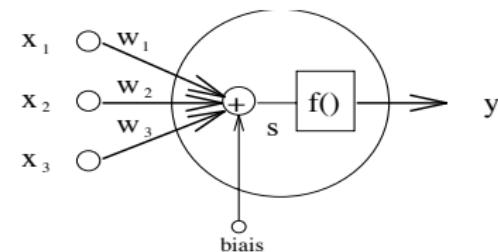


# Artificial neural networks

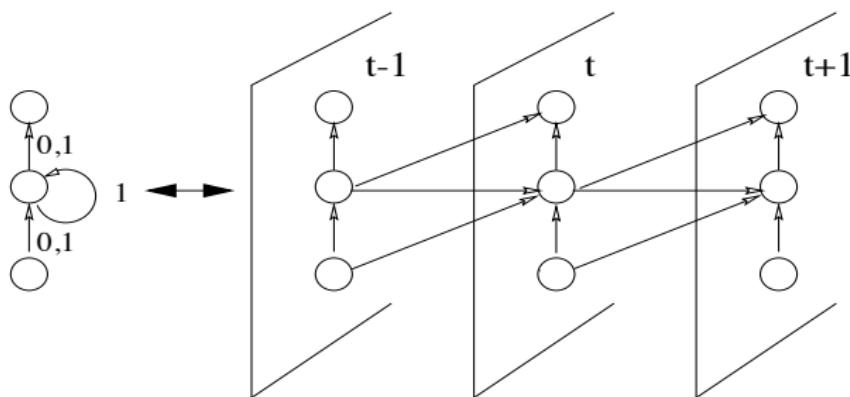
- The 'standard' **neuron**:

$$y = f\left(\sum_{j=0}^N w_j x_j\right)$$

- **Linear combination** of the input signals.
- $w_j$  are **ajustable parameters**.
- $y$  is the neuron **activation**.
- $x_0$  is clamped to 1.0,  $w_0$  is the bias term.
- The **activation function**  $f$  is typically non-linear, bounded and derivable over  $]-\infty, \infty[$ .



# Time-delayed networks



- TDNN unfolded in time.

# Lorenz attractor

- Lorenz equations:

$$\begin{aligned}\frac{dx}{dt} &= -\sigma x + \sigma y \\ \frac{dy}{dt} &= -xz + rx - y \\ \frac{dz}{dt} &= xy - bz\end{aligned}$$

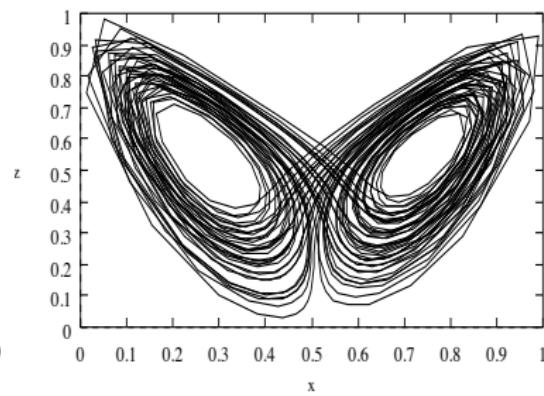
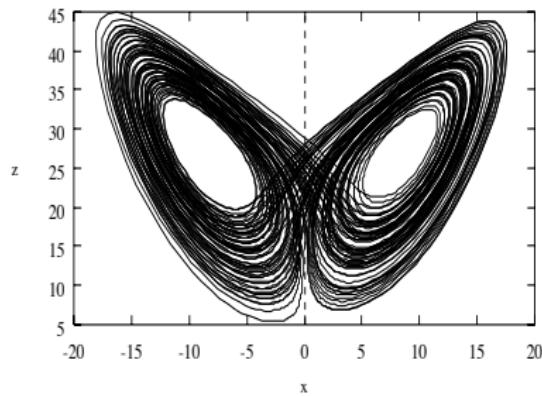
- Configuration:

$$[\hat{x}(k), \hat{z}(k)] = \mathcal{N}_k(\mathbf{w}, x(k-1))$$

- The TDNN has dimension  $1 \times 10 \times 2$  with delays  $(30, 1, 0)$  and training is performed over 4000 points.



# Lorenz attractor



- True vs. reconstructed attractor (i.e.  $[x(k), z(k)]$  vs.  $[\hat{x}(k), \hat{z}(k)]$ ).

