

Statistical learning

Gérard Biau

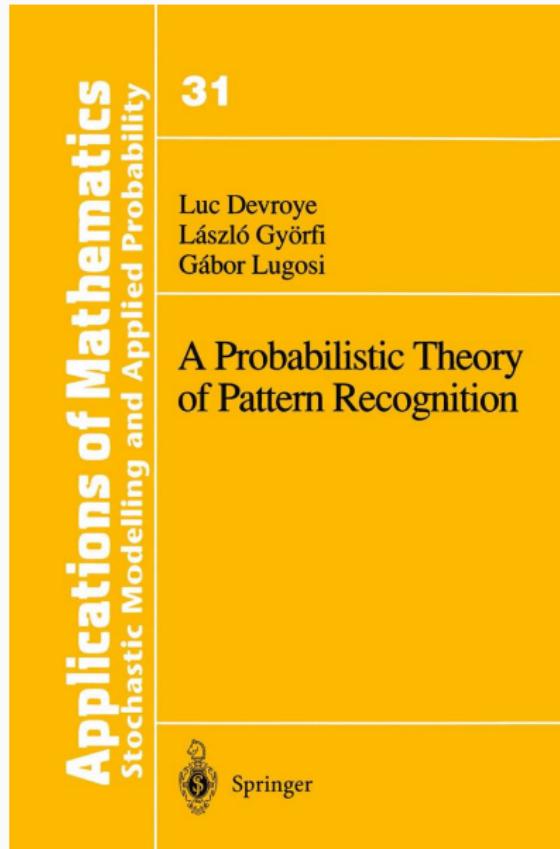
Semester 1

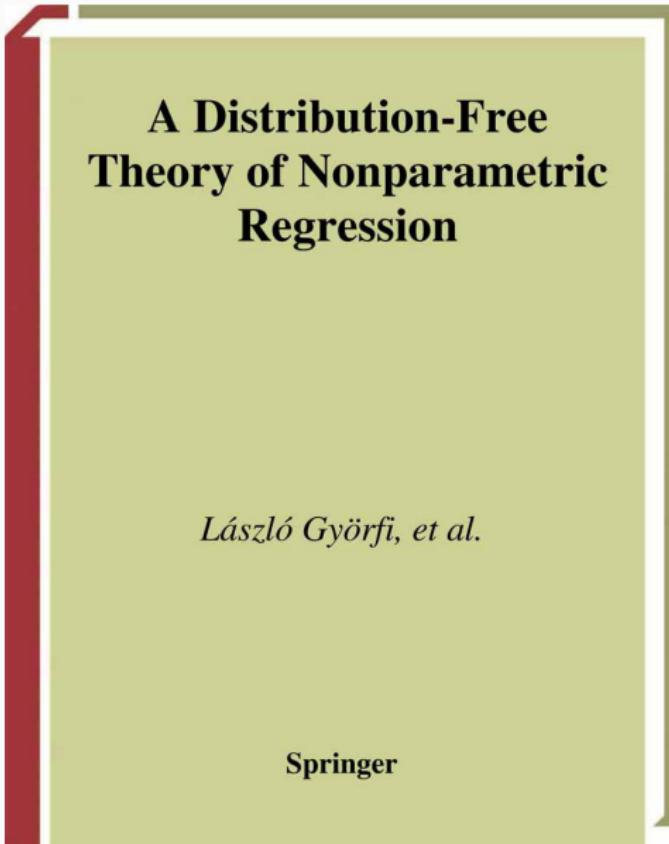


Overview

- **Prerequisites:** basics of probability and statistics.
- **Objectives:**
 - ▷ Learn and apply the fundamental concepts of statistical learning;
 - ▷ Understand the basic theory underlying data science, machine learning, and AI;
 - ▷ Be able to read current research books and papers.
- **Keywords:** classification, prediction, performance guarantees.
- **Resources:** slides and proof textbook.
- **Teacher:** Gérard Biau, gerard.biau@sorbonne-universite.fr.
- **Office:** ☎ 15-25, 2nd floor, ☎ 01 44 27 85 63.

References



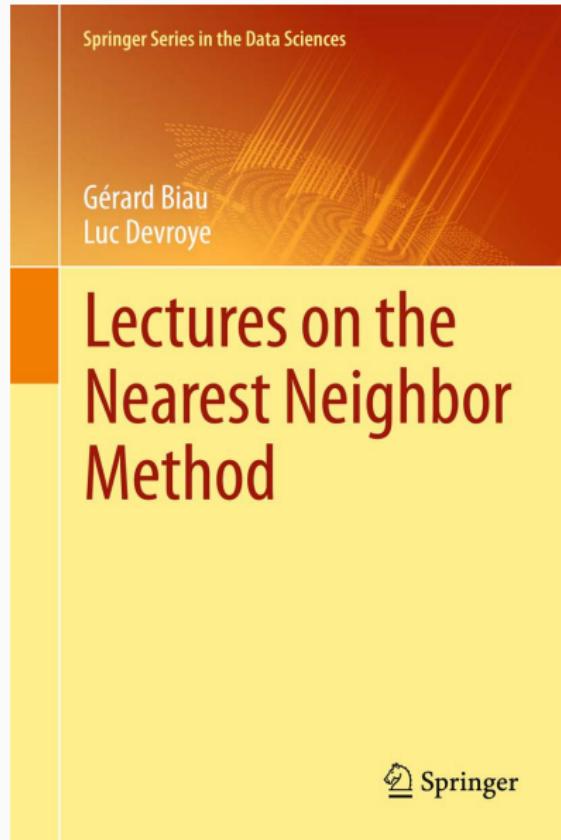


A Distribution-Free
Theory of Nonparametric
Regression

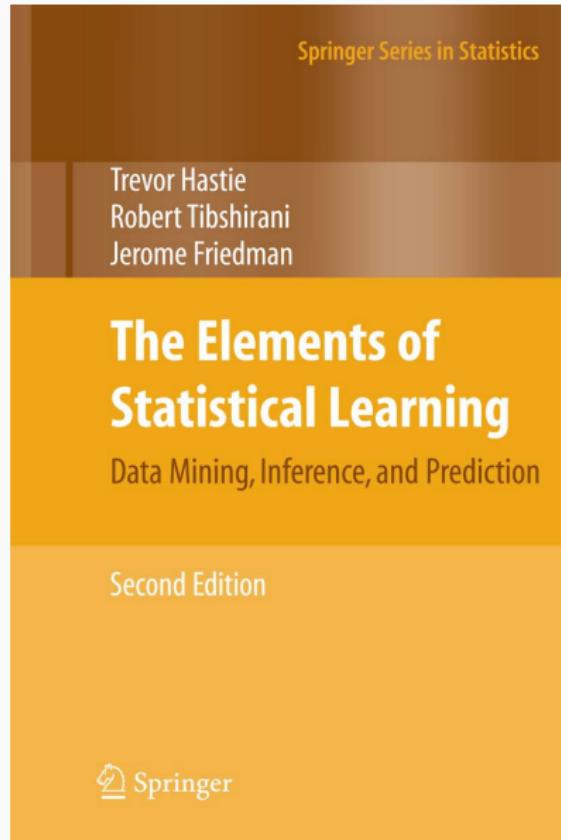
László Györfi, et al.

Springer

References



References



References

ESAIM: PS
November 2005, Vol. 9, p. 323–375
DOI: 10.1051/ps:2005018

ESAIM: Probability and Statistics

THEORY OF CLASSIFICATION: A SURVEY OF SOME RECENT ADVANCES*

STÉPHANE BOUCHERON¹, OLIVIER BOUSQUET² AND GÁBOR LUGOSI³

Abstract. The last few years have witnessed important new developments in the theory and practice of pattern classification. We intend to survey some of the main new ideas that have led to these recent results.

Mathematics Subject Classification. 62G08, 60E15, 68Q32.

Received June 18, 2004. Accepted September 12, 2005.

1. INTRODUCTION

The last few years have witnessed important new developments in the theory and practice of pattern classification. The introduction of new and effective techniques of handling high-dimensional problems – such as boosting and support vector machines – have revolutionized the practice of pattern recognition. At the same time, the better understanding of the application of empirical process theory and concentration inequalities have led to effective new ways of studying these methods and provided a statistical explanation for their success. These new tools have also helped develop new model selection methods that are at the heart of many classification algorithms.

The purpose of this survey is to offer an overview of some of these theoretical tools and give the main idea of the analysis of some of the important algorithms. This survey does not attempt to be exhaustive. The selection of the topics is largely biased by the personal taste of the authors. We also limit ourselves to describing the key ideas in a simple way, often sacrificing generality. In these cases the reader is pointed to the references for the sharpest and more general results available. References and bibliographical remarks are given at the end of each section, in an attempt to avoid interruptions in the arguments.

2. BASIC MODEL

The problem of pattern classification is about guessing or predicting the unknown class of an observation. An observation is often a collection of numerical and/or categorical measurements represented by a d -dimensional vector x but in some cases it may even be a curve or an image. In our model we simply assume that $x \in \mathcal{X}$

Keywords and phrases. Pattern recognition, statistical learning theory, concentration inequalities, empirical processes, model selection.

* The authors acknowledge support by the PASCAL Network of Excellence under EC grant no. 506778. The work of the third author was supported by the Spanish Ministry of Science and Technology and FEDER, grant BFM2003-0324.

¹ Laboratoire Probabilités et Modèles Aléatoires, CNRS & Université Paris VII, Paris, France.

² Pertinax SA, 32 rue des Jeûneurs, 75002 Paris, France.

³ Department of Economics, Pompeu Fabra University, Ramon Trias Fargas 25-27, 08005 Barcelona, Spain; lugosi@upf.es

© EDP Sciences, SMAI 2005

References

+ This is page 3
Printer: Opaque this

1

Pattern classification and learning theory

Gábor Lugosi

1.1 A binary classification problem

Pattern recognition (or *classification* or *discrimination*) is about guessing or predicting the unknown class of an observation. An *observation* is a collection of numerical measurements, represented by a d -dimensional vector x . The unknown nature of the observation is called a *class*. It is denoted by y and takes values in the set $\{0, 1\}$. (For simplicity, we restrict our attention to binary classification.) In pattern recognition, one creates a function $g(x) : \mathbb{R}^d \rightarrow \{0, 1\}$ which represents one's guess of y given x . The mapping g is called a *classifier*. A classifier *errs* on x if $g(x) \neq y$.

To model the learning problem, we introduce a probabilistic setting, and let (X, Y) be an $\mathbb{R}^d \times \{0, 1\}$ -valued random pair.

The random pair (X, Y) may be described in a variety of ways; for example, it is defined by the pair (μ, η) , where μ is the probability measure for X and η is the regression of Y on X . More precisely, for a Borel-measurable set $A \subseteq \mathbb{R}^d$,

$$\mu(A) = \mathbb{P}\{X \in A\},$$

and for any $x \in \mathbb{R}^d$,

$$\eta(x) = \mathbb{P}\{Y = 1 | X = x\} = \mathbb{E}\{Y | X = x\}.$$

Thus, $\eta(x)$ is the conditional probability that Y is 1 given $X = x$. The distribution of (X, Y) is determined by (μ, η) . The function η is called the *a posteriori* probability.

Any function $g : \mathbb{R}^d \rightarrow \{0, 1\}$ defines a classifier. An error occurs if $g(X) \neq Y$, and the probability of error for a classifier g is

$$L(g) := \mathbb{P}\{g(X) \neq Y\}.$$

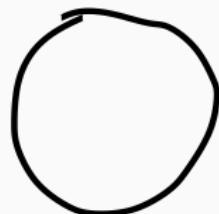
The Bayes classifier given by

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

Motivating examples

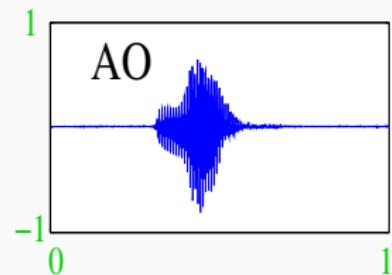
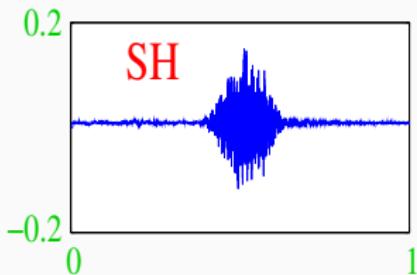
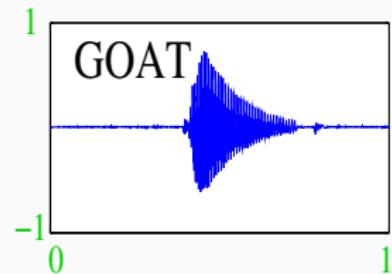
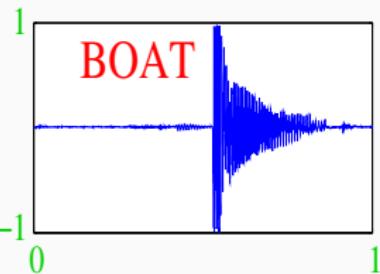
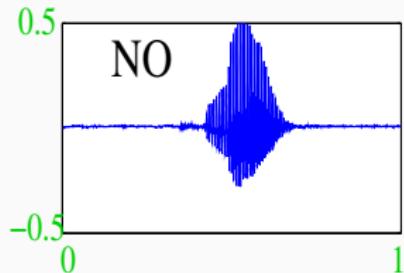
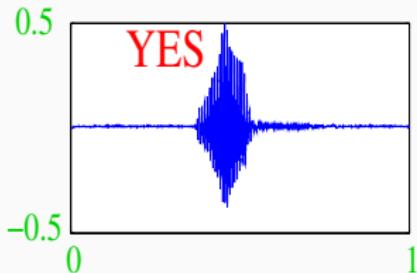
Handwritten digit recognition

0	0	0	0	0
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9



Which digit is this? 0, 1, 2,...?

Voice recognition



Weather forecast



Temperature prediction

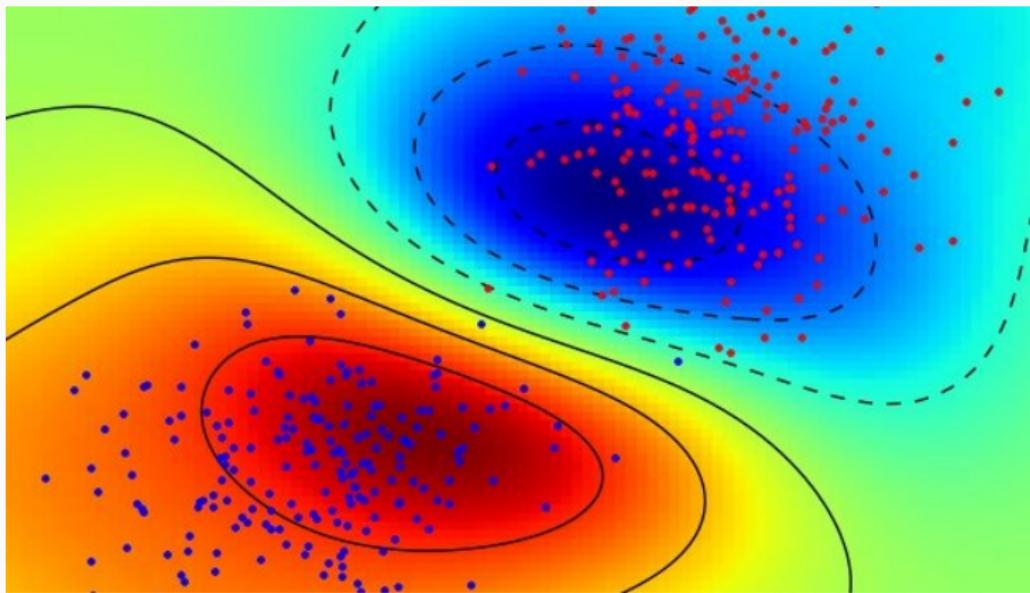


Image segmentation



Outline

1. Supervised classification
2. Empirical risk minimization
3. Vapnik-Chervonenkis theory
4. Margin-based bounds
5. Convex loss functions
6. Model selection
7. Stone's theorem
8. k -nearest neighbor classifiers
9. Partitioning classifiers and trees
10. Neural networks
11. Quantization and clustering

Supervised classification

Basics

- **Random pair:** $(X, Y) \in \mathbb{R}^d \times \{0, 1\}$.
- X is the **observation** and Y is the **label** (or **class**).
- **Variants:** $X \in \mathcal{X}$, $Y \in \{-1, 1\}$, $Y \in \{1, \dots, M\}$, etc.
- The distribution of (X, Y) is described by:
 - ▷ The **distribution of X :** $\mu(A) = \mathbb{P}(X \in A)$, $A \in \mathcal{B}(\mathbb{R}^d)$;
 - ▷ The **a posteriori probability:** $\eta(x) = \mathbb{P}(Y = 1|X = x)$.
- **Note:** $\eta(x) = \mathbb{E}(Y|X = x) \stackrel{\text{def}}{=} r(x)$.
- Are X and Y independent? **Not necessarily**.
- Do we have $Y = \varphi(X)$? **Not necessarily**.
- **Objective:** find a **classifier** $g : \mathbb{R}^d \rightarrow \{0, 1\}$ such that $g(X) \approx Y$.
- **Error probability:** $L(g) = \mathbb{P}(g(X) \neq Y)$.
- **0-1 loss:** $L(g) = \mathbb{E} \mathbb{1}_{[g(X) \neq Y]}$.

Bayes classifier

Definition

$$g^*(x) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) > \mathbb{P}(Y = 0|X = x) \\ 0 & \text{otherwise,} \end{cases}$$

i.e.,

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

Lemma

For any classifier g , $L(g^*) \leq L(g)$. Thus, g^* is the optimal decision.

▷ $L^* \stackrel{\text{def}}{=} L(g^*)$ is the Bayes risk (or Bayes error).

▷ One has

$$L^* = \inf_{g: \mathbb{R}^d \rightarrow \{0,1\}} \mathbb{P}(g(X) \neq Y).$$

Bayes risk

- Exercise 1:

$$L^* = 1 - \mathbb{E}[\mathbf{1}_{[\eta(X) > 1/2]}\eta(X) + \mathbf{1}_{[\eta(X) \leq 1/2]}(1 - \eta(X))].$$

- Exercise 2:

$$L^* = \mathbb{E}[\min(\eta(X), 1 - \eta(X))] = \frac{1}{2} - \frac{1}{2}\mathbb{E}|2\eta(X) - 1|.$$

- Exercise 3:

$$L^* = 0 \Leftrightarrow Y = \varphi(X) \text{ with probability one.}$$

Statistical learning

- **Sample:** $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ i.i.d. copies of (X, Y) .
- (X, Y) and \mathcal{D}_n are **independent**.
- A **classifier**: $g_n(x) = g_n(x; \mathcal{D}_n)$ taking values in $\{0, 1\}$.
- **Error probability**: $L(g_n) = \mathbb{P}(g_n(X) \neq Y | \mathcal{D}_n)$.
- **Objective**: construct g_n such that $L(g_n) \approx L^*$.



Statistics = inference. Learning = prediction.

Consistency

Definition (Consistency)

A classifier g_n is **consistent** for a certain distribution of (X, Y) if

$$\mathbb{E}L(g_n) = \mathbb{P}(g_n(X) \neq Y) \rightarrow L^* \quad \text{as } n \rightarrow \infty,$$

and **strongly consistent** if

$$L(g_n) \rightarrow L^* \quad \text{almost surely.}$$

- ▷ **Exercise:** consistency $\Leftrightarrow L(g_n) \xrightarrow{L^1} L^* \Leftrightarrow L(g_n) \xrightarrow{\mathbb{P}} L^*$.

Definition (Universal consistency)

A classifier is called **universally** (strongly) **consistent** if it is (strongly) consistent for **any** distribution of (X, Y) .

Empirical risk minimization

Principle

- A class \mathcal{C} of classifiers $g : \mathbb{R}^d \rightarrow \{0, 1\}$.
- Objective: find $g_n^* \in \mathcal{C}$ such that $L(g_n^*) \approx \inf_{g \in \mathcal{C}} L(g)$.
- Natural choice: $g_n^* \in \arg \min_{g \in \mathcal{C}} L_n(g)$, where

$$L_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[g(X_i) \neq Y_i]} \quad (\text{empirical risk of } g).$$

- Error probability: $L(g_n^*) = \mathbb{P}(g_n^*(X) \neq Y | \mathcal{D}_n)$.
- Decomposition:

$$\begin{aligned} L(g_n^*) - L^* &= [L(g_n^*) - \inf_{g \in \mathcal{C}} L(g)] + [\inf_{g \in \mathcal{C}} L(g) - L^*] \\ &= \text{estimation error} + \text{approximation error}. \end{aligned}$$

DON'T FORGET!

Estimation = random. Approximation = deterministic.

Overfitting

- Estimation vs. approximation.
- Small \mathcal{C} : restrictive. Large \mathcal{C} : overfitting.
- Example:

- ▷ $\mathcal{C} = \text{all measurable functions};$
- ▷ Approximation error = 0;
- ▷ But $L_n(g_n^*) = 0$ with

$$g_n^*(x) = \begin{cases} Y_i & \text{for } x = X_i, \ i = 1, \dots, n \\ 0 & \text{otherwise;} \end{cases}$$

- ▷ No generalization ability.
- Learning strategy: keep \mathcal{C} under control.

Starting point

Lemma

One has

$$(i) \quad L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \leq 2 \sup_{g \in \mathcal{C}} |L_n(g) - L(g)|,$$

and

$$(ii) \quad |L_n(g_n^*) - L(g_n^*)| \leq \sup_{g \in \mathcal{C}} |L_n(g) - L(g)|.$$

Take-home message: $\sup_{g \in \mathcal{C}} |L_n(g) - L(g)|$ small guarantees:

- ▷ That $L(g_n^*)$ is not much larger than the best error probability in \mathcal{C} ;
- ▷ That the empirical estimate $L_n(g_n^*)$ of $L(g_n^*)$ is also good.

Concentration

Fact: $nL_n(g) \stackrel{\mathcal{D}}{=} \text{Bin}(n, L(g))$.

☞ We need **uniform deviations** of binomial r.v. from their means.

Theorem (Hoeffding's inequality)

Let X_1, \dots, X_n be **independent** real-valued random variables. Assume that each X_i takes its values in $[a_i, b_i]$ ($a_i < b_i$) with probability one. Then, for all $t > 0$,

$$\mathbb{P}\left(\sum_{i=1}^n X_i - \mathbb{E} \sum_{i=1}^n X_i \geq t\right) \leq e^{-2t^2/\sum_{i=1}^n (b_i - a_i)^2}$$

and

$$\mathbb{P}\left(\sum_{i=1}^n X_i - \mathbb{E} \sum_{i=1}^n X_i \leq -t\right) \leq e^{-2t^2/\sum_{i=1}^n (b_i - a_i)^2}.$$

In particular,

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i - \mathbb{E} \sum_{i=1}^n X_i\right| \geq t\right) \leq 2e^{-2t^2/\sum_{i=1}^n (b_i - a_i)^2}.$$

Chernoff's bounding method

Lemma

Let X be a real-valued random variable with $\mathbb{E}X = 0$ and $X \in [a, b]$ ($a < b$) with probability one. Then, for all $s \geq 0$,

$$\mathbb{E}e^{sX} \leq e^{s^2(b-a)^2/8}.$$

Proof of Hoeffding's inequality:

$$\begin{aligned}\mathbb{P}\left(\sum_{i=1}^n X_i - \mathbb{E}\sum_{i=1}^n X_i \geq t\right) &\leq e^{-st} \prod_{i=1}^n e^{s^2(b_i-a_i)^2/8} \\ &= e^{-st} e^{s^2 \sum_{i=1}^n (b_i-a_i)^2/8} \\ &= e^{-2t^2 / \sum_{i=1}^n (b_i-a_i)^2},\end{aligned}$$

by choosing $s = 4t / \sum_{i=1}^n (b_i - a_i)^2$.

