

Regressions Linéaire

Contexte: Relation entre

- une variable unique target y

- une ou plusieurs explicatives
plurielles

$$X_1, \dots, X_p$$

Hypothèse d'une relation linéaire .
Caractéristiques de cette rel. Validité?

Relation linéaire entre y et x_1, \dots, x_p :

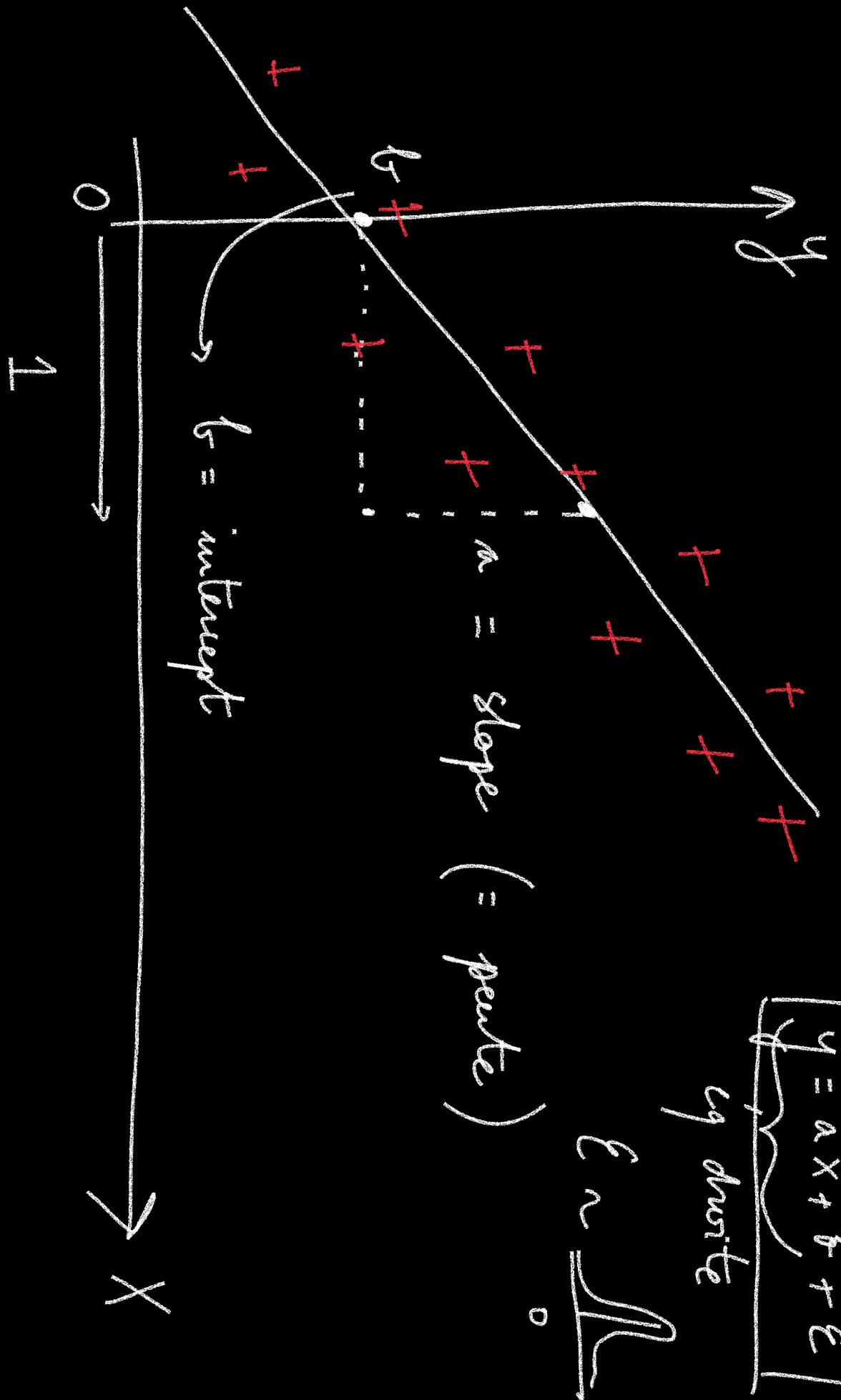
$$y = d_1 \cdot x_1 + d_2 \cdot x_2 + \dots + d_p \cdot x_p + d_{p+1} + \varepsilon$$

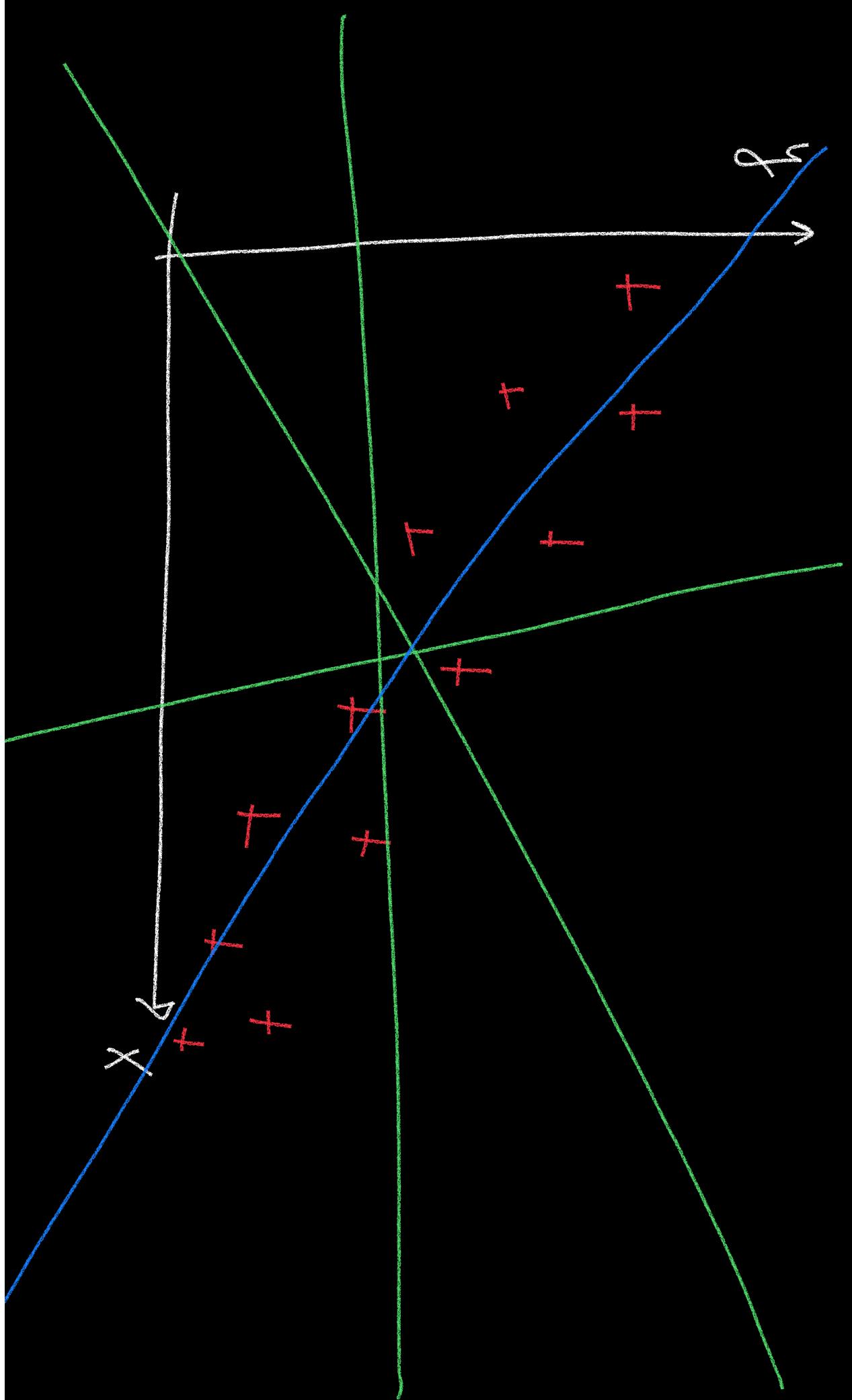
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y = d_1 x_1 + d_2 + \varepsilon$$

suit une loi:

équation d'une droite





Régression Linéaire avec Moindres Carrés

chercher la droite qui minimise