# Ikeda attractor

- **Ikeda** equations:

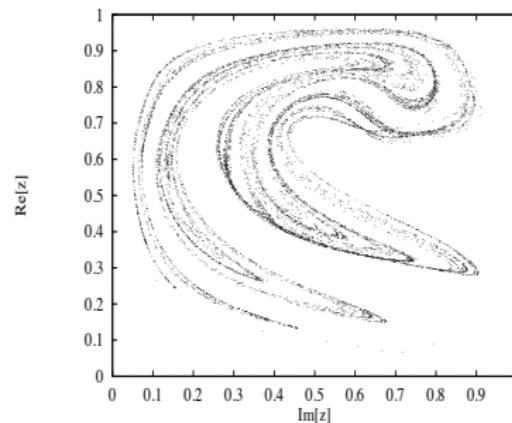
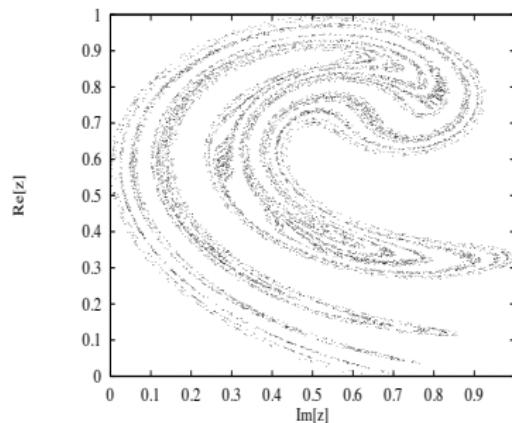
$$z(k+1) = a + bz(k)\exp(i[\phi - c/(1 + |z(k)|^2)])$$

- Configuration:

$$(\text{Re}[\hat{z}(k+1)], \text{Im}[\hat{z}(k+1)]) = \mathcal{N}_k(\mathbf{w}, \text{Re}[z(k)])$$

- The TDNN has dimension  $1 \times 15 \times 2$  with delays  $(5, 1, 0)$  and training is performed over 10,000 points.

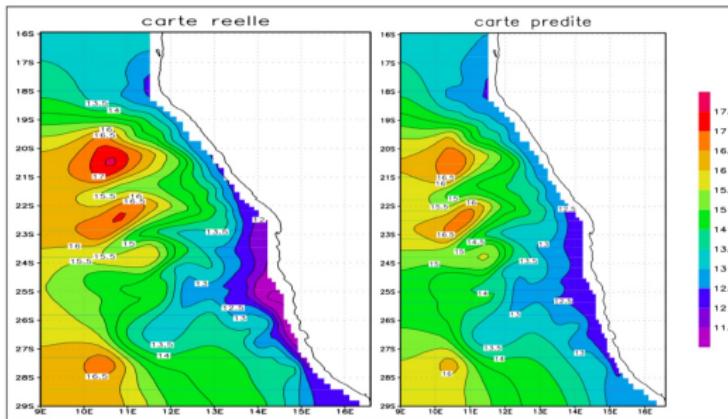
# Ikeda attractor



- True vs. reconstructed attractor (i.e.  $[\text{Re}[z(k)], \text{Im}[z(k)]$  vs.  $[\text{Re}[\hat{z}(k)], \text{Im}[\hat{z}(k)]$ ).

# Weather forecasts

- Reconstruction of 2D sea temperature maps
- A recurrent NN is trained at each point in the grid. Predictions are combined by cubic interpolation.



# Convolutional neural networks

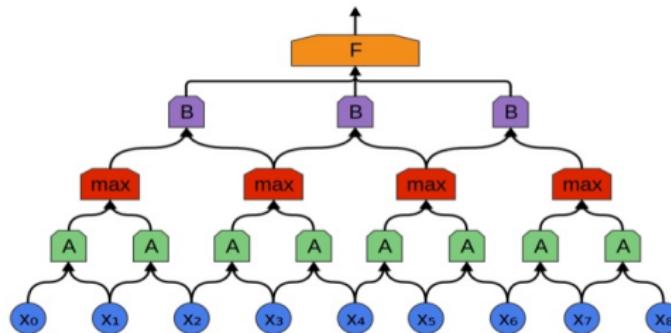


Figure from colah.github.io "Understanding LSTM Networks"

- A convolutional layer  $A$  look at small time segments of our data, computing certain features.
- A max-pooling layer takes the maximum of features over small blocks of a previous layer (kind of zoom out).
- $B$  is used to create another convolutional layer stacked on top of the previous one.  $F$  is a fully-connected layer.