The case $|\mathcal{C}| < \infty$

Corollary

Let $X \stackrel{\mathcal{D}}{=} \text{Bin}(n, p)$, $n \geq 1$ and $p \in [0, 1]$. Then, for all $t > 0$,

$$\mathbb{P}(|X - np| \geq t) \leq 2e^{-2t^2/n}.$$

Theorem

Assume that $|\mathcal{C}|$ is finite, with $|\mathcal{C}| \leq N$. Then, for all $t > 0$,

$$\mathbb{P}\left(\sup_{g \in \mathcal{C}} |L_n(g) - L(g)| \geq t\right) \leq 2Ne^{-2nt^2}.$$

- ▷ The bound is universal.
- ▷ Borel-Cantelli: $\sup_{g \in \mathcal{C}} |L_n(g) - L(g)| \rightarrow 0$ almost surely.
- ▷ Consequence: $L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \rightarrow 0$ almost surely.
- ▷ Bound on $\mathbb{E}L(g_n^*) - \inf_{g \in \mathcal{C}} L(g)$?

From \mathbb{P} to \mathbb{E}

Lemma

Let X be a random variable taking values in \mathbb{R}_+ . Assume that there exists a constant $C \geq 1$ such that, for all $t > 0$, $\mathbb{P}(X \geq t) \leq Ce^{-2nt^2}$. Then

$$\mathbb{E}X \leq \sqrt{\frac{\log(Ce)}{2n}}.$$

▷ Consequence:

$$\mathbb{E} \sup_{g \in \mathcal{C}} |L_n(g) - L(g)| \leq \sqrt{\frac{\log(2eN)}{2n}}$$

and

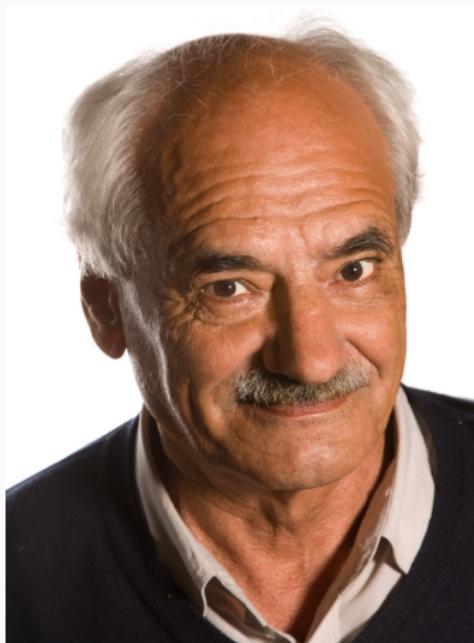
$$\mathbb{E}L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \leq 2\sqrt{\frac{\log(2eN)}{2n}}.$$

▷ Take-home message: estimation error = $O(\sqrt{\frac{\log N}{n}})$.

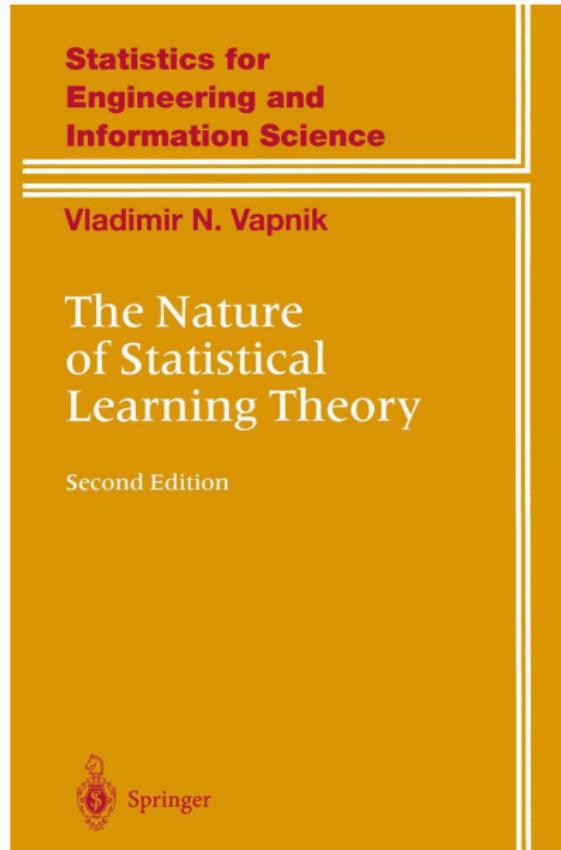
▷ Next objective: handle more complex classes of functions.

Vapnik-Chervonenkis theory

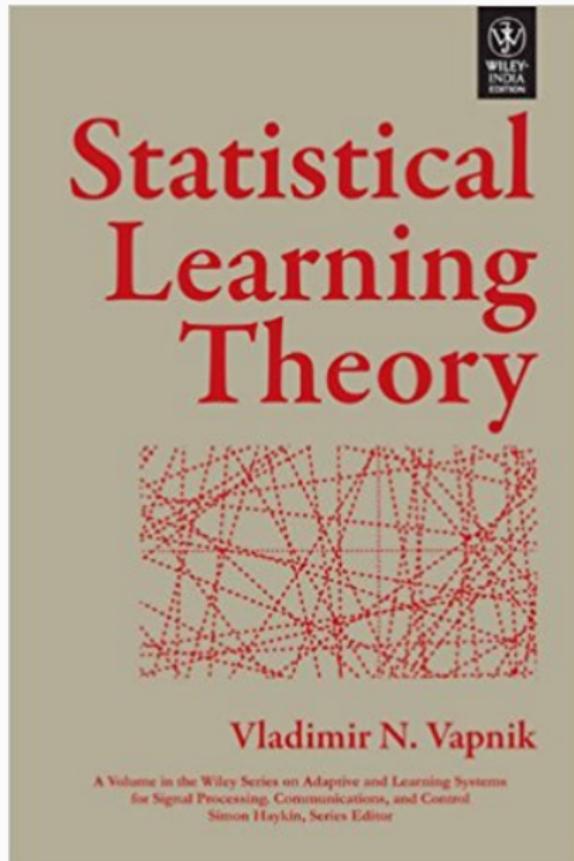
Vladimir Vapnik and Alexey Chervonenkis



References



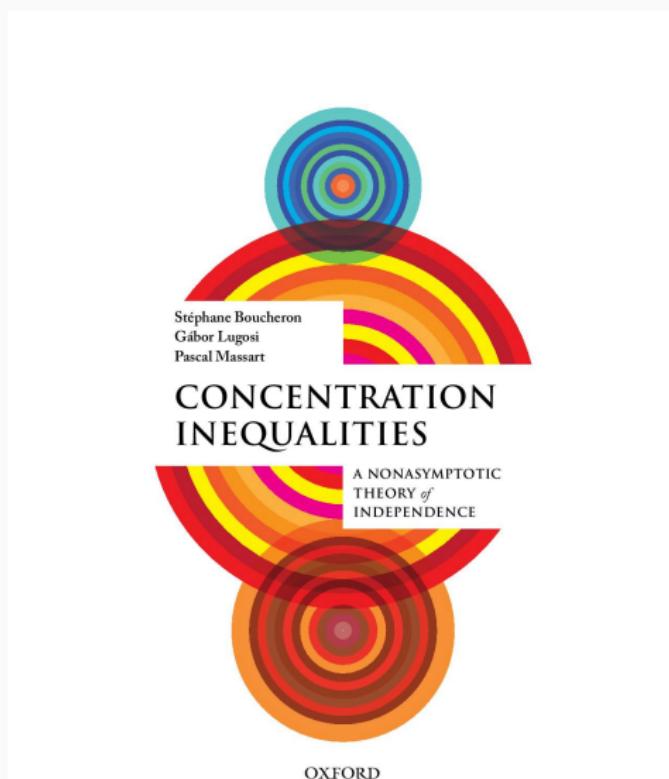
References



From finite to infinite classes

- **Learning:** functions of the form $f = \mathbb{1}_{[(x,y):g(x)\neq y]}, g \in \mathcal{C}$.
- **General context:**
 - ▷ X_1, \dots, X_n = i.i.d. random variables taking values in a set \mathcal{X} ;
 - ▷ \mathcal{F} = a class of bounded functions $f : \mathcal{X} \rightarrow [0, 1]$ (or $[-1, 1]$).
- **Notation:** $Pf = \mathbb{E}f(X_1)$ and $P_n f = \frac{1}{n} \sum_{i=1}^n f(X_i)$.
- **Objective:** bound $\sup_{f \in \mathcal{F}} |P_n f - Pf|$.
- **Basic tool:** concentration inequalities.

Reference



Bounded difference inequality

Theorem (Bounded difference inequality)

Let X_1, \dots, X_n be independent random variables taking values in a set \mathcal{X} . Assume that $g : \mathcal{X}^n \rightarrow \mathbb{R}$ is Borel measurable and satisfies

$$\sup_{\substack{(x_1, \dots, x_n) \in \mathcal{X}^n \\ x_i' \in \mathcal{X}}} |g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)| \leq c_i, \quad 1 \leq i \leq n,$$

for some positive constants c_1, \dots, c_n (*bounded difference assumption*).

Then, for all $t > 0$,

$$\mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}g(X_1, \dots, X_n) \geq t) \leq e^{-2t^2 / \sum_{i=1}^n c_i^2}$$

and

$$\mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}g(X_1, \dots, X_n) \leq -t) \leq e^{-2t^2 / \sum_{i=1}^n c_i^2}.$$

In particular,

$$\mathbb{P}(|g(X_1, \dots, X_n) - \mathbb{E}g(X_1, \dots, X_n)| \geq t) \leq 2e^{-2t^2 / \sum_{i=1}^n c_i^2}.$$

Concentration

- Natural choice: $g(X_1, \dots, X_n) = \sup_{f \in \mathcal{F}} |P_n f - P f|$.
- The bounded difference assumption is satisfied with $c_i = 1/n$.
- Concentration: for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} |P_n f - P f| \leq \mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f| + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

- Consequence: focus on the expected value $\mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f|$.
- Key: Rademacher averages and symmetrization principle.

Rademacher averages

- A random variable σ is **Rademacher** if $\mathbb{P}(\sigma = \pm 1) = 1/2$.
- **Definition:**
 - ▷ Let $A \subseteq \mathbb{R}^n$ be a **bounded** set of vectors $a = (a_1, \dots, a_n)$;
 - ▷ Let $\sigma_1, \dots, \sigma_n$ be **i.i.d.** Rademacher random variables;
 - ▷ **Rademacher average** associated with A :

$$R_n(A) = \mathbb{E} \sup_{a \in A} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i a_i \right|.$$

- $\mathcal{F}(x_1^n)$ = the class of n -vectors $(f(x_1), \dots, f(x_n))$, $f \in \mathcal{F}$.
- **Key quantity:**

$$R_n(\mathcal{F}(X_1^n)) = \mathbb{E} \left(\sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i f(X_i) \right| \mid X_1, \dots, X_n \right),$$

where $\sigma_1, \dots, \sigma_n$ are **independent** of the X_i .

DON'T FORGET!

$R_n(\mathcal{F}(X_1^n)) = \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) \right|$ is **data-dependent**.

Fundamental inequalities

Theorem

One has

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f| \leq 2 \mathbb{E} R_n(\mathcal{F}(X_1^n)).$$

In addition, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} |P_n f - P f| \leq 2 \mathbb{E} R_n(\mathcal{F}(X_1^n)) + \sqrt{\frac{\log(1/\delta)}{2n}},$$

and, with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} |P_n f - P f| \leq 2 R_n(\mathcal{F}(X_1^n)) + 3 \sqrt{\frac{\log(2/\delta)}{2n}}.$$

- ▷ The last inequality is a **data-dependent** performance bound.
- ▷ **Sharp** estimates of the maximal deviations.
- ▷ In fact,

$$\frac{1}{2} \mathbb{E} R_n(\mathcal{F}(X_1^n)) - \frac{1}{2\sqrt{n}} \leq \mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f| \leq 2 \mathbb{E} R_n(\mathcal{F}(X_1^n)).$$

Properties of Rademacher averages

Theorem

Let A, B be bounded subsets of \mathbb{R}^n , and let $c \in \mathbb{R}$ be a constant. Then

$$R_n(A \cup B) \leq R_n(A) + R_n(B), \quad R_n(c \cdot A) = |c|R_n(A),$$

and

$$R_n(A \oplus B) \leq R_n(A) + R_n(B),$$

where $c \cdot A = \{ca : a \in A\}$ and $A \oplus B = \{a + b : a \in A, b \in B\}$. In addition, if $A = \{a^{(1)}, \dots, a^{(N)}\} \subseteq \mathbb{R}^n$ is a **finite** set, then

$$R_n(A) \leq \max_{1 \leq j \leq N} \|a^{(j)}\| \frac{\sqrt{2 \log(2N)}}{n},$$

where $\|\cdot\|$ denotes **Euclidean norm**.

Class of indicators

- **Learning:** functions of the form $f = \mathbb{1}_{[(x,y):g(x)\neq y]}, g \in \mathcal{C}$.
- **Binary framework:** $\mathcal{F} = \{\mathbb{1}_A : A \in \mathcal{A}\}$, with $|\mathcal{A}| \geq 2$.
- In this case, for $x_1^n = (x_1, \dots, x_n)$, $\mathcal{F}(x_1^n)$ is a **finite subset** of \mathbb{R}^n .
- **Combinatorial quantity:** $|\mathcal{F}(x_1^n)|$.
- **Shatter coefficient:** $S_{\mathcal{A}}(n) = \max_{x_1^n} |\mathcal{F}(x_1^n)|$.
- **Properties:** $S_{\mathcal{A}}(1) = 2$, $2 \leq S_{\mathcal{A}}(n) \leq 2^n$, and

$$S_{\mathcal{A}}(k) < 2^k \text{ for some } k > 1 \Leftrightarrow S_{\mathcal{A}}(n) < 2^n \text{ for all } n \geq k.$$

- For all x_1^n ,

$$R_n(\mathcal{F}(x_1^n)) \leq \sqrt{\frac{2 \log(2S_{\mathcal{A}}(n))}{n}} \leq 2 \sqrt{\frac{\log S_{\mathcal{A}}(n)}{n}}.$$

- $\log S_{\mathcal{A}}(n)$ may be bounded in terms of the **VC dimension**.

VC dimension

Definition (VC dimension)

The VC dimension $V_{\mathcal{A}}$ of \mathcal{A} is the largest integer $n_0 \geq 1$ for which $S_{\mathcal{A}}(n_0) = 2^{n_0}$. If $S_{\mathcal{A}}(n) = 2^n$ for all $n \geq 1$, then $V_{\mathcal{A}} = \infty$.

- ▷ **Exercise 1:** if $|\mathcal{A}| < \infty$, then $S_{\mathcal{A}}(n) \leq |\mathcal{A}|$ and $V_{\mathcal{A}} \leq \log_2 |\mathcal{A}|$.
- ▷ **Exercise 2:** if $\mathcal{A} = \{(-\infty, a] : a \in \mathbb{R}\}$, then $V_{\mathcal{A}} = 1$. If $\mathcal{A} = \{[a, b] : (a, b) \in \mathbb{R}^2\}$, then $V_{\mathcal{A}} = 2$.
- ▷ **Generalization:** if $\mathcal{A} = \{(-\infty, a_1] \times \cdots \times (-\infty, a_p] : (a_1, \dots, a_p) \in \mathbb{R}^p\}$, then $V_{\mathcal{A}} = p$. If $\mathcal{A} = \{\text{rectangles of } \mathbb{R}^p\}$, then $V_{\mathcal{A}} = 2p$. If $\mathcal{A} = \{\text{convex polygons of } \mathbb{R}^2\}$, then $V_{\mathcal{A}} = \infty$.
- ▷ **Important:** for \mathcal{G} a finite-dimensional vector space of functions $\mathbb{R}^p \rightarrow \mathbb{R}$, and for

$$\mathcal{A} = \{\{x \in \mathbb{R}^p : g(x) \geq 0\} : g \in \mathcal{G}\},$$

then $V_{\mathcal{A}} \leq \dim \mathcal{G}$.

- ▷ **Consequence:** if $\mathcal{A} = \text{subsets of } \mathbb{R}^p$ of the form $\{x \in \mathbb{R}^p : a^\top x + b \geq 0 : a \in \mathbb{R}^p, b \in \mathbb{R}\}$, then $V_{\mathcal{A}} \leq p + 1$.

Sauer's lemma

Theorem (Sauer's lemma)

If $V_{\mathcal{A}} < \infty$, then, for all $n \geq 1$, $S_{\mathcal{A}}(n) \leq \sum_{i=1}^{V_{\mathcal{A}}} \binom{n}{i}$.

- ▷ Exercise: $S_{\mathcal{A}}(n) \leq (n + 1)^{V_{\mathcal{A}}}$.
- ▷ Only two cases:
 - (i) Either $V_{\mathcal{A}} = \infty \rightarrow S_{\mathcal{A}}(n) = 2^n$ for all $n \geq 1$;
 - (ii) Either $V_{\mathcal{A}} < \infty \rightarrow S_{\mathcal{A}}(n) \leq (n + 1)^{V_{\mathcal{A}}}$ for all $n \geq 1$.
- ▷ Never $S_{\mathcal{A}}(n) \sim 2^{\sqrt{n}}$.
- ▷ Important consequence: $\log S_{\mathcal{A}}(n) \leq V_{\mathcal{A}} \log(n + 1)$ and

$$R_n(\mathcal{F}(x_1^n)) \leq 2 \sqrt{\frac{V_{\mathcal{A}} \log(n + 1)}{n}}.$$

Vapnik-Chervonenkis inequality

Theorem (Vapnik-Chervonenkis inequality)

For *all* distributions, one has

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \leq 4 \sqrt{\frac{V_{\mathcal{A}} \log(n+1)}{n}}.$$

- ▷ **Chaining arguments:** $\mathbb{E} \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \leq c \sqrt{\frac{V_{\mathcal{A}}}{n}}.$
- ▷ **Variants:** $\mathbb{P}(\sup_{f \in \mathcal{F}} |P_n(f) - P(f)| \geq t) \leq c_1 S_{\mathcal{A}}(n) e^{-c_2 n \varepsilon^2}.$
- ▷ **Take-home message:** $\mathbb{E} \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| = O\left(\sqrt{\frac{V_{\mathcal{A}} \log n}{n}}\right).$

DON'T FORGET!

The VC inequality is valid for **all distributions!**

Back to learning

- **Reminder:** $L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \leq 2 \sup_{g \in \mathcal{C}} |L_n(g) - L(g)|.$
- **Strategy:**
 - ▷ $L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \leq 2 \sup_{f \in \mathcal{F}} |P_n f - Pf|;$
 - ▷ $\mathcal{F} = \text{indicator functions}$ of the form $f = \mathbb{1}_{[(x,y):g(x)\neq y]}, g \in \mathcal{C};$
 - ▷ $\mathcal{A} = \{A_g : g \in \mathcal{C}\},$ where $A_g = \{(x, y) \in \mathbb{R}^d \times \{0, 1\} : g(x) \neq y\}.$

Proposition

For all $n \geq 1, S_{\bar{\mathcal{A}}}(n) = S_{\mathcal{A}}(n),$ where

$$\bar{\mathcal{A}} = \left\{ \{x \in \mathbb{R}^d : g(x) = 1\} : g \in \mathcal{C} \right\}.$$

In particular, $V_{\bar{\mathcal{A}}} = V_{\mathcal{A}}.$

- ▷ **Notation:** $V_{\mathcal{C}}$ instead of $V_{\bar{\mathcal{A}}}.$
- ▷ **Important consequence:**

$$\mathbb{E} L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) = O\left(\sqrt{\frac{V_{\mathcal{C}} \log n}{n}}\right).$$

Examples

- Linear classification:

$$g(x) = \begin{cases} 1 & \text{if } \sum_{j=1}^d a_j x^{(j)} + a_0 > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $(a_0, a_1, \dots, a_d) \in \mathbb{R}^{d+1}$. Thus,

$$\bar{\mathcal{A}} \subseteq \left\{ \{x \in \mathbb{R}^d : a^\top x + a_0 \geq 0\} : a \in \mathbb{R}^d, a_0 \in \mathbb{R} \right\},$$

and $V_{\mathcal{C}} \leq d + 1$.

- Exercise: $\eta(x) = 0.1\mathbb{1}_{[x^{(1)} \leq 1/2]} + 0.9\mathbb{1}_{[x^{(1)} > 1/2]}$. Compute g^* , L^* , and $\inf_{g \in \mathcal{C}} L(g) - L^*$.
- Closed balls:

$$\bar{\mathcal{A}} = \left\{ \{x \in \mathbb{R}^d : \sum_{j=1}^d (x^{(j)} - a_j)^2 \leq a_0\} : (a_0, a_1, \dots, a_d) \in \mathbb{R}^{d+1} \right\}.$$

We see that $\bar{\mathcal{A}} \subseteq \{ \{x \in \mathbb{R}^d : g(x) \geq 0\} : g \in \mathcal{G} \}$, where \mathcal{G} is a vector space of dimension $d + 2 \rightarrow V_{\mathcal{C}} \leq d + 2$.

Examples, cont.

- Convex polygons of \mathbb{R}^2 : bad idea since $V_{\mathcal{A}} = \infty$.
- Generalized linear classification: $\psi_1, \dots, \psi_{d^*}$ functions from $\mathbb{R}^d \rightarrow \mathbb{R}$ and

$$g(x) = \begin{cases} 1 & \text{if } \sum_{j=1}^{d^*} a_j \psi_j(x) + a_0 > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $(a_0, a_1, \dots, a_{d^*}) \in \mathbb{R}^{d^*+1}$.

- ▷ $\psi_j(x) = x^{(j)}$: linear classifiers ✓
- ▷ ψ_j = coordinates + products of coordinates:

$$\mathcal{A} \subseteq \left\{ a_0 + \sum_{j=1}^d a_j x^{(j)} + \sum_{j=1}^d b_j (x^{(j)})^2 + \sum_{1 \leq j_1 < j_2 \leq d} c_{j_1 j_2} x^{(j_1)} x^{(j_2)} \geq 0 \right\}.$$

Conclusion: $V_{\mathcal{C}} \leq d^* + 1$, where $d^* = 2d + \frac{d(d-1)}{2}$.

Margin-based bounds

Drawbacks of empirical risk minimization

- For classes with finite VC dimension:

$$\mathbb{E}L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) = O\left(\sqrt{\frac{V_{\mathcal{C}} \log n}{n}}\right).$$

- However, such classes are nearly always too small.
- $L(g_n^*)$ is typically far from the Bayes risk L^* .
- Example: for any class \mathcal{C} with finite VC dimension and for all $\varepsilon \in (0, 1/2)$, there exists (X, Y) such that

$$\inf_{g \in \mathcal{C}} L(g) - L^* > 1/2 - \varepsilon.$$

- Besides, minimizing $L_n(g)$ is often a computationally difficult problem.
- Solution: modify the empirical functional to be minimized.

A different framework

- We consider **± 1 -classifiers** of the form

$$g_f(x) = \begin{cases} 1 & \text{if } f(x) > 0 \\ -1 & \text{otherwise,} \end{cases}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $f \in \mathcal{F}$ bounded.

- Error probability:** $\mathbb{P}(Yf(X) < 0) \leq L(g_f) \leq \mathbb{P}(Yf(X) \leq 0)$.
- Notation:** $L(f) = L(g_f)$.
- Empirical error probability:**

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[Y_i f(X_i) < 0]} \leq L_n(f) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[Y_i f(X_i) \leq 0]}.$$

- The product $Yf(X)$ is called the **margin**.



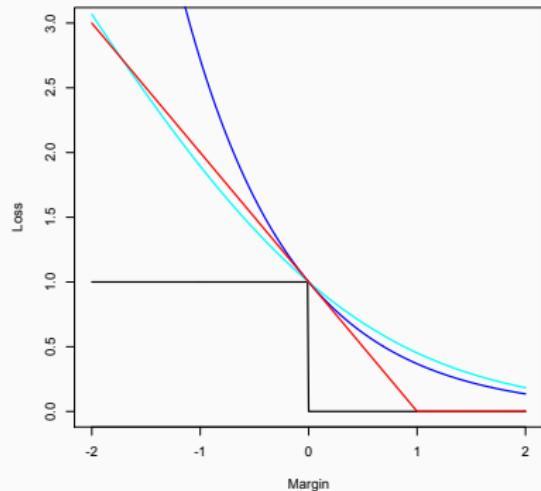
Large margin = good confidence.

Loss functions

- $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ a **loss function** such that
 $\phi(0) = 1, \quad \mathbb{1}_{[x \leq 0]} \leq \phi(x), \quad \phi$ is L_ϕ -Lipschitz on the range of \mathcal{F} .

