

# Statistical Learning and Applications

Laurent Jacob

September 15, 2014

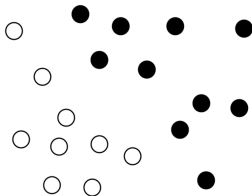
- 8 three hour classes.
- Assessment: 2/3 project, 1/3 homeworks.
- Projects: study article, either methods (implementation), or theoretical. You are free to suggest articles, or pick one from the website.
- End of November: preliminary report (25% of the grade). January: final (short) report.
- 3 homeworks along the semester, due within three weeks.
- Website:  
<http://lear.inrialpes.fr/people/mairal/teaching/2014-2015/M2ENS/>
- Scribe: For each course, a duo of students commit to turn their notes into latex format.

- Hastie, Tibshirani, Friedman. The Elements of Statistical Learning, 2001. (free online)
- Theoretical statistics class by P. Bartlett:  
<http://www.stat.berkeley.edu/~bartlett/courses/2013spring-stat210b/>.
- Theoretical statistics class by S. Arlot and F. Bach (in French):  
<http://www.di.ens.fr/~arlot/2013orsay.htm>.
- Boyd and Vandenberghe. Convex Optimization, 2004. (free online)
- The matrix cookbook.

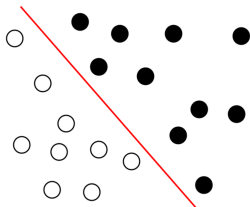
# Outline of this class

- ① A few examples.
- ② Bias/variance trade-off and how to deal with it.
- ③ Supervised learning.
- ④ Unsupervised learning.
- ⑤ Statistical learning theory.

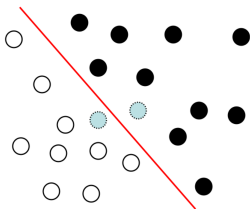
- This class is concerned with learning from data. Essentially:



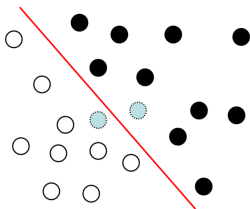
- This class is concerned with learning from data. Essentially:



- This class is concerned with learning from data. Essentially:



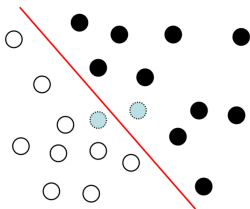
- This class is concerned with learning from data. Essentially:



- Also: multi-class, regression, unsupervised...

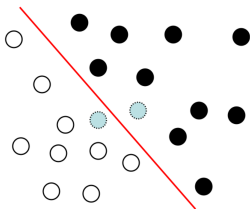


- This class is concerned with learning from data. Essentially:



- Also: multi-class, regression, unsupervised...
- We start with a few examples to make things concrete.

- This class is concerned with learning from data. Essentially:



- Also: multi-class, regression, unsupervised...
- We start with a few examples to make things concrete.
- These examples highlight a general problem which we will discuss right after.

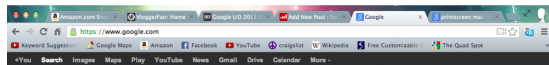
# Part I

## A few examples

# Recommender systems



Given a user and the movies he liked, what should he watch next?



Given a query what are the most relevant webpages?

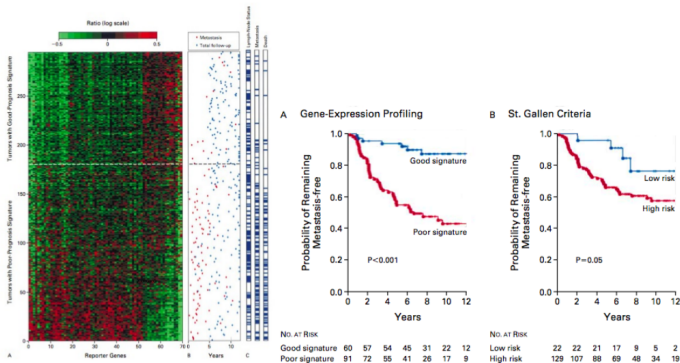
- Given a text, predict its topic.
- Given an email, predict whether it is a spam.
- Given a text, predict its translation in another language.

Modern technologies in molecular biology provide descriptions of individuals through thousands/millions of descriptors:

- Gene expression (arrays, sequencing),
- SNPs,
- Methylations,
- ...

Potential to allow better understanding/prediction of complex phenomena.

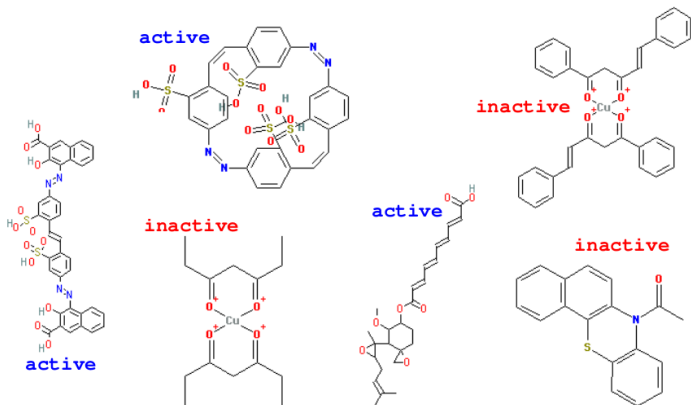
# Tumor classification for prognosis



- Given the expression of the genes in a new tumor, predict the occurrence of a metastasis in the next 5 years.
- Similarly: diagnosis.

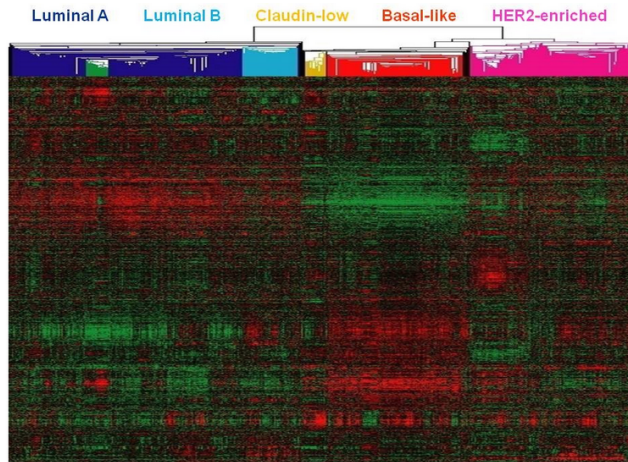


# Molecule classification for drug design



Given a candidate molecule, is it active against a therapeutical target.

# Gene expression clustering



(from C. Perou's website)

Are there groups of breast tumors with similar gene expression profile?