# Convolutional neural networks

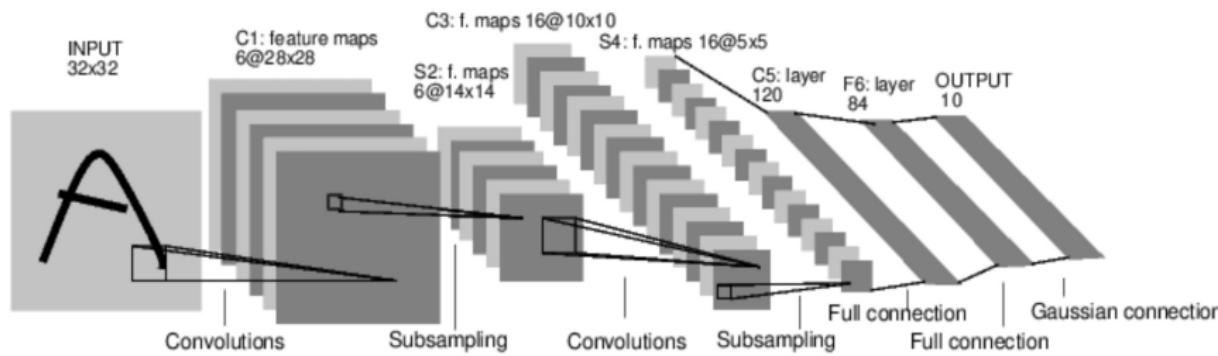


Figure from Pytorch tutorials

- For example, this network that classifies digit images.

# Python and Pytorch

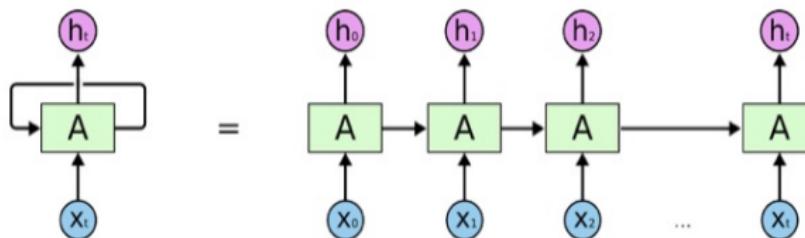
```
import torch.nn as nn
import torch.nn.functional as F

class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        # 1 input image channel, 6 output channels, 5x5 square conv. kernel
        self.conv1 = nn.Conv2d(1, 6, 5)
        self.conv2 = nn.Conv2d(6, 16, 5)
        # an affine operation: y = Wx + b
        self.fc1 = nn.Linear(16 * 5 * 5, 120)
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        # Max pooling over a (2, 2) window
        x = F.max_pool2d(F.relu(self.conv1(x)), (2, 2))
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(-1, self.num_flat_features(x))
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x

    def num_flat_features(self, x):
        size = x.size()[1:] # all dimensions except the batch dimension
        num_features = 1
        for s in size:
            num_features *= s
        return num_features
```

## Sequence learning with recurrent networks



An unrolled recurrent neural network.

Figure from colah.github.io "Understanding LSTM Networks"

- A loop allows information to be passed from one step of the network to the next, allowing information to persist.

## Sequence learning with recurrent networks

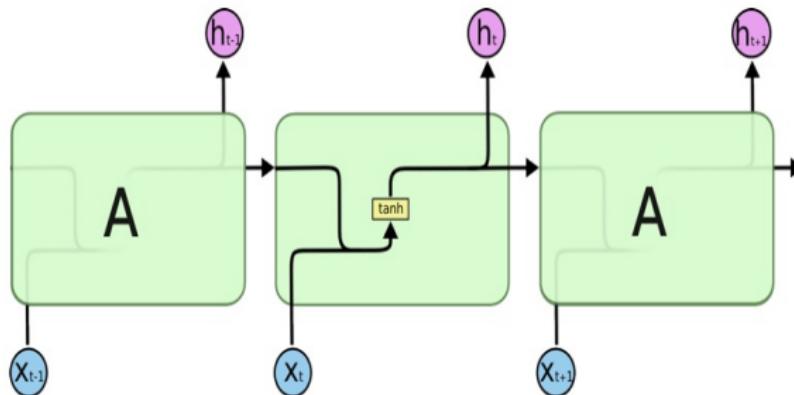


Figure from colah.github.io "Understanding LSTM Networks"

- Such RNN models are not always capable of handling long-term dependencies (i.e. when the time-gap between the relevant information and the moment where it is needed becomes very large).