- **Typical choices:**

- ▷ **Exponential:** $\phi(x) = e^{-x}$
- ▷ **Logit:** $\phi(x) = \log_2(1 + e^{-x})$
- ▷ **Hinge:** $\phi(x) = \max(1 - x, 0) !!!$



- **Risk functional:** $A(f) = \mathbb{E}\phi(Yf(X))$.
- **Empirical risk functional:** $A_n(f) = \frac{1}{n} \sum_{i=1}^n \phi(Y_i f(X_i))$.
- **Clearly:** $L(f) \leq A(f)$ and $L_n(f) \leq A_n(f)$.

Contraction principle

Theorem

Let A be a bounded subset of \mathbb{R}^n . If

$$\bar{A} = \left\{ \sum_{j=1}^N c_j a^{(j)} : N \in \mathbb{N}, a^{(j)} \in A, \sum_{j=1}^N |c_j| = 1 \right\},$$

then $R_n(\bar{A}) = R_n(A)$. In addition, if $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is a function with $\psi(0) = 0$ and Lipschitz constant L_ψ , then

$$R_n(\psi \circ A) \leq 2L_\psi R_n(A) \quad (\text{contraction principle}),$$

where $\psi \circ A$ is the set of vectors of the form $(\psi(a_1), \dots, \psi(a_n)) \in \mathbb{R}^n$, $a = (a_1, \dots, a_n) \in A$.

Rademacher averages and performance bounds

Theorem

Let $f_n^* \in \arg \min_{f \in \mathcal{F}} A_n(f)$, let B denote a uniform upper bound on $\phi(yf(x))$, and let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$,

$$L(f_n^*) \leq A_n(f_n^*) + 4L_\phi \mathbb{E} R_n(\mathcal{F}(X_1^n)) + B \sqrt{\frac{\log(1/\delta)}{2n}}.$$

- ▷ Warning: L is a probability and A_n is empirical!
- ▷ The Rademacher average of \mathcal{F} bounds the performance of f_n^* .

Illustration 1: Weighted voting schemes

- Class:

$$\mathcal{F}_\lambda = \left\{ f = \sum_{j=1}^N c_j g_j : N \in \mathbb{N}, g_1, \dots, g_N \in \mathcal{C}, \sum_{j=1}^N |c_j| = \lambda \right\},$$

where \mathcal{C} = class of ± 1 base classifiers.

- Rademacher averages:

$$R_n(\mathcal{F}_\lambda(x_1^n)) = \lambda R_n(\mathcal{C}(x_1^n)) \leq 2\lambda \sqrt{\frac{V_{\mathcal{C}} \log(n+1)}{n}}.$$

- Example:

▷ \mathcal{C} = classifiers of the form $g(x) = 2\mathbb{1}_{[x \geq a]} - 1$, $a \in \mathbb{R}$;

▷ $V_{\mathcal{C}} = 1$ and $\overline{\mathcal{F}_\lambda}$ = all functions of total variation bounded by 2λ .

- To summarize: with probability $1 - \delta$,

$$L(f_n^*) \leq A_n(f_n^*) + 8L_\phi \lambda \sqrt{\frac{V_{\mathcal{C}} \log(n+1)}{n}} + B \sqrt{\frac{\log(1/\delta)}{2n}}.$$

- Advantage: only $V_{\mathcal{C}}$ in the bound. Inconvenient: A_n instead of L_n .

Weighted voting schemes, cont.

- γ a positive parameter and

$$\phi(x) = \begin{cases} 1 & \text{if } x \leq 0 \\ 0 & \text{if } x \geq \gamma \\ 1 - x/\gamma & \text{otherwise.} \end{cases}$$

- Quick check: $B = 1$, $L_\phi = 1/\gamma$, $\mathbb{1}_{[x \leq 0]} \leq \phi(x) \leq \mathbb{1}_{[x < \gamma]}$.
- Consequence: $A_n(f) \leq L_n^\gamma(f)$, where

$$L_n^\gamma(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[Y_i f(X_i) < \gamma]} \quad \text{is the margin error.}$$

- Properties: $L_n(f) \leq L_n^\gamma(f)$ and $L_n^\gamma(f)$ is nondecreasing in γ .
- Interpretation: $L_n^\gamma(f)$ counts the number of misclassified pairs + those which are well classified but only with a small margin.
- Margin-based bound: with probability at least $1 - \delta$,

$$L(f_n^*) \leq L_n^\gamma(f_n^*) + \frac{8\lambda}{\gamma} \sqrt{\frac{V_C \log(n+1)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

Illustration 2: Kernel methods

- Examples: Support Vector Machines and kernel Fisher discriminant.
- Necessitates a detour through the theory of Hilbert spaces.

Definition (Reproducing kernel)

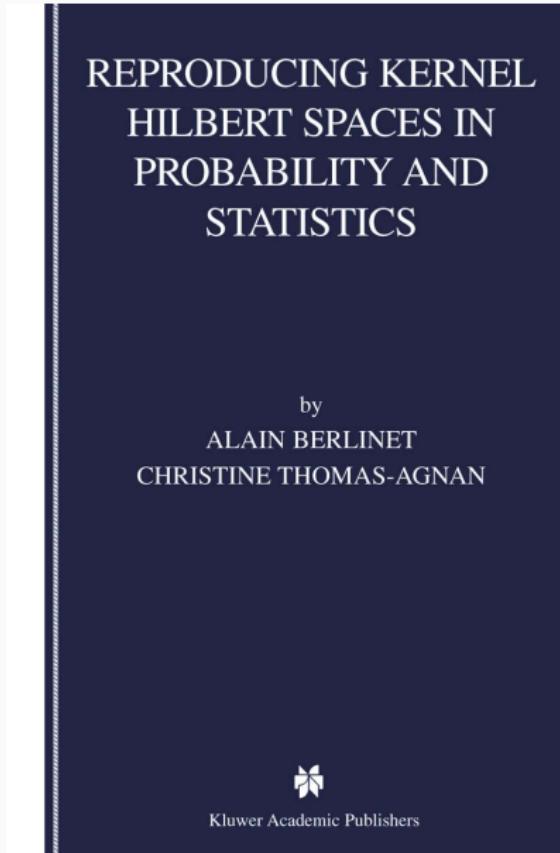
Let \mathcal{F} be a Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ defined on a non-empty set \mathcal{X} . A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a reproducing kernel of \mathcal{F} if

- (i) $\forall x \in \mathcal{X}, k(\cdot, x) \in \mathcal{F};$
- (ii) $\forall f \in \mathcal{F}, \langle f, k(\cdot, x) \rangle = f(x)$ (reproducing property).

- ▷ $k(x, x') = k(x', x) = \langle k(\cdot, x), k(\cdot, x') \rangle.$
- ▷ Vocabulary: \mathcal{F} is a reproducing kernel Hilbert space (RKHS).

DON'T FORGET!

RKHS \Rightarrow Hilbert space, but Hilbert space $\not\Rightarrow$ RKHS.



Kernel methods, cont.

- **Example 1:** any finite-dimensional Hilbert space of functions is a RKHS, with $k(x, x') = \sum_{i=1}^n e_i(x)e_i(x')$.
- **Example 2:** the space $L^2(\mathbb{R})$ is **not** a RKHS.
- **Example 3:** the space

$$\mathcal{F} = \{f : f(0) = 0, f \text{ is absolutely continuous, } f, f' \in L^2(\mathbb{R})\}$$

is a RKHS, with $k(x, x') = \frac{1}{2}e^{-|x-x'|}$.

Kernel methods, cont.

- Key: any reproducing kernel is **symmetric** and **positive definite**, i.e.,

$$\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0 \quad \forall n \geq 1, \forall (x_1, \dots, x_n) \in \mathcal{X}^n, \forall (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n.$$

- Reciprocal property? Yes.

Theorem (Moore-Aronszajn theorem)

Let k be a **positive definite** function on $\mathcal{X} \times \mathcal{X}$. There exists only one Hilbert space \mathcal{F} of functions on \mathcal{X} with k as **reproducing kernel**. The subspace \mathcal{F}_0 of \mathcal{F} spanned by the functions $(k(\cdot, x))_{x \in \mathcal{X}}$ is dense in \mathcal{F} , and \mathcal{F} is the set of functions on \mathcal{X} that are pointwise limits of Cauchy sequences in \mathcal{F}_0 with the inner product

$$\langle f, g \rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x'_j),$$

where $f = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$ and $g = \sum_{j=1}^m \beta_j k(\cdot, x'_j)$.

Statistical consequences

- Non-linear embedding:

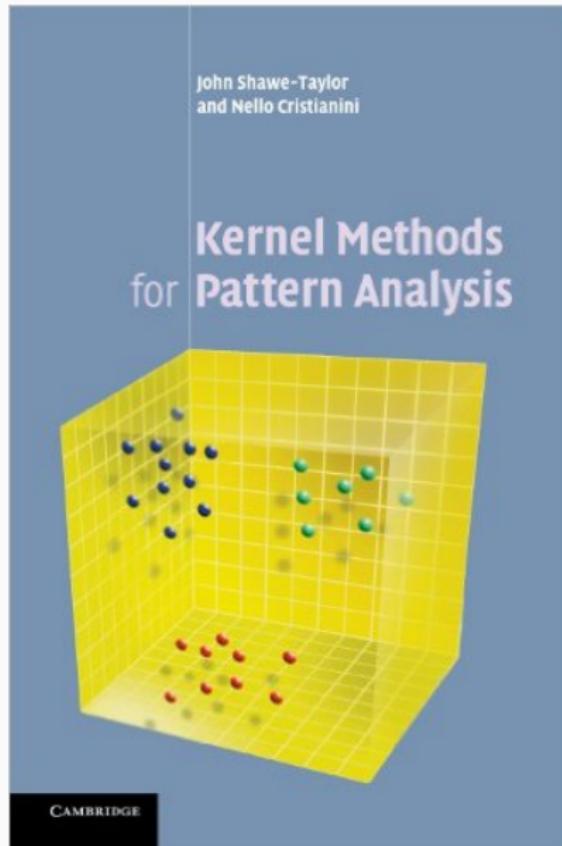
$$\begin{aligned}\varphi &: \mathcal{X} \rightarrow \mathcal{F} \\ x &\mapsto k(\cdot, x).\end{aligned}$$

- Rationale: $\forall (x, x') \in \mathcal{X}^2, \langle \varphi(x), \varphi(x') \rangle = k(x, x')$.
- Advantage: any linear algorithm based on computing inner products can be extended into a non-linear version by replacing the inner products by a kernel function → kernel trick.
- Even though the algorithm remains of low complexity, it works in a big class of functions.
- Science of kernel design (vectors, strings, graphs, etc.).

DON'T FORGET!

φ is never computed.

Reference



Optimization in a RKHS

- Kernel algorithms often perform minimization over the **ball of a RKHS**:

$$\mathcal{F}_\lambda = \{f \in \mathcal{F} : \|f\| \leq \lambda\}$$

$$\approx \left\{ f = \sum_{i=1}^N c_i k(\cdot, x_i) : N \in \mathbb{N}, x_1, \dots, x_N \in \mathcal{X}, \sum_{i,j=1}^N c_i c_j k(x_i, x_j) \leq \lambda^2 \right\}.$$

- Example:** Support Vector Machines (SVM).

Theorem

$$\frac{\lambda}{n\sqrt{2}} \sqrt{\sum_{i=1}^n k(X_i, X_i)} \leq R_n(\mathcal{F}_\lambda(X_1^n)) \leq \frac{\lambda}{n} \sqrt{\sum_{i=1}^n k(X_i, X_i)}.$$

- Nice bound**, which can be computed very easily from the data.
- Performance bound in the RKHS**: with probability $1 - \delta$,

$$L(f_n^*) \leq L_n^\gamma(f_n^*) + \frac{4\lambda}{\gamma n} \sqrt{\sum_{i=1}^n k(X_i, X_i)} + 3\sqrt{\frac{\log(2/\delta)}{2n}}.$$

References

The Annals of Statistics
2008, Vol. 36, No. 2, 489–511
DOI: 10.1214/060720006002000309
© Institute of Mathematical Statistics, 2008

STATISTICAL PERFORMANCE OF SUPPORT VECTOR MACHINES

BY GILLES BLANCHARD,¹ OLIVIER BOUSQUET AND PASCAL MASSART

Fraunhofer-Institute FIRST, Google and Université Paris-Sud

The support vector machine (SVM) algorithm is well known to the computer learning community for its very good practical results. The goal of the present paper is to study this algorithm from a statistical perspective, using tools of concentration theory and empirical processes.

Our main result builds on an observation made by other authors that the SVM can be viewed as a statistical regularization procedure. From this point of view, it can also be interpreted as a model selection principle using a penalized criterion. It is thus possible to adapt some of the results related to model selection in the framework to study two important points: (1) what is the minimum penalty and how does it compare to the penalty actually used in the SVM algorithm; (2) is it possible to obtain “oracle inequalities” in that setting, for the specific loss function used in the SVM algorithm? We show that the answer to the latter question is positive and provides relevant insight to the former. Our result shows that it is possible to obtain fast rates of convergence for SVMs.

1. Introduction. The success of the support vector machine (SVM) algorithm for pattern recognition is probably mainly due to the number of remarkable experimental results that have been obtained in very diverse domains of application. The algorithm itself can be written as a nice convex optimization problem for which there exists a unique optimum except in rare degenerate cases. It can also be expressed as the minimization of a regularized functional where the regularizer is the squared norm in a Hilbert space of functions on the input space. Although these are nice mathematical formulations, quite amenable to analysis, the statistical behavior of this algorithm remains only partially understood. Our goal in this work is to investigate the properties of the SVM algorithm in a statistical setting.

1.1. The abstract classification problem and convex loss approximation. We consider a generic (binary) classification problem, defined by the following setting: assume that the product $\mathcal{X} \times \mathcal{Y}$ is a measurable space endowed with an unknown probability measure P , where $\mathcal{Y} = \{-1, 1\}$ and \mathcal{X} is called the input space. The pair (X, Y) denotes a random variable with values in $\mathcal{X} \times \mathcal{Y}$ distributed according to P . We will denote P_X the marginal distribution of variable X . We observe a

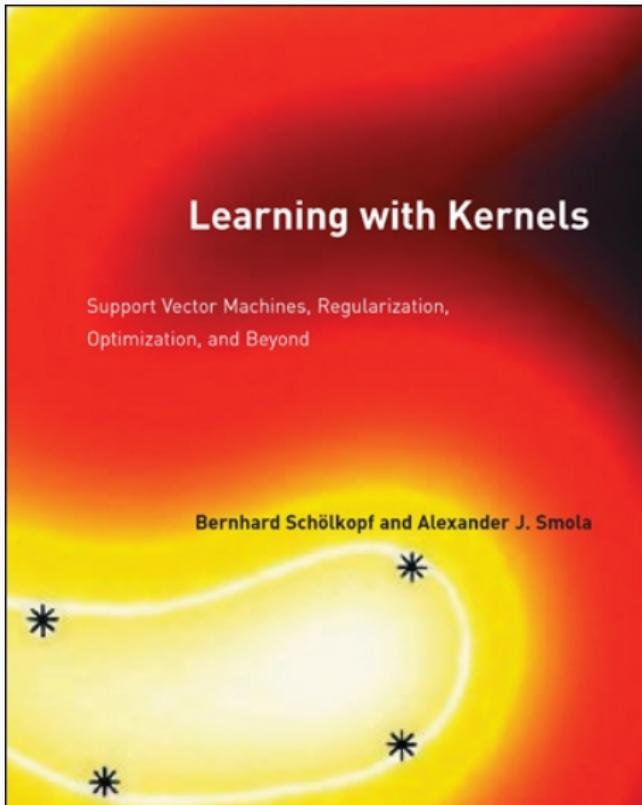
Received October 2006; revised April 2007.

¹Supported in part by a grant of the Humboldt Foundation.

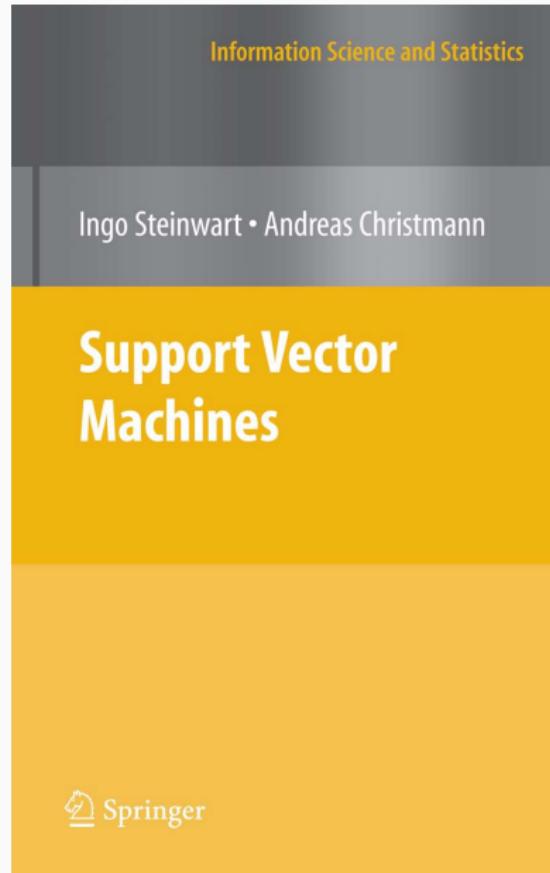
AMS 2000 subject classifications. 62G05, 62G20.

Key words and phrases. Classification, support vector machine, model selection, oracle inequality.

References



References

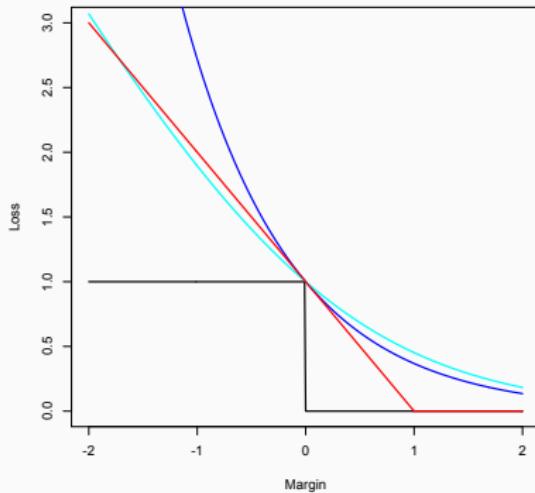


Convex loss functions

Convex surrogates

- **Error probability:** $\mathbb{E} \mathbb{1}_{[Yf(X) < 0]} \leq L(g_f) \leq \mathbb{E} \mathbb{1}_{[Yf(X) \leq 0]}$.
- **Idea:** smooth the indicator by a convex function of $Yf(X)$.
- $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ a loss function such that ϕ is **strictly decreasing, strictly convex**, of class C^1 , $\phi(0) = 1$, and $\lim_{x \rightarrow \infty} \phi(x) = 0$.
- **Risk functionals:** $A(f) = \mathbb{E}\phi(Yf(X))$ and $A_n(f) = \frac{1}{n} \sum_{i=1}^n \phi(Y_i f(X_i))$.
- **Typical choices:**

- ▷ **Exponential:** $\phi(x) = e^{-x}$
- ▷ **Logit:** $\phi(x) = \log_2(1 + e^{-x})$
- ▷ **Hinge:** $\phi(x) = \max(1 - x, 0)$!!!!



Convex surrogates, cont.

- What is $f^* = \arg \min_f A(f)$ for ϕ strictly convex and differentiable?
- Clearly, $\mathbb{E}(\phi(Yf(X)) | X = x) = \eta(x)\phi(f(x)) + (1 - \eta(x))\phi(-f(x))$.
- Consequence: $f^*(x) = \arg \min_\alpha h_{\eta(x)}(\alpha)$, where

$$h_\eta(\alpha) = \eta\phi(\alpha) + (1 - \eta)\phi(-\alpha), \quad \eta \in [0, 1].$$

- Note: h_η is strictly convex and therefore f^* is well defined.

Convex surrogates, cont.

- The minimum is achieved for $h'_\eta(\alpha) = 0$, that is, when

$$\frac{\eta}{1-\eta} = \frac{\phi'(-\alpha)}{\phi'(\alpha)}.$$

- Since ϕ' is strictly increasing, the solution is **positive** if and only if $\eta > 1/2$.
- Conclusion:** $g^*(x) \stackrel{\text{def}}{=} 2\mathbb{1}_{[f^*(x)>0]} - 1$ is the **Bayes classifier**!
- Examples:**
 - $\triangleright \phi(x) = e^{-x} \rightarrow f^*(x) = \frac{1}{2} \log\left(\frac{\eta(x)}{1-\eta(x)}\right);$
 - $\triangleright \phi(x) = \log_2(1 + e^{-x}) \rightarrow f^*(x) = \log\left(\frac{\eta(x)}{1-\eta(x)}\right).$
- For the **hinge loss**: $f^*(x) = 2\mathbb{1}_{[\eta(x)>1/2]} - 1 = \text{Bayes classifier itself!}$

Connection between 0-1-risk and ϕ -risk

- **Objective:** connect $L(f) - L^*$ with $A(f) - A^*$.
- **Tool:** $H : [0, 1] \rightarrow \mathbb{R}$ defined by $H(\eta) = \inf_{\alpha} h_\eta(\alpha)$.

Lemma

Let ϕ be a **convex** loss function such that the following hold:

- (i) $f^*(x) > 0$ if and only if $\eta(x) > 1/2$;
- (ii) There exist constants $c \geq 0$ and $s \geq 1$ satisfying

$$\left| \frac{1}{2} - \eta \right|^s \leq c^s (1 - H(\eta)), \quad \eta \in [0, 1].$$

Then, for **any** function $f : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$L(f) - L^* \leq 2c(A(f) - A^*)^{1/s}.$$

- ▷ **Exponential:** $H(\eta) = 2\sqrt{\eta(1-\eta)}$.
- ▷ **Logit:** $H(\eta) = -\eta \log_2 \eta - (1-\eta) \log_2 (1-\eta)$.
- ▷ **In both cases:** $c = 1/\sqrt{2}$ and $s = 2$.
- ▷ **Hinge:** $H(\eta) = 2 \min(\eta, 1-\eta) \rightarrow c = 1/2$ and $s = 1$.

Excess risk of convex minimizers

Context: \mathcal{C} = a class of ± 1 **base** classifiers, and

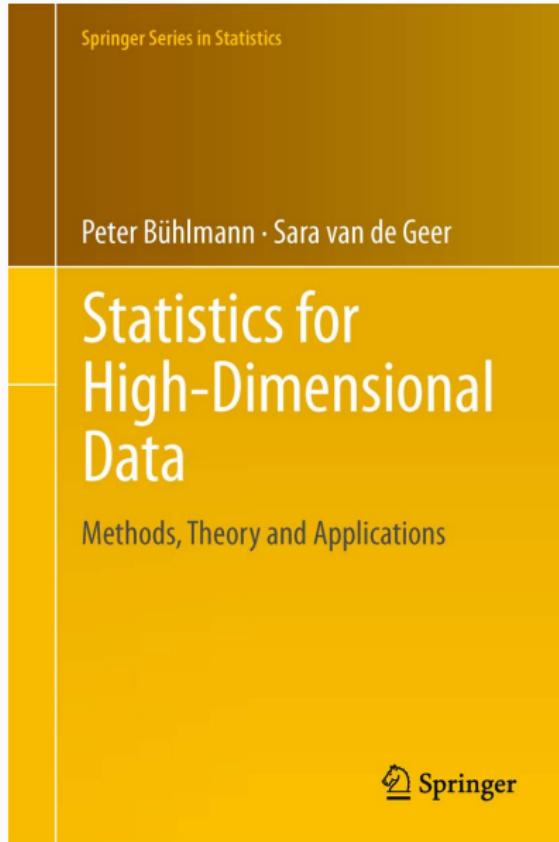
$$\mathcal{F}_\lambda = \left\{ f = \sum_{j=1}^N c_j g_j : N \in \mathbb{N}, g_1, \dots, g_N \in \mathcal{C}, \sum_{j=1}^N |c_j| = \lambda \right\}.$$

Theorem

Let $f_n^* \in \arg \min_{f \in \mathcal{F}_\lambda} A_n(f)$, using either the **exponential** or the **logit** loss function, and let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$,

$$L(f_n^*) - L^* \leq 2 \left(8L_\phi \lambda \sqrt{\frac{V_{\mathcal{C}} \log(n+1)}{n}} + B \sqrt{\frac{\log(1/\delta)}{2n}} \right)^{1/2} + \sqrt{2} \left(\inf_{f \in \mathcal{F}_\lambda} A(f) - A^* \right)^{1/2}.$$

- ▷ **Exponential:** $B = e^\lambda$ and $L_\phi = e^\lambda$.
- ▷ **Logit:** $B = \log_2(1 + e^\lambda)$ and $L_\phi = 1/\log 2$.
- ▷ If $\inf_{f \in \mathcal{F}_\lambda} A(f) - A^* = 0$, then $L(f_n^*) - L^* = O\left(\left(\frac{\log n}{n}\right)^{1/4}\right)$.
- ▷ The exponent in the rate of convergence is **dimension-free**!
- ▷ Convex optimization → **boosting algorithms**.



