

Machine Learning

Master 2 Computer Science

Aurélien Garivier

2018-2019



Table of contents

1. Before we start
2. What is Machine Learning?
3. The Learning Models
4. Machine Learning Methodology
5. Supervised Classification
6. Nearest-Neighbor Classification
7. Deviation Bound for Bernoulli Variables

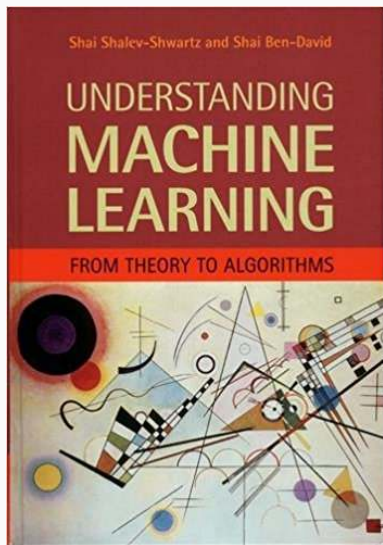
Before we start

Outline (1/2)

- 1. 09.10 Introduction, nearest-neighbor classification
 - 09.17 Pot du DI at IFE Descartes (meet your tutor)
- 2. 09.24 ML methodology, k-nearest neighbors, decision trees
- 3. 10.01 PAC Learning Theory, no-free-lunch theorem
- 4. 10.8 VC dimension, empirical risk minimization
- 5. 10.15 Linear separators, Support Vector Machines
- 6. 10.22 Kernels, regularization
 - 10.29 holidays

Outline (2/2)

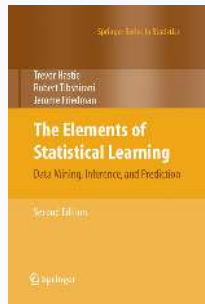
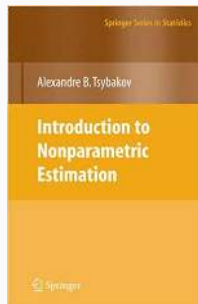
- 7. 11.5 Boosting, Bagging, Random Forests
 - 11.12 winter school
- 8. 11.19 Neural networks and stochastic gradient descent
 - 11.26 winter school
 - 12.03 no lecture
- 9. 12.10 Regression, model selection
- 10. 12.17 Dimension reduction
- 11. 01.07 Clustering
- 12. 01.14 Online learning



General introduction to Machine Learning theory, by two leading researchers of the field.

Covers most of the content of this course (the converse is also almost true).

Additional References



Homework and in-class exercises

and

- analysis and review of a research article (report + oral presentation)
- or participation in a student ML challenge

(you choose)

What is Machine Learning?

Machine Learning (ML): Definition

Arthur Samuel (1959)

Field of study that gives computers the ability to learn without being explicitly programmed

Tom M. Mitchell (1997)

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .

ML: Learn from and make predictions on data

- Algorithms operate by building a model from **example** inputs in order to make data-driven **predictions or decisions**...
- ...rather than following strictly static program instructions: useful when designing and programming explicit algorithms is unfeasible or poorly efficient.

Within Artificial Intelligence

- evolved from the study of pattern recognition and computational learning theory in artificial intelligence.
- AI: emulate cognitive capabilities of humans (big data: humans learn from abundant and diverse sources of data).
- a machine mimics "cognitive" functions that humans associate with other human minds, such as "learning" and "problem solving".

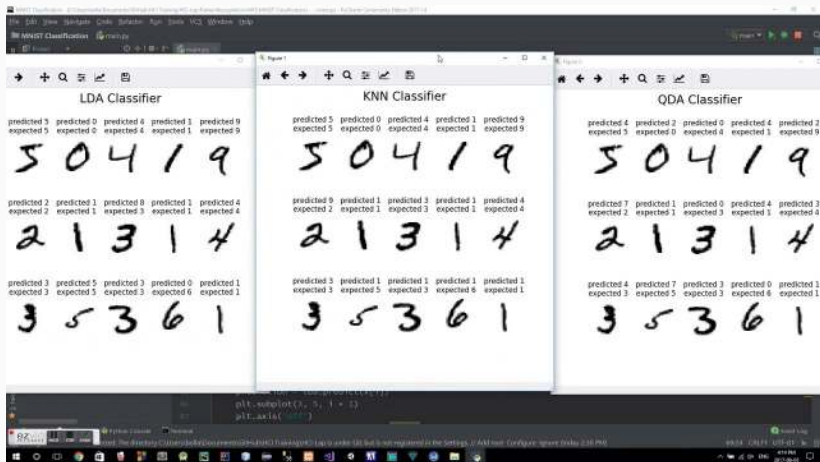
Machine Learning: Typical Problems

- spam filtering, text classification
- optical character recognition (OCR)
- search engines
- recommendation platforms
- speech recognition software
- computer vision
- bio-informatics, DNA analysis, medicine
- etc.