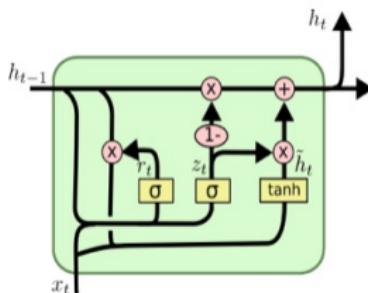
# Elman RNN

- Applying a multi-layer Elman RNN with tanh or ReLU non-linearity on an input sequence.
- For each element in the input sequence, each layer computes the following function:

$$h_t = \tanh(w_{ih} * x_t + b_{ih} + w_{hh} * h_{(t-1)} + b_{hh})$$

$h_t$  is the hidden state at time  $t$ , where  $h_t$  is the hidden state at time  $t$ , and  $x_t$  is the hidden state of the previous layer at time  $t$  or input  $t$  for the first layer.

# Gated Recurrent Unit



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Figure from colah.github.io "Understanding LSTM Networks"

- A standard variant of the LSTM is the multi-layer Gated Recurrent Unit (GRU, Cho, et al. 2014). The forget and input gates are merged into a single “update gate.” It also merges the cell state and hidden state.
- The resulting model is simpler than standard LSTM models, and has been increasingly popular.

# LSTM

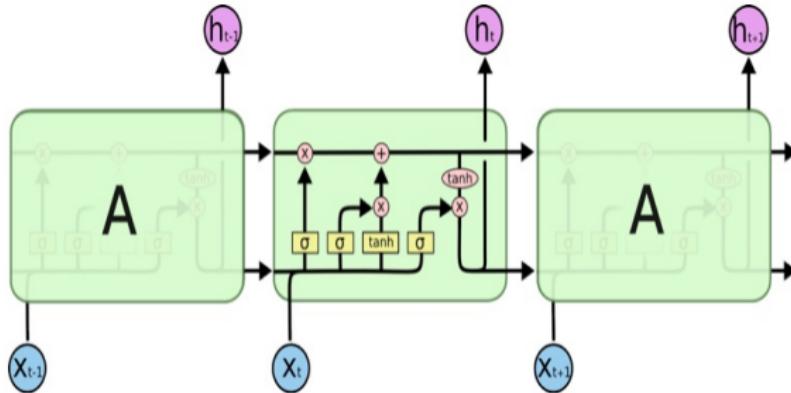


Figure from colah.github.io "Understanding LSTM Networks"

- Long Short Term Memory networks (LSTM) are a special kind of RNN, capable of learning long-term dependencies (Hochreiter et al. 1997). There are two gates: the forget and input gates that controls the information flow and the memory.

- LSTM RNN applied to an input sequence. For each element in the input sequence, each layer computes the following function:

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{(t-1)} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{(t-1)} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{(t-1)} + b_{hg}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{(t-1)} + b_{ho}) \\ c_t &= f_t * c_{(t-1)} + i_t * g_t \\ h_t &= o_t * \tanh(c_t) \end{aligned}$$

$h_t$  is the hidden state,  $c_t$  is the cell state,  $x_t$  is the hidden state of the previous layer at time  $t$ , and  $i_t, f_t, g_t, o_t$  are the input, forget, cell, and out gates, respectively.