Model selection

Oracle inequalities

- Choosing the right set \mathcal{C} of possible classifiers is a **key to success**.
- If \mathcal{C} is **too large**: overfitting risk.
- If \mathcal{C} is **too small**: unable to approximate complex decisions.
- Welcome to the domain of **model selection**.
- A possibly infinite collection of **models** $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_k, \dots$
- For **each** model: empirical risk minimizer $g_{n,k}^* \in \mathcal{C}_k$.
- **Model selection problem**: select the best among the $(g_{n,k}^*)_k$.
- **Oracle inequality**:

$$\mathbb{E}L(g_{n,\hat{k}}^*) - L^* \leq C \inf_k \left(\mathbb{E}L(g_{n,k}^*) - L_k^* + L_k^* - L^* + \gamma(k, n) \right),$$

where $L_k^* = \inf_{g \in \mathcal{C}_k} L(g)$.

Data-driven penalization

- **Penalty-based** model selection: choose \hat{k} that minimizes

$$L_n(g_{n,k}^*) + \text{pen}(n, k).$$

- The penalty should estimate the **amount of overfitting**

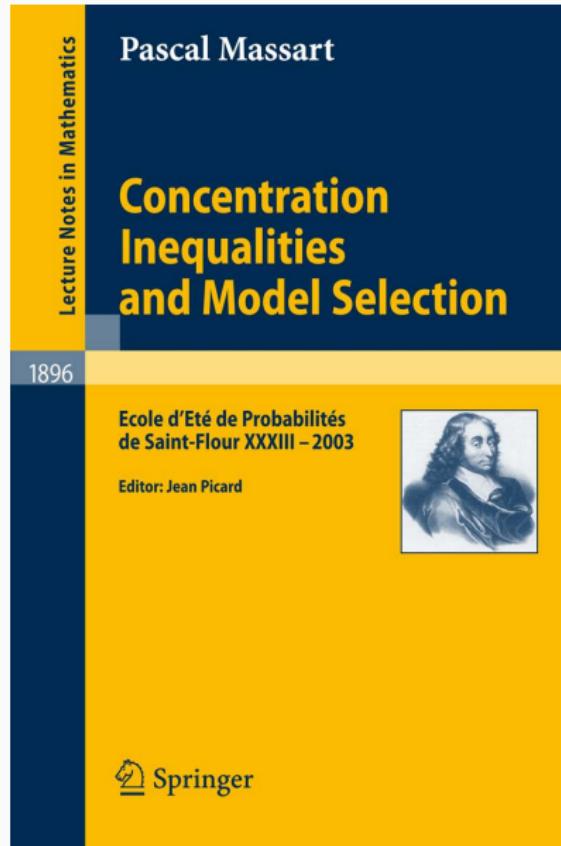
$$\mathbb{E} \sup_{g \in \mathcal{C}_k} |L_n(g) - L(g)|.$$

- **Distribution-free** penalties \rightarrow VC dimension $V_{\mathcal{C}_k} \rightarrow$ too conservative.
- **Data-driven** penalties \rightarrow Rademacher averages.

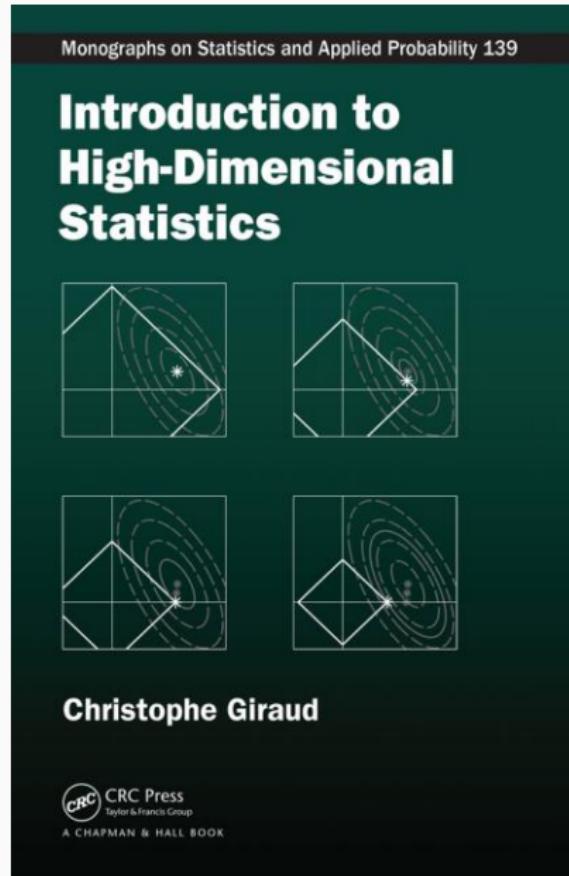


Model selection is a **science**.

References



References



Oracle inequality

Theorem

Let $\text{pen}(n, k)$ be defined by

$$\text{pen}(n, k) = 2R_n(\mathcal{F}_k(X_1^n)) + 4\sqrt{\frac{\log k}{n}}.$$

Then

$$\mathbb{E}L(g_{n,\hat{k}}^*) - L^* \leq \inf_k \left(4\mathbb{E}R_n(\mathcal{F}_k(X_1^n)) + L_k^* - L^* + 4\sqrt{\frac{\log k}{n}} \right) + \frac{\sqrt{2}\pi^{5/2}}{3\sqrt{n}}.$$

- ▷ Optimal constant in front of the infimum.
- ▷ Satisfactory in the general case.

Stone's theorem

Plug-in principle

- Starting point:

$$g^*(x) = \begin{cases} 1 & \text{if } r(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

- Idea: estimate $r(x)$ from the training data $\mathcal{D}_n \rightsquigarrow r_n(x)$.
- Plug-in classifier:

$$g_n(x) = \begin{cases} 1 & \text{if } r_n(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

- Question 1: connection $r_n \leftrightarrow L(g_n)$?
- Question 2: which choice for r_n ?

DON'T FORGET!

Plug-in = regression estimation problem.

Connection $r_n \leftrightarrow L(g_n)$

Theorem (Classification and regression)

Let r_n be a *regression function estimate* of r , and let g_n be the *corresponding plug-in classifier*. Then

$$0 \leq L(g_n) - L^* \leq 2 \int_{\mathbb{R}^d} |r_n(x) - r(x)| \mu(dx).$$

In particular, for *all* $p \geq 1$,

$$0 \leq L(g_n) - L^* \leq 2 \left(\int_{\mathbb{R}^d} |r_n(x) - r(x)|^p \mu(dx) \right)^{1/p},$$

and

$$0 \leq \mathbb{E}L(g_n) - L^* \leq 2 \mathbb{E}^{1/p} |r_n(X) - r(X)|^p.$$

▷ Take-home message:

$$\mathbb{E} \int_{\mathbb{R}^d} |r_n(x) - r(x)|^2 \mu(dx) \rightarrow 0$$

implies that the corresponding plug-in classifier g_n is *consistent*.

Local averaging estimates

- **Definition:** $r_n(x) = \sum_{i=1}^n W_{ni}(x) Y_i$.
- **Weight vector:** $(W_{n1}(x), \dots, W_{nn}(x))$.
- **Interpretation:** X_i “close” to x should provide more information.
- **Often** (but not always) $(W_{n1}(x), \dots, W_{nn}(x))$ is a **probability vector**.
- **Important:** each $W_{ni}(x)$ is a function of x and X_1, \dots, X_n (**not** Y_1, \dots, Y_n).
- **Equivalently:** $r_n(x) = \sum_{i=1}^n W_{ni}(x) \mathbb{1}_{[Y_i=1]}$.
- Companion **plug-in classifier**:

$$g_n(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^n W_{ni}(x) Y_i > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

- Whenever $\sum_{i=1}^n W_{ni}(x) = 1$:

$$g_n(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^n W_{ni}(x) \mathbb{1}_{[Y_i=1]} > \sum_{i=1}^n W_{ni}(x) \mathbb{1}_{[Y_i=0]} \\ 0 & \text{otherwise.} \end{cases}$$

Example 1: kernel estimate

- Definition:

$$r_n(x) = \frac{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i}{\sum_{j=1}^n K\left(\frac{x-X_j}{h}\right)}.$$

- Kernel K : a nonnegative real-valued function on \mathbb{R}^d .
- Bandwidth h : a positive number (may depend on n).
- Weights:

$$W_{ni}(x) = \frac{K\left(\frac{x-X_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{x-X_j}{h}\right)}.$$

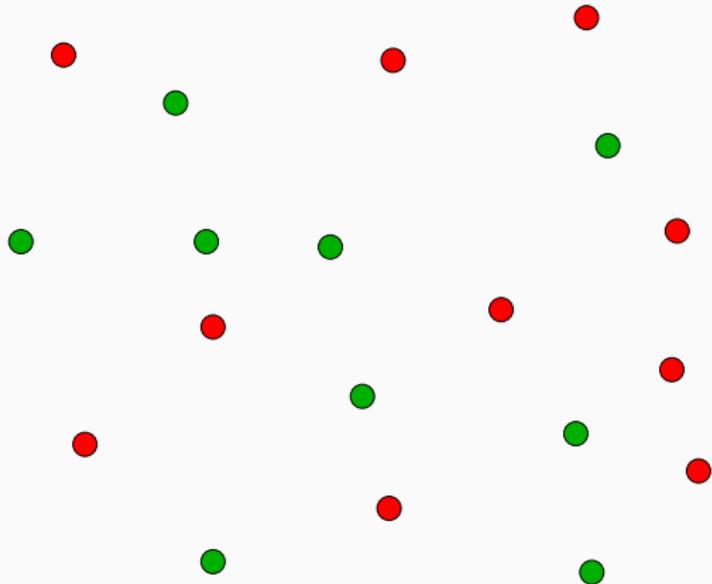
- If both denominator and numerator are zero: $r_n(x) = \frac{1}{n} \sum_{i=1}^n Y_i$.
- Kernels:

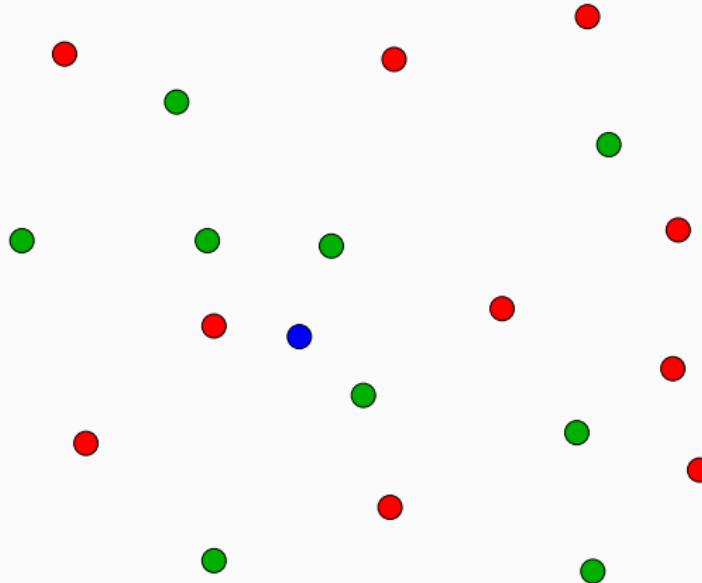
▷ Naive: $K(z) = \mathbb{1}_{[\|z\| \leq 1]}$,

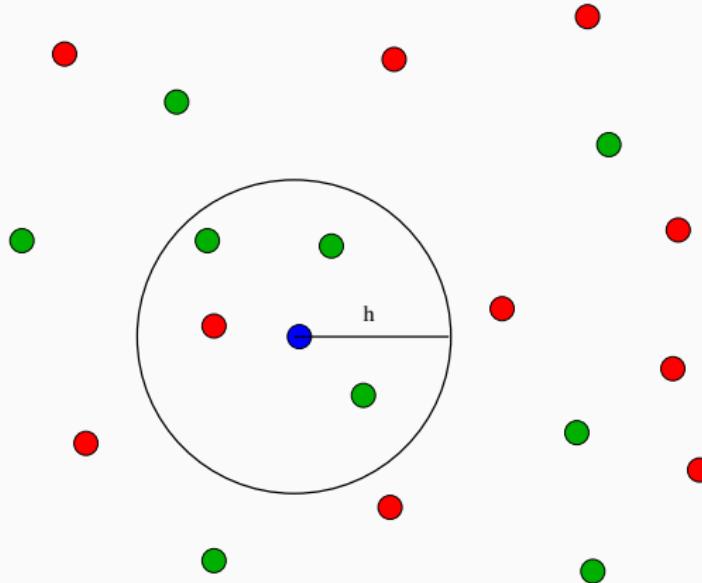
$$r_n(x) = \frac{\sum_{i=1}^n \mathbb{1}_{[\|x-X_i\| \leq h]} Y_i}{\sum_{j=1}^n \mathbb{1}_{[\|x-X_j\| \leq h]}};$$

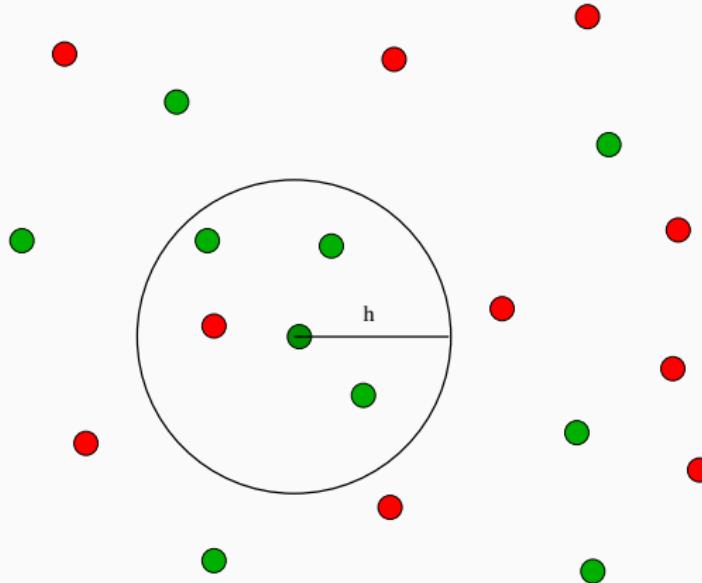
▷ Epanechnikov: $K(z) = (1 - \|z\|^2) \mathbb{1}_{[\|z\| \leq 1]}$;

▷ Gaussian: $K(z) = e^{-\|z\|^2}$.







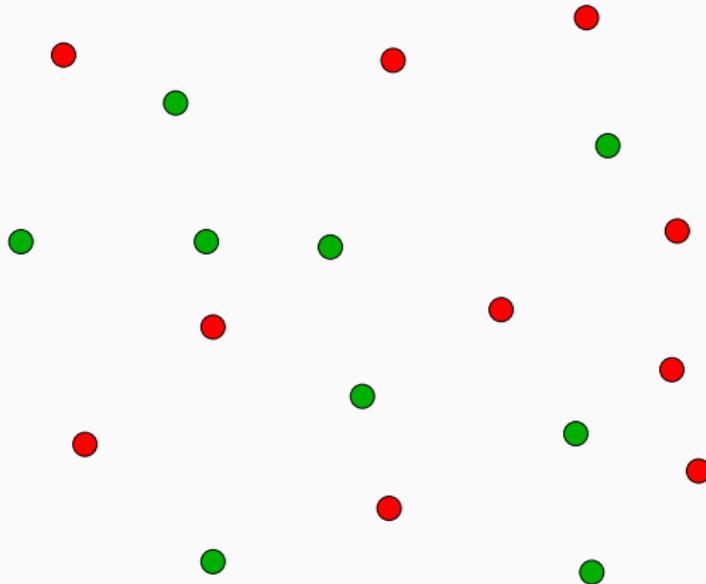


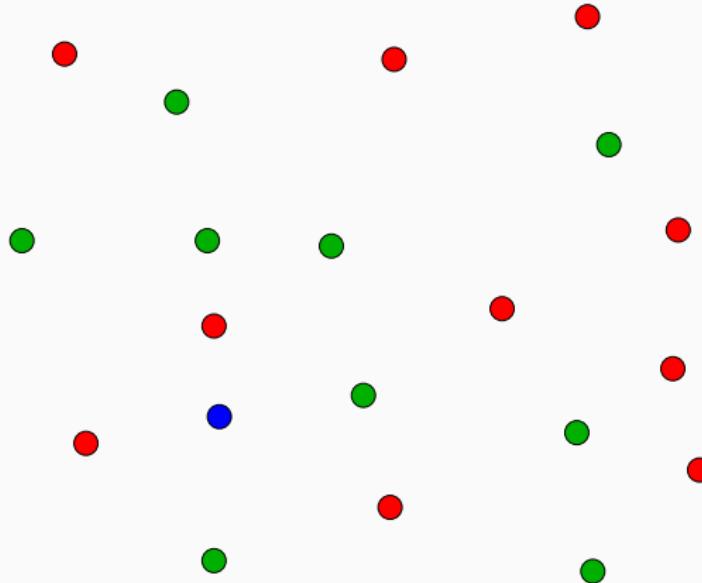
Example 2: nearest neighbor (NN) estimate

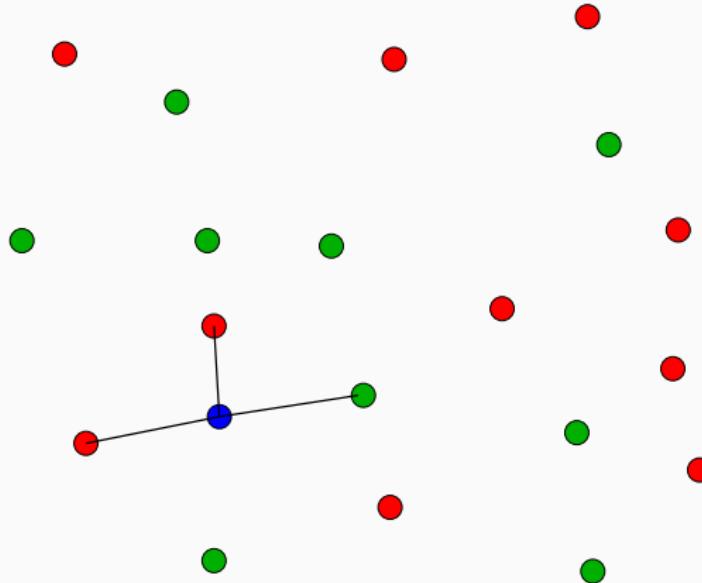
- Definition:
 - ▷ $(X_{(1)}(x), Y_{(1)}(x)), \dots, (X_{(n)}(x), Y_{(n)}(x))$ reordering of \mathcal{D}_n according to
$$\|X_{(1)}(x) - x\| \leq \dots \leq \|X_{(n)}(x) - x\|;$$
 - ▷ Whenever $\|X_i - x\| = \|X_j - x\|$ and $i < j$, we declare X_i closer to x ;
 - ▷ NN estimate: $r_n(x) = \sum_{i=1}^n v_{ni} Y_{(i)}(x)$, where $\sum_{i=1}^n v_{ni} = 1$.
- $(\Sigma_1, \dots, \Sigma_n)$: permutation of $(1, \dots, n)$ such that X_i is the Σ_i -th nearest neighbor of x for all i .
- Local averaging: $r_n(x) = \sum_{i=1}^n W_{ni}(x) Y_i$, where $W_{ni}(x) = v_{n\Sigma_i}$.
- k -NN estimate:

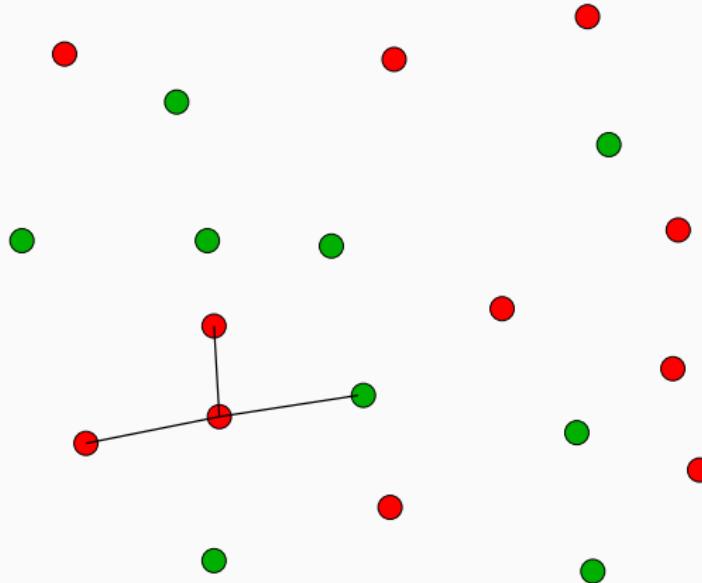
$$v_{ni} = \begin{cases} \frac{1}{k} & \text{for } 1 \leq i \leq k \\ 0 & \text{for } k < i \leq n. \end{cases}$$

- To keep in mind: $r_n(x) = \frac{1}{k} \sum_{i=1}^k Y_{(i)}(x)$.









Stone's theorem

$r_n(x) = \sum_{i=1}^n W_{ni}(x) Y_i$, with $(W_{n1}(x), \dots, W_{nn}(x))$ a probability vector.

Theorem (Stone's theorem)

Assume that for any distribution of X , the weights satisfy the following conditions:

- (i) There is a constant C such that, for every Borel measurable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with $\mathbb{E}|f(X)| < \infty$,

$$\mathbb{E}\left(\sum_{i=1}^n W_{ni}(X)|f(X_i)|\right) \leq C\mathbb{E}|f(X)| \quad \text{for all } n \geq 1.$$

- (ii) For all $a > 0$,

$$\mathbb{E}\left(\sum_{i=1}^n W_{ni}(X)\mathbb{1}_{[\|X_i - X\| > a]}\right) \rightarrow 0.$$

- (iii) One has

$$\mathbb{E} \max_{1 \leq i \leq n} W_{ni}(X) \rightarrow 0.$$

Then the corresponding plug-in classifier g_n is universally consistent, i.e., $\mathbb{E}L(g_n) \rightarrow L^*$ for all distributions of (X, Y) .

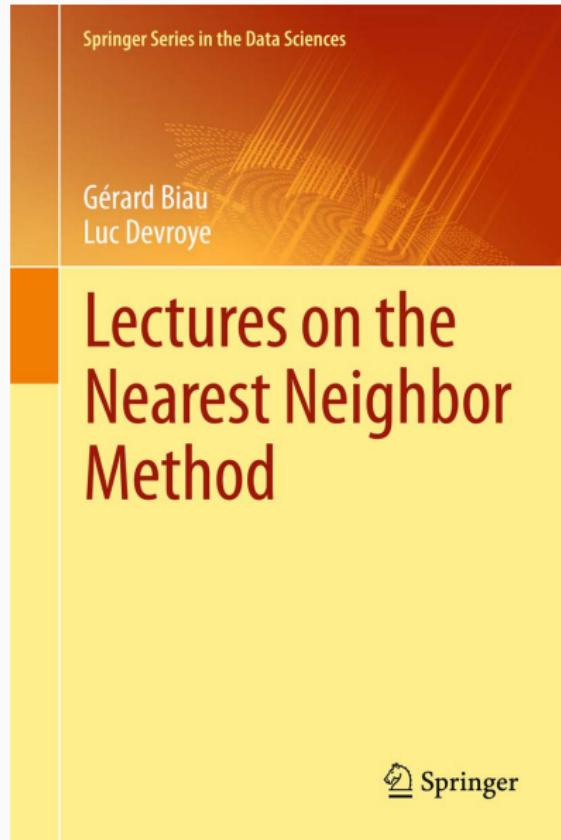
Discussion

- Condition (i) is merely technical.
- Condition (ii) ensures that $r_n(X)$ is asymptotically mostly influenced by the data points close to X .
- Condition (iii) states that asymptotically all weights become small.
- No single observation has a too large contribution to the estimate.
- The number of points in the averaging must tend to infinity.



