

## Module 2: Statistical learning

---

TRAB268  
15.01.2018

### Topics for statistical learning:

i) aim: inference vs prediction  
          ↑  
          understand  
          interpret

ii) type of problem: regression vs classification

regression: quantitative output

- bus arrival: predict lateness
- power grid: predict consumption

classification: qualitative output

- prison: classify inmate as violate or not violate
- customer: good or bad (paying debt)

iii) type of set-up (data/sim/method)

supervised vs unsupervised learning

↑  
response available  
as in regression and  
classification

↑ want to detect  
unknown patterns in data

## iv) type of method

parametric vs non-parametric

- |                         |                        |
|-------------------------|------------------------|
| - simple to use         | - flexible             |
| - but constrained       | - can overfit          |
| - computationally cheap | ← need more data       |
|                         | ← comp. more expensive |

## v) over- and underfitting

Ex: truth 2nd poly (black - end fitted = orange)  
↓  
see figure below  
1st poly (red)  
2d poly (purple)

a) best: orange - purple - red  
2 2d 1 ← students

b) red: not hitting the "correct" curve, but less variable than 2d

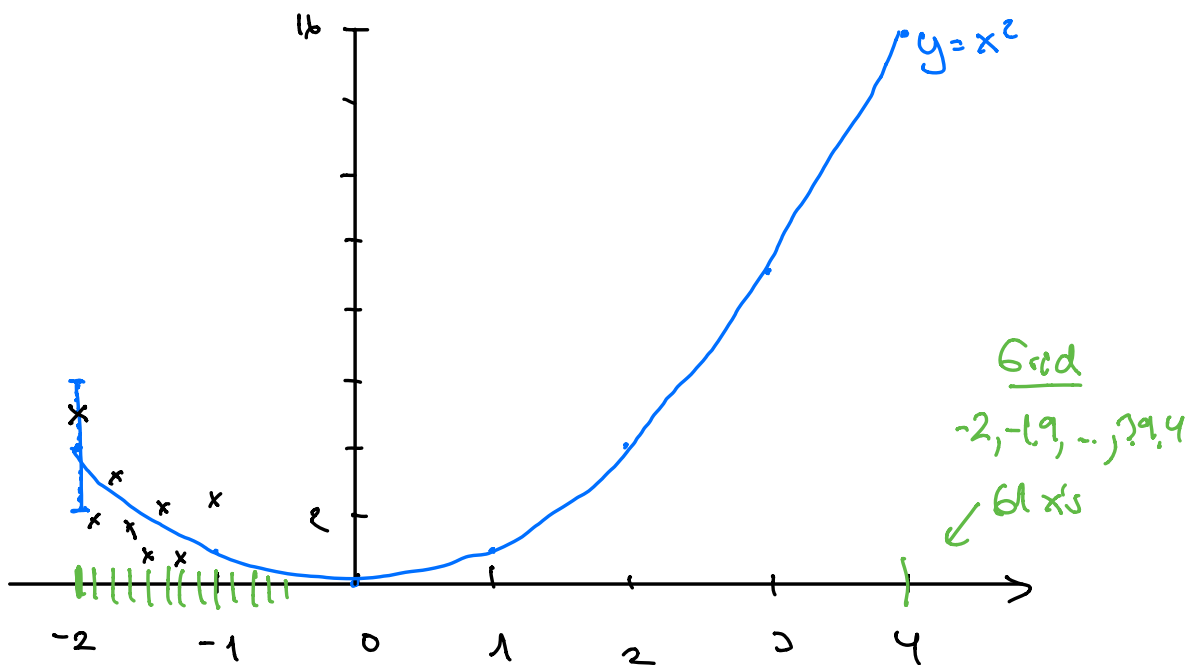
2d  
purple: on average hit correct, but rather variable

Training MSE vs Test MSE

if we choose the model with minimum on training MSE this may lead to high test MSE if we have used a too flexible method.

Polynomial example:

$$Y = \underset{\hat{f}(x)}{x^2} + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$



our choice of parametric model

- fit:
- poly 1:  $f(x) = \beta_0 + \beta_1 x$
  - 2:  $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$
  - $\vdots$
  - 20:  $f(x) = \beta_0 + \dots + \beta_{20} x^{20}$

First: generate one data set

→ fit model

→ make prediction at each  $x$  (61)

→ plot predictions (line) ← compare with observed obs.

→ generate one test set → calc predictions & compare observation

## Bias-variance trade-off

$$Y = f(x) + \varepsilon \quad \text{where } E(\varepsilon) = 0$$

Training data  $(x_i, y_i)$ , independent pairs,  
used to fit model and give  $\hat{f}$ .  
↑  
given model

Now, want to predict new observation at  $x_0$ .  
Use  $\hat{f}(x_0)$  as estimator  
↑  
function of  $(x_i, y_i)$   
↑  
RV's

Quadratic loss at new observation  $Y$  at  $x_0$

$$(Y - \hat{f}(x_0))^2$$

We want the "long term average" = the  
expected value ← over  $Y$

$$\begin{aligned} E[(Y - \hat{f}(x_0))^2] &= E \left\{ Y^2 - 2Y\hat{f}(x_0) + \hat{f}(x_0)^2 \right\} \\ &= \underline{E(Y^2)} - 2E(Y) \cdot E(\hat{f}(x_0)) + \underline{E(\hat{f}(x_0)^2)} \end{aligned}$$

↑                    ↑  
new  $Y$        based  
                  on  
                   $(y_1, \dots, y_n)$

$$\left\{ \begin{array}{l} \text{Remember: } \text{Var}(Y) = E(Y^2) - E(Y)^2 \\ \Leftrightarrow E(Y^2) = \text{Var}(Y) + E(Y)^2 \\ \text{Same for } \hat{f}(x_0): E(\hat{f}(x_0)^2) = \text{Var}(\hat{f}(x_0)) + E(\hat{f}(x_0))^2 \end{array} \right.$$

$$= \text{Var}(Y) + E(Y)^2 + \text{Var}(\hat{f}(x_0)) + E(\hat{f}(x_0))^2 - 2E(Y) \cdot E(\hat{f}(x_0))$$

$$Y = \underbrace{f(x_0)}_{\text{not random}} + \underbrace{\varepsilon}_{\text{since } E(\varepsilon) = 0}$$

$$= \text{Var}(\varepsilon) + \underbrace{f(x_0)^2} + \text{Var}(\hat{f}(x_0)) + \underbrace{E(\hat{f}(x_0))^2} - \underbrace{2f(x_0) \cdot E(\hat{f}(x_0))}^2$$

$$= \text{Var}(\varepsilon) + \text{Var}(\hat{f}(x_0)) + \underbrace{(E(\hat{f}(x_0)) - f(x_0))^2}$$

irreducible error

Variance of prediction

expected value of  $\hat{f}(x_0)$  true value at  $x_0$

Squared bias of prediction

This is at one  $x_0$ , we might average over all  $x_0$ 's.