The screenshot shows the article page for "The half-life of DNA in bone: measuring decay kinetics in 158 dated fossils" in the journal Proceedings of the Royal Society B: Biological Sciences. The page includes the journal logo, a navigation bar with links like Home, Current issue, Past issues, Submit, Subscribe, and Alerts, and a CrossMark logo. The article title is prominently displayed, followed by the authors' names: Morten E. Allentoft, Matthew Collins, David Harker, James Haile, Charlotte L. Oskam, Marie L. Hale, Paula F. Campos, Jose A. Samaniego, M. Thomas P. Gilbert, Eske Willerslev, Guojie Zhang, R. Paul Scofield, Richard N. Holdaway, and Michael Bunce. Below the authors, there are links for "Author Affiliations" and "Abstract". The abstract text is visible, starting with "Claims of extreme survival of DNA have emphasized the need for reliable models of DNA degradation through time. By analysing mitochondrial DNA (mtDNA)". On the right side, there is a "Return To Issue" button and a "This Article" section containing publication details (published online 10 October 2012, doi: 10.1098/rspb.2012.1745) and a list of services such as "Abstract Free", "Full Text Free", and "Email this article to a friend".

**PROCEEDINGS OF THE ROYAL SOCIETY B** | BIOLOGICAL SCIENCES

Home | Current issue | Past issues | Submit | Subscribe | Alerts

CrossMark  
click for updates

**The half-life of DNA in bone: measuring decay kinetics in 158 dated fossils**

Morten E. Allentoft<sup>1,2,3,\*</sup>, Matthew Collins<sup>4</sup>, David Harker<sup>4</sup>, James Haile<sup>1</sup>, Charlotte L. Oskam<sup>1</sup>, Marie L. Hale<sup>2</sup>, Paula F. Campos<sup>3,5</sup>, Jose A. Samaniego<sup>3</sup>, M. Thomas P. Gilbert<sup>1,3</sup>, Eske Willerslev<sup>3</sup>, Guojie Zhang<sup>6</sup>, R. Paul Scofield<sup>7</sup>, Richard N. Holdaway<sup>2,8</sup> and Michael Bunce<sup>1,\*</sup>

[+ Author Affiliations](#)

\* Authors for correspondence (morten.allentoft@gmail.com; m.bunce@murdoch.edu.au).

**Abstract**

Claims of extreme survival of DNA have emphasized the need for reliable models of DNA degradation through time. By analysing mitochondrial DNA (mtDNA)

[Return To Issue](#)

**This Article**

Published online before print 10 October 2012  
doi: 10.1098/rspb.2012.1745  
Proc. R. Soc. B  
rspb20121745

- **Abstract** Free
- Full Text **Free**
- Full Text (PDF) **Free**
- Data Supplement

➤ All Versions of this Article:  
rspb.2012.1745v1  
279/1748/4724 *most recent*

**Classifications**

- Research article

**Services**

- Email this article to a friend
- Alert me when this article

Decay of DNA molecules.

# Ancestral genome reconstruction

**WIRED** GEAR SCIENCE ENTERTAINMENT BUSINESS SECURITY DESIGN OPINION MAGAZINE

bon appétit SANDWICH

Get your favorite magazines on your tablet!

SMART CAR CLICK HERE

SCIENCE | Biology | Dinosaurs | DNA | Half Life

## Jurassic Park Impossible Because of Stupid Laws of Physics

BY WIRED UK 10.10.12 | 12:15 PM | PERMALINK

Share | Tweet | +1 | 122 | Share | Print



By Ian Steadman, Wired UK

The lesson of the *Jurassic Park* tragedy was clear — man and dinosaur were not meant to coexist. It's lucky then that dinosaur fossils are far too old to contain any genetic material that could be used for cloning. DNA breaks down over time, even when kept in ideal conditions, and a study of extinct moa bones has revealed an estimate of the half-life for our genes.

LE FIGARO.fr ACTUALITE ECONOMIE SPORT CULTURE LIFESTYLE MAISON LE FIGARO BUSINESS

IN BREVET | TOPS | PLUS | REPORTAGE | LE FIGARO | LE FIGARO BUSINESS | iRobot do you? avec technologie AeroForce

LE FLASH ACTU 11/118 Rythmes scolaires: Gaudin en concert

### L'ADN est trop fragile pour envisager un «Jurassic Park»

ACTUALITE | SCIENCE & ENVIRONNEMENT | 10/10/2012 à 19:38 | Publié le 10/10/2012 à 18:04

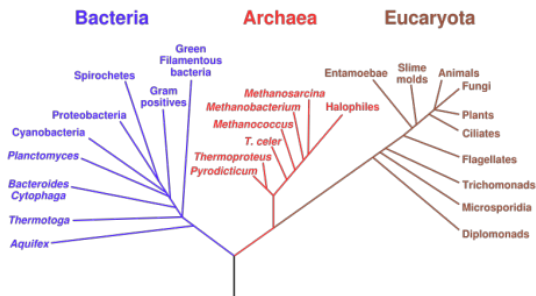
UNIVERS | SAN ET BIEN-ÊTRE | RECHERCHES | PARTAGER | INTRODUCE

**Wired** Apprenez une langue avec Babbel et Lingua111

Des scientifiques viennent de montrer qu'un brin d'ADN devient illisible au bout de quelques millions d'années seulement. Un délai bien plus court que les 66 millions d'années qui nous séparent de la disparition des dinosaures.

Does it make Jurassic Park unrealistic?

## Phylogenetic Tree of Life



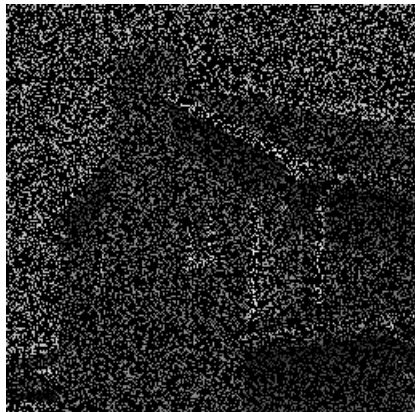
Actually it does. But given enough descendants, we can infer the genome of extinct ancestors (black death, LUCA).

# Image inpainting



Complete an image with missing parts.

# Image inpainting



Estimation problem: predict each image patch, as a linear combination of dictionary elements.