For each of this task, it is possible but very inefficient to write an explicit program reaching the prescribed goal.

It proves much more succesful to have a machine infer what the good decision rules are.

Example: MNIST dataset



Related Fields

- **Computational Statistics**: focuses in prediction-making through the use of computers together with statistical models (ex: Bayesian methods).
- **Statistical Learning**: ML by statistical methods, with statistical point of view (probabilistic guarantees: consistency, oracle inequalities, minimax)
- **Data Mining** (unsupervised learning) focuses more on exploratory data analysis: discovery of (previously) unknown properties in the data. This is the analysis step of Knowledge Discovery in Databases.
- Importance of **probability**- and **statistics**-based methods → **Data Science** (Michael Jordan)
- Strong ties to **Mathematical Optimization**, which delivers methods, theory and application domains to the field

Machine Learning and Statistics

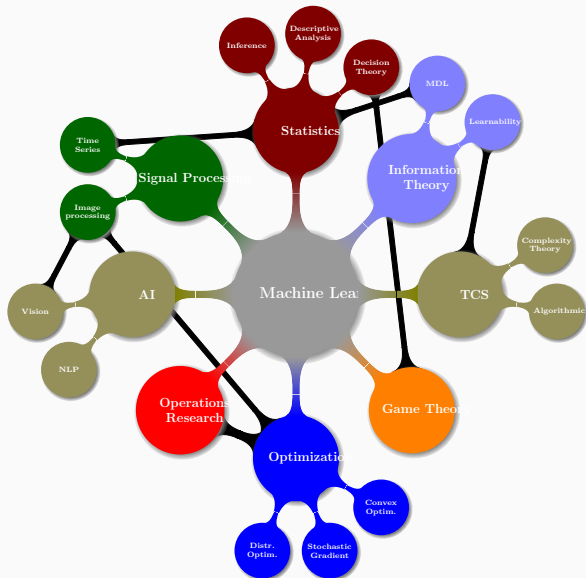
- Data analysis (inference, description) is the goal of statistics for long.
- Machine Learning has more **operational** goals (ex: consistency is important the statistics literature, but often makes little sense in ML).

Models (if any) are *instrumental*.

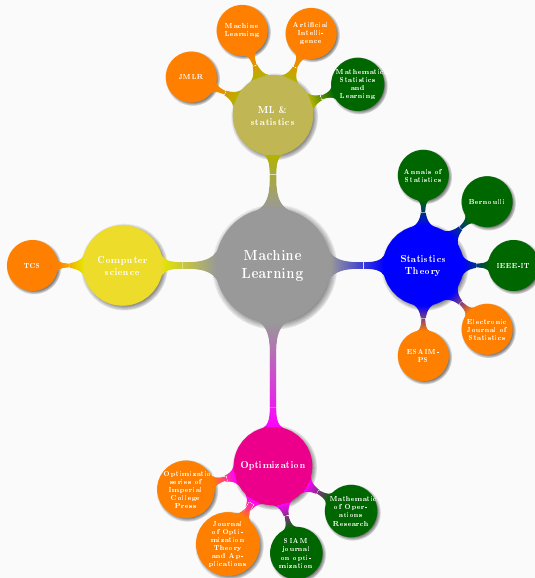
Ex: linear model (nice mathematical theory) vs Random Forests.

- Machine Learning/big data: no separation between statistical modelling and optimization (in contrast to the statistics tradition).
- In ML, data is often here before (unfortunately).
- ML more focused on *correlation*, less on *causality*.
- Algorithmic considerations play a major role in ML.
- No clear separation (statistics evolves as well), but different hypotheses focus of interest. Ex: model-free versus model-based, asymptotic consistency versus finite sample bounds.