# Neural-Transfer

- The Neural-Style (or Neural-Transfer) algorithm takes as input a content-image (e.g. a turtle), a style-image (e.g. artistic waves) and return the transformed content-image as if it was ‘painted’ using the artistic style of the style-image:

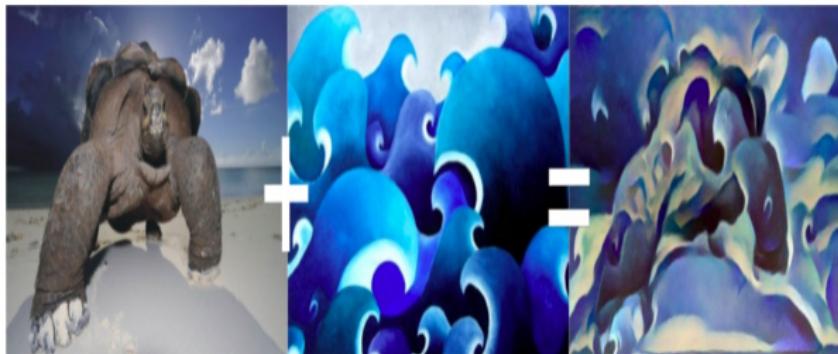


Figure from (Gatys, Ecker and Bethge 2015)

# Neural-Transfer

- Principle: define two distances, one for the content (DC) and one for the style (DS).
- DC measures how different the content is between two images, while DS measures how different the style is between two images.
- Then, we take a third image, the input with some noise, and transform it in order to both minimize its content-distance with the content-image and its style-distance with the style-image.
- The distance is based on the feature maps at all depth layers of a pre-trained deep convolutional neural network applied on the input.

# Deep convolutional generative adversarial networks

- Learning reusable feature representations from large unlabeled datasets has been an area of active research.
- In computer vision, one can leverage the practically unlimited amount of unlabeled images and videos to learn good intermediate representations, which can then be used on a variety of supervised learning tasks such as image classification.
- Good image representations can be built by (GANs) (Goodfellow et al., 2014), and later reusing parts of the generator and discriminator networks as feature extractors for supervised tasks

## Walking in the latent space

- Walking in the latent space results in semantic changes to the image generations (such as objects being added and removed), the model has learned relevant and interesting representations

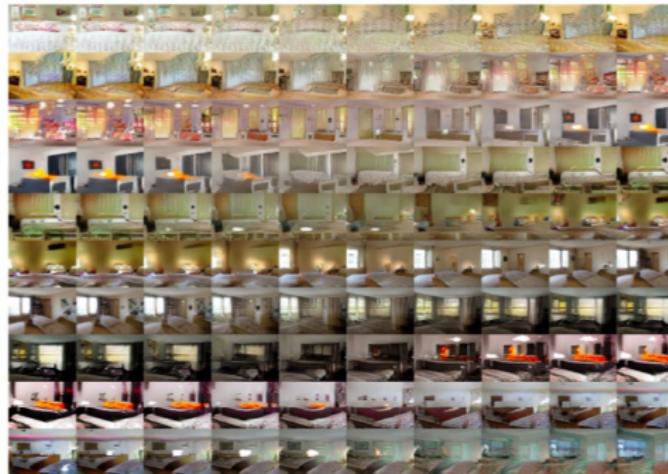


Figure from (Radford, Metz, Chintala, ICLR 2016)

Université Claude Bernard Lyon 1

- Every image in the space plausibly looking like a bedroom (e.g. TV slowly being transformed into a window)

## Vector arithmetic on face samples

- In the context of evaluating learned representations of words simple arithmetic operations reveal rich linear structure in representation space (Mikolov et al., 2013).

$$\text{Vector("King")} - \text{vector("Man")} + \text{vector("Woman")} = \text{vector("Queen")}$$

- A similar structure emerges in the latent representation of the generators.

# Vector arithmetic on face samples

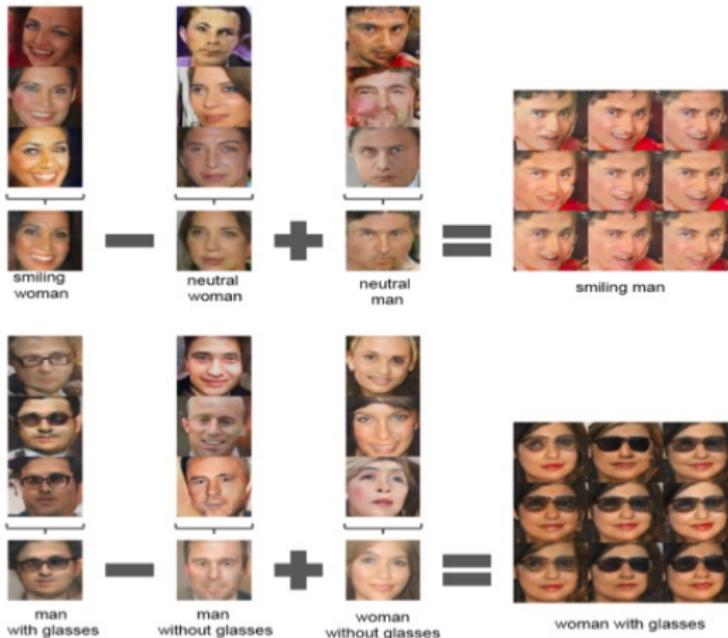


Figure from (Radford, Metz, Chintala, ICLR 2016)