Charles Stone

k-nearest neighbor classifiers



Reminder

- Reordering $(X_{(1)}(x), Y_{(1)}(x)), \dots, (X_{(n)}(x), Y_{(n)}(x))$ according to

$$\|X_{(1)}(x) - x\| \leq \dots \leq \|X_{(n)}(x) - x\|.$$

- Whenever $\|X_i - x\| = \|X_j - x\|$ and $i < j$, we declare X_i closer to x .
- k -NN regression function estimate: $r_n(x) = \frac{1}{k} \sum_{i=1}^k Y_{(i)}(x)$.
- k -NN classifier:

$$g_n(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^k \mathbb{1}_{[Y_{(i)}(x)=1]} > \sum_{i=1}^k \mathbb{1}_{[Y_{(i)}(x)=0]} \\ 0 & \text{otherwise.} \end{cases}$$

DON'T FORGET!

If X has a density, then there is no distance tie.

Universal consistency

Theorem

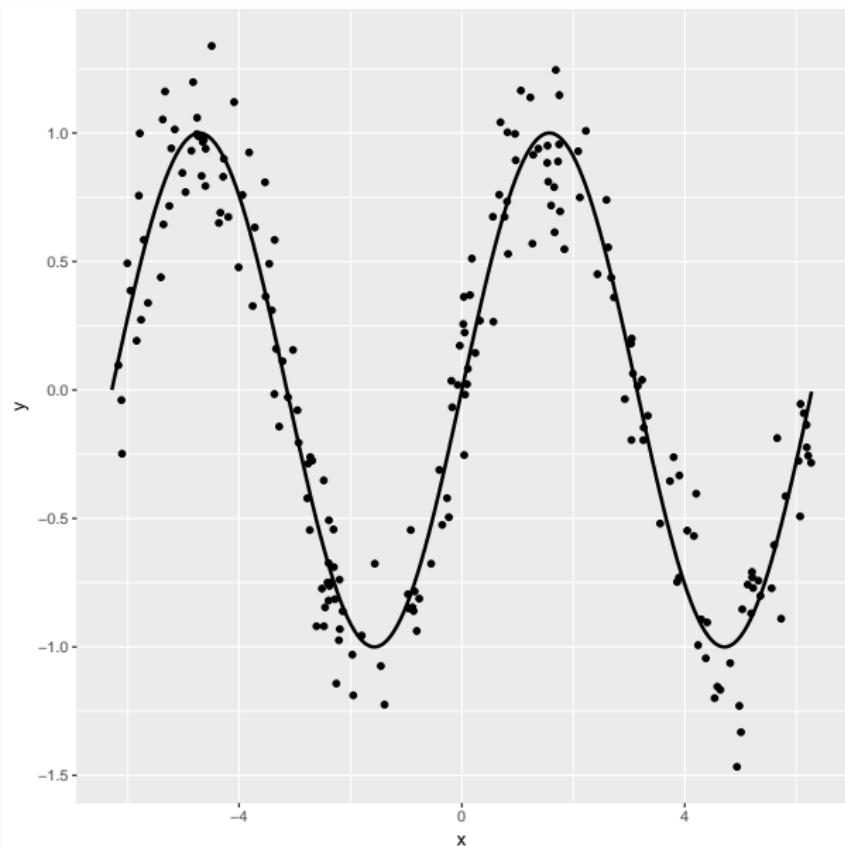
Assume that $k \rightarrow \infty$ and $k/n \rightarrow 0$. Then the k -NN classifier is **universally consistent**, i.e., $\mathbb{E}L(g_n) \rightarrow L^*$ for **all** distributions of (X, Y) .

- ▷ Proof's agenda: verify Stone's conditions (i)-(iii).
- ▷ Simplification: distance ties $\|X_i - X\| = \|X_j - X\|$ occur with zero probability.

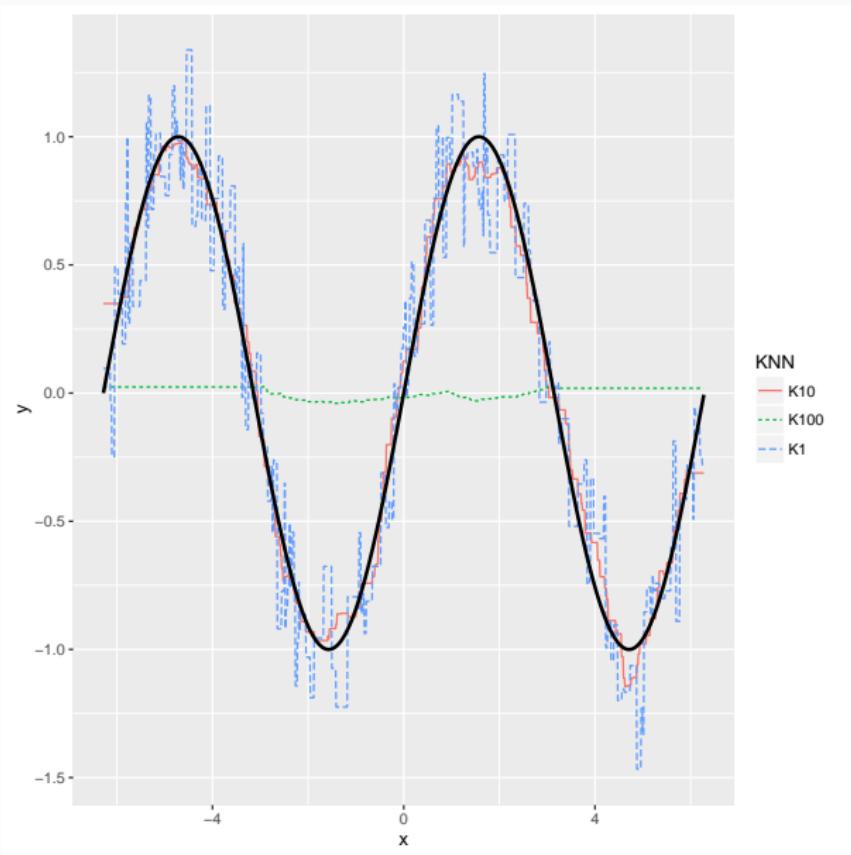


k is large but small with respect to n : bias/variance compromise.

Choice of k



Choice of k



Properties of the k th nearest neighbor

- The **support** of μ is defined by

$$\text{supp}(\mu) = \{x \in \mathbb{R}^d : \mu(B(x, \rho)) > 0 \text{ for all } \rho > 0\}.$$

- **Properties:**

- ▷ $\text{supp}(\mu)$ is a **closed** set;
- ▷ $\text{supp}(\mu)$ is the **smallest** closed subset of \mathbb{R}^d of μ -measure one;
- ▷ One has $\mathbb{P}(X \in \text{supp}(\mu)) = 1$.

Lemma

If $x \in \text{supp}(\mu)$ and $k/n \rightarrow 0$, then

$$\|X_{(k)}(x) - x\| \rightarrow 0 \quad \text{almost surely.}$$

Combinatorial lemmas

Lemma

Let ν be a probability measure on \mathbb{R}^d . Fix $x' \in \mathbb{R}^d$ and let, for $a \geq 0$,

$$B_a(x') = \left\{ x \in \mathbb{R}^d : \nu(B(x, \|x' - x\|)) \leq a \right\}.$$

Then

$$\nu(B_a(x')) \leq \gamma_d a,$$

where γ_d is a positive constant depending only upon d .

Corollary

If distance ties occur with zero probability, then

$$\sum_{i=1}^n \mathbb{1}_{[\text{ }X \text{ is among the } k\text{-NN of } X_i \text{ in } \{X_1, \dots, X_{i-1}, X, X_{i+1}, \dots, X_n\}]} \leq k\gamma_d,$$

with probability one.

Stone's lemma

Lemma (Stone's lemma)

Assume that distance ties occur with zero probability. Then, for every Borel measurable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathbb{E}|f(X)| < \infty$,

$$\sum_{i=1}^k \mathbb{E}|f(X_{(i)}(X))| \leq k\gamma_d \mathbb{E}|f(X)|,$$

where γ_d is a positive constant depending only upon d .

Consistency of the k -NN classifier

- To do: verify Stone's conditions with $W_{ni}(x) = 1/k$ if X_i is among the k nearest neighbors of x and $W_{ni}(x) = 0$ otherwise.
- Condition (iii) is clear since $k \rightarrow \infty$.
- Condition (ii): note that

$$\mathbb{E}\left(\sum_{i=1}^n W_{ni}(X) \mathbf{1}_{[\|X_i - X\|] > a}\right) = \mathbb{E}\left(\frac{1}{k} \sum_{i=1}^k \mathbf{1}_{[\|X_{(i)}(X) - X\|] > a}\right).$$

So,

$$\mathbb{P}(\|X_{(k)}(X) - X\| > a) \rightarrow 0 \Rightarrow \mathbb{E}\left(\sum_{i=1}^n W_{ni}(X) \mathbf{1}_{[\|X_i - X\|] > a}\right) \rightarrow 0.$$

But, for all $a > 0$,

$$\mathbb{P}(\|X_{(k)}(X) - X\| > a) = \int_{\mathbb{R}^d} \mathbb{P}(\|X_{(k)}(x) - x\| > a) \mu(dx).$$

Assuming that $k/n \rightarrow 0$, the conclusion follows by the Lebesgue dominated convergence theorem.

Consistency of the k -NN classifier, cont.

- Condition (i): take f such that $\mathbb{E}|f(X)| < \infty$; we have to show that for some constant C

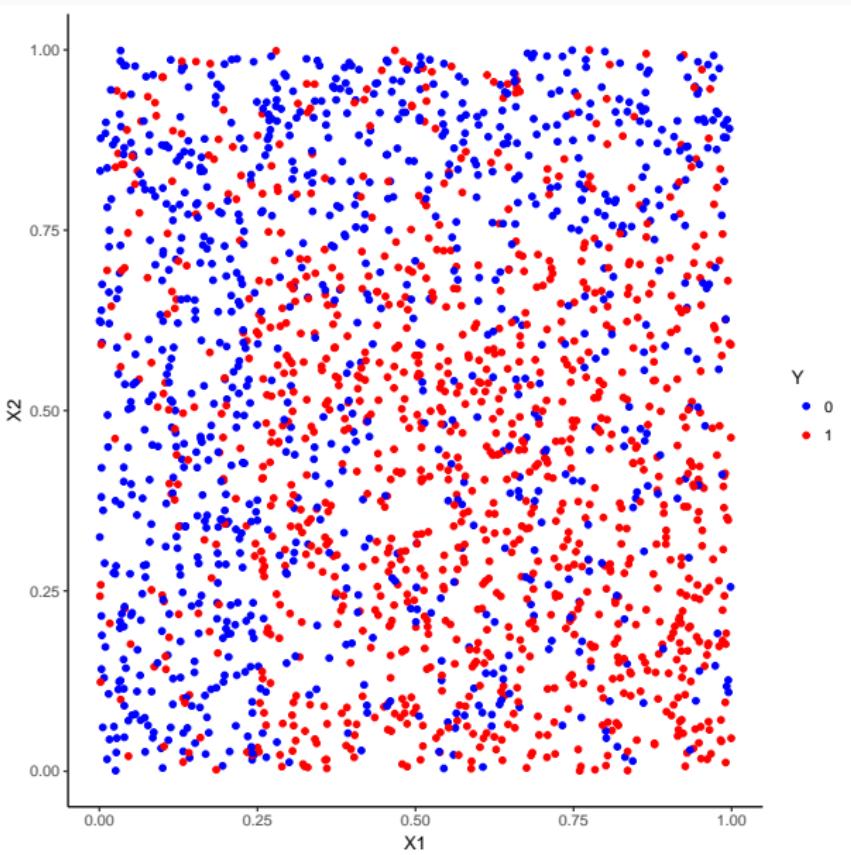
$$\mathbb{E}\left(\frac{1}{k} \sum_{i=1}^n |f(X_i)| \mathbb{1}_{[X_i \text{ is among the } k\text{-NN of } X]}\right) \leq C \mathbb{E}|f(X)|.$$

Since

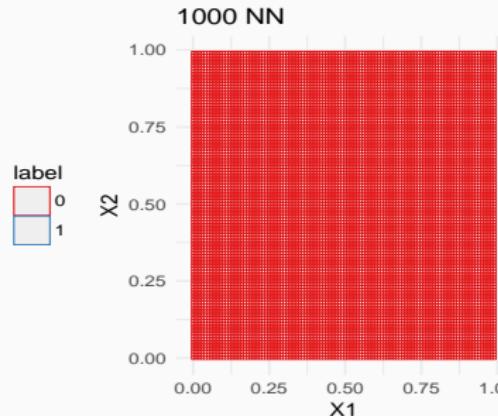
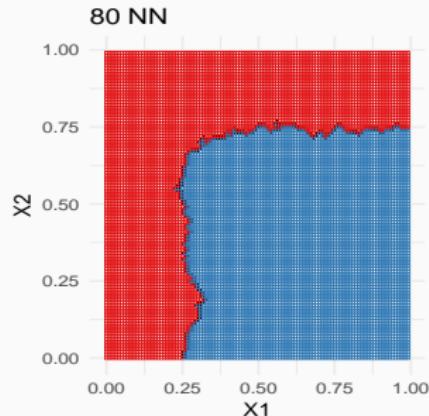
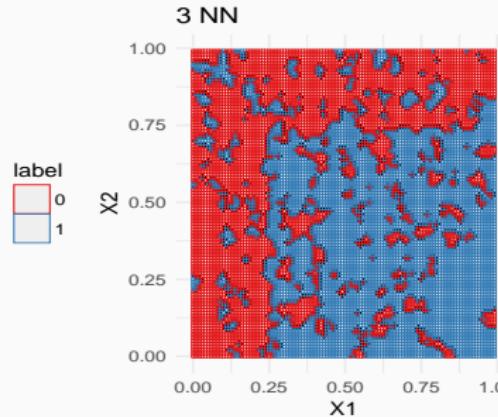
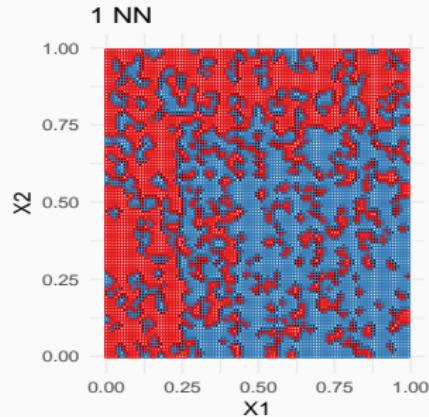
$$\mathbb{E}\left(\frac{1}{k} \sum_{i=1}^n |f(X_i)| \mathbb{1}_{[X_i \text{ is among the } k\text{-NN of } X]}\right) = \mathbb{E}\left(\frac{1}{k} \sum_{i=1}^k |f(X_{(i)}(X))|\right),$$

this is precisely the statement of Stone's lemma. ✓

Example



Bias/Variance compromise



Choice of k by data splitting

- Choosing k : minimizing the empirical error is **not** a good idea.
- Data splitting:
 - ▷ A **training** set $\mathcal{D}_m = \{(X_1, Y_1), \dots, (X_m, Y_m)\}$;
 - ▷ A **test** set $\mathcal{D}_\ell = \{(X_{m+1}, Y_{m+1}), \dots, (X_n, Y_n)\}$, with $m + \ell = n$.
- Candidates: $\mathcal{C}_m = \{g_k : 1 \leq k \leq m\} \rightarrow k\text{-NN classifiers using } \mathcal{D}_m$.
- Strategy: choose $g_n^* \in \mathcal{C}_m$ such that

$$g_n^* \in \arg \min_{g_k \in \mathcal{C}_m} \frac{1}{\ell} \sum_{i=m+1}^n \mathbb{1}_{[g_k(X_i) \neq Y_i]}.$$

Theorem

One has

$$\mathbb{E}(L(g_n^*) - \inf_{g_k \in \mathcal{C}_m} L(g_k)) \leq 2 \sqrt{\frac{\log(2em)}{2\ell}}.$$

- ▷ The classifier g_n^* is **universally consistent** provided

$$\lim_{n \rightarrow \infty} m = \infty \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\ell}{\log m} = \infty.$$

Partitioning classifiers and trees

Partitioning classifiers

- **Principle:** partition \mathbb{R}^d into disjoint cells A_1, A_2, \dots
- Classification by a **majority vote** in each cell.
- **Classifier:**

$$g_n(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^n \mathbb{1}_{[X_i \in A(x)]} \mathbb{1}_{[Y_i = 1]} > \sum_{i=1}^n \mathbb{1}_{[X_i \in A(x)]} \mathbb{1}_{[Y_i = 0]} \\ 0 & \text{otherwise,} \end{cases}$$

where $A(x)$ = cell containing x .

- **X-property:** the partitions depend **only** on X_1, \dots, X_n (**not** Y_1, \dots, Y_n).
- **Notation:** $\text{diam}(A) = \sup_{(x,y) \in A^2} \|x - y\|$ and $N(x) = \sum_{i=1}^n \mathbb{1}_{[X_i \in A(x)]}$.

Theorem

Let g_n be a partitioning classifier with the **X-property**. If

(i) $\text{diam}(A(X)) \rightarrow 0$ in probability,

and

(ii) $N(X) \rightarrow \infty$ in probability,

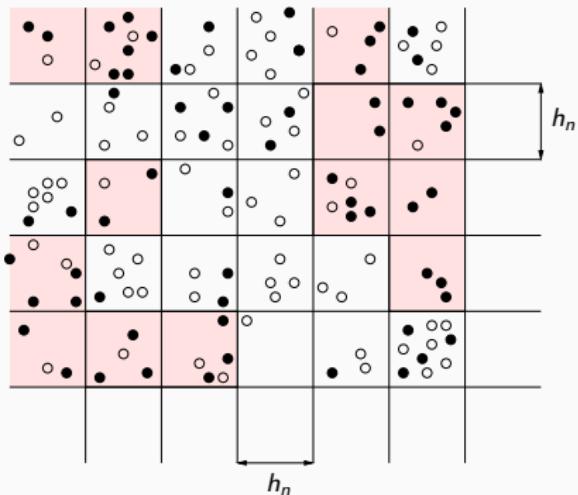
then $\mathbb{E}L(g_n) \rightarrow L^*$.

Example 1: cubic histogram classifier

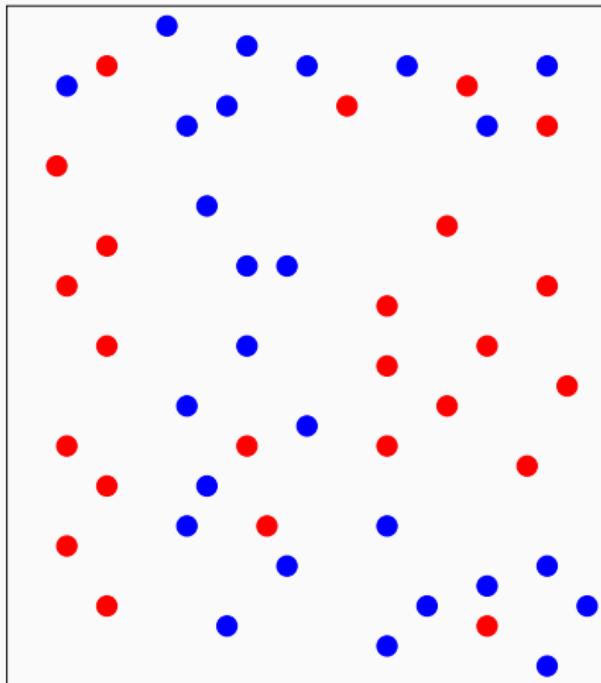
- **Definition:** A_{n1}, A_{n2}, \dots a partition of \mathbb{R}^d into **cubes** of size h .
- So, each cell = $\prod_{j=1}^d [k_j h, (k_j + 1)h]$, where the k_j are integers.

Theorem

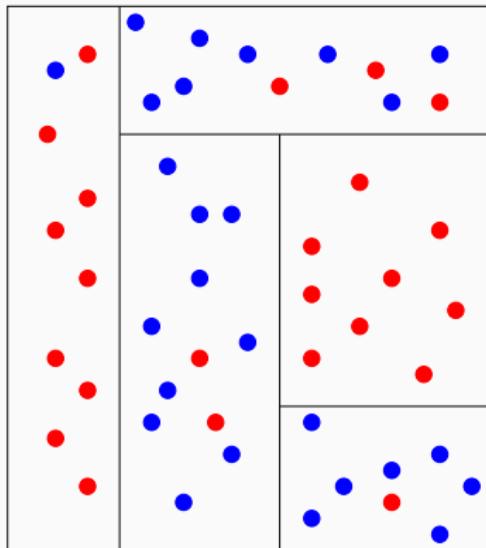
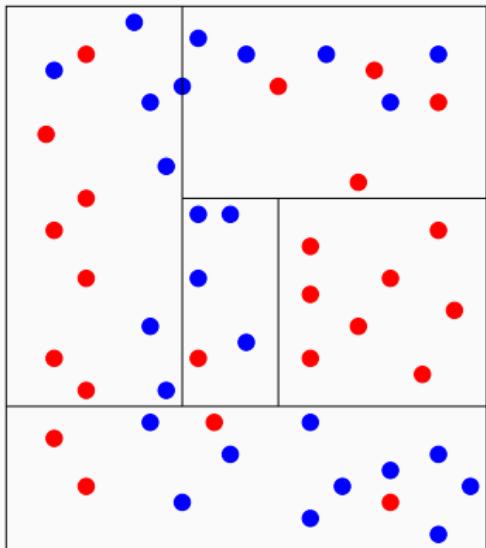
Assume that $h \rightarrow 0$ and $nh^d \rightarrow \infty$. Then the cubic histogram classifier is **universally consistent**, i.e., $\mathbb{E}L(g_n) \rightarrow L^*$ for **all** distributions of (X, Y) .



