# 4/ Cross Validation

$\longrightarrow$ most robust method that can be used to select a
model / estimator

$\longrightarrow$ computationnally intensive

n iid observations $X = \left( X_1, \ldots, X_n \right)$
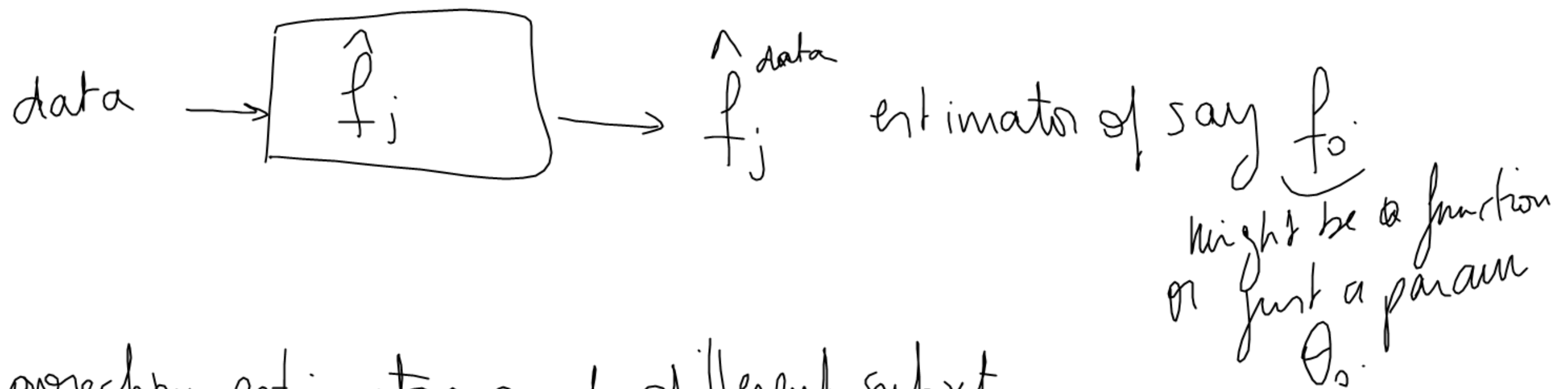
you split it into two sets

$\longrightarrow$ a learning sample / estimation sample / training sample

$\longrightarrow$ a validation sample / transfer sample / testing sample

you also need x a contrast

* d different estimators that may depend
on d different models but not necessary

estimator $\longrightarrow$ black box

$$data \longrightarrow \boxed{\hat{f}_j} \longrightarrow \hat{f}_j^{\,data} \quad \text{estimator of say } \underline{f_0}$$

might be a function
or just a param
$\theta_0$.

ex : — projection estimators on d different subsets
— density estimators ( kernel with different bandwidth )

a) Hold-out

$\boxed{S_L \text{ learning} \mid S_T \text{ Transfer}}$ ✗

* for each $\hat{f}_j$, you compute them only with $S_L$

$$\longrightarrow \hat{f}_1^{S_L}, \ldots, \hat{f}_d^{S_L}$$

* for each of them, you compute $C\left(\hat{f}_j^{S_L}, S_T\right)$ where $C$ is a contrast designed for our target $f_0$

$$\left(\text{It means that } \mathbb{E}_{f_0}\left(C(f, S_T)\right) \text{ is minimal when } f = f_0\right)$$

the estimator that you choose is given by $\boxed{\hat{j} = \underset{j = 1 \ldots d}{\arg\min}\, C\left(\hat{f}_j^{S_L}, S_T\right)}$

Computed with the whole sample

* $\hat{f}_{\hat{j}}^{\hat{X}}$

ex    $X_1, \ldots, X_n$  iid  with  density $f_0$.

$$\hat{f}_j^X(x) = \frac{1}{n h_j} \sum_{i=1}^{n} K\left(\frac{x - X_i}{h_j}\right)$$

with $h_j$ from $j = 1 \ldots d$
$d$ different bandwidth.

$S_L = X_1, \ldots, X_{n/2}$

$S_T = X_{n/2+1}, \ldots, X_n$

$$\hat{f}_j^{S_L}(x) = \frac{1}{n/2 \, h_j} \sum_{i=1}^{n/2} K\left(\frac{x - X_i}{h_j}\right)$$

$$C\left(\hat{f}_j^{S_L}, S_T\right) = -\frac{2}{n/2} \sum_{i=n/2+1}^{n} \hat{f}_j^{S_L}(X_i) + \int \left(\hat{f}_j^{S_L}(x)\right)^2 dx$$

when $h_j$ is large

when $h_j$ is very small

h in between

$X_i$

$\underline{n \text{ is even}}$

$$\hat{j} = \underset{j=1\dots d}{\text{argmin}} \quad C\left(\hat{f}_j^{S_L}, S_T\right)$$

$$= \underset{j=1,\dots,d}{\text{argmin}} - \frac{2}{n/2} \textcolor{red}{\sum_{i=n/2+1}^{n}} \hat{f}_j^{S_L}(x_i) + \int\left[\hat{f}_j^{S_L}(x)\right]^2 dx$$