Since 1699, when French explorers landed at the great bend of the Mississippi River and celebrated the first Mardi Gras in North America, New Orleans has brewed a fascinating melange of cultures. It was French, then Spanish, then French again, then sold to the United States. Through all these years, and even into the 1900s, others arrived from everywhere: Acadians (Cajuns), Africans, indige-



# Image inpainting



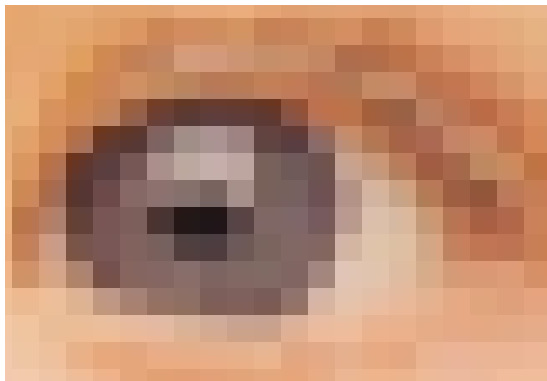


Improve the quality of an image.

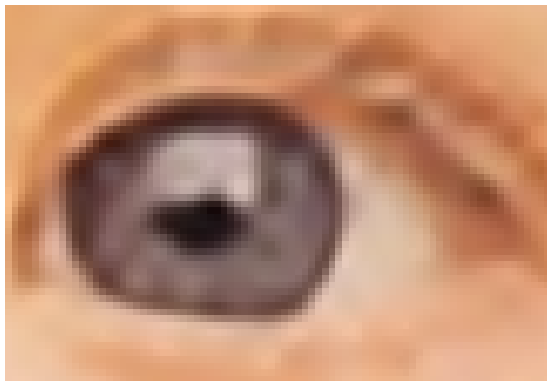
# Image up-scaling



Improve the quality of an image.

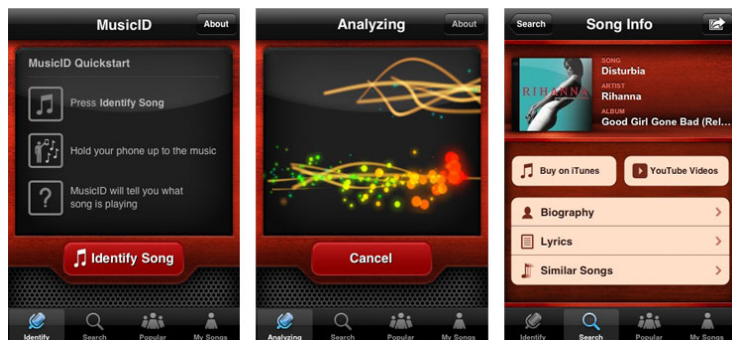


Improve the quality of an image.



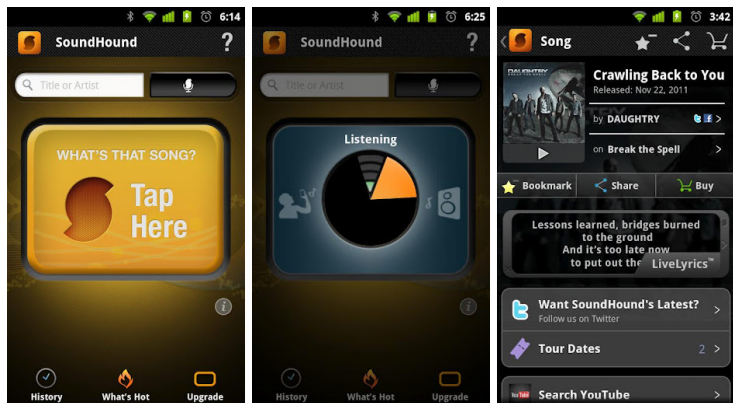
Improve the quality of an image.

# Music recognition



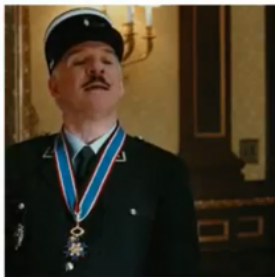
Guess which tune is being played.

# Music recognition



Guess which tune is being tapped/hummed.

Presented clip



Clip reconstructed  
from brain activity



- Collect fMRI data of people watching videos.
- Reconstitute new video they are watching based on fMRI measurements (“brain reading”).



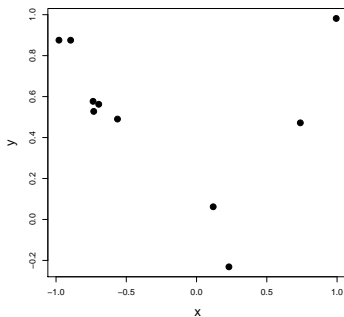
- Each of these examples involves **complex objects**/**large numbers of features** for a **restricted number of samples**.
- Intuitively, observing all these characteristics should allow us to predict or understand complex mechanisms.
- We now discuss why this wealth of features can cause trouble in statistical learning.
- Understanding this problem should give more perspective to the tools we will present later.

## Part II

Overfitting, bias-variance tradeoff: what is the problem?

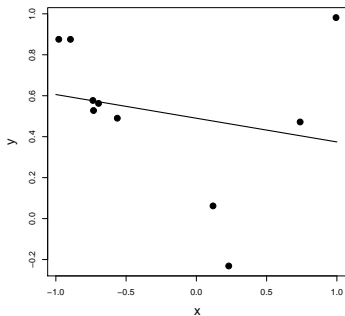
- We start with an informal example.
- We will formalize what we observe later.

# Bias-variance tradeoff: intuition



- We observe 10 couples  $(x_i, y_i)$ .
- We want to estimate  $y$  from  $x$ .
- Strategy: find  $f$  such that  $f(x_i)$  is close to  $y_i$ .