Université Claude Bernard



# Cycle-Consistent Adversarial Networks

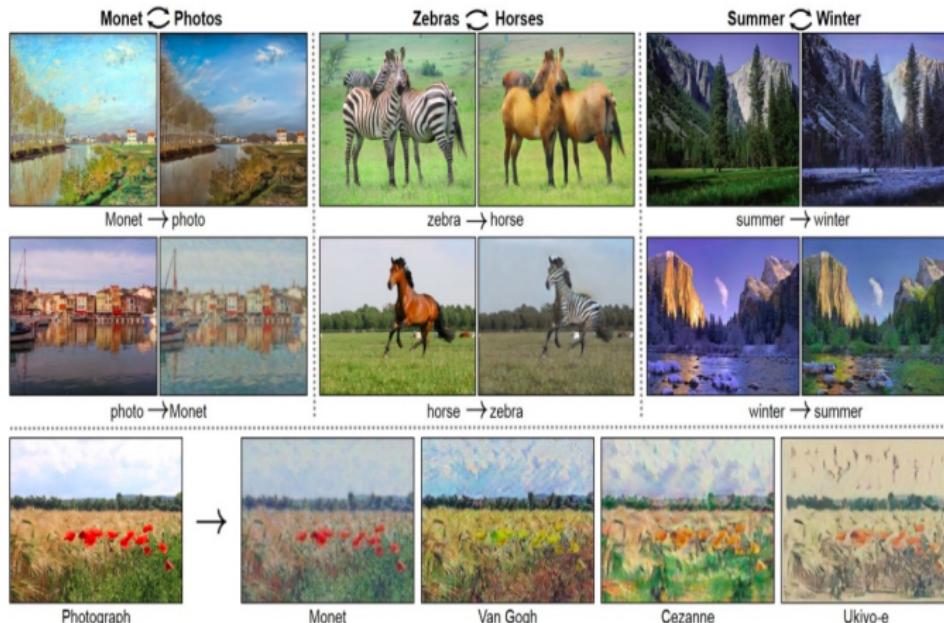


Figure from (Zhu, Park, Isola, arxiv. 2017)

## References: Deep in Image

Jun-Yan Zhu, Taesung Park, Phillip Isola, Alexei A. Efros. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In arxiv, 2017.

Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, Alexei A. Efros. Image-to-Image Translation with Conditional Adversarial Networks. In CVPR 2017.

Leon A. Gatys, Alexander S. Ecker, Matthias Bethge. A Neural Algorithm of Artistic Style. In arxiv, 2015.

# Statistical Machine Translation

- A neural machine translation system is a neural network that directly models the conditional probability  $p(y|x)$  of translating a source sentence,  $x_1, \dots, x_n$ , to a target sentence,  $y_1, \dots, y_m$ .
- A basic form of NMT consists of two components:
  - 1 an encoder which computes a representation  $s$  for each source sentence,
  - 2 a decoder which generates one target word at a time and hence decomposes the conditional probability as:

$$\log p(y|x) = \sum_{j=1}^m \log(p(y_j|y_{<j}, s))$$

- A natural choice to model such a decomposition in the decoder is to use a recurrent neural network (RNN) architecture (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Cho et al., 2014; Bahdanau et al., 2015; Luong et al., 2015; Jean et al., 2015).



# Statistical Machine Translation

- The probability of decoding each word  $y_j$  is parameterized as:

$$p(y_j | y_{<j}, s) = \text{softmax}(g(h_j))$$

- $g$  is the transformation function that outputs a vocabulary-sized vector,
- $h_j$  is the RNN hidden unit, abstractly computed as:

$$h_j = f(h_{j-1}, s)$$

- $f$  computes the current hidden state given the previous hidden state and can be either a standard RNN, a GRU, or an LSTM unit.
- In earlier works,  $s$  was used once to initialize the decoder hidden state. When  $s$  depends selectively on the encoder hidden states, the approach is referred to as an **attention mechanism** and was proposed in (Bahdanau et al., 2015; Jean et al., 2015).