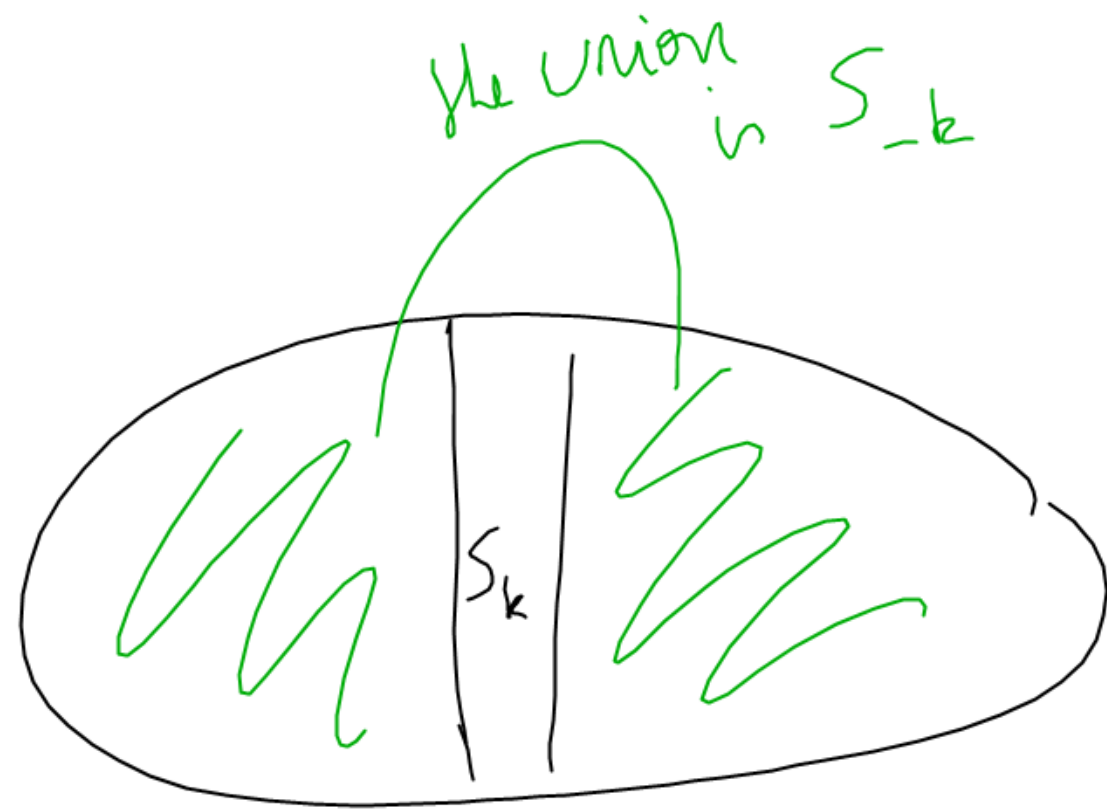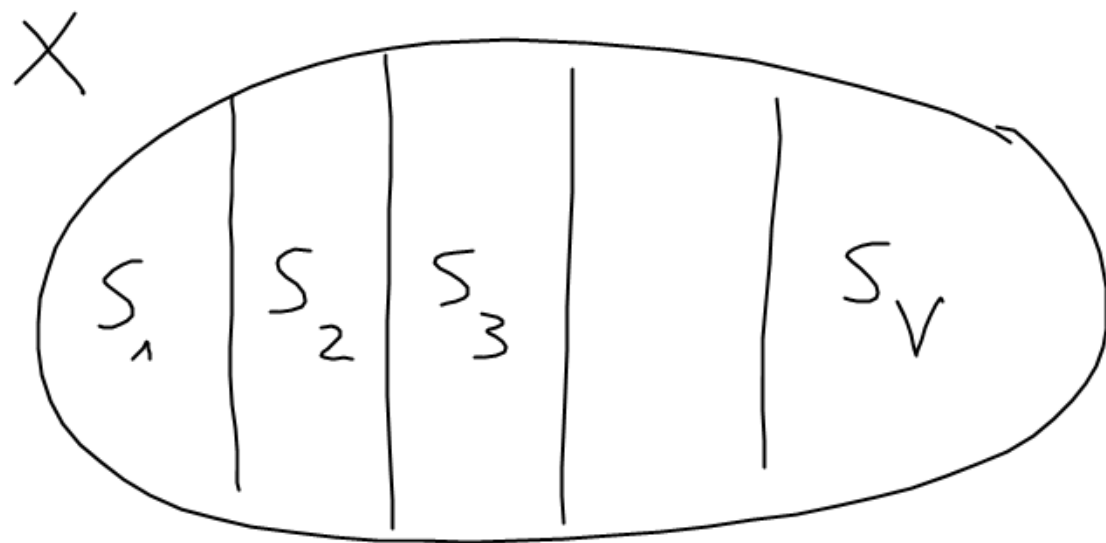<span style="color:red">this is in $S_T$<br>the transfer sample</span>

The good estimator is therefore $\hat{f}_{\hat{j}}$ 

<span style="color:blue">$\hat{\times}$ ⟶ I compute it with the whole sample.</span>

$\hat{\hat{j}}$ ⟶ I selected the good method by Holdout

b) K-fold or V-fold

X



the union in $S_{-k}$

$S_{-k}$ = in the sample when we remove $S_k$ , $k = 1..V$

for each $k$, $\hat{f}_j^{S_{-k}}$ for $j = 1...d$ is computed with $S_{-k}$.