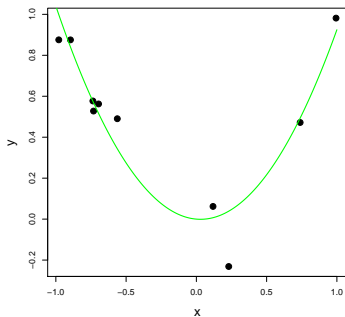
# Bias-variance tradeoff: intuition



Find  $f$  as a line

$$\min_{f(x)=ax+b} \|Y - f(X)\|^2$$

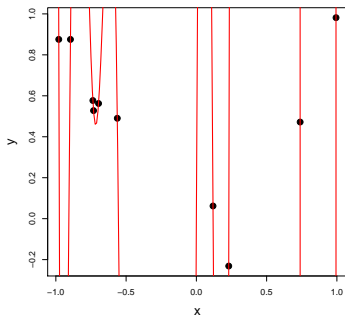
# Bias-variance tradeoff: intuition



Find  $f$  as a quadratic function

$$\min_{f(x)=ax+bx^2} \|Y - f(X)\|^2$$

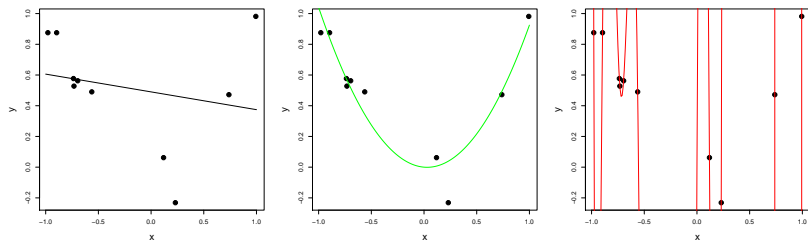
# Bias-variance tradeoff: intuition



Find  $f$  as a polynomial of degree 10

$$\min_{f(x)=\sum_{j=1}^{10} a_j x^j} \|Y - f(X)\|^2$$

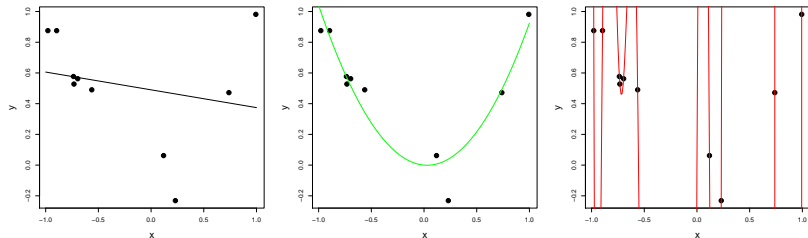
# Bias-variance tradeoff: intuition



Which function would you trust to predict  $y$  corresponding to  $x = 0.5$ ?

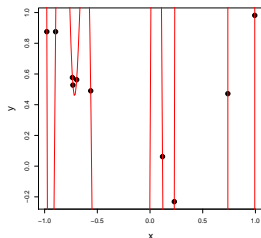


# Bias-variance tradeoff: intuition



- Reminder: we aim at “finding  $f$  such that  $f(x_i)$  is close to  $y_i$ ”.
- With the polynomial of degree 10,  $f(x_i) - y_i = 0$  for all 10 points.
- There is something wrong with our objective.

# Bias-variance tradeoff: intuition

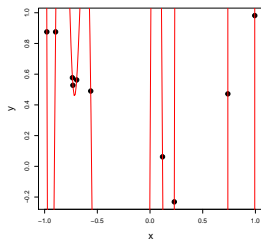


More precisely:

- If we allow any function  $f$ , we can find **a lot** of perfect solutions.
- Our actual goal is to estimate  $y$  for new points  $x$  from the same population :

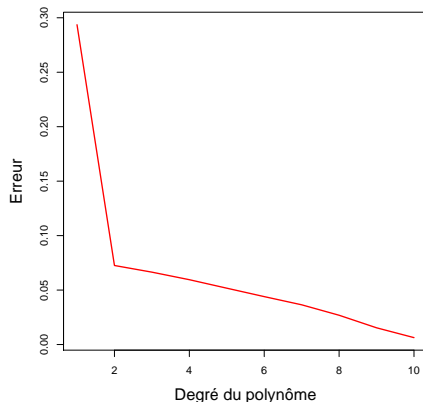
$$\min_f \mathbb{E}_{(X,Y)} \|Y - f(X)\|^2$$

# Biais-variance tradeoff: intuition

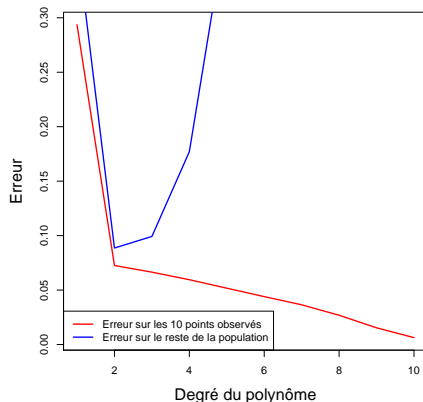


Even more precisely :

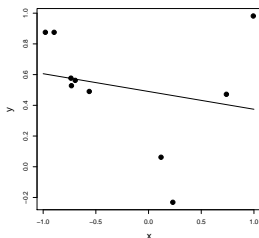
- We did not take into account the fact that our 10 points are a subsample from the population.
- If we sample 10 new points from the same population, the complex functions are likely to change more than the simple ones.
- Consequence: these fonctions will probably generalize less well to the rest of the population.



- When the degree increases, the error  $\|y - f(x)\|^2$  over the 10 observations always decreases.
- Over the rest of the population, the error decreases, then increases.

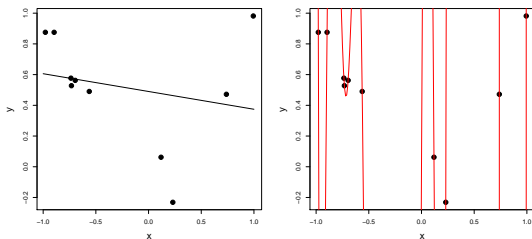


- When the degree increases, the error  $\|y - f(x)\|^2$  over the 10 observations always decreases.
- Over the rest of the population, the error decreases, **then increases**.



This suggests the existence of a **tradeoff** between two types of errors:

- Sets of functions which are too simple cannot contain functions which explain the data well enough.
- Sets of functions which are too rich may contain functions which are too specific to the observed sample.



This suggests the existence of a **tradeoff** between two types of errors:

- Sets of functions which are too simple cannot contain functions which explain the data well enough.
- Sets of functions which are too rich may contain functions which are too specific to the observed sample.

## Parenthesis: complexity vs dimension (1/3)

- Our introductory examples had a **large number of descriptors**.
- This case involves increasingly **complex** functions of a single variable.



- In fact, the two notions are related: here in particular, the three functions are linear in different representations.
- Reminder (linear regression):  
$$\arg \min_{\theta \in \mathbb{R}^p} \|Y - X\theta\|^2 = (X^\top X)^{-1} X^\top Y \text{ (if } X^\top X \text{ is invertible).}$$
- How can we use this fact to compute  
$$\arg \min_{f(x) = \sum_{j=1}^p a_j x^j} \|Y - f(X)\|^2?$$

## Parenthesis : complexity vs dimension (3/3)

- We could have illustrated the same principle using linear functions involving more and more variables.
- Example : predicting a phenotype using the expression of an increasing number of genes.
- We stuck to polynomials, which allow for better visual representations.
- Along this class, the notion of complexity of a set of functions will become more and more precise.
- Complexity is what causes problems for inference, not just dimension.