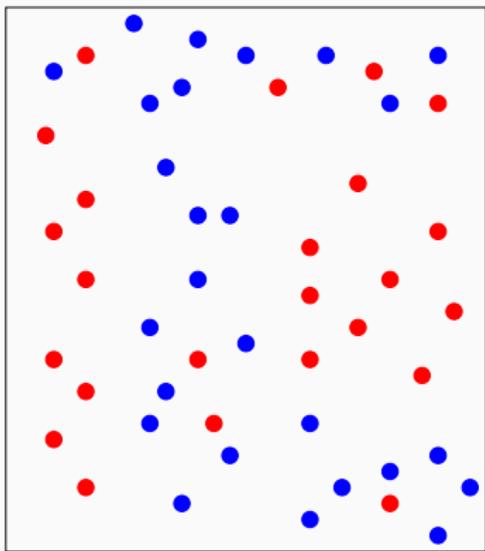
Example 2: tree classifiers



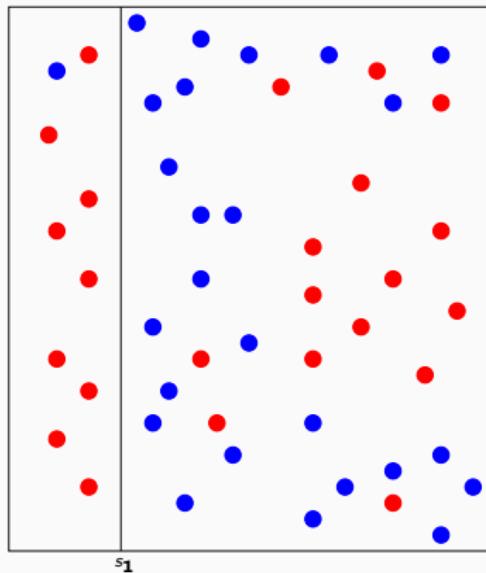
Example 2: tree classifiers



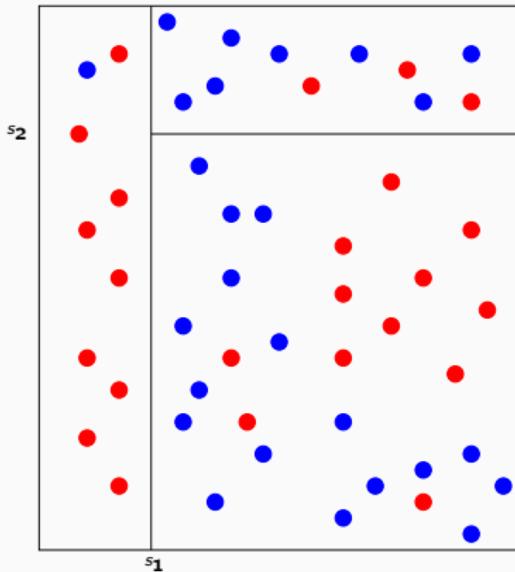
Example 2: tree classifiers



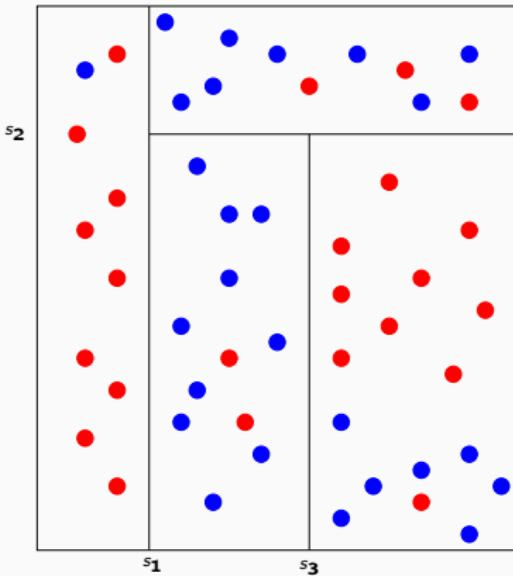
Example 2: tree classifiers



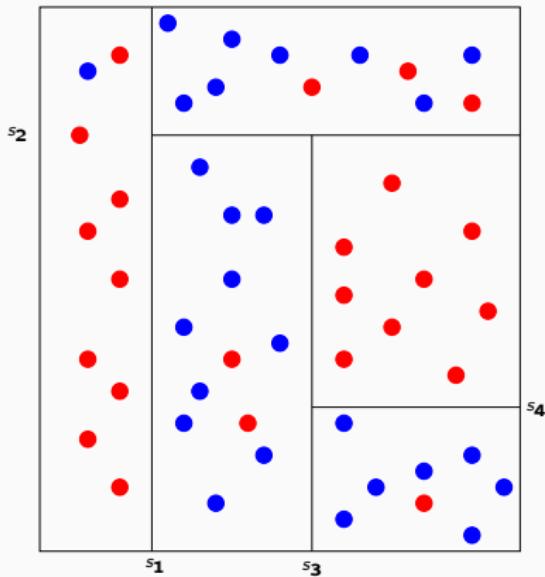
Example 2: tree classifiers



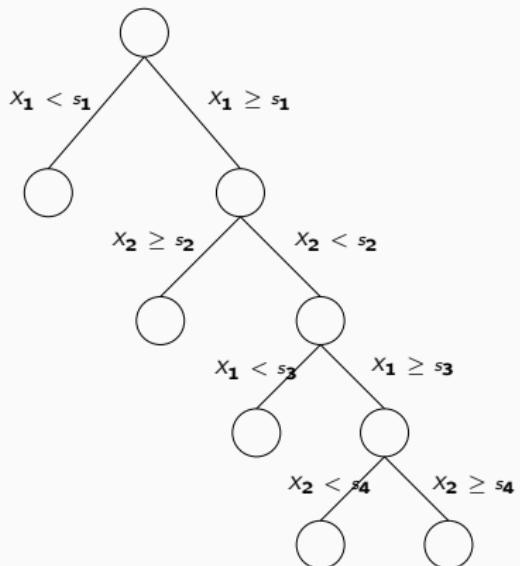
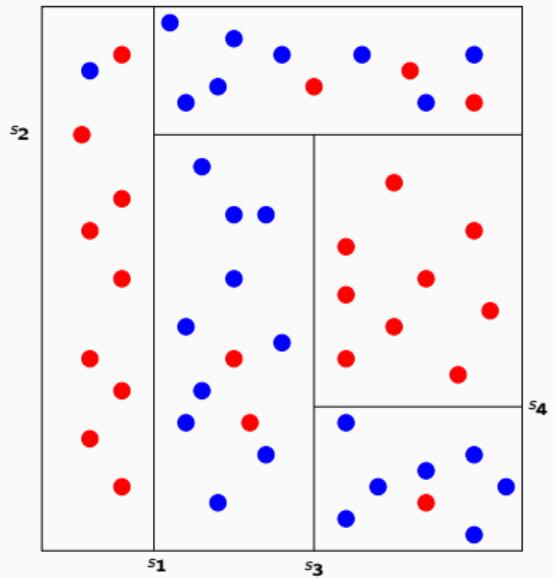
Example 2: tree classifiers



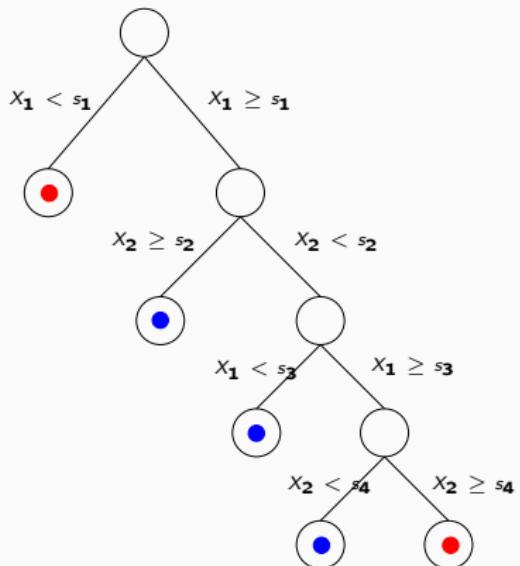
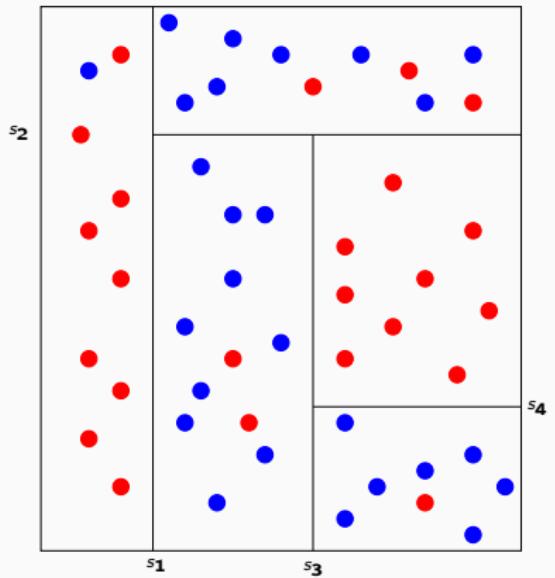
Example 2: tree classifiers



Example 2: tree classifiers



Example 2: tree classifiers



Binary trees

- **Definition:** recursive binary partitioning of \mathbb{R}^d , represented by a tree.
- A node has exactly either zero or two **children**.
- A node with zero children is called a **leaf**.
- If $u \leftrightarrow A$ and $u_L, u_R \leftrightarrow A_L, A_R$, then $A = A_L \cup A_R$ and $A_L \cap A_R = \emptyset$.
- The **root** $\leftrightarrow \mathbb{R}^d$ and the **leaves** \leftrightarrow a partition of \mathbb{R}^d .
- We pass from A to A_L and A_R by **answering a question** on x :

"Is $x^{(j)} \geq \alpha?$ ", for some coordinate j and some α .

- \mathbb{R}^d is partitioned into **hyperrectangles**.
- **Principle:** x is passed into the root and then **iteratively transmitted** to the child nodes. This is repeated until a leaf is reached.

Tree classifiers

- **Classifier:** for $x \in A$,

$$g_n(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^n \mathbb{1}_{[X_i \in A]} \mathbb{1}_{[Y_i=1]} > \sum_{i=1}^n \mathbb{1}_{[X_i \in A]} \mathbb{1}_{[Y_i=0]} \\ 0 & \text{otherwise.} \end{cases}$$

- **Two questions:**
 - ▷ Do we cut?
 - ▷ In the affirmative, where do we cut?
- Many **tree species** (median, centered, CART, C_{4.5}, etc.).
- A data-modeling technique that proved itself to be **effective**.
- Trees are **fast**, **simple**, and able to explain **complex** data.
- Many **modern** predictive algorithms rely on tree principles.
- **Impact** in geometry, statistics, machine learning, and AI.

Median tree classifier

- Principle:
 - ▷ At each node: find the median according to one coordinate;
 - ▷ n points → two children with sizes $\lfloor(n-1)/2\rfloor$ and $\lceil(n-1)/2\rceil$;
 - ▷ The median itself stays behind and is not sent down to the subtrees;
 - ▷ Repeat this for k levels of nodes, in a rotational manner.
- 2^k leaf regions, each having at least $n/2^k - 2$ and at most $n/2^k$ points.

Theorem

Assume that X has a density. If $k \rightarrow \infty$ and $\frac{n}{k2^k} \rightarrow \infty$, then the median tree classifier is consistent, i.e., $\mathbb{E}L(g_n) \rightarrow L^*$. (Note: the conditions on k are fulfilled if $k \leq \log_2 n - 2 \log_2 \log_2 n$, $k \rightarrow \infty$.)

DON'T FORGET!

A median tree is invariant under monotone transformations of the axes.

Label-dependent cuts

- **Strategy:** at each step, choose (j, s) with

$$R_1(j, s) = \{x : x^{(j)} < s\} \quad \text{and} \quad R_2(j, s) = \{x : x^{(j)} \geq s\}.$$

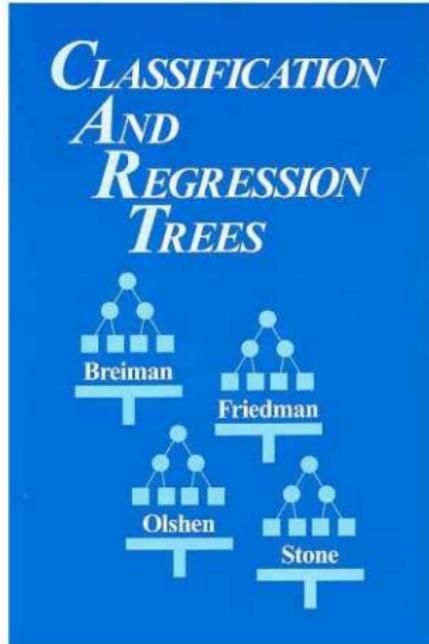
- The pair (j, s) is selected by minimizing the **impurity** of the children.
- **Impurity:** low when the labels are homogeneous, high otherwise.
- **Example:** choose (j, s) that minimizes

$$\sum_{X_i \in R_1(j, s)} \mathbb{1}_{[Y_i \neq \hat{Y}_1]} + \sum_{X_i \in R_2(j, s)} \mathbb{1}_{[Y_i \neq \hat{Y}_2]},$$

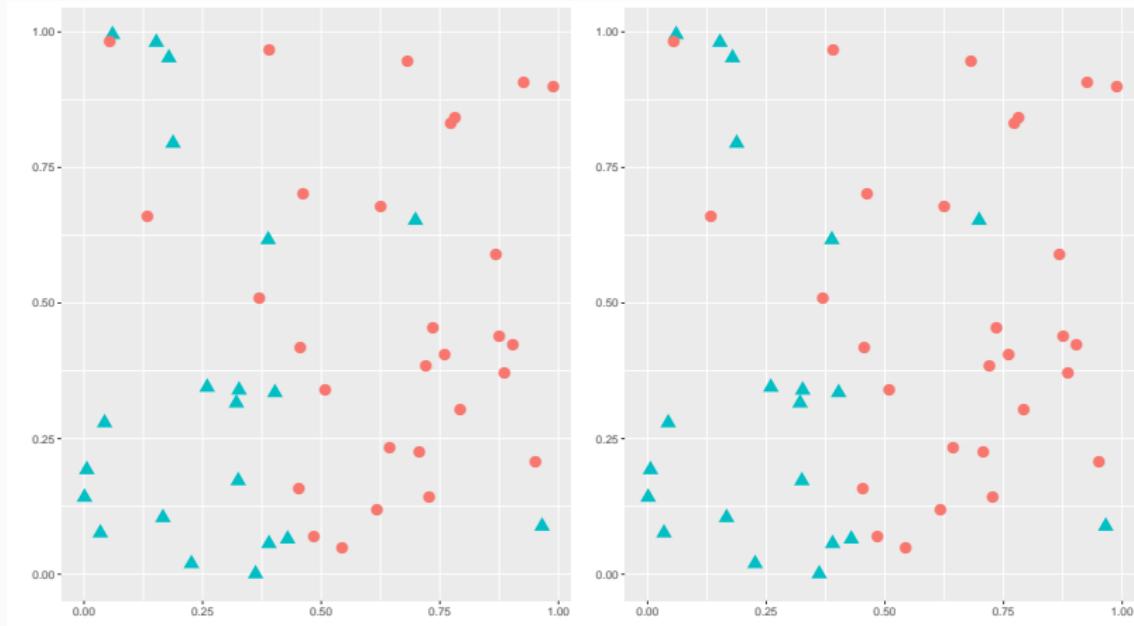
where \hat{Y}_k is the majority label in the node $R_k(j, s)$.

- There are **other criteria:** Gini impurity, entropy, etc.
- **Example:** CART algorithm.

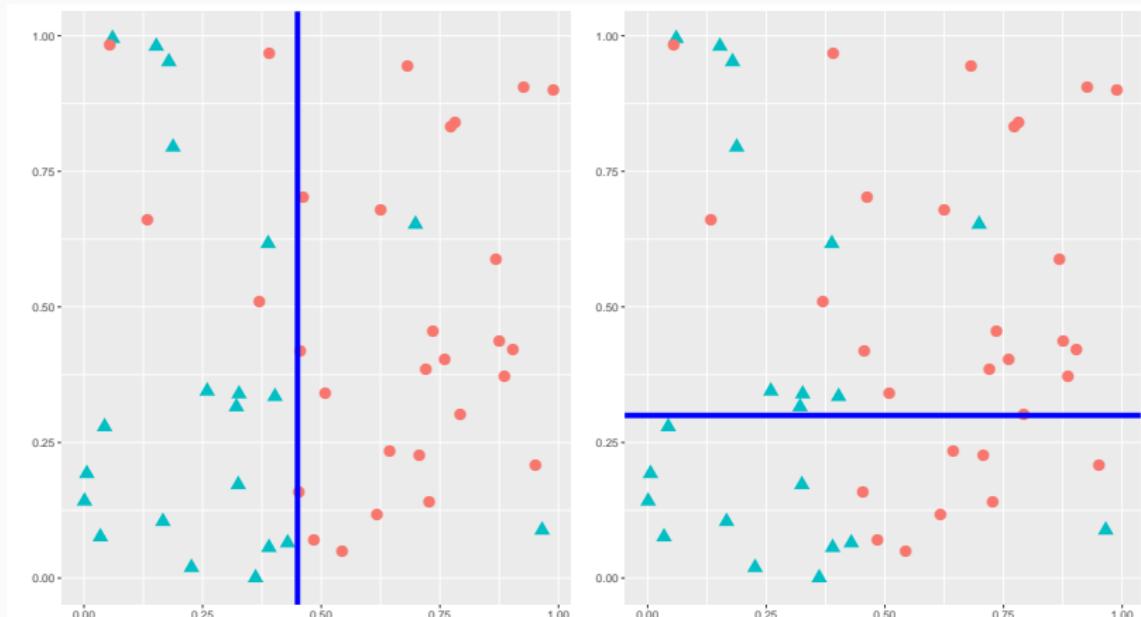
Leo Breiman



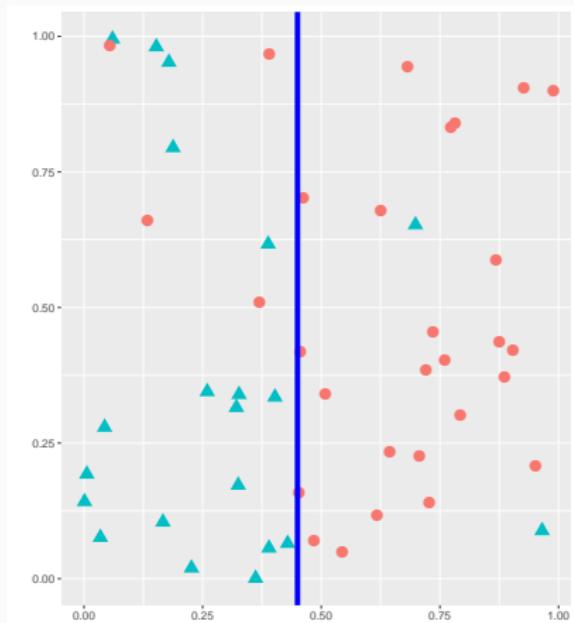
Example



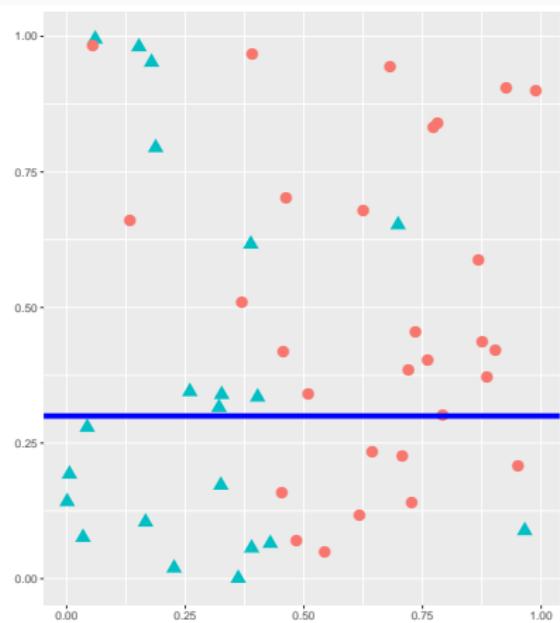
Example



Example

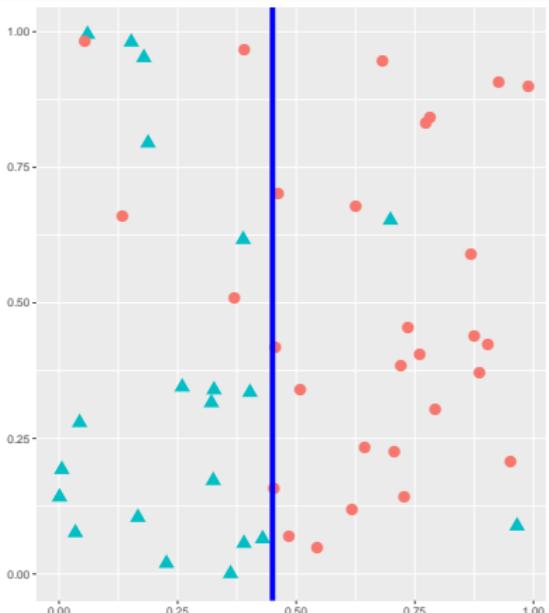


4+2 errors

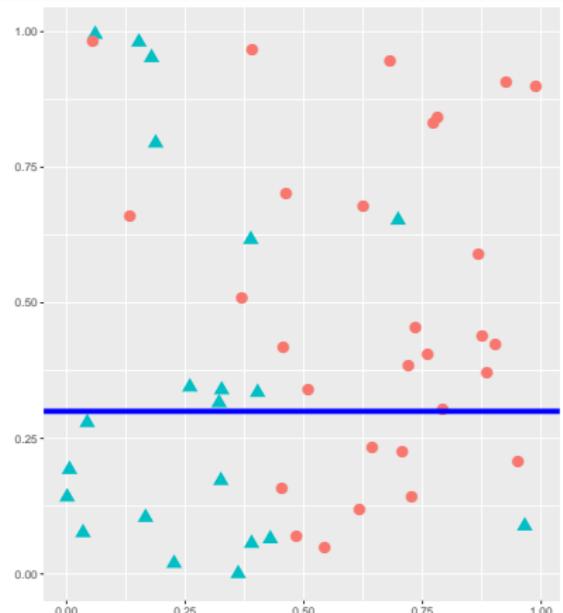


8+10 errors

Example



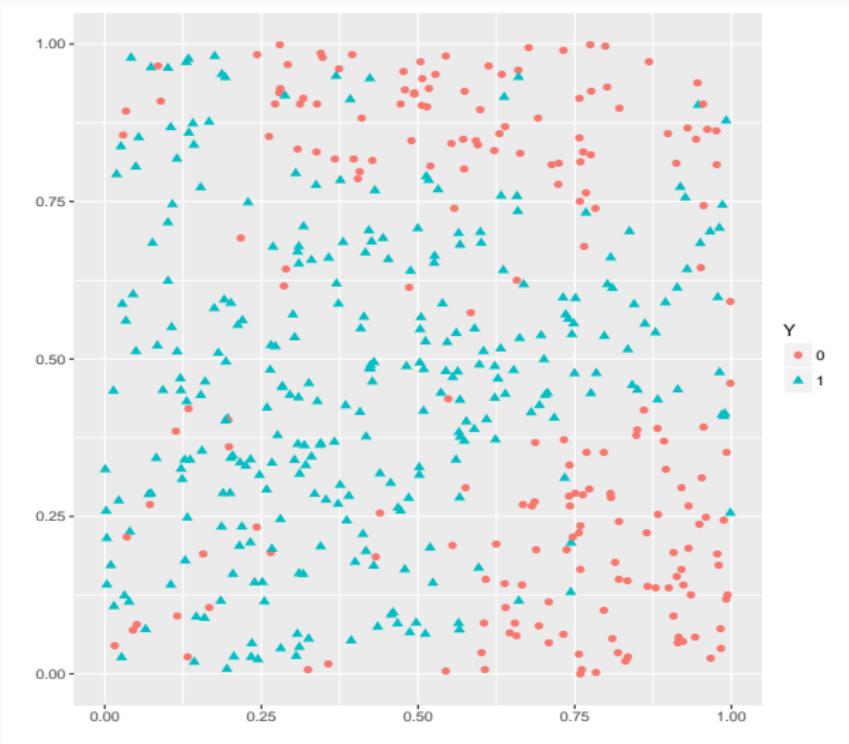
4+2 errors



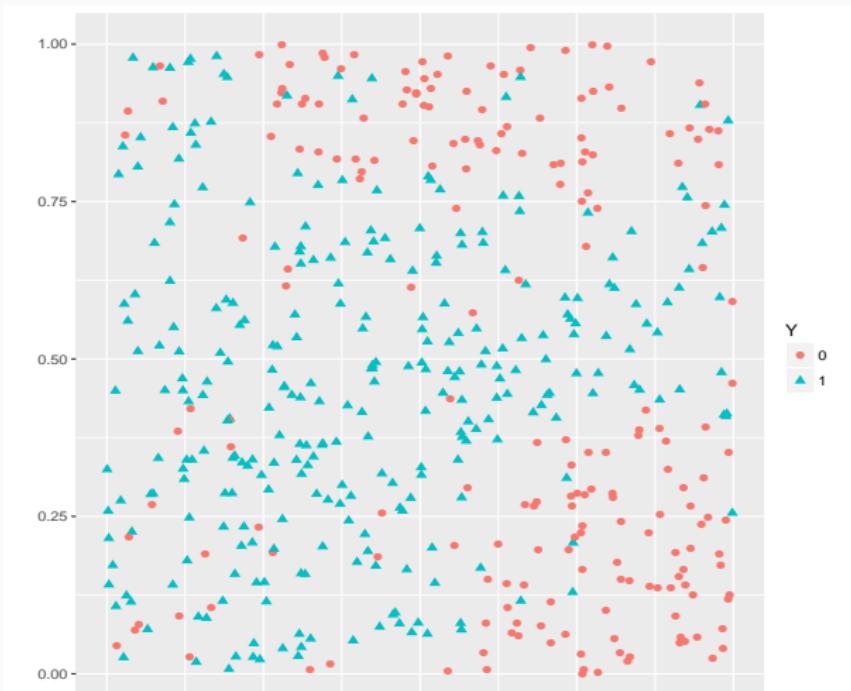
8+10 errors

Best cut: **left**.