$\quad \hookrightarrow \, C(\hat{f}_j^{S_{-k}}, S_k)$

So far it looks like Holdout
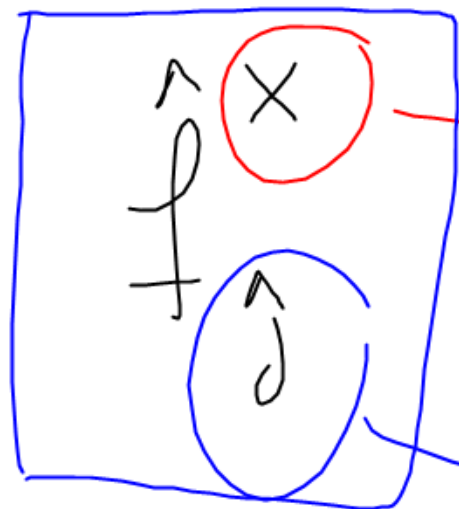$S_L = S_{-k}$ ; $S_T = S_k$

I select

$$\hat{j} = \underset{j = 1 \dots d}{\text{argmin}} \left[ \frac{1}{V} \sum_{k=1}^{V} C\left( \hat{f}_j^{S_{-k}}, S_k \right) \right]$$

score for each method $j$ should give an idea of how close method $j$ is to $f_0$.

here in turn a data is in the learning sample and the transfer sample — average

then I use as an estimator of $f_0$

$$\hat{f}_{\hat{j}}$$

$\hat{f}_{\times}$ → I use the whole sample to compute it.

$\hat{j}$ → the method I selected by V-fold cross-validation

V-fold is "an average" of V Hold-out. It is more stable.

V=5 or 10 are the best choices... (V=n : leave-one-out method)

## other method

you can also use $\quad \dfrac{1}{V} \sum\limits_{k=1}^{V} \hat{\rho}_{\hat{j}}^{S_{-k}}$

But this has problems. especially because depending on the problem computing averages of estimators do not make sense.

# IV. What about testing?

Testing and model selection are different
in the sense that they do not answer to the same question

$$m = 1, \ldots, M \qquad M \text{ different models}$$

model selection $\longrightarrow$ you will always get a model $\hat{m}$
$\hookrightarrow$ this is the one which is the "best"
in the sense of bias/variance equilibrium

$\longrightarrow$ It does not select necessarily the true model!!
But only one which is not too far and
which has a reasonable nb of parameters
given your data.

goodness-of-fit tests:

they are testing $H_0$: my model $m$ is true

vs $H_1$: ———————— false

e.g. Shapiro and Wilk's test of gaussianity

If you know the models well enough, you can compute one for each model. Do not forget to correct for multiplicity with Bonferroni.

the answer will be

→ this $m_1$ is plausible     eg $\mathcal{N}(m, \sigma^2)$

↝ this other $m_2$ is plausible eg $\mathcal{E}(\lambda)$.

↘ this $m_3$ is not likely    e.g. Uniform ...

→ So this may select no model at all or models that are not compatible!!

When you have $\Delta_1, \ldots, \Delta_K$ tests of level $\alpha$ $\longrightarrow$ <span style="color:red">Bonferroni $\alpha/K$</span>

then $P(\exists$ one test which wrongly rejects$)$

$$\leq \sum_{k=1}^{K} P(\Delta_i \text{ wrongly rejects})$$

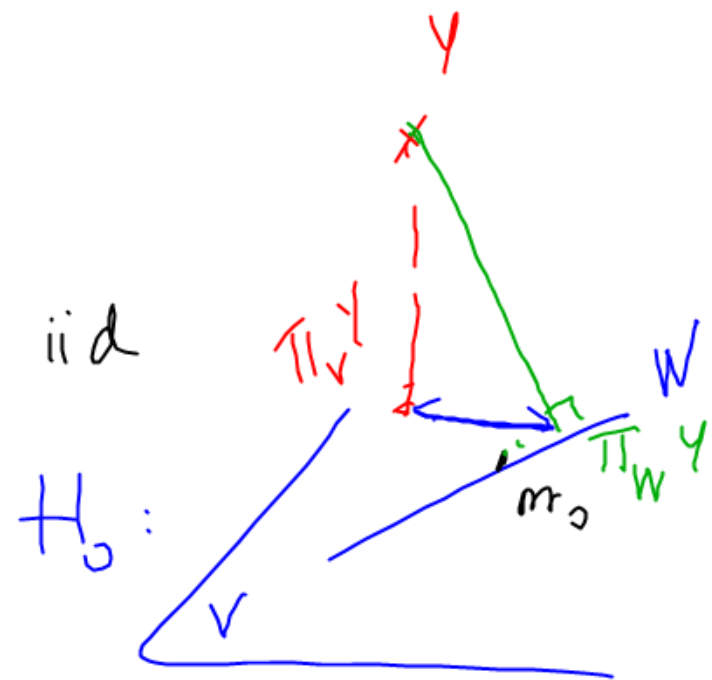$$\leq \sum_{k=1}^{K} \alpha = K\alpha \qquad \longrightarrow \quad \leq \alpha.$$

with $K=15 \longrightarrow$ one chance over 2 to make a mistake

## 1/ Fisher test

In linear gaussian models.

$$Y = m + \varepsilon \qquad \cdot Y \in \mathbb{R}^n, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \varepsilon_i \text{ iid}$$
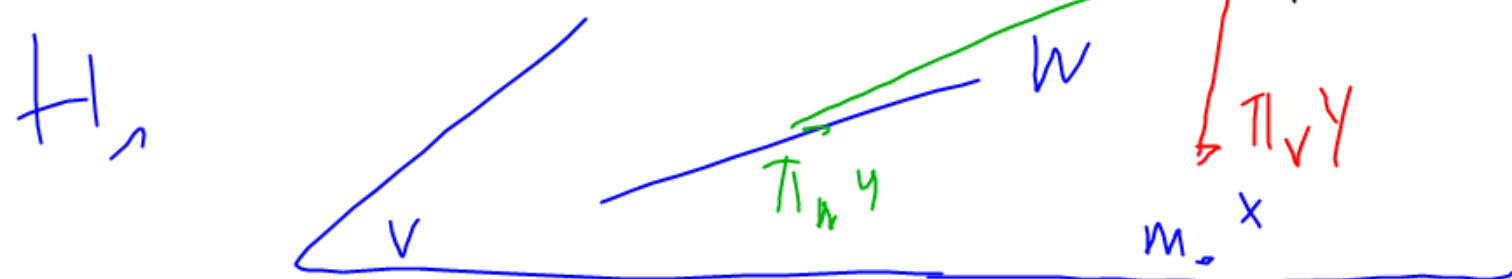
$$\circ \; m \in V \not\subseteq \mathbb{R}^n$$

$H_0:$

$H_0: m \in W \qquad vs \qquad H_1: m \in V \setminus W \qquad \text{where } W \not\subseteq V.$

For instance $\qquad W = Vect \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

$$V = Vect \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} \cos u_1 \\ \vdots \\ \cos u_n \end{pmatrix}, \cdots \begin{pmatrix} \cos(d u_1) \\ \vdots \\ \cos(d u_n) \end{pmatrix}$$

$H_1$

$\Pi_V Y$

$\Pi_\Lambda Y$

$m_\circ$

the Fisher test is based on the statistic

$$T = \frac{\|\Pi_V Y - \Pi_W Y\|^2}{\|Y - \Pi_V Y\|^2} \times \frac{n - \dim V}{\dim V - \dim W}$$

under $H_0$, $T$ obeys a Fisher distribution $\mathcal{F}(\dim V - \dim W, n - \dim V)$

$\longrightarrow$ you reject when $T$ is larger than the corresponding quantile $1 - \alpha$.

$\longrightarrow$ you transform that into p-values.

NB: in R, lm()

pval one per variable

$\longrightarrow$ this is the pvalue of the Fisher test for $W = \mathrm{Vect}\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

(pval)

## 2) Wilk's theorem and likelihood ratio test

You have two models for your data $X = (X_1 \cdots, X_n)$ iid.

model 1 $\qquad\qquad\qquad\qquad\qquad$ model 2

$$\{ P_\theta, \theta \in \Theta_0 \} \quad C \quad \{ P_K, K \in \mathcal{K} \}$$

ex: model 1 is $\{ \mathcal{N}(m, 1), \underset{\theta}{\underline{m}} \text{ is unknown} \}$ $\Theta_0 = \mathbb{R}$.

model 2 is $\{ \mathcal{N}(m, \sigma^2), m \text{ and } \sigma^2 \text{ unknown} \}$. $\mathcal{K} = \mathbb{R} \times \mathbb{R}^+$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad K = (m, \sigma^2)$

eg Transfer model

$$G C M \qquad P\left(\underset{x}{\text{object is put in class}} A\right) = \frac{\sum\limits_{y \in A} S(x, y)}{\sum\limits_{y \text{ are learned}} S(x, y)}$$

$$S(x, y) = \exp\left(-c \underbrace{d(x, y)}_{= \sum\limits_{i=1}^{3} w_i |x_i - y_i|}\right)$$

$i = 1 \hookleftarrow$ color
$i = 2 \hookleftarrow$ shape
$i = 3 \hookleftarrow$ size

model 1

$$\{ P_\theta, \theta \in \Theta_0 \}$$

$$\theta = (c, w_1)$$

$c \in \mathbb{R}_+$

$w_1 \in [0,1]$ $\quad w_2 = 1 - w_1$ $\qquad w_3 = 0$

model 2.

$$\{ P_k, k \in K \}.$$

$$k = (c, w_1, w_2)$$

$c \in \mathbb{R}_+$.

$(w_1, w_2) \in [0,1]^2$ st $w_1 + w_2 < 1$

$w_3 = 1 - w_1 - w_2$ ".

$\longrightarrow$ So the test will here answer to the question " is the size important for the categorization ? "

The likelihood ratio statistic is given by

$$T = \frac{\overbrace{\max_{K \in \mathcal{K}} L_K(X)}^{\text{model 2}}}{\underbrace{\max_{\theta \in \Theta_0} L_\theta(X)}_{\text{model 1}}}$$

$\longrightarrow$ how plausible $X$ is under model 2.

$\longrightarrow$ "   " model 1

if $T$ is large you reject $H_0$: model 1 holds because model 2 seems more plausible

$$T = \frac{\overset{\text{model 2}}{L_{\hat{K}}(X)}}{\underset{\text{model 1}}{L_{\hat{\theta}}(X)}}$$

where $\hat{K}$ is the MLE in model 2

$\hat{\theta}$ —————— 1

equivalently you reject when

$$W = 2 \left( \underbrace{\ell_{\hat{K}}(x)}_{\substack{\text{log likelihood} \\ \text{in model 2}}} - \underbrace{\ell_{\hat{\Theta}}(x)}_{\text{log likelihood in model 1}} \right)$$

Wilk's thm says that under $H_0$, $W \xrightarrow[n \to +\infty]{\text{distrib}} \chi^2(d)$

We reject when $W$ is larger than the corresponding quantile

$\longrightarrow$ pvalues.

$d =$ nb of param in model 2

minus nb of param in model 1

$d = \dim(\text{model 2}) - \dim(\text{model 1})$

In practice, for intricate models, being sure of the nb of parameters can be intricate

ex    you could have parametrized the transfer model with.

$$(c, w_1, w_2) \text{ in model } \underline{1} \quad (\text{forgetting that } w_1 + w_2 = 1)$$

$$(c, w_1, w_2, w_3) \quad \text{---} \quad 2 \quad (\text{------} \quad w_1 + w_2 + w_3 = 1)$$

$$\exp\left(-r\left[w_1 d_1(x,y) + w_2 d_2(x,y) + w_3 d_3(x,y)\right]\right)$$

$$\exp\left(-c_1 d_1(x,y) + c_2 d_2(x,y) + c_3 d_3(x,y)\right)$$

My advise before using this test

$\longrightarrow$ perform simulation

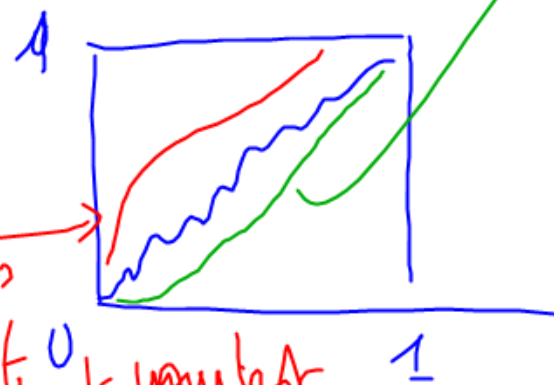$\longrightarrow$ verify that under $H_0$, pvalues are uniform

$\rightsquigarrow$ ecdf should be diagonal
(or under the diagonal) to guarantee
the level of the test.

simulate Nsimu times $(X_1..X_n)$ under model 1. (under $H_0$)

compute each time $W^{(X_1..X_n)} \longrightarrow$ pvalue.

$\longrightarrow$ Nsimu pvalues $\longrightarrow$ ecdf

you have too much
small pval under $H_0$
$\rightarrow$ you cannot trust your test

you are smaller than uniform
$\rightarrow$ you still control the level of the test

(under $H_0$ pval are uniform generally)

eg $X_1 \ldots X_n$ simulate Vas $\overset{iid}{\sim}$ $\mathcal{N}(\overset{\checkmark}{m}, 1)$ model 1

1 param that you don't know

$$W = 2\left( \ell_{\hat{k}}^{model\,2}(X) - \ell_{\hat{\theta}}^{model\,1}(X) \right)$$

$$\hat{\theta} = \hat{m} = \overline{X} \qquad \hat{\sigma}^2$$

2 param unknown

$$\hat{k} = \left( \overline{X}, \frac{1}{n}\sum_{i=1}^{n}(X_i - \overline{X})^2 \right) \qquad \left( model\,2: \mathcal{N}(\widehat{m, \sigma^2}) \right)$$

$$\ell_{\hat{k}}(X) = -\underbrace{\frac{\sum_{i=1}^{n}(X_i-\overline{X})^2}{2\,\hat{\sigma}^2}} - \frac{1}{2}\log(2\pi\hat{\sigma}^2) = -\frac{n}{2} - \frac{1}{2}\log\left(2\frac{\pi}{n}\sum_{i=1}^{n}(X_i-\overline{X})^2\right)$$

$$\ell_{\hat{\theta}}(X) = -\underbrace{\frac{\sum_{i=1}^{n}(X_i-\overline{X})^2}{2\times 1}} - \frac{1}{2}\log(2\pi\times 1) \qquad \left(\text{Here I know } \sigma = 1 \text{ in model 1}\right)$$

# 3/ Bootstrap

## example without bootstrap

$$X \sim \mathcal{E}(\theta_0)$$

I'm interested in the distribution of $\left| \frac{1}{X} - \theta_0 \right|$ where $X \sim \mathcal{E}(\theta_0)$

Nsimu $\quad X_1 \ldots, X_{Nsimu} \sim \mathcal{E}(\theta_0)$

$$T_i = \left| \frac{1}{X_i} - \theta_0 \right| \quad \text{for } i = 1, \ldots, Nsimu.$$

$\longrightarrow$ histograms etc to have an idea of the density

$\longrightarrow$ ecdf $\longrightarrow$ cdf

$\longrightarrow$ empirical quantile $\longrightarrow$ quantile

Why would I need that?

if I observe $X_1 \sim \mathcal{E}(\theta_0)$

then my estimator of $\theta_0$ would be $\frac{1}{X_1} =$

and now you want to know $\longrightarrow$ Confidence interval on $\theta_0$

$\longrightarrow$ make test

$\leadsto$ you need a distribution for how far

is $\hat{\theta}$ from $\theta_0$

$\leadsto$ But you don't know $\theta_0$, so what would you do?

$\longrightarrow$ Bootstrap But you need $n$ observation to do that

a) parametric bootstrap

ex $\qquad X_1, \ldots, X_n \sim \mathcal{E}(\theta_0)$ observations

$\longrightarrow \hat{\theta} = \dfrac{1}{\overline{X}}$

$\longrightarrow$ you would like to know the distribution of $|\hat{\theta} - \theta_0|$ to compute for instance CI.

$\mathcal{E}$ the $1 \cdot \alpha$ quantile of this dist

and my CI would be $[\hat{\theta} \pm \varepsilon]$

$\longrightarrow$ I cannot do that because my dist° depends on $\theta_0$ and I don't know it $\longrightarrow$

$\longrightarrow$ OK so I simulate Nsimu times

typical notation for bootstrap sample

$X^{*}_{1}, ..., X^{*}_{n}$ iid $\mathcal{E}(\hat{\theta})$
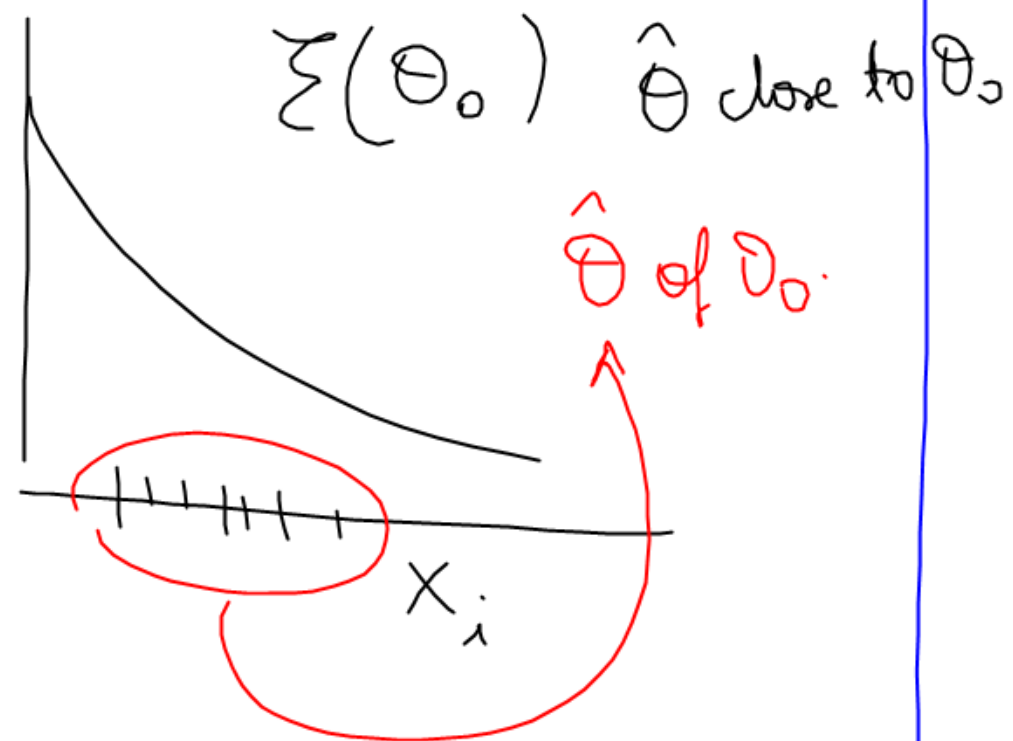
the one that is computed with the original observations

$\longrightarrow$ I compute the bootstrap version of $\hat{\theta}$ : $\hat{\theta}^{*}_{i}$ the $i^{th}$ simulation

$\longrightarrow$ $T^{*}_{i} := \left| \hat{\theta}^{*}_{i} - \hat{\theta} \right|$ is the surrogate of $\left| \hat{\theta} - \theta_{0} \right|$, that you cannot access

$\longrightarrow$ you can use the Nsimu $T^{*}_{i}$ to get cdf, quantiles etc ...

$\longrightarrow$ for instance with the $1-\alpha$ quantile $\xi^{*}$ of this bootstrapped distribution you get the bootstrap confidence interval" $[\hat{\theta} \pm \xi^{*}]$
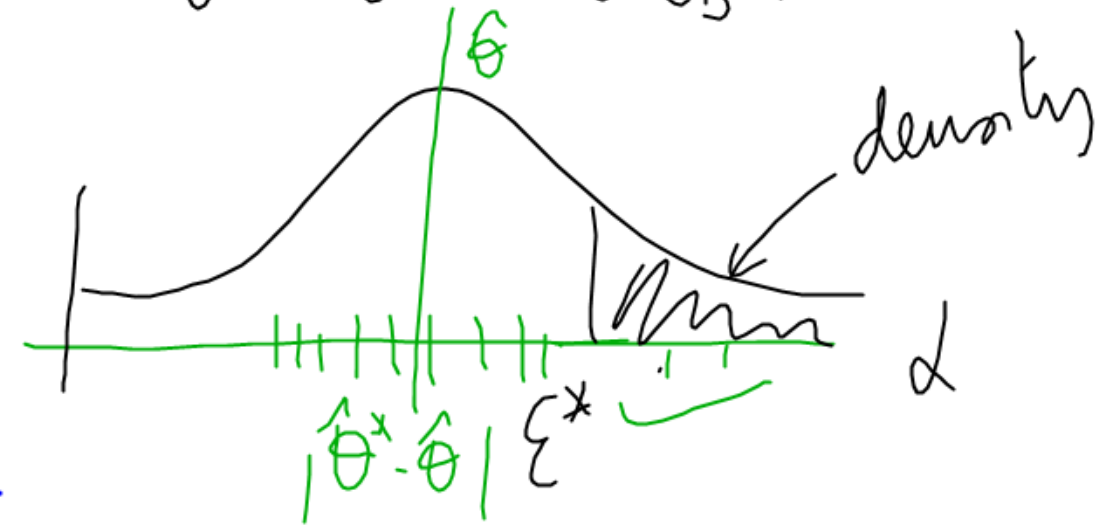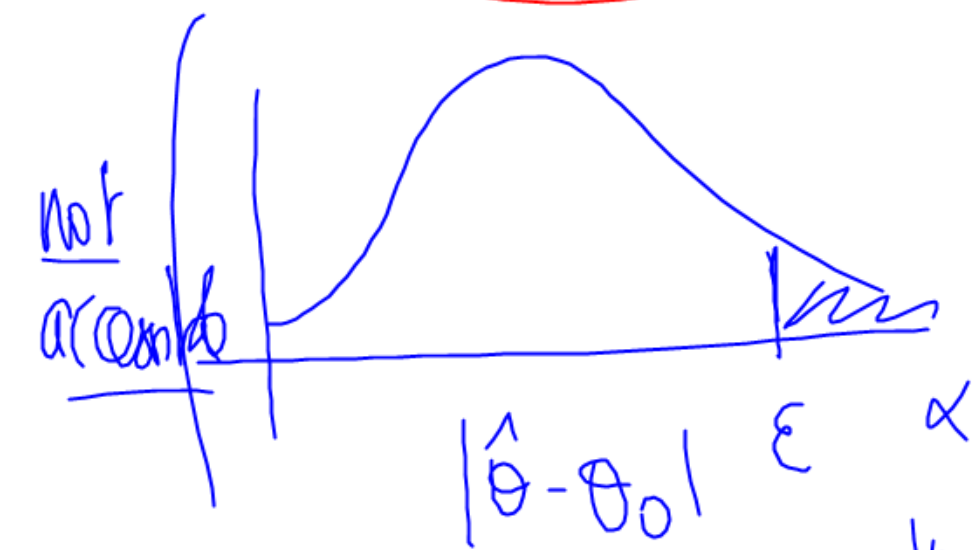
Original world

$\xi(\theta_0)$  $\hat{\theta}$ close to $\theta_0$



$\hat{\theta}$ of $\theta_0$.

$X_i$

not
accept

$|\hat{\theta} - \theta_0| \; \xi \; \alpha$

Bootstrap world

$\xi(\hat{\theta})$

$\times$ Nsimu

$X_i^*$ $\longrightarrow$ $\hat{\theta}^*$

you hope that $\hat{\theta}^*$ is close to $\hat{\theta}$ i the same way as
$\hat{\theta}$ is close to $\theta_0$.

$\longrightarrow$ Nsimu times $\hat{\theta}^*$

But if n is large, $\xi^*$ should be close to $\xi$.

$\hat{\theta}$

density

$|\hat{\theta}^* - \hat{\theta}| \; \xi^* \quad \alpha$

In general,

given a model and data $X_1, ..., X_n$ observed and thought to
with parameter $\theta$
come from this model with unknown parameter $\theta_0$.

$\longrightarrow$ propose an estimator $\hat{\theta}$ of $\theta_0$ ( MLE, least square, empirical mean)

you need $\hat{\theta} \underset{n \to +\infty}{\longrightarrow} \theta_0$ )

$\longrightarrow$ simulate $X_1^*, ..., X_n^*$ with parameter $\hat{\theta}$.

N simu times

$\rightsquigarrow$ compute each time $\hat{\theta}^* - \hat{\theta}$ (and then up to you to do distance

absolute value etc )

$\longrightarrow$ arrive at the empirical distribution of the quantity that you want

Bootstrap

$\rightarrow$ use this empirical bootstrap dist° as if it was the one

of $\hat{\theta} - \theta_0$.

$\left( \text{to build CI, test etc...} \right)$

⚠ DO NOT FORGET the centering

$$- \theta_0 \longrightarrow - \hat{\theta}.$$

Theories exist to show that it works but it combines 2 things. $N_{sim} \rightarrow +\infty$

if $n$ is not big enough, you will pay the fact that $\hat{\theta}$ is far from $\theta_0$. $n \rightarrow +\infty$

## b) non parametric bootstrap

When you don't have a model, you have at least data and you can always pick again in the sample to create new bootstrap sample.

$\longrightarrow$ bootstrap of the mean
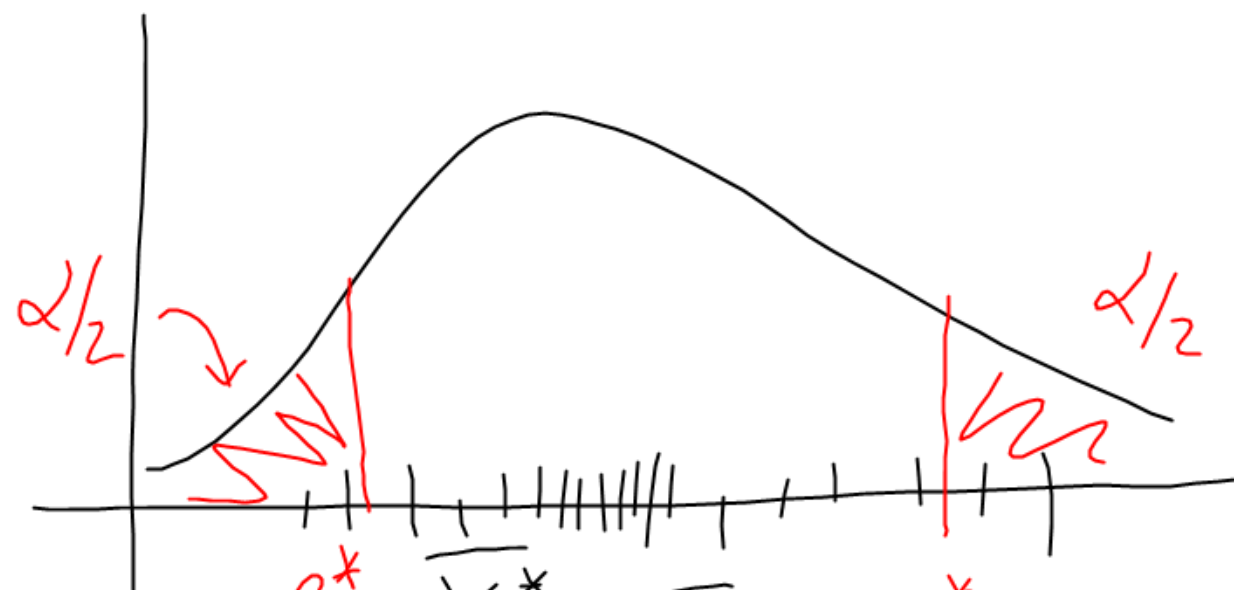
$X_1, \ldots, X_n$ iid with unknown mean $m = E(X)$

$\longrightarrow$ you can estimate $m$ by $\overline{X}$

$\longrightarrow$ N sim times, you pick uniformly at random in $\{X_1, \ldots, X_n\}$ (with replacement)

to get $X_1^*, \ldots, X_n^* \longrightarrow \overline{X^*}$

$\longrightarrow$ you get Nsim $\overline{X}_i^* - \overline{X}$

to approximate the distribution of $\overline{X} - m$.

for instance

in R command
quantile will do
the job.

$\alpha/2$ $\qquad$ $\alpha/2$



$q_{\alpha/2}^*$ $\qquad$ $\overline{X}_i^* - \overline{X}$ $\qquad$ $q_{1-\alpha/2}^* \Leftrightarrow$ the data $W_i^*$
such that a fraction $\alpha/2$
of the data are bigger
than that

the data $W_i^*$
such that a fraction
$\alpha/2$ is smaller than that.

$W_i^*$

bootstrap

$\longrightarrow$ the CI at confidence level $1-\alpha$ is $\left[ \overline{X} + q_{\alpha/2}^* , \overline{X} + q_{1-\alpha/2}^* \right]$

# V Other methods with independance

## 1) Supervised classification

⤳ Deep Learning

## 2) Unsupervised classification / clustering

* k-means    needs k = the nb of clusters

* to estimate k  ⟶ Bic criterion
                 ⟶ hierarchical clustering