# Statistical Machine Translation

- Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation
- A RNN Encoder-Decoder consists of two recurrent neural networks (RNN). One RNN encodes a sequence of symbols into a fixed-length vector representation, and the other decodes the representation into another sequence of symbols.
- The encoder and decoder of the proposed model are jointly trained to maximize the conditional probability of a target sequence given a source sequence.
- The conditional probabilities of phrase pairs are computed by the RNN Encoder-Decoder
- Qualitatively, the model learns a semantically and syntactically meaningful representation of linguistic phrases.

# Learning to Align and Translate

- The use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture,
- The model can be extended by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly.

# Neural Conversational Model

- Conversational modeling is an important task in natural language understanding and machine intelligence. Although previous approaches exist, they are often restricted to specific domains (e.g., booking an airline ticket) and require hand-crafted rules.
- a simple approach for this task uses the sequence to sequence framework. Our model converses by predicting the next sentence given the previous sentence or sentences in a conversation.
- it can be trained end-to-end and thus requires much fewer hand-crafted rules. We find that this straightforward model can generate simple conversations given a large conversational training dataset.
- despite optimizing the wrong objective function, the model is able to converse well.
- On a domain-specific IT helpdesk dataset, the model can find a solution to a technical problem via conversations.
- On a noisy open-domain dataset, the model can perform simple forms of common sense reasoning, but the lack of consistency is a common failure mode



# References: Deep in Text

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. In arxiv, 2017.

Ilya Sutskever, Oriol Vinyals, Quoc V. Le. Sequence to Sequence Learning with Neural Networks. NIPS 2014.

Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. ICLR 2015.

Oriol Vinyals, Quoc Le. A Neural Conversational Model. ICML 2015.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. In Proceedings of Workshop at ICLR, 2013.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed Representations of Words and Phrases and their Compositionality. In Proceedings of NIPS, 2013.

Minh-Thang Luong, Hieu Pham, Christopher D. Manning. Effective Approaches to Attention-based Neural Machine Translation. In arXiv 2015.

# Outline

## 1 Introduction

- Machine Learning: An overview
- Typical problems
- The Supervised Learning setting
- The curse of dimensionality
- Polynomial curve fitting
- The Bias-Variance decomposition

## 2 Classifier Evaluation

- Cross-validation
- Evaluation of binary classifiers

## 3 Multi-class Classification

- Error-correcting output code

## 4 Multi-Label Classification

## 5 MLC loss functions

## 6 MLC learning methods

## 5 Classification and Regression Models

- Linear regression
- Logistic regression
- Shrinkage methods
- Naive Bayes classifiers
- Support vector machines
- Decision and regression trees
- Ensemble Methods

## 6 Neural networks

## 7 References & Further Material



# Useful MOOCs

- French:

- FUN : <https://www.fun-mooc.fr/>
- EPFL : <http://moocs.epfl.ch/> (FR+EN)

- English:

- Coursera : <http://coursera.org>
- Stanford Online : <http://lagunita.stanford.edu>
- MIT OpenCourseWare : <https://ocw.mit.edu>
- edX : <https://www.edx.org/>
- Udacity : <http://udacity.com>



# References

- C. M. Bishop. *Pattern Recognition and Machine Learning*, Springer, 2006.
- G. Bontempi and S. Ben Taieb. *Statistical foundations of machine learning*, OTexts: Melbourne, Australia, 2016.
- I. Goodfellow, Y. Bengio and A. Courville. *Deep Learning*, MIT Press, 2016.
- G. Obozinski and F. Bach. *Introduction to Graphical Models*, Master "Mathématiques Appliquées", Parcours "Mathématiques, Vision et Apprentissage", ENS Cachan 2017.
- A. Aussem *Dynamical recurrent neural networks towards prediction and modeling of dynamical systems*. Neurocomputing, 1999.
- W. Waegeman, K. Dembczynski, E. Hüllermeier. *Multi-target Prediction, International Conference on Machine Learning*, Atlanta, USA, June 2013.