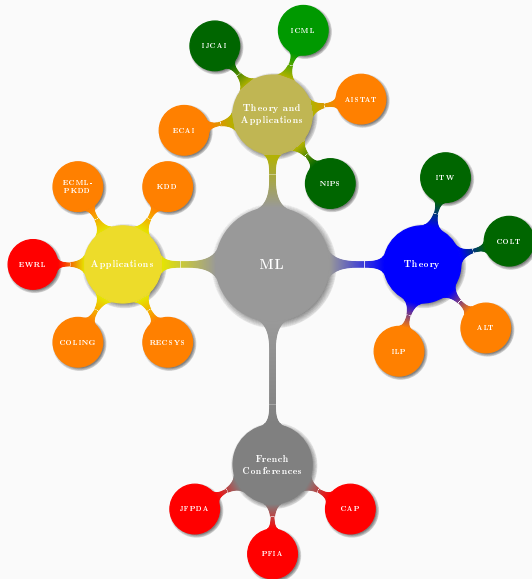
ML and its neighbors



ML journals

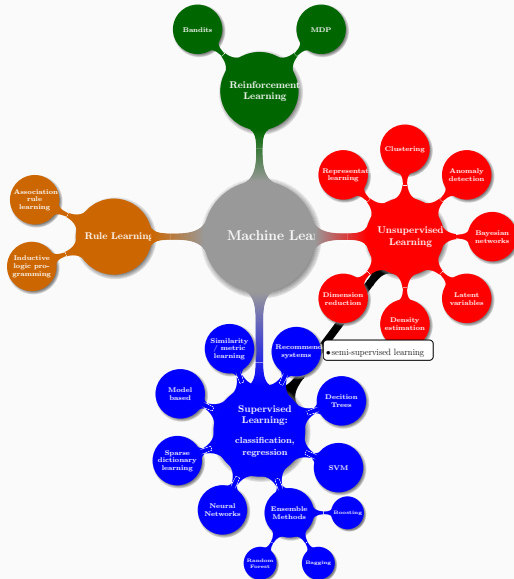


ML conferences



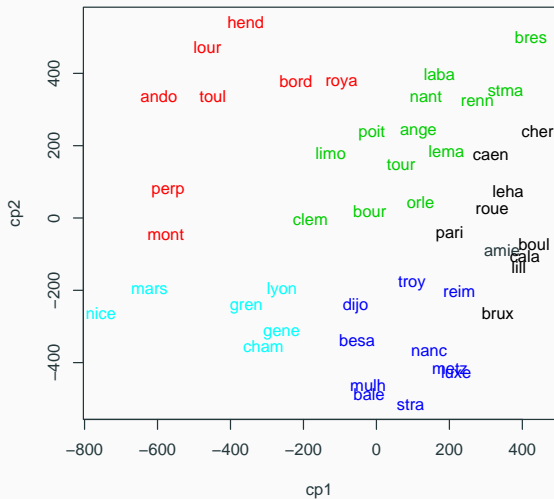
The Learning Models

What ML is composed of



- (many) observations on (many) individuals
- need to have a simplified, structured overview of the data
- *taxonomy*: untargeted search for *homogeneous clusters* emerging from the data
- Examples:
 - customer segmentation
 - image analysis (recognizing different zones)
 - exploration of data

Example: representing the climate of cities



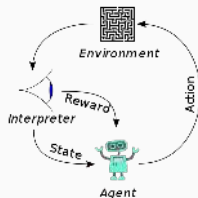
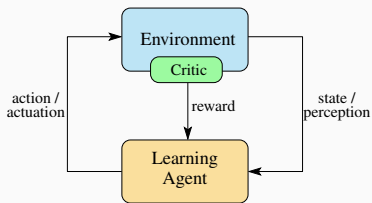
Supervised Learning

- Observations = pairs (X_i, Y_i)
- Goal = learn to *predict* Y_i given X_i
- Regression (when Y is continuous)
- Classification (when Y is discrete)

Examples:

- Spam filtering / text categorization
- Image recognition
- Credit risk ranking

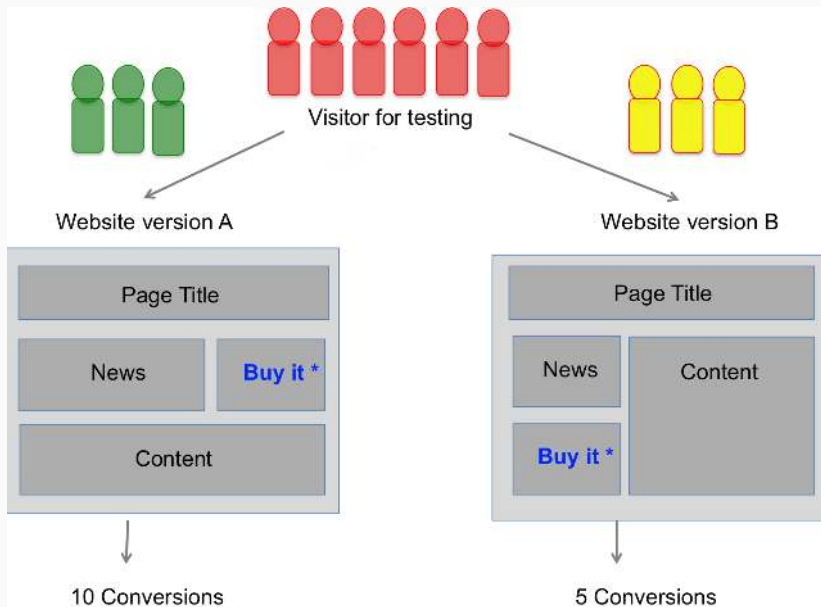
Reinforcement Learning



[Src: https://en.wikipedia.org/wiki/Reinforcement_learning]

- area of machine learning inspired by behaviourist psychology
- how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.
- Model: random system (typically : Markov Decision Process)
 - agent
 - state
 - actions
 - rewards
- sometimes called approximate dynamic programming, or neuro-dynamic programming

Example: A/B testing



Machine Learning Methodology

n -by- p matrix X

- n examples = points of observations
- p features = characteristics measured for each example

Questions to consider:

- Are the features centered?
- Are the features normalized? bounded?

In `scikitlearn`, all methods expect a 2D array of shape (n, p) often called

`X (n_samples, n_features)`

- Inside R: package datasets
- Inside scikitlearn: package sklearn.datasets
- UCI Machine Learning Repository



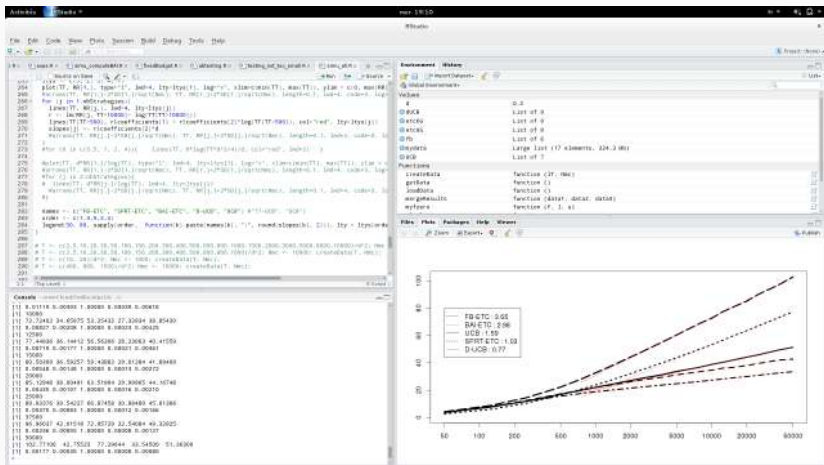
- Challenges: Kaggle, etc.

The big steps of data analysis

1. Extracting the data to expected format
2. Exploring the data
 - detection of outliers, of inconsistencies
 - descriptive exploration of the distributions, of correlations
 - data transformations

 - learning sample
 - validation sample
 - test sample
3. For each algorithm: parameter estimation using training and validation samples
4. Choice of final algorithm using testing sample, risk estimation

Machine Learning tools: R



The screenshot shows the scikit-learn website homepage. At the top, there is a navigation bar with the 'learn' logo and links for Home, Tutorials, User guides, and Examples. A search bar is also present. The main header features a grid of colorful images and the text 'scikit-learn Machine Learning in Python'. Below this, there are three columns of content: Classification, Regression, and Clustering. Each column has a title, a brief description, application examples, and a 'More' link. At the bottom, there are three sections: News, Community, and Who uses scikit-learn?, with the latter featuring the AWeber logo.