- Until now, we did not need to introduce a **model** for the data, *i.e.*, a distribution over  $\mathcal{X} \times \mathcal{Y}$  :
  - Data could come from any population.
  - The functions we used to predict  $y$  can be derived from particular probabilistic models, but this is not necessary (they were in fact historically introduced without a model).
- The objective is not to criticize the use of models, but to show that the tradeoff problem we introduced goes beyond probabilistic models.
- We now show how using a model can give a better insight into the problem.

## A little more formally: biais-variance decomposition

- We now assume that the data follow:

$$y = f(x) + \varepsilon, \quad (1)$$

and  $\mathbf{E}[\varepsilon] = 0$ .

- Without loss of generality, we consider an estimator  $\hat{f}$  of  $f$ , fonction function of the data  $\mathcal{D} = (x_i, y_i)_{(i=1, \dots, n)}$  generated under (1) (so don't forget:  $\hat{f}$  is a random quantity).
- We consider the mean **quadratic error**  $\mathbf{E}[(y - \hat{f}(x))^2]$  incurred when using  $\hat{f}$  to estimate  $y$  from  $x$ , generated under (1) but independent from  $\mathcal{D}$ .
- Expectation is taken over the  $(n + 1)$   $(x, y)$  pairs :  $n$  to build  $\hat{f}$ , plus the one over which we compute the error.

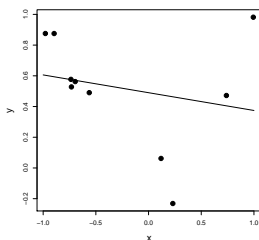
## Proposition

*Under the previous hypotheses,*

$$\begin{aligned} \mathbf{E}[(y - \hat{f}(x))^2] &= \left( \mathbf{E}[\hat{f}(x)] - f(x) \right)^2 + \mathbf{E} \left[ \left( \mathbf{E}[\hat{f}(x)] - \hat{f}(x) \right)^2 \right] \\ &\quad + \mathbf{E}[(y - f(x))^2] \end{aligned}$$

- The first term is the squared **bias** of  $\hat{f}$ : the difference between its mean (over the sample of  $\mathcal{D}$ ) and the true  $f$ .
- The second term is the **variance** of  $\hat{f}$ : how much  $\hat{f}$  varies around its average when the data change.
- The third term is the **Bayes error**, and does not depend on the estimator. The actual quantity of interest is the **excess of risk**  $\mathbf{E}[(y - \hat{f}(x))^2] - \mathbf{E}[(y - f(x))^2]$ .

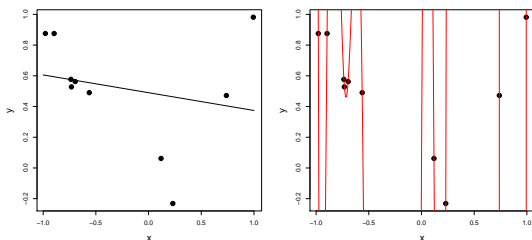
## Back to our example



### Tradeoff between two types of error:

- Sets of functions which are too simple cannot contain functions which explain the data well enough:  
these sets lead to estimators with a large **bias**.
- Sets of functions which are too rich may contain functions which are too specific to the observed sample:  
these sets lead to estimators with a large **variance**.

# Back to our example



## Tradeoff between two types of error:

- Sets of functions which are too simple cannot contain functions which explain the data well enough:  
these sets lead to estimators with a large **bias**.
- Sets of functions which are too rich may contain functions which are too specific to the observed sample:  
these sets lead to estimators with a large **variance**.

## Reminder (König-Huygens)

For any real random variable  $Z$ ,  $\mathbf{E}[(Z - \mathbf{E}[Z])^2] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$

$$\mathbf{E}[(y - \hat{f}(x))^2] = \mathbf{E}[y^2 - 2y\hat{f}(x) + \hat{f}(x)^2]$$



## Reminder (König-Huygens)

For any real random variable  $Z$ ,  $\mathbf{E} \left[ (Z - \mathbf{E}[Z])^2 \right] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$

$$\begin{aligned} \mathbf{E}[(y - \hat{f}(x))^2] &= \mathbf{E}[y^2 - 2y\hat{f}(x) + \hat{f}(x)^2] \\ &= \mathbf{E}[y^2] - \mathbf{E}[2y\hat{f}(x)] + \mathbf{E}[\hat{f}(x)^2] \end{aligned}$$

## Reminder (König-Huygens)

For any real random variable  $Z$ ,  $\mathbf{E}[(Z - \mathbf{E}[Z])^2] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$

$$\begin{aligned}\mathbf{E}[(y - \hat{f}(x))^2] &= \mathbf{E}[y^2 - 2y\hat{f}(x) + \hat{f}(x)^2] \\ &= \mathbf{E}[y^2] - \mathbf{E}[2y\hat{f}(x)] + \mathbf{E}[\hat{f}(x)^2] \\ &= \mathbf{E}[y^2] + \mathbf{E}[(y - \mathbf{E}[y])^2] \\ &\quad - 2\mathbf{E}[y]\mathbf{E}[\hat{f}(x)] \\ &\quad + \mathbf{E}[\hat{f}(x)]^2 + \mathbf{E}[(\hat{f}(x) - \mathbf{E}[\hat{f}(x)])^2]\end{aligned}$$

## Reminder (König-Huygens)

For any real random variable  $Z$ ,  $\mathbf{E}[(Z - \mathbf{E}[Z])^2] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$

$$\begin{aligned}\mathbf{E}[(y - \hat{f}(x))^2] &= \mathbf{E}[y^2 - 2y\hat{f}(x) + \hat{f}(x)^2] \\ &= \mathbf{E}[y^2] - \mathbf{E}[2y\hat{f}(x)] + \mathbf{E}[\hat{f}(x)^2] \\ &= f(x)^2 + \mathbf{E}[(y - f(x))^2] \\ &\quad - 2f(x)\mathbf{E}[\hat{f}(x)] \\ &\quad + \mathbf{E}[\hat{f}(x)]^2 + \mathbf{E}[(\hat{f}(x) - \mathbf{E}[\hat{f}(x)])^2]\end{aligned}$$

## Reminder (König-Huygens)

For any real random variable  $Z$ ,  $\mathbf{E}[(Z - \mathbf{E}[Z])^2] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$

$$\begin{aligned}\mathbf{E}[(y - \hat{f}(x))^2] &= \mathbf{E}[y^2 - 2y\hat{f}(x) + \hat{f}(x)^2] \\ &= \mathbf{E}[y^2] - \mathbf{E}[2y\hat{f}(x)] + \mathbf{E}[\hat{f}(x)^2] \\ &= f(x)^2 + \mathbf{E}[(y - f(x))^2] \\ &\quad - 2f(x)\mathbf{E}[\hat{f}(x)] \\ &\quad + \mathbf{E}[\hat{f}(x)]^2 + \mathbf{E}[(\hat{f}(x) - \mathbf{E}[\hat{f}(x)])^2] \\ &= \mathbf{E}[(y - f(x))^2] + \mathbf{E}[(\hat{f}(x) - \mathbf{E}[\hat{f}(x)])^2] \\ &\quad + \left(\mathbf{E}[\hat{f}(x)] - f(x)\right)^2\end{aligned}$$

- Using a (rather general) model, we managed to start formalizing the tradeoff introduced with our example.
- We now generalize this formalization.