When do we stop?

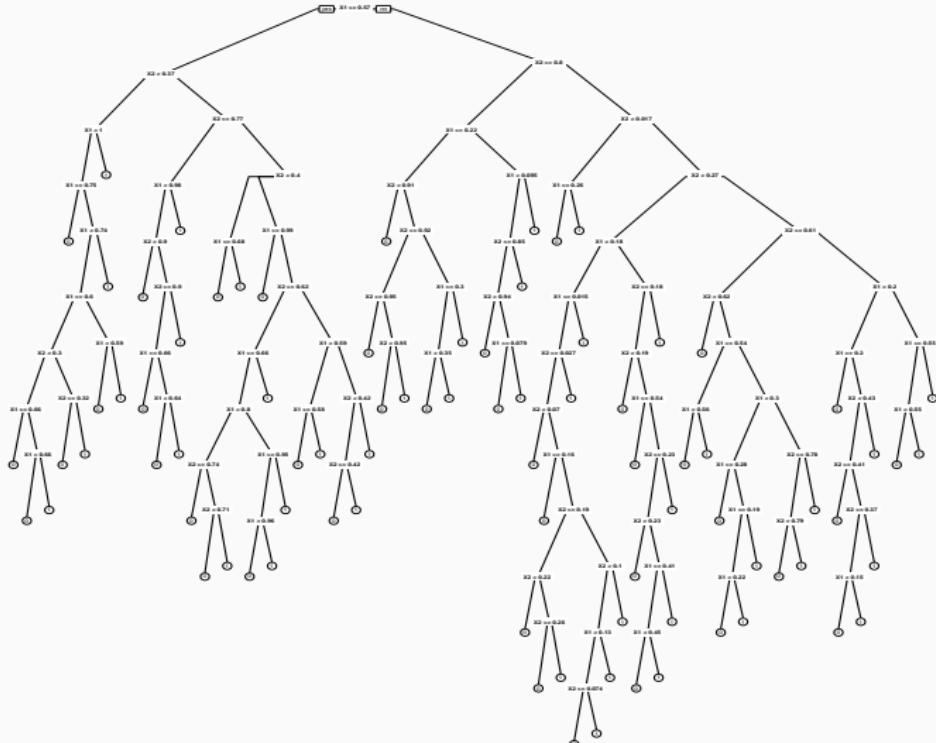


When do we stop?

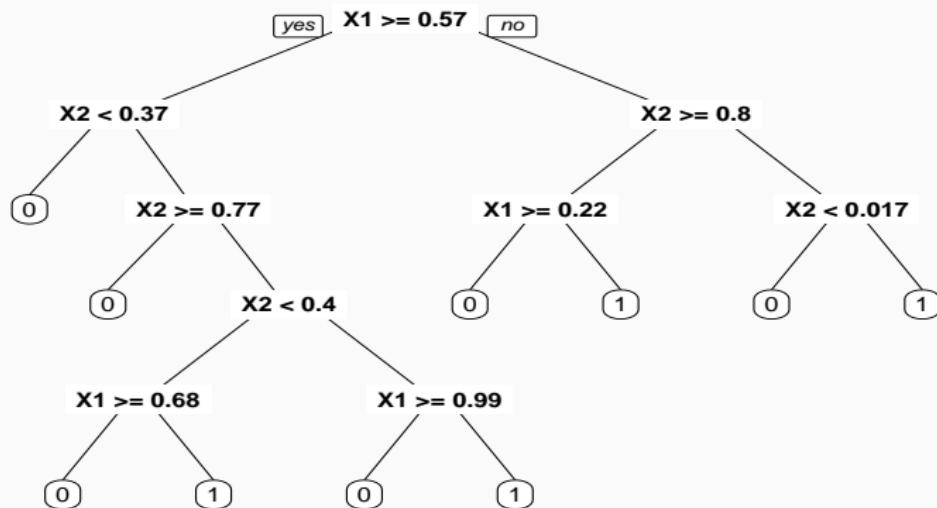


Guess: about 5 cells.

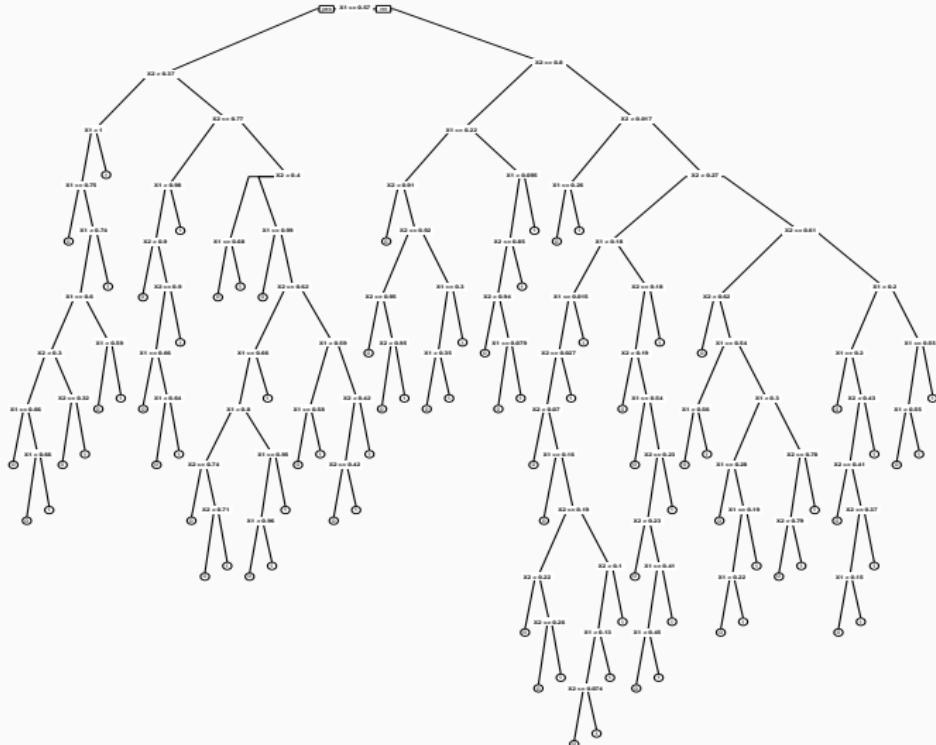
Maximal tree?



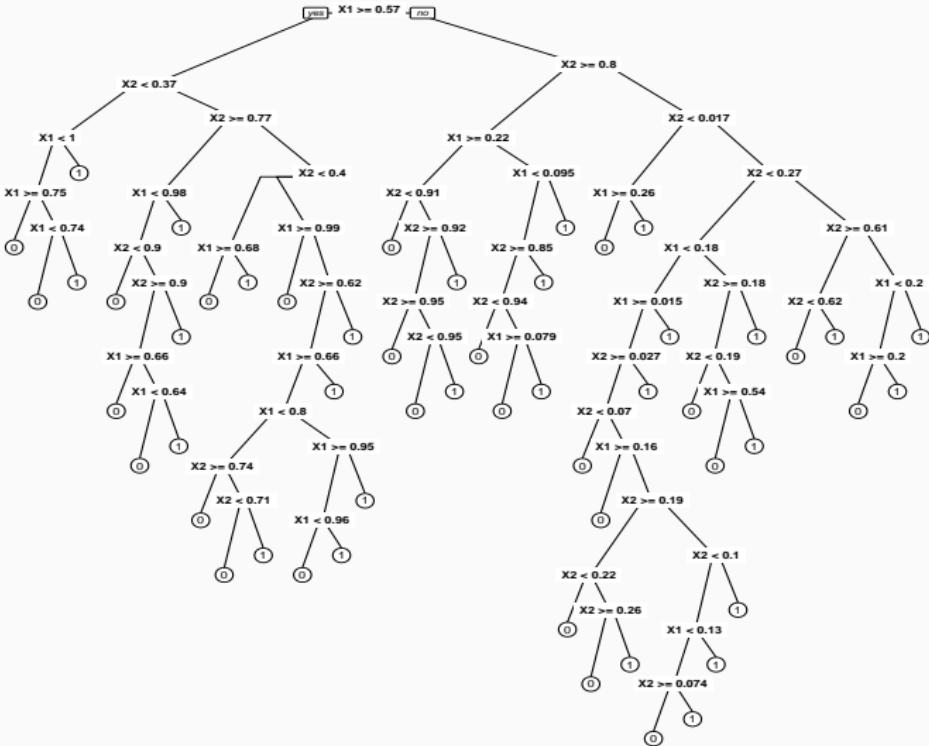
Smaller tree?



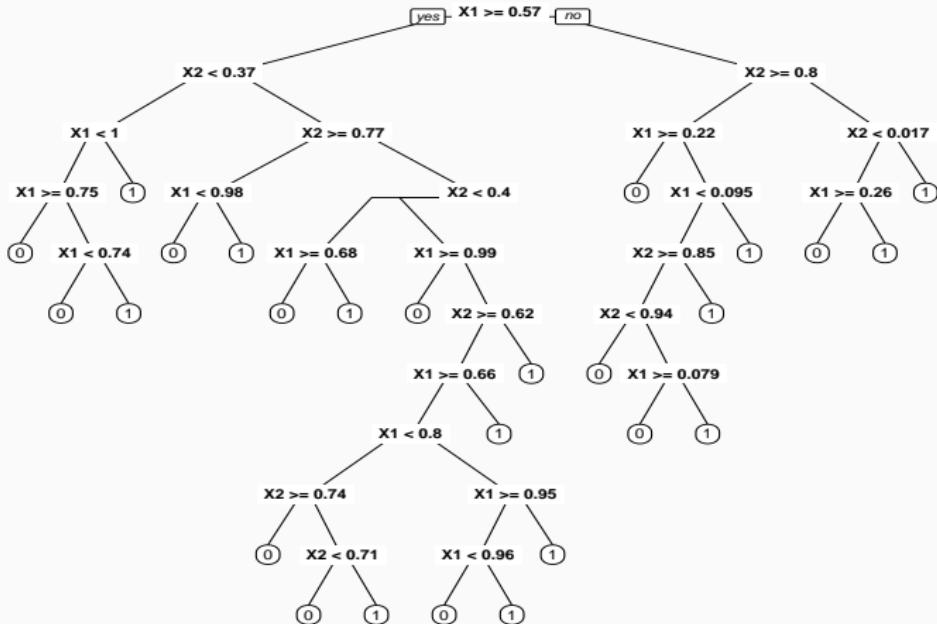
Nested trees



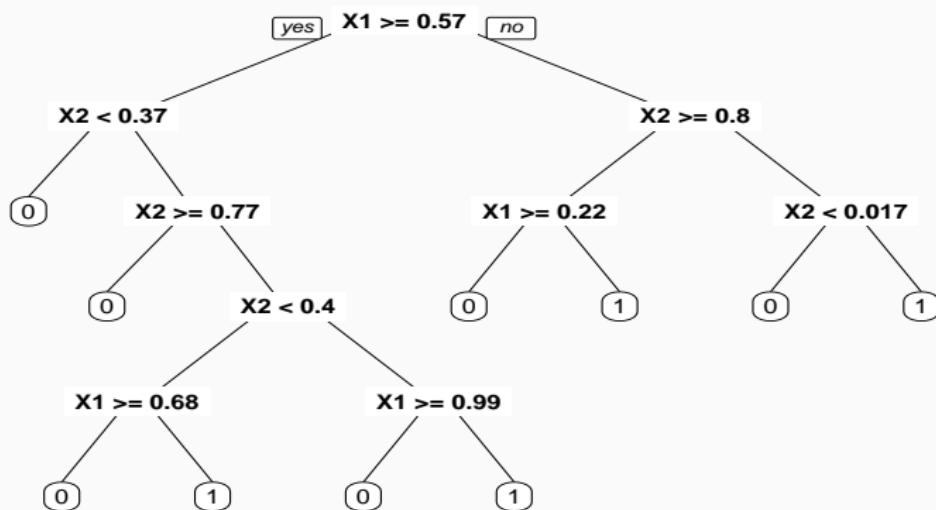
Nested trees



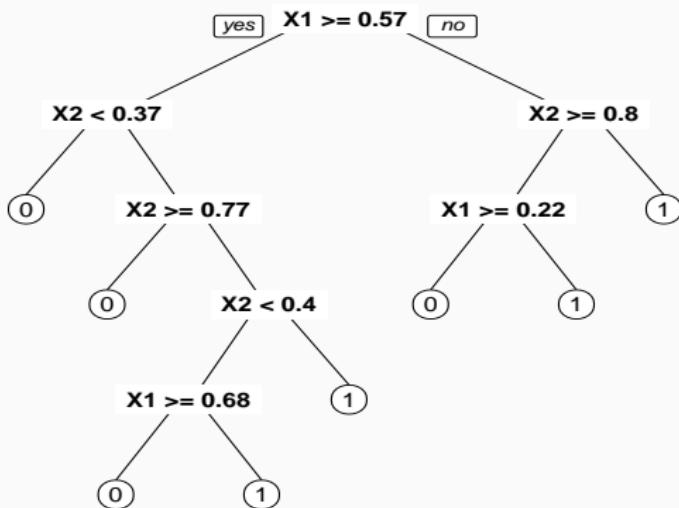
Nested trees



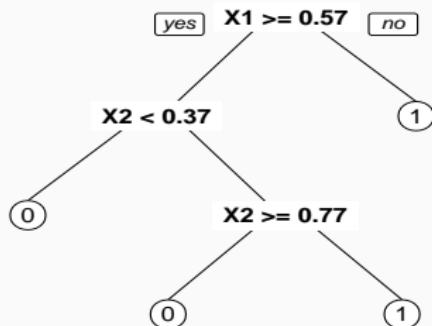
Nested trees



Nested trees



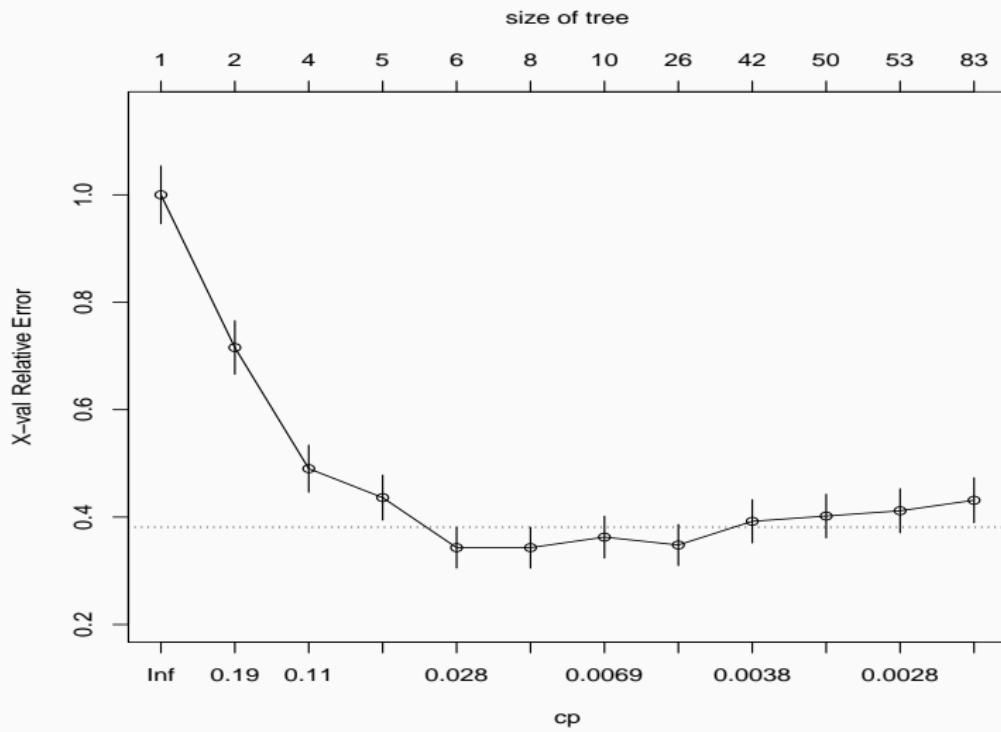
Nested trees



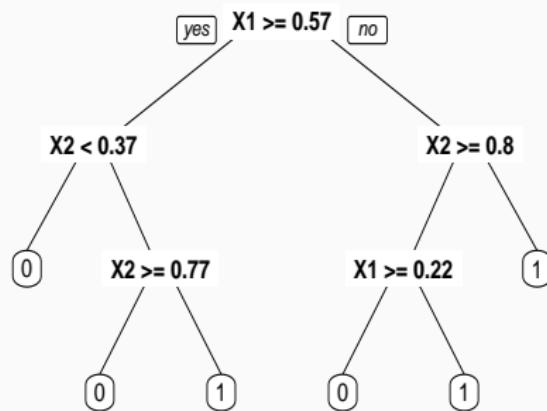
Nested trees

(1)

... and pruning

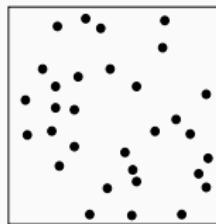
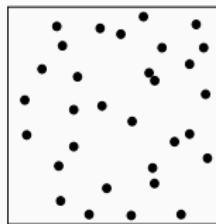
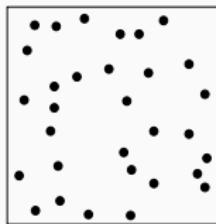


Final tree

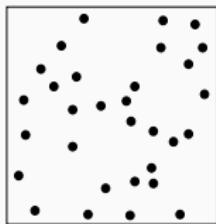


... Breiman's random forests

1. Data **resampling** (bootstrap or subsampling).
2. At each cell, randomly preselect a **small number** of features.
3. Construct the trees by a **CART-type** procedure.
4. **Stopping rule**: leave one point in the leaves.
5. **Final decision**: majority vote.



...



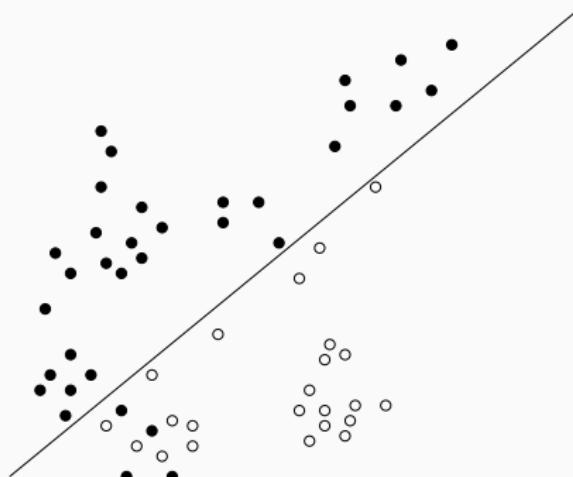
Neural networks

Rosenblatt's perceptron (1958)

- Linear discriminant or **perceptron**:

$$g(x) = \begin{cases} 1 & \text{if } \psi(x) > 1/2 \\ 0 & \text{otherwise,} \end{cases}$$

where $\psi(x) = \sum_{j=1}^d c_j x^{(j)} + c_0 = c^\top x + c_0$.



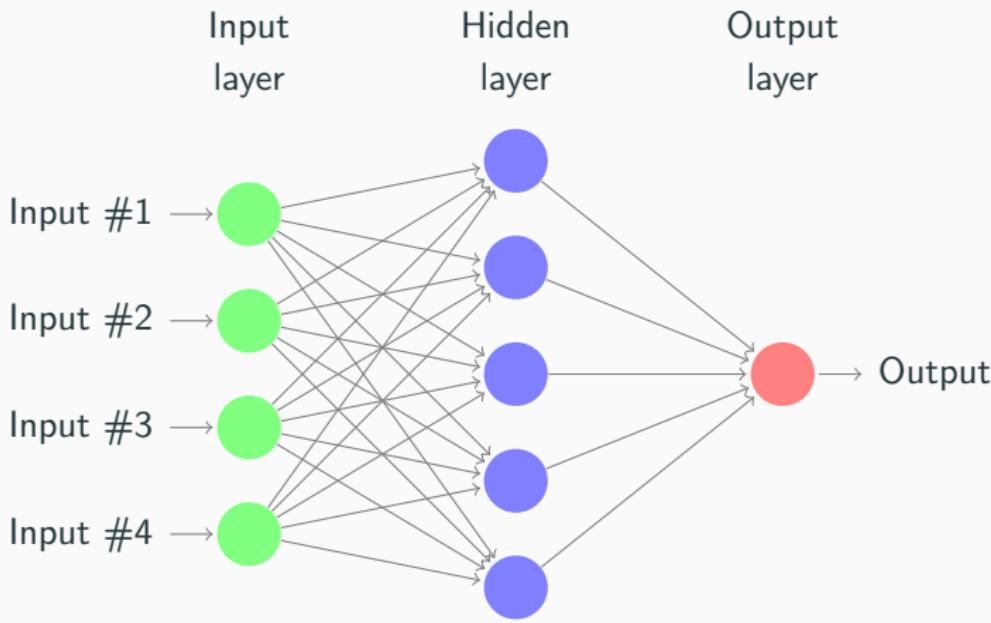
- This is a **neural network without hidden layers**.

Neural networks

Neural network with **one hidden layer**:

$$\psi(x) = \sum_{i=1}^k c_i \sigma(\psi_i(x)) + c_0,$$

where $\psi_i(x) = \sum_{j=1}^d a_{ij}x^{(j)} + a_{i0}$ and σ is the **activation function**.