Classification
Trying to predict category or class, belongs to.
Applications: Spam detection, image recognition.
Also there: Soft-margin SVM, ensemble methods, etc.
More »

Regression
Predicting a continuous value (stock price, house value, etc.).
Applications: Drug response, stock prices.
Algorithms: RF, Ridge, Lasso, etc.
More »

Clustering
Automatic grouping of data into clusters.
Applications: Customer segmentation, anomaly detection, etc.
Also there: K-Means, hierarchical clustering, etc.
More »

Dimensionality reduction
Reducing the number of features while retaining most of the information.
Applications: SVM, k-NN, etc. on high-dimensional data.
Algorithms: PCA, LDA, t-SNE, etc.
More »

Model selection
Choosing the best model from a set of candidates.
Goal: Minimize overfitting and maximize generalization.
Algorithms: Grid search, cross-validation, etc.
More »

Preprocessing
Preparing data for machine learning.
Applications: Feature scaling, feature selection, etc.
Modules: StandardScaler, MinMaxScaler, etc.
More »

News
Getting development: What's new?
2014-01-12

Community
About us: We're looking for contributors!
Note: Machine Learning is a skill.

Who uses scikit-learn?
AWeber

Knime, Weka and co: integrated environments

The screenshot displays the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'J48 -C 0.25 -M 2'. The 'Test options' section includes 'Use training set' (selected), 'Supplied test set' (Set...), 'Cross-validation' (Folds: 10), and 'Percentage split' (%: 66). The 'Classifier output' section shows a summary of stratified cross-validation results:

==== Stratified cross-validation ====		
==== Summary ====		
Correctly Classified Instances	144	96 %
Incorrectly Classified Instances	6	4 %
Kappa statistic	0.94	
Mean absolute error	0.035	
Root mean square		

The 'Tree View' window shows a decision tree structure:

```
graph TD
    Root((petalwidth)) -- "<= 0.6" --> Node1[Iris-setosa (50.0)]
    Root -- "> 0.6" --> Node2((petalwidth))
    Node2 -- "<= 1.7" --> Node3((petalength))
    Node2 -- "> 1.7" --> Node4[Iris-virginica (46.0/1.0)]
    Node3 -- "<= 4.9" --> Node5[Iris-versicolor (48.0/1.0)]
    Node3 -- "> 4.9" --> Node6((petalwidth))
    Node6 -- "<= 1.5" --> Node7[Iris-virginica (3.0)]
    Node6 -- "> 1.5" --> Node8[Iris-versicolor (3.0/1.0)]
```

The interface also includes a 'Result list' on the left, a 'Start' button, and a 'Visualize' section on the right with a scatter plot and a 'Jitter' slider. The scatter plot shows data points for Iris-versicolor (red) and Iris-virginica (green).

Supervised Classification

- Domain set \mathcal{X}
- Label set \mathcal{Y}
- Statistical Model: $\{D \text{ probability over } \mathcal{X} \times \mathcal{Y}\}$
- Training data: pairs $(X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$, $1 \leq i \leq m$
m = sample size
- Learner's output: $\hat{h} : \mathcal{X} \rightarrow \mathcal{Y}$. Possibly $\hat{h} \in \mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$.
- Measures of success: risk measure

$$L_{\mathcal{D}}(h) = \mathbb{P}_{(X, Y) \sim \mathcal{D}}(h(X) \neq Y) = D\left(\{(x, y) : h(x) \neq y\}\right).$$

Example: Character Recognition

Domain set \mathcal{X} Label set \mathcal{Y} Joint distribution \mathcal{D}	64×64 images $\{0, 1, \dots, 9\}$?
Prediction function $h \in \mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ Risk $R(h) = P_{\mathcal{X}, \mathcal{Y}}(h(X) \neq Y)$	
Sample $S = \{(x_i, y_i)\}_{i=1}^m$ Empirical risk $L_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(x_i) \neq y_i\}$	MNIST dataset
Learning algorithm $\mathcal{A} = (\mathcal{A}_n)_n, \mathcal{A}_n : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$ Expected risk $R_n(\mathcal{A}) = \mathbb{E}_n[R(\mathcal{A}_n(D_n))]$	neural nets, boosting...
Empirical risk minimizer $\hat{h}_n = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ Regularized empirical risk minimizer $\hat{h}_n = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h) + \lambda C(h)$	

Realizable case vs agnostic learning

One usually distinguishes

- the *realizable case*: there exists $h : \mathcal{X} \rightarrow \mathcal{Y}$ such that $\mathbb{P}_{(X,Y) \sim \mathcal{D}}(h(X) = Y) = 1$,
- and the *agnostic case* otherwise (x does not permit to predict y with certainty).