- We now suppose more generally that the observations are sampled from a joint distribution  $\mathbb{P}(x, y)$ .
- This does not necessarily mean that we assume a particular probabilistic model: given a deterministic set of couples  $(x, y)$ ,  $\mathbb{P}$  can be their empirical distribution.
- We also consider a **loss function**

$$L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

$L(y, y')$  quantifies the cost of the error made by predicting  $y'$  when the true value is  $y$ .

- Special case (our example):  $L(y, y') = (y - y')^2$ .

We look for an estimator  $f : \mathcal{X} \rightarrow \mathcal{Y}$  minimizing

$$R(f) = \int_{\mathcal{X} \times \mathcal{Y}} L(y, f(x)) d\mathbb{P} = \mathbf{E}[L(y, f(x))]. \quad (2)$$

$R$  is the **risk** of  $f$  : the average cost of using  $f$  to predict  $y$  from  $x$  over the joint distribution.

- In practice, we cannot compute  $R(f)$  because the distribution  $\mathbb{P}$  is unknown.
- We therefore use a training set ( $\mathcal{D}$  in the previous example) to estimate  $R$ , for example through the **empirical risk**:

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)). \quad (3)$$

- **Empirical risk minimization** : choose  $f$  minimizing  $\hat{R}$ .
- We saw in our example that minimizing the empirical risk was not enough to obtain a low risk  $R$



## A little more generally : structural risk minimization

- More generally, we can minimize the risk over a function space  $\mathcal{H}$  (polynomials of a certain degree in our example).
- If  $R^*$  is the Bayes risk, we can decompose the **Bayes regret** :

$$R(f) - R^* = \left( R(f) - \inf_{g \in \mathcal{H}} R(g) \right) + \left( \inf_{g \in \mathcal{H}} R(g) - R^* \right). \quad (4)$$

- The second term is the approximation error: the smallest excess of risk we can reach using a function of  $\mathcal{H}$ .
- This is a **bias** term, which does not depend on the data but only on the size of  $\mathcal{H}$ .
- The first term is the excess of risk of  $f$  with respect to the best function in  $\mathcal{H}$ .

- We consider  $\hat{f}$  obtained by minimization of the empirical risk over  $\mathcal{H}$ :

$$\hat{f} \in \arg \min_{g \in \mathcal{H}} \hat{R}(g)$$

- We want to bound the excess of risk  $R(\hat{f}) - \inf_{g \in \mathcal{H}} R(g) \geq 0$
- This term (estimation error) can be decomposed:

$$\begin{aligned} R(\hat{f}) - \inf_{g \in \mathcal{H}} R(g) &\stackrel{\Delta}{=} R(\hat{f}) - R(f_{\mathcal{H}}^*) \\ &= R(\hat{f}) - \hat{R}(\hat{f}) \\ &\quad + \hat{R}(\hat{f}) - \hat{R}(f_{\mathcal{H}}^*) \\ &\quad + \hat{R}(f_{\mathcal{H}}^*) - R(f_{\mathcal{H}}^*). \end{aligned}$$

$$\begin{aligned} R(\hat{f}) - \inf_{g \in \mathcal{H}} R(g) &= R(\hat{f}) - R(f_{\mathcal{H}}^*) \\ &= R(\hat{f}) - \hat{R}(\hat{f}) \\ &\quad + \hat{R}(\hat{f}) - \hat{R}(f_{\mathcal{H}}^*) \\ &\quad + \hat{R}(f_{\mathcal{H}}^*) - R(f_{\mathcal{H}}^*). \end{aligned}$$

- Reminder :
  - $R$  is the **population** risk,  $\hat{R}$  the **empirical** risk, an estimator.
  - $\hat{f}$  is the estimator minimizing  $\hat{R}$  over  $\mathcal{H}$ ,  $f_{\mathcal{H}}^*$  the one obtained by minimizing  $R$  over  $\mathcal{H}$ .
  - We therefore estimate at two levels: the function  $f$  and the risk  $R$ .

$$\begin{aligned} R(\hat{f}) - \inf_{g \in \mathcal{H}} R(g) &= R(\hat{f}) - \hat{R}(\hat{f}) \\ &\quad + \hat{R}(\hat{f}) - \hat{R}(f_{\mathcal{H}}^*) \\ &\quad + \hat{R}(f_{\mathcal{H}}^*) - R(f_{\mathcal{H}}^*). \end{aligned}$$

- The first term is the difference between the true risk and the estimated risk, for our estimator of  $f$ .
- This is a complex object to study. **Statistical learning theory** (Vapnik and Chervonenkis) aims at bounding this quantity as a function of  $n$  and the complexity of  $\mathcal{H}$ .
- The second term is nonpositive by construction.
- The third one is easier to control as it involves a deterministic function and the law of large numbers applies.

## A little more generally : structural risk minimization

We can however bound the first term:

$$R(\hat{f}) - \hat{R}(\hat{f}) \leq \sup_{f \in \mathcal{H}} \left| \mathbf{E}[L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right|,$$

and since this quantity also bounds the third term, we get

$$R(\hat{f}) - \inf_{g \in \mathcal{H}} R(g) \leq 2 \sup_{f \in \mathcal{H}} \left| \mathbf{E}[L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right|.$$

- This bound of the estimation error suggests that it corresponds to a **variance** term, which increases with the size of  $\mathcal{H}$ .
- The more complex  $\mathcal{H}$  is, the more likely it is to contain a function for which the empirical risk and the population risk are very different.

## A little more generally : structural risk minimization

We can make this notion of size more precise by introducing the **Rademacher complexity** of  $\mathcal{H}$ :

### Definition

Let  $\epsilon_i, i = 1, \dots, n$  i.i.d such that  $\mathbb{P}(\epsilon_i = 1) = \mathbb{P}(\epsilon_i = -1) = 1/2$ ,  $Z_i, i = 1, \dots, n$  i.i.d data and  $\mathcal{H}$  a space of functions defined over this data, then

$$\mathfrak{R}(\mathcal{H}) = \mathbf{E}_{\epsilon_1^n, Z_1^n} \left[ \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Z_i) \right| \right]$$

is the Rademacher complexity of  $\mathcal{H}$ .