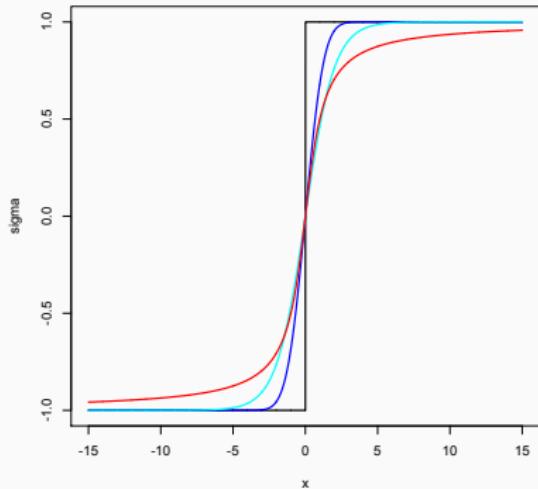
Activation function

- **Definition:** nondecreasing function with

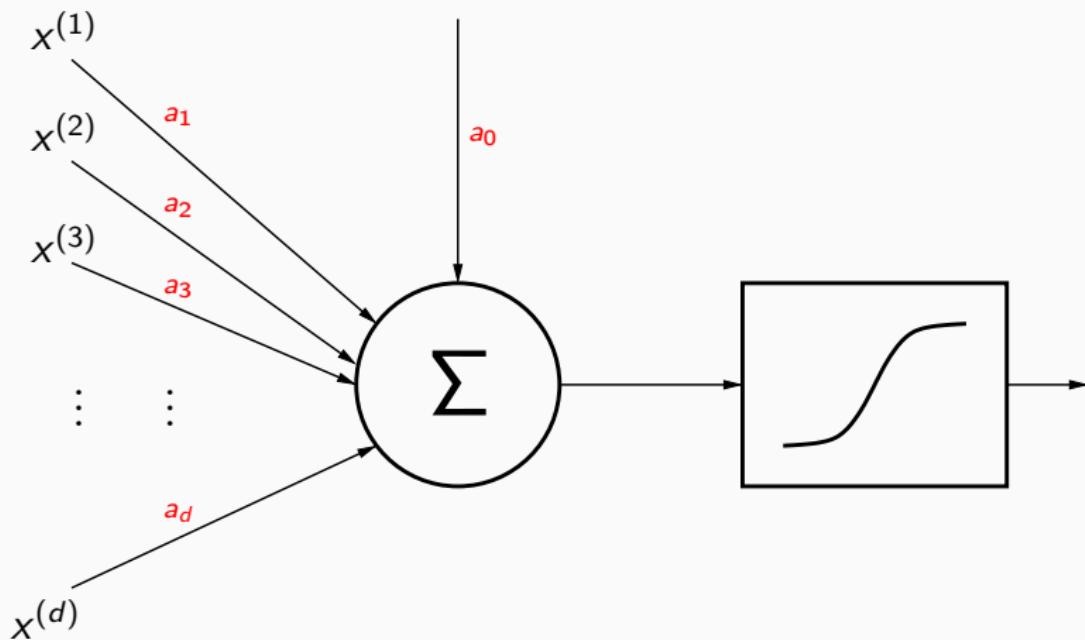
$$\lim_{x \rightarrow -\infty} \sigma(x) = -1 \quad \text{and} \quad \lim_{x \rightarrow \infty} \sigma(x) = 1.$$

- **Examples:**

- ▷ **Threshold:** $\sigma(x) = 2\mathbb{1}_{[x \geq 0]} - 1$
- ▷ **Gaussian:** $\sigma(x) = 2 \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2} dt - 1$
- ▷ **Logistic:** $\sigma(x) = \frac{1-e^{-x}}{1+e^{-x}}$
- ▷ **Arctan:** $\sigma(x) = \frac{2}{\pi} \arctan(x)$



A neuron



More layers, more neurons

- A one-hidden-layer neural network has k hidden neurons.
- Output of the i -th neuron: $u_i = \sigma(\psi_i(x))$. So,

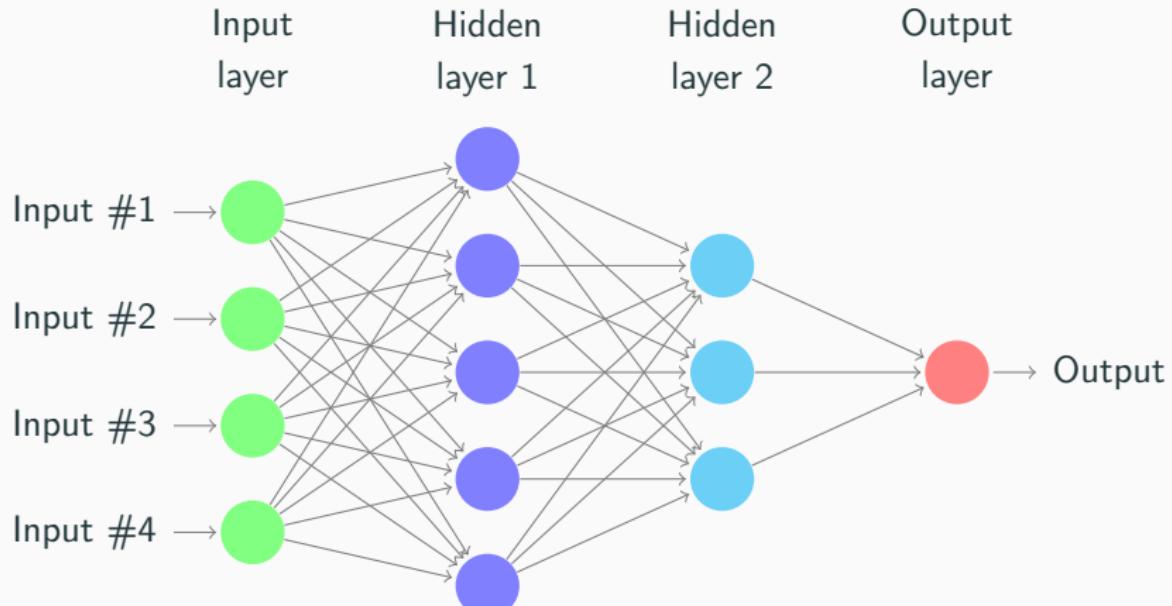
$$\psi(x) = \sum_{i=1}^k c_i u_i + c_0.$$

- Two-hidden-layer neural network: $\psi(x) = \sum_{i=1}^\ell c_i z_i + c_0$, where

$$z_i = \sigma\left(\sum_{j=1}^k b_{ij} u_j + b_{i0}\right) \quad \text{and} \quad u_j = \sigma\left(\sum_{i=1}^d a_{ji} x^{(i)} + a_{j0}\right).$$

- First hidden layer: k neurons. Second hidden layer: ℓ neurons.
- And so on... Deep learning: more than 100 layers!

Two-hidden-layer neural network



Theoretical roots

- A theorem by **Kolmogorov** (1957) and **Lorentz** (1976): for every continuous function f on $[0, 1]^d$,

$$f(x) = \sum_{i=1}^{2d+1} F_i \left(\sum_{j=1}^d G_{ij}(x^{(j)}) \right),$$

where the G_{ij} and the F_i are continuous functions depending on f .

- **Example:** with $d = 2$ and $f(x) = x^{(1)}x^{(2)}$,

$$f(x) = \frac{1}{4} \left((x^{(1)} + x^{(2)})^2 - (x^{(1)} - x^{(2)})^2 \right).$$

- **Next step:** approximate the G_{ij} and F_i by functions $\sigma(a^\top x + a_0)$.
- One-hidden-layer: **universally consistent** classifiers provided $k_n \rightarrow \infty$.
- **Unrealistic**, since the hardware is fixed beforehand.

Approximation by neural networks

Theorem

If \mathcal{C}_k is the class of all neural network classifiers with the **threshold activation function** and k neurons in **two hidden layers**, then

$$\lim_{k \rightarrow \infty} \inf_{g \in \mathcal{C}_k} L(g) - L^* = 0$$

for **all** distributions of (X, Y) .

Theorem

If \mathcal{C}_k is the class of all neural network classifiers with **one hidden layer** of k nodes and an **arbitrary activation function**, then

$$\lim_{k \rightarrow \infty} \inf_{g \in \mathcal{C}_k} L(g) - L^* = 0$$

for **all** distributions of (X, Y) .

Approximation results

Lemma

Let $(\mathcal{F}_k)_k$ be a sequence of classes of functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and let $(\mathcal{C}_k)_k$ be the *companion sequence* of classes of classifiers. Assume that for every $a, b \in \mathbb{R}^d$ and every *continuous* function h on $[a, b]^d$,

$$\lim_{k \rightarrow \infty} \inf_{f \in \mathcal{F}_k} \sup_{x \in [a, b]^d} |h(x) - f(x)| = 0.$$

Then, for *any* distribution of (X, Y) ,

$$\lim_{k \rightarrow \infty} \inf_{g \in \mathcal{C}_k} L(g) - L^* = 0.$$

Approximation results, cont.

Theorem

For every *continuous* function $h : [a, b]^d \rightarrow \mathbb{R}$ and for every $\varepsilon > 0$, there exists a neural network with *one hidden layer*, of the form

$$\psi(x) = \sum_{i=1}^k c_i \sigma(\psi_i(x)) + c_0,$$

such that

$$\sup_{x \in [a, b]^d} |h(x) - \psi(x)| \leq \varepsilon.$$

Consistency results

- Assumptions:
 - ▷ \mathcal{C}_k = one-hidden-layer neural network classifiers with k nodes;
 - ▷ σ = threshold activation function.
- Shatter coefficient:
$$S_{\mathcal{A}}(n) \leq (ne)^{kd+2k+1}.$$
- VC dimension:
$$2\lfloor k/2 \rfloor d \leq V_{\mathcal{C}_k} \leq (2kd + 4k + 2) \log_2(e(kd + 2k + 1)).$$
- Consequence: if $k \rightarrow \infty$ and $k \log n/n \rightarrow 0$, then $\mathbb{E}L(g_n^*) \rightarrow L^*$ for all distributions of (X, Y) .
- Optimization: gradient descent (backpropagation) \rightarrow smooth σ needed.
- But careful: there exists an activation σ that is monotone increasing, continuous, concave on $(0, \infty)$ and convex on $(-\infty, 0)$, such that $V_{\mathcal{C}_k} = \infty$ for each $k \geq 8\dots$

Quantization and clustering

References

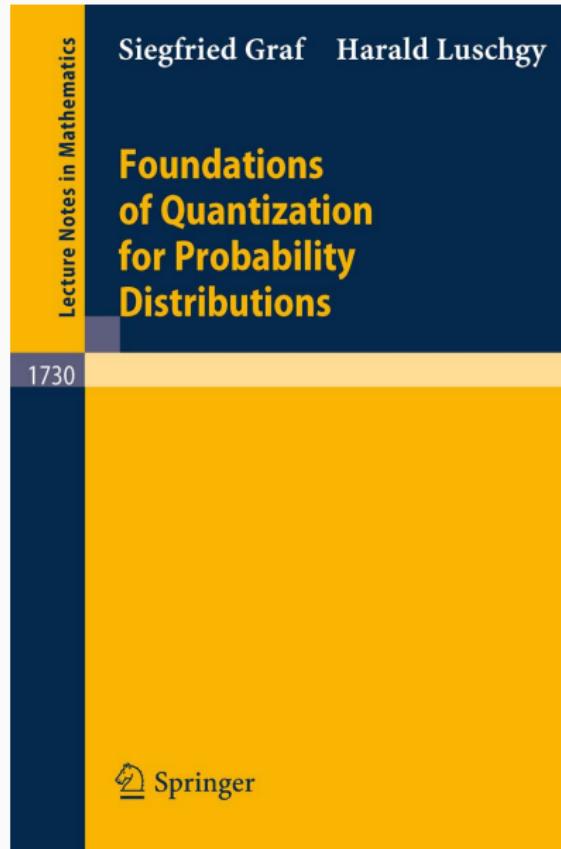
Learning-Theoretic Methods in Vector Quantization

Lecture Notes for the
Advanced School on the Principles of Nonparametric Learning
Udine, Italy, July 9-13, 2001.

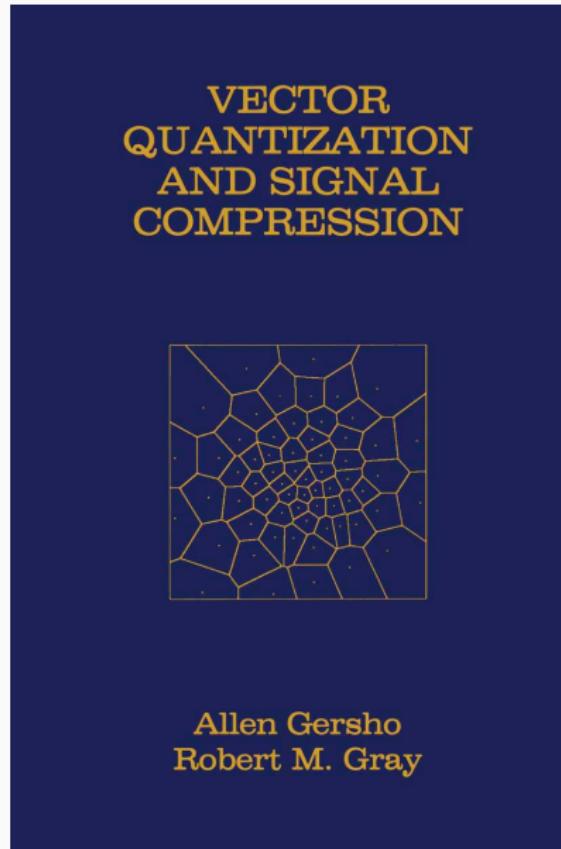
To appear in:
Principles of Nonparametric Learning
L. Györfi, editor, CISM Lecture Notes, Wien, New York: Springer 2001.

Tamás Linder
Department of Mathematics & Statistics
and
Department of Electrical & Computer Engineering
Queen's University
Kingston, Ontario, Canada, K7L 3N6
email: linder@mast.queensu.ca

References



References



Quantization

- **Quantization**: probabilistic principle to compress information.
- **Context**: a random variable X taking values in $(\mathbb{R}^d, \|\cdot\|)$.
- $\|\cdot\| =$ **Euclidean norm**.
- **Assumption**: $\mathbb{E}\|X\|^2 < \infty \Leftrightarrow \int_{\mathbb{R}^d} \|x\|^2 \mu(dx) < \infty$.

Definition (Quantizer)

Let $k \geq 1$ be an integer. A **quantizer q of order k** is a Borel measurable function $q : \mathbb{R}^d \rightarrow \mathcal{C} \subseteq \mathbb{R}^d$, with $|\mathcal{C}| \leq k$.

- A quantizer q of order k is characterized by:
 - ▷ A **codebook** $\mathcal{C} = \{c_1, \dots, c_k\}$;
 - ▷ A **partition** $\mathcal{P} = \{A_1, \dots, A_k\}$ of \mathbb{R}^d , with $q(x) = c_j \Leftrightarrow x \in A_j$.
- **Notation**: $q = (\mathcal{C}, \mathcal{P})$.

Quality measure of compression

Definition (Distortion)

The **distortion** (for X or μ) of a quantizer $q = (\mathcal{C}, \mathcal{P})$ of order k is

$$D(\mu, q) = \mathbb{E} \|X - q(X)\|^2 = \int_{\mathbb{R}^d} \|x - q(x)\|^2 \mu(dx).$$

The **minimal distortion** at the order k is $D_k^*(\mu) = \inf_q D(\mu, q)$, where the infimum is taken over all quantizers of order k .

- ▷ The **smaller** the distortion, the **better** the compression.
- ▷ The compression quality **improves** with k .

Lemma

One has $D_k^*(\mu) \downarrow 0$ as $k \rightarrow \infty$.

Nearest neighbor (NN) quantizers

- **Context:** quantizers of order k .
- **Voronoi partition:** for $\mathcal{C} = \{c_1, \dots, c_k\}$, the Voronoi partition is

$$A_1 = \{x \in \mathbb{R}^d : \|x - c_1\| \leq \|x - c_\ell\|, \forall \ell = 1, \dots, k\}, \text{ and}$$

$$A_j = \{x \in \mathbb{R}^d : \|x - c_j\| \leq \|x - c_\ell\|, \forall \ell = 1, \dots, k\} \setminus \bigcup_{t=1}^{j-1} A_t,$$

for $2 \leq j \leq k$.

Definition (NN quantizer)

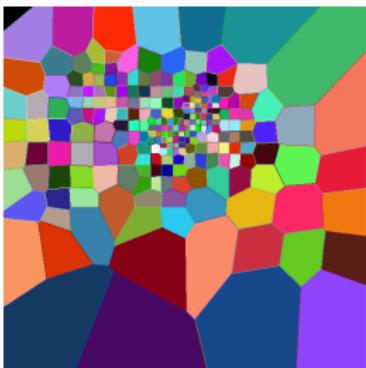
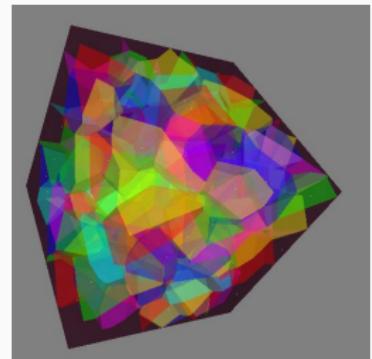
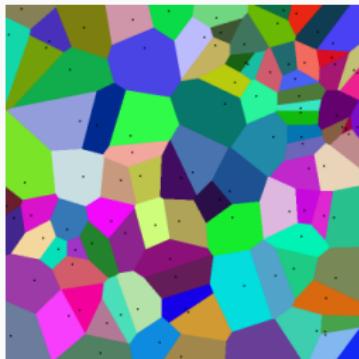
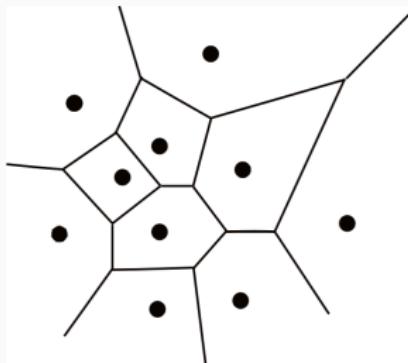
A quantizer of order k is a NN quantizer if its partition is the **Voronoi partition** associated with its codebook. Thus, a NN quantizer takes the form $q = (\mathcal{C}, \mathcal{P}_V(\mathcal{C}))$, where $|\mathcal{C}| \leq k$.

- ▷ A NN quantizer is entirely **characterized** by its codebook, via the rule

$$\|x - q(x)\| = \min_{c_j \in \mathcal{C}} \|x - c_j\|.$$

- ▷ **Vocabulary:** the c_j are the **centers** or the **centroids**.

Voronoi diagrams



Bryant park



Properties of NN quantizers

Proposition

Let q_{NN} be a **NN quantizer** with codebook $\mathcal{C} = \{c_1, \dots, c_k\}$. Then

$$D(\mu, q_{\text{NN}}) = \mathbb{E} \min_{1 \leq j \leq k} \|X - c_j\|^2 = \int_{\mathbb{R}^d} \min_{1 \leq j \leq k} \|x - c_j\|^2 \mu(dx).$$

In addition, for **any** quantizer $q = (\mathcal{C}, \mathcal{P})$, $D(\mu, q_{\text{NN}}) \leq D(\mu, q)$.

- ▷ If quantizers with **minimal distortion** exist, they are **NN quantizers**.
- ▷ **Notation:** $q_{\text{NN}} = (c, \mathcal{P}_V(c))$, with $c = (c_1, \dots, c_k) \in \mathbb{R}^{dk}$ and distortion

$$W(\mu, c) \stackrel{\text{def}}{=} D(\mu, q_{\text{NN}}).$$

Theorem

There exists a quantizer with **minimal** distortion.

Empirical quantization

- In practice, the distribution of X is **unknown**.
- **Sample**: X_1, \dots, X_n i.i.d., distributed as (and independent of) X .
- **Objective**: construct a **good** $q_n(\cdot) = q_n(\cdot; X_1, \dots, X_n)$.
- The **distortion** of q_n is naturally defined by

$$D(\mu, q_n) = \mathbb{E}(\|X - q_n(X)\|^2 | X_1, \dots, X_n) = \int_{\mathbb{R}^d} \|x - q_n(x)\|^2 \mu(dx).$$

- **Empirical measure**: $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$.
- **Empirical distortion**:

$$D(\mu_n, q) = \int_{\mathbb{R}^d} \|x - q(x)\|^2 \mu_n(dx) = \frac{1}{n} \sum_{i=1}^n \|X_i - q(X_i)\|^2.$$

- For $q_{NN} = (c, \mathcal{P}_V(c))$, with $c = (c_1, \dots, c_k) \in \mathbb{R}^{dk}$,

$$D(\mu_n, q_{NN}) = W(\mu_n, c) = \frac{1}{n} \sum_{i=1}^n \min_{1 \leq j \leq k} \|X_i - c_j\|^2.$$

Consistent quantization

Definition (Consistency)

A quantizer is **consistent** if

$$\mathbb{E}D(\mu, q_n) \rightarrow D^*(\mu) \quad \text{as } n \rightarrow \infty.$$

It has **rate of convergence** $(v_n)_n$ if

$$\mathbb{E}D(\mu, q_n) - D^*(\mu) = O(1/v_n), \quad \text{with } v_n \rightarrow \infty.$$

- ▷ **Exercise:** consistency $\Leftrightarrow D(\mu, q_n) \xrightarrow{L^1} D^*(\mu)$.
- ▷ **Natural choice:** q_n^* that **minimizes** the empirical distortion over all NN quantizers.
- ▷ **Definition:** $c_n^* = (c_{n,1}^*, \dots, c_{n,k}^*)$ such that

$$W(\mu_n, c_n^*) = \inf_{c \in \mathbb{R}^{dk}} W(\mu_n, c).$$

So,

$$q_n^* = (c_n^*, \mathcal{P}_V(c_n^*)).$$

Quantization and clustering

- q_n^* allows a **clustering** of X_1, \dots, X_n into k groups.
- **Principle:** X_i is assigned to group j if $q_n^*(X_i) = j$.
- **Cluster** $\#j$ = the X_i such that $\|X_i - c_{n,j}^*\| \leq \|X_i - c_{n,\ell}^*\|, \forall \ell = 1, \dots, k$.

DON'T FORGET!

Clustering = unsupervised learning.

Towards an algorithm

- Computation of q_n^* is often a NP hard problem \rightarrow **k-means algorithm**.
- **Basic idea:** for $\mathcal{C} = \{c_1, \dots, c_k\}$ and $\mathcal{P} = \{A_1, \dots, A_k\}$, let $q = (\mathcal{C}, \mathcal{P})$ and $q_n = (\mathcal{C}_n, \mathcal{P})$, with $\mathcal{C}_n = \{c_{n,1}, \dots, c_{n,k}\}$ such that

$$c_{n,j} = \arg \min_{y \in \mathbb{R}^d} \sum_{i=1}^n \|X_i - y\|^2 \mathbb{1}_{[X_i \in A_j]} = \frac{\sum_{i=1}^n X_i \mathbb{1}_{[X_i \in A_j]}}{\sum_{i=1}^n \mathbb{1}_{[X_i \in A_j]}}, \quad 1 \leq j \leq k.$$

Then

$$\begin{aligned} D(\mu_n, q) &= \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^n \|X_i - c_j\|^2 \mathbb{1}_{[X_i \in A_j]} \\ &\geq \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^n \|X_i - c_{n,j}\|^2 \mathbb{1}_{[X_i \in A_j]} \\ &= D(\mu_n, q_n). \end{aligned}$$

k -means algorithm

1. **Initialization:** $\mathcal{C}^{(1)} = \{c_1^{(1)}, \dots, c_k^{(1)}\}$ and $\mathcal{P}_V^{(1)} = \{A_1^{(1)}, \dots, A_k^{(1)}\}$.
2. **Lloyd's iteration:** compute $\mathcal{C}^{(\ell+1)} = \{c_1^{(\ell+1)}, \dots, c_k^{(\ell+1)}\}$ from $\mathcal{C}^{(\ell)} = \{c_1^{(\ell)}, \dots, c_k^{(\ell)}\}$ via the iteration

$$c_j^{(\ell+1)} = \frac{\sum_{i=1}^n X_i \mathbb{1}_{[X_i \in A_j^{(\ell)}]}}{\sum_{i=1}^n \mathbb{1}_{[X_i \in A_j^{(\ell)}]}}, \quad 1 \leq j \leq k,$$

where $\{A_1^{(\ell)}, \dots, A_k^{(\ell)}\}$ is the **Voronoi partition** associated with $\mathcal{C}^{(\ell)}$.

3. The algorithm **stops after a finite number of iterations.**
4. **Warning:** the output codebook is **not** c_n^* .

Consistency of q_n^*

- Reminder: $c_n^* = (c_{n,1}^*, \dots, c_{n,k}^*)$ such that

$$W(\mu_n, c_n^*) = \inf_{c \in \mathbb{R}^{dk}} W(\mu_n, c).$$

So,

$$q_n^* = (c_n^*, \mathcal{P}_V(c_n^*)).$$

- Next step: prove the consistency of q_n^* .

Definition (Wasserstein distance)

Let ν_1 and ν_2 be probability measures on \mathbb{R}^d with finite second moment. The **Wasserstein distance** ρ_W between ν_1 and ν_2 is

$$\rho_W(\nu_1, \nu_2) = \inf_{\substack{X \stackrel{\mathcal{D}}{\equiv} \nu_1, Y \stackrel{\mathcal{D}}{\equiv} \nu_2}} \sqrt{\mathbb{E}\|X - Y\|^2}.$$

Wasserstein distance

- **Property 1:** There exists (X_0, Y_0) such that $X_0 \stackrel{\mathcal{D}}{=} \nu_1$, $Y_0 \stackrel{\mathcal{D}}{=} \nu_2$, and

$$\rho_W(\nu_1, \nu_2) = \sqrt{\mathbb{E}\|X_0 - Y_0\|^2}.$$

- **Property 2:** One has $\rho_W(\nu_n, \nu) \rightarrow 0$ if

$$\nu_n \Rightarrow \nu \quad \text{and} \quad \int_{\mathbb{R}^d} \|x\|^2 \nu_n(dx) \rightarrow \int_{\mathbb{R}^d} \|x\|^2 \nu(dx).$$

Proposition

Let ν_1 and ν_2 be probability measures on \mathbb{R}^d with finite second moment. If q is a **NN quantizer**, then

$$|D(\nu_1, q)^{1/2} - D(\nu_2, q)^{1/2}| \leq \rho_W(\nu_1, \nu_2).$$

Consistency of q_n^*

Theorem

One has $D(\mu, q_n^*) \rightarrow D^*(\mu)$ almost surely, and $\mathbb{E}D(\mu, q_n^*) \rightarrow D^*(\mu)$.

- ▷ Conclusion: the quantizer q_n^* is consistent.

Theorem

If $\|X\| \leq R$ with probability one, then

$$\mathbb{E}D(\mu, q_n^*) - D^*(\mu) \leq \frac{12kR^2}{\sqrt{n}}.$$

- ▷ $\|X\| \leq R$ is called the peak power constraint.
- ▷ Take-home message: the rate of convergence does not depend on d !