Examples:

- spam filtering, character recognition
- credit risk, heart disease prediction

We generally focus on the agnostic case.

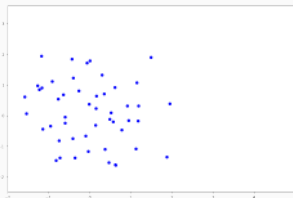
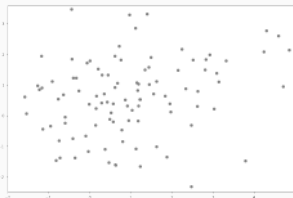
One can have 2 visions of D :

As a pair (D_x, k) , where

- for $A \subset \mathcal{X}$, $D_x(A) = D(A \times \mathcal{Y})$ is the marginal distribution of X ,
- and for $x \in \mathcal{X}$ and $B \subset \mathcal{Y}$,
 $k(B|x) = \mathbb{P}(Y \in B|X = x)$ is (a version of) the conditional distribution of Y given X .

As a pair $(D_y, (D(\cdot|y))_y)$, where

- for $y \in \mathcal{Y}$, $D_y(y) = D(\mathcal{X} \times y)$ is the marginal distribution of Y ,
- and for $A \subset \mathcal{X}$ and $y \in \mathcal{Y}$,
 $D(A|y) = \mathbb{P}(X \in A|Y = y)$ is the conditional distribution of X given $Y = y$.



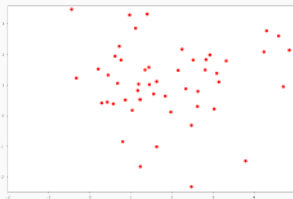
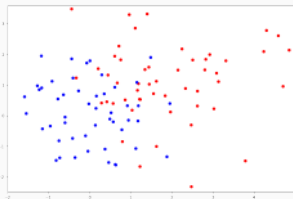
One can have 2 visions of D :

As a pair (D_x, k) , where

- for $A \subset \mathcal{X}$, $D_x(A) = D(A \times \mathcal{Y})$ is the marginal distribution of X ,
- and for $x \in \mathcal{X}$ and $B \subset \mathcal{Y}$, $k(B|x) = \mathbb{P}(Y \in B|X = x)$ is (a version of) the conditional distribution of Y given X .

As a pair $(D_y, (D(\cdot|y))_y)$, where

- for $y \in \mathcal{Y}$, $D_y(y) = D(\mathcal{X} \times y)$ is the marginal distribution of Y ,
- and for $A \subset \mathcal{X}$ and $y \in \mathcal{Y}$, $D(A|y) = \mathbb{P}(X \in A|Y = y)$ is the conditional distribution of X given $Y = y$.



Bayes Classifier

Consider binary classification $\mathcal{Y} = \{0, 1\}$.

Theorem

The Bayes classifier is defined by

$$h^*(x) = \mathbb{1}\{\eta(x) \geq 1/2\} = \mathbb{1}\{\eta(x) \geq 1 - \eta(x)\} = \mathbb{1}\{2\eta(x) - 1 \geq 0\}.$$

For every classifier $h : \mathcal{X} \rightarrow \mathcal{Y} = \{0, 1\}$,

$$L_{\mathcal{D}}(h) \geq L_{\mathcal{D}}(h^*) = \mathbb{E}\left[\min(\eta(X), 1 - \eta(X))\right].$$

The Bayes risk $L_{\mathcal{D}}^ = L_{\mathcal{D}}(h^*)$ is called the **noise** of the problem.*

More precisely,

$$L_{\mathcal{D}}(h) - L_{\mathcal{D}}(h^*) = \mathbb{E}\left[|2\eta(X) - 1| \mathbb{1}\{h(X) \neq h^*(X)\}\right].$$

Extends to $|\mathcal{Y}| > 2$.