Intuition:  $\mathfrak{R}$  measures the capacity of  $\mathcal{H}$  to provide functions which align with noise.

## A little more generally : structural risk minimization

We can make this notion of size more precise by introducing the **Rademacher complexity** of  $\mathcal{H}$ :

### Definition

Let  $\epsilon_i, i = 1, \dots, n$  i.i.d such that  $\mathbb{P}(\epsilon_i = 1) = \mathbb{P}(\epsilon_i = -1) = 1/2$ ,  $Z_i, i = 1, \dots, n$  i.i.d data and  $\mathcal{H}$  a space of functions defined over this data, then

$$\mathfrak{R}(\mathcal{H}) = \mathbf{E}_{\epsilon_1^n, Z_1^n} \left[ \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Z_i) \right| \right]$$

is the Rademacher complexity of  $\mathcal{H}$ .

This complexity increases with the size of  $\mathcal{H}$  and decreases with the size  $n$  of the sample.

## A little more generally : structural risk minimization

We can bound the mean estimation error in terms of the Rademacher complexity of  $\mathcal{H}$ .

### Proposition

$$\mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \leq 2\mathfrak{R}(\mathcal{H}).$$

Therefore,

$$\mathbf{E}_{(x,y)_1^n} [R(\hat{f}) - R^*] \leq \left( \min_{g \in \mathcal{H}} R(g) - R^* \right) + 4\mathfrak{R}(\mathcal{H}).$$



Therefore

$$\mathbf{E}_{(x,y)_1^n} \left[ R(\hat{f}) - R^* \right] \leq \left( \min_{g \in \mathcal{H}} R(g) - R^* \right) + 4\mathfrak{R}(\mathcal{H}),$$

- This result illustrates a little more generally the bias variance tradeoff for risk minimization.
- It makes explicit the link between complexity and sample size: lots of points are needed to estimate in large  $\mathcal{H}$  (otherwise  $\mathfrak{R}(\mathcal{H})$  is large).

Therefore

$$\mathbf{E}_{(x,y)_1^n} \left[ R(\hat{f}) - R^* \right] \leq \left( \min_{g \in \mathcal{H}} R(g) - R^* \right) + 4\mathfrak{R}(\mathcal{H}),$$

Concretely, this analysis is at the core of two major elements of statistical learning (Vapnik and Chervonenkis, late 60's):

- It is used in learning theory to establish consistency of empirical risk minimization: only families with bounded complexity allow to learn by ERM (are consistent).
- **It also suggests a strategy to design estimators:** build small classes  $\mathcal{H}$  which we think contain good approximations.

$$\mathbf{E}_{(x,y)_1^n} \left[ R(\hat{f}) - R^* \right] \leq \left( \min_{g \in \mathcal{H}} R(g) - R^* \right) + 4\mathfrak{R}(\mathcal{H}),$$

Practical procedure proposed by Vapnik and Chervonenkis: **structural risk minimization**:

- 1 Define nested function sets of increasing complexity.
- 2 Minimize the empirical risk over each family.
- 3 Choose the solution giving the best generalization performances.

## Structural risk minimization:

- 1 Define nested function sets of increasing complexity.
- 2 Minimize the empirical risk over each family.
- 3 Choose the solution giving the best generalization performances.

We will study practical instances of this strategy later in this class.

## A little more generally : structural risk minimization

Proof of the previous bound (inspired from Peter Bartlett's slides)

$$\mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right|$$

## A little more generally : structural risk minimization

Proof of the previous bound (inspired from Peter Bartlett's slides)

$$\begin{aligned} & \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) \right] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \end{aligned}$$

## A little more generally : structural risk minimization

Proof of the previous bound (inspired from Peter Bartlett's slides)

$$\begin{aligned} & \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) \right] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right] \right| \end{aligned}$$

## A little more generally : structural risk minimization

Proof of the previous bound (inspired from Peter Bartlett's slides)

$$\begin{aligned} & \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) \right] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right] \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - L(y_i, f(x_i)) \right] \right| \end{aligned}$$



## A little more generally : structural risk minimization

Proof of the previous bound (inspired from Peter Bartlett's slides)

$$\begin{aligned} & \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x,y)} [L(y, f(x))] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) \right] - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \right] \right| \\ &= \mathbf{E}_{(x,y)_1^n} \sup_{f \in \mathcal{H}} \left| \mathbf{E}_{(x',y')_1^n} \left[ \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - L(y_i, f(x_i)) \right] \right| \\ &\leq \mathbf{E}_{(x,y)_1^n} \mathbf{E}_{(x',y')_1^n} \left[ \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - L(y_i, f(x_i)) \right| \right] \end{aligned}$$

## A little more generally : structural risk minimization

We now introduce  $\epsilon_i, i = 1, \dots, n \in \{-1, 1\}$ . Notice that

$$\begin{aligned} & \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n L(y'_i, f(x'_i)) - L(y_i, f(x_i)) \right| \\ &= \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (L(y'_i, f(x'_i)) - L(y_i, f(x_i))) \right|, \end{aligned}$$

since the data is i.i.d, switching the two terms does not affect the distribution of the sup.

The equality holds for any choice of  $\epsilon_i$ , so we can take the expectation over a uniform i.i.d choice.

Finally,

$$\begin{aligned} & \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (L(y'_i, f(x'_i)) - L(y_i, f(x_i))) \right| \\ & \leq \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i L(y'_i, f(x'_i)) \right| + \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i L(y_i, f(x_i)) \right| \\ & = 2 \mathbf{E} \sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i L(y_i, f(x_i)) \right| = 2 \mathfrak{R}(\mathcal{H}). \end{aligned}$$

This proof technique is called **symmetrization**.

# More intuition about the complexity of a set of functions: VC dimension

- In practice, we sometimes use VC dimension of a set of functions to bound the Rademacher complexity.
- We restrict ourselves to the sets  $\mathcal{H}$  of binary valued functions (useful for classification).
- We say a set  $Z = (Z_1, \dots, Z_n)$  is **shattered** by  $\mathcal{H}$  if  $\text{Card} \{f(Z_1), \dots, f(Z_n) \mid f \in \mathcal{H}\} = 2^n$ .
- Interpretation: we can find an  $f \in \mathcal{H}$  assigning 0 to any subset of  $Z$  and 1 to its complement.
- The VC dimension  $\nu(\mathcal{H})$  of  $\mathcal{H}$  is the largest integer  $n$  such that there exists a set  $(Z_1, \dots, Z_n)$  shattered by  $\mathcal{H}$ .

