

Validation of the model

Model:

$$\forall i \in \{1, \dots, n\}, y_i = \beta_0 + \beta_1 x_i^{(1)} + \dots + \beta_p x_i^{(p)} + \epsilon_i$$

$$Y = X\beta + U$$

assumptions:

• X : a full rank matrix

• $U \sim \mathcal{N}(0, \sigma^2 I_n)$

Two scopes:

Estimation

Validation

Main difficulty: to verify that
 $U \sim \mathcal{N}(0, \sigma^2 I_n)$ because no
observation!

Because no observation of the ϵ_i ,
we are considering the $\hat{\epsilon}_i$. (residual)

RL: $\hat{\epsilon}_i$ are estimations of the ϵ_i .

P_X : matrix associated to the orthogonal
projection on \mathcal{E}_X

$$\begin{matrix} \uparrow & \text{---} & \times & \cdot & \left(\begin{matrix} \times & \times \end{matrix} \right)^{\top} & \text{---} & \times \\ & \text{---} & & & \text{---} & & \\ & \text{---} & & & \text{---} & & \end{matrix}$$

$$E[\varepsilon_i] = 0$$

$$V[u] = \sigma^2 I_n$$

$$E[\hat{\beta}] = 0$$

$$V[\hat{u}] = \sigma^2 (I_n - P)$$

We deduce that $V(\hat{\epsilon}_i)$ depends
on i .

\Rightarrow no homoscedasticity for the
variance of the residual

To suppress this non homogeneity,
we make some normalization

notation: h_{ij} the term at position
(i, j) into the T_x matrix

Let introduce:

$$r_i = \frac{\sum_{j=1}^n h_{ij}}{n}$$

$$\Delta \sqrt{1 - h_{ii}}$$

$$\rightarrow V[r_i] = 1 \quad \forall i \in \{1, \dots, n\}$$

Since σ^2 is unknown, we can not
compute the r_i in practice

→ we prefer the standardized
residuals whose definition is:

$$t_i = \frac{\hat{\epsilon}_i}{\hat{\sigma}_n \sqrt{1 - h_{ii}}}$$

RR:

$$\hat{\sigma}_n^2 = \frac{1}{n - rk} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

$$\rightarrow (n - rk) \frac{\hat{\sigma}_n^2}{\sigma^2} \sim \chi^2(n - rk)$$

$$t_i = \frac{\epsilon_i}{\sigma_i} \sim \mathcal{N}(0, 1)$$

$$\sqrt{\frac{(n-rk) \hat{\sigma}_n^2}{\sigma^2}} / (n-rk)$$

$$\sim \chi^2(n-rk)$$

with $\epsilon_i \sim \mathcal{D}(0, \sigma^2(1-h_{i,i}))$

$$\frac{\epsilon_i}{\sigma \sqrt{1-h_{i,i}}}$$

We do not know the distribution
of the t_i because we do
not have the independance between
 $\hat{\beta}_i$ and $\hat{\sigma}_n^2$

This explains why we introduce
the studentized residuals:

$$t_i^* = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}}$$

where $\hat{\sigma}_{(i)}$ is an estimation
of σ^2 without using the
individual i .

In practice, how to construct $\hat{\sigma}_{(i)}$?

$\underline{r}_{(i)}$: The vector \underline{r} where

we delete the row i .

$\mathbb{X}_{(i)}$: The matrix \mathbb{X} where

we delete the row i .

$U_{(i)}$: under U where we delete the row i .

$$\Rightarrow Y_{(i)} = X_{(i)} \beta + \epsilon_{(i)}$$

We are able to estimate β and σ^2
with these $(n-1)$ observations

$$\Rightarrow \hat{\beta}_{(i)} = \left(\begin{matrix} 1 & X_{(i)} & X_{(i)} \end{matrix} \right)^{-1} \begin{matrix} 1 \\ X_{(i)} \\ Y_{(i)} \end{matrix}$$

$$\hat{Y}_{-(i)} = \hat{X}_{(i)} \hat{\beta}_{(i)}$$

$$\hat{\sigma}_{(i)}^2 = \frac{1}{n - rk \hat{X}_{(i)}} \times \left(\hat{Y}_{-(i)} - \hat{X}_{-(i)} \hat{\beta}_{(i)} \right)^T \left(\hat{Y}_{-(i)} - \hat{X}_{-(i)} \hat{\beta}_{(i)} \right)$$

$$\hat{\epsilon}_i = y_i - \hat{y}_i$$

with $\hat{y}_i = (1 \ x_i^{(1)} \ \dots \ x_i^{(p)}) \hat{\beta}_{(1)}$

RL: Thus, $\hat{\epsilon}_i$ can be viewed as

The prediction error

$$l_i^* = \frac{y_i - \hat{y}_i}{\hat{\sigma}_i \sqrt{1 + \hat{\gamma}_i (x_i^{\text{tr}})^T x_i}}$$

$$\sim T(n - k - 1)$$


$$x_i = \begin{pmatrix} 1 & x_i^{(1)} & \dots & x_i^{(p)} \end{pmatrix}$$

R_k : The computation of $\hat{\sigma}_{(i)}$

↳ in some sense a cross-validation method.

What is cross-validation?

Q: Learning sample

We split the learning sample into 

fold

integer
between 1 and n

→ we create V folds denoted

$\mathcal{L}_1, \dots, \mathcal{L}_V$

1st step:

we take $\mathcal{L}_1 \cup \mathcal{L}_2 \cup \dots \cup \mathcal{L}_V$ as the
new learning sample

→ we build our model thanks to
the new learning sample
for example: for linear model
B.

We take Z_i which is for this
step a test sample.

→ $\hat{\varepsilon}_i = \hat{y}_i - \bar{y}_i$ for the
individuals that are into Z_i

Step 2: we take $\mathcal{L}_1 \cup \mathcal{L}_3 \cup \dots \cup \mathcal{L}_V$

as the new learning sample

→ we build a new model

for example in linear model: $\hat{\beta}_2$

→ \mathcal{S}_2 is a test sample for this model

→ $\hat{\epsilon}_i = \hat{y}_i - y_i$ for all the individuals into \mathcal{S}_2

We do the same for step 3 until step V.

→ cross-validation error:

$$:= \frac{1}{V} \sum_{i=1}^V \hat{\epsilon}_i^2$$

This error is a way to evaluate
the performance of the model
constructed. Thank to the whole
learning sample!

Word:

1) Simulated observation

like this:

$$n = 1000$$

$$p = \frac{1}{2}$$

$$X^{(1)} \sim \mathcal{E}(0.2)$$

$$X^{(2)} \sim \mathcal{U}([-2, 2])$$

$$X^{(3)} \sim \text{T}(7)$$

$$X^{(4)} \sim \text{B}(10, 0.5)$$

$$Y = 3 - 2 \times X^{(1)} + 3 \times X^{(2)} \\ + X^{(3)} - 5 \times X^{(4)} + U$$

with $U \sim \mathcal{D}(0, 1)$

2) Construct the estimation

of the linear model

3) Evaluate the performance
of your model with cross-validation.