Nearest-Neighbor Classification

The Nearest-Neighbor Classifier

We assume that \mathcal{X} is a metric space with distance d .

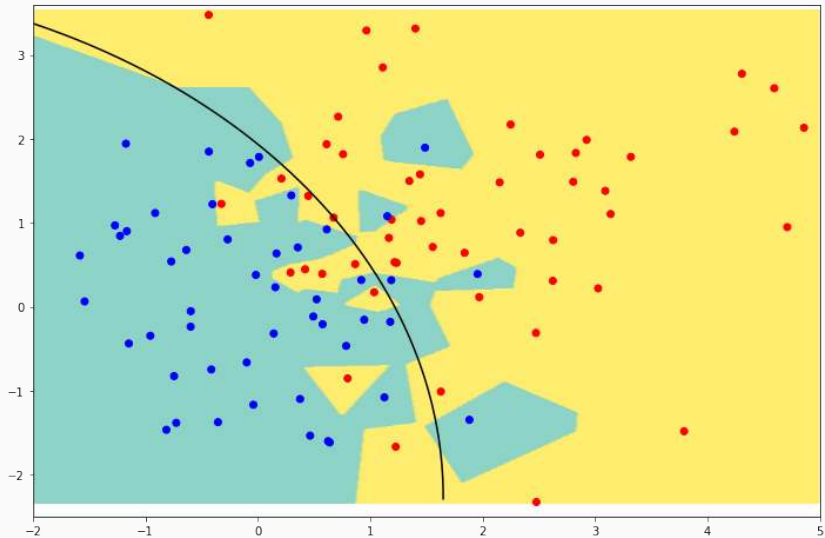
The nearest-neighbor classifier $\hat{h}_m^{NN} : \mathcal{X} \rightarrow \mathcal{Y}$ is defined as

$$\hat{h}_m^{NN}(x) = Y_I \text{ where } I \in \arg \min_{1 \leq i \leq m} d(x - X_i).$$

Typical distance: L^2 norm on \mathbb{R}^d $\|x - x'\| = \sqrt{\sum_{j=1}^d (x_j - x'_j)^2}$.

Buts many other possibilities: Hamming distance on $\{0, 1\}^d$, etc.

Numerically



- A1.** $\mathcal{Y} = \{0, 1\}$.
- A2.** $\mathcal{X} = [0, 1]^d$.
- A3.** η is c -Lipschitz continuous:

$$\forall x, x' \in \mathcal{X}, |\eta(x) - \eta(x')| \leq c \|x - x'\| .$$

Theorem

Under the previous assumptions, for all distributions D and all $m \geq 1$

$$L_D(\hat{h}_m^{NN}) \leq 2L_D^* + \frac{3c\sqrt{d}}{m^{1/(d+1)}}$$

Proof Outline

- Conditioning: as $l(x) = \arg \min_{1 \leq i \leq m} \|x - X_i\|$,

$$L_D(\hat{h}_n^{NN}) = \mathbb{E} \left[\mathbb{E} [\mathbb{1}\{Y \neq Y_{l(X)}\} | X, X_1, \dots, X_m] \right].$$

- $Y \sim \mathcal{B}(p)$, $Y' \sim \mathcal{B}(q) \implies \mathbb{P}(Y \neq Y') \leq 2 \min(p, 1-p) + |p - q|$,

$$\mathbb{E} [\mathbb{1}\{Y \neq Y_{l(X)}\} | X, X_1, \dots, X_m] \leq 2 \min(\eta(X), 1-\eta(X)) + c \|X - X_{l(X)}\|.$$

- Partition \mathcal{X} into $|\mathcal{C}| = T^d$ cells of diameter \sqrt{d}/T :

$$\mathcal{C} = \left\{ \left[\frac{j_1 - 1}{T}, \frac{j_1}{T} \right] \times \dots \times \left[\frac{j_d - 1}{T}, \frac{j_d}{T} \right], \quad 1 \leq j_1, \dots, j_d \leq T \right\}.$$

- 2 cases: either the cell of X is occupied by a sample point, or not:

$$\|X - X_{l(X)}\| \leq \sum_{c \in \mathcal{C}} \mathbb{1}\{X \in c\} \left(\frac{\sqrt{d}}{T} \mathbb{1} \bigcup_{i=1}^m \{X_i \in c\} + \sqrt{d} \mathbb{1} \bigcap_{i=1}^m \{X_i \notin c\} \right).$$

- $\implies \mathbb{E}[\|X - X_{l(X)}\|] \leq \frac{\sqrt{d}}{T} + \frac{\sqrt{d}T^d}{em}$ and choose $T = \lfloor m^{\frac{1}{d+1}} \rfloor$.

What does the analysis say?

- Is it loose? (sanity check: uniform \mathcal{D}_X)
- Non-asymptotic (finite sample bound)
- The second term $\frac{3c\sqrt{d}}{m^{1/(d+1)}}$ is distribution independent
- Does not give the trajectorial decrease of risk
- Exponential bound d (cannot be avoided...)
 \implies *curse of dimensionality*

- How to improve the classifier?

Deviation Bound for Bernoulli Variables

Remember: Jensen's Inequality

Basic version: if $\phi : \mathcal{X} \rightarrow \mathbb{R}$ is convex and $t \in (0, 1)$ then for all $x, x' \in \mathcal{X}$, $f(tx + (1 - t)x') \leq tf(x) + (1 - t)f(x')$.

Probabilistic version: If $\phi : \mathcal{X} \rightarrow \mathbb{R}$ is convex and if X is a random variable with range in \mathcal{X} , then $\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)]$.

Example: For a real-valued random variable X $\mathbb{E}[X^2] \geq \mathbb{E}[X]^2$ and thus $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \geq 0$.

Think about equality case.

Chernoff's Bound

Theorem (Chernoff-Hoeffding Deviation Bound)

Let $\mu \in (0, 1)$. $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{B}(\mu)$, and let $x \in (\mu, 1]$.

(i) Chernoffs' bound for Bernoulli variables:

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-n \text{kl}(x, \mu)), \quad (1)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1 - p) \log \frac{1 - p}{1 - q}$.

(ii) Hoeffding's bound for Bernoulli variables: since $\text{kl}(p, q) \geq 2(p - q)^2$,

$$\mathbb{P}(\bar{X}_n \geq x) \leq \exp(-2n(x, \mu)^2). \quad (2)$$

(iii) Inequalities (1) and (2) hold for arbitrary independent random variables with range $[0, 1]$ and expectation μ .

Reason: $\exp(\lambda x) \leq (1 - x) \exp(0) + x \exp(\lambda)$.

Kullback-Leibler Divergence

Definition

Let P and Q be two probability distributions on a set \mathcal{X} . The Kullback-Leibler divergence from Q to P is defined as follows:

- if P is not absolutely continuous with respect to Q , then $\text{KL}(P, Q) = +\infty$;
- otherwise, let $\frac{dP}{dQ}$ be the Radon-Nikodym derivative of P with respect to Q . Then

$$\text{KL}(P, Q) = \int_{\mathcal{X}} \log \frac{dP}{dQ} dP = \int_{\mathcal{X}} \frac{dP}{dQ} \log \frac{dP}{dQ} dQ .$$

The integral always exists but may be equal to $+\infty$.

Examples:

- $\text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = \text{kl}(p, q)$,
- $\text{KL}(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) = \frac{(\mu_1 - \mu_2)^2}{2\sigma^2}$.

Tensorization of entropy: If $P = P_1 \otimes P_2$ and $Q = Q_1 \otimes Q_2$, then

$$\text{KL}(P, Q) = \text{KL}(P_1, Q_1) + \text{KL}(P_2, Q_2) .$$

Contraction of entropy aka data-processing inequality:

Let P and Q be probability distributions on \mathcal{X} , and let $X \sim P$ and $Y \sim Q$. If $f : \mathcal{X} \rightarrow \mathcal{X}'$ is a measurable function and if \tilde{P} (resp. \tilde{Q}) is the distribution of $f(X)$ (resp. $f(Y)$), then

$$\text{KL}(\tilde{P}, \tilde{Q}) \leq \text{KL}(P, Q) .$$

Application: Lower bound

Let $\mu \in (0, 1)$. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{B}(\mu)$, and let $x \in (\mu, 1]$. Then

$$\liminf_m \frac{1}{m} \log \mathbb{P}(\bar{X}_m > x) \geq -\text{kl}(x, \mu)$$

"Chernoff's bound is asymptotically almost tight"