# More intuition about the complexity of a set of functions: VC dimension

- We extend the VC dimension to real valued functions by thresholding functions at 0.
- Linear functions in  $p$  dimensions:  $\mathcal{H}_L = \{f_\theta(x) = \text{sign}(\theta^\top x), \theta \in \mathbb{R}^p\}$ .
- Includes linear functions and polynomials in our introduction.
- We can show that  $\nu(\mathcal{H}_L) = p$ .

# More intuition about the complexity of a set of functions: VC dimension

- Proof of  $\nu(\mathcal{H}_L) \geq p$ : we build a set of  $p$  points in  $p$  dimensions shattered by a function of  $\mathcal{H}_L$ . Let  $\mathcal{E}_p$  be the canonical basis of  $\mathbb{R}^p$ . For any set  $y \in \{0, 1\}^p$  and any  $i = 1, \dots, p$ ,  $f_\theta(e_i) = y_i$  by choosing  $\theta_i = y_i$ .
- Proof of  $\nu(\mathcal{H}_L) < p + 1$ : no set of  $p + 1$  points in  $p$  dimensions can be shattered by a linear function.

# More intuition about the complexity of a set of functions: VC dimension

- Let  $x_1, \dots, x_{p+1} \in \mathbb{R}^p$ . One of the points can necessarily be written as a linear combination of the  $p$  others.

# More intuition about the complexity of a set of functions: VC dimension

- Let  $x_1, \dots, x_{p+1} \in \mathbb{R}^p$ . One of the points can necessarily be written as a linear combination of the  $p$  others.
- Without loss of generality, let us write  $x_{p+1} = \sum_{i=1}^p \alpha_i x_i$  and  $f_\theta(x_{p+1}) = \sum_{i=1}^p \alpha_i \theta^\top x_i$ .



# More intuition about the complexity of a set of functions: VC dimension

- Let  $x_1, \dots, x_{p+1} \in \mathbb{R}^p$ . One of the points can necessarily be written as a linear combination of the  $p$  others.
- Without loss of generality, let us write  $x_{p+1} = \sum_{i=1}^p \alpha_i x_i$  and  $f_\theta(x_{p+1}) = \sum_{i=1}^p \alpha_i \theta^\top x_i$ .
- Let  $y = (\text{sign}(\alpha_1), \dots, \text{sign}(\alpha_p), -1)$ , and assume there exists  $\theta \in \mathbb{R}^p$  such that  $\text{sign}(\theta^\top x_i) = y_i, i = 1, \dots, p$ .

# More intuition about the complexity of a set of functions: VC dimension

- Let  $x_1, \dots, x_{p+1} \in \mathbb{R}^p$ . One of the points can necessarily be written as a linear combination of the  $p$  others.
- Without loss of generality, let us write  $x_{p+1} = \sum_{i=1}^p \alpha_i x_i$  and  $f_\theta(x_{p+1}) = \sum_{i=1}^p \alpha_i \theta^\top x_i$ .
- Let  $y = (\text{sign}(\alpha_1), \dots, \text{sign}(\alpha_p), -1)$ , and assume there exists  $\theta \in \mathbb{R}^p$  such that  $\text{sign}(\theta^\top x_i) = y_i, i = 1, \dots, p$ .
- Then necessarily  $\text{sign}(\theta^\top x_{p+1}) = \text{sign}(\sum_{i=1}^p \alpha_i \theta^\top x_i) = 1$  since  $\text{sign}(\theta^\top x_i) = \text{sign}(\alpha_i), i = 1, \dots, p$ .

# More intuition about the complexity of a set of functions: VC dimension

- Let  $x_1, \dots, x_{p+1} \in \mathbb{R}^p$ . One of the points can necessarily be written as a linear combination of the  $p$  others.
- Without loss of generality, let us write  $x_{p+1} = \sum_{i=1}^p \alpha_i x_i$  and  $f_\theta(x_{p+1}) = \sum_{i=1}^p \alpha_i \theta^\top x_i$ .
- Let  $y = (\text{sign}(\alpha_1), \dots, \text{sign}(\alpha_p), -1)$ , and assume there exists  $\theta \in \mathbb{R}^p$  such that  $\text{sign}(\theta^\top x_i) = y_i, i = 1, \dots, p$ .
- Then necessarily  $\text{sign}(\theta^\top x_{p+1}) = \text{sign}(\sum_{i=1}^p \alpha_i \theta^\top x_i) = 1$  since  $\text{sign}(\theta^\top x_i) = \text{sign}(\alpha_i), i = 1, \dots, p$ .
- $y$  can therefore not be obtained by any function of  $\mathcal{H}_L$ , and no set of  $p + 1$  vectors in  $\mathbb{R}^p$  is shattered by  $\mathcal{H}_L$ .

- We saw how the risk could generally be decomposed as a term of bias/approximation and a term of variance/estimation.
- This decomposition highlights the tradeoff that needs to be dealt with in inference. This tradeoff is related to the complexity of the set of functions under consideration:
  - Sets too simple lead to a large approximation error.
  - Sets too large lead to a large estimation error.
- We defined this notion of complexity more precisely (Rademacher, VC), and saw it also depended on the number of samples.
- These ideas are crucial in modern applications, where we sometimes have few samples in high dimension.

- We saw how the risk could generally be decomposed as a term of bias/approximation and a term of variance/estimation.
- This decomposition highlights the tradeoff that needs to be dealt with in inference. This tradeoff is related to the complexity of the set of functions under consideration:
  - Sets too simple lead to a large approximation error.
  - Sets too large lead to a large estimation error.
- We defined this notion of complexity more precisely (Rademacher, VC), and saw it also depended on the number of samples.
- These ideas are crucial in modern applications, where we sometimes have few samples in high dimension.

- We saw how the risk could generally be decomposed as a term of bias/approximation and a term of variance/estimation.
- This decomposition highlights the tradeoff that needs to be dealt with in inference. This tradeoff is related to the complexity of the set of functions under consideration:
  - Sets too simple lead to a large approximation error.
  - Sets too large lead to a large estimation error.
- We defined this notion of complexity more precisely (Rademacher, VC), and saw it also depended on the number of samples.
- These ideas are crucial in modern applications, where we sometimes have few samples in high dimension.

- We saw how the risk could generally be decomposed as a term of bias/approximation and a term of variance/estimation.
- This decomposition highlights the tradeoff that needs to be dealt with in inference. This tradeoff is related to the complexity of the set of functions under consideration:
  - Sets too simple lead to a large approximation error.
  - Sets too large lead to a large estimation error.
- We defined this notion of complexity more precisely (Rademacher, VC), and saw it also depended on the number of samples.
- These ideas are crucial in modern applications, where we sometimes have few samples in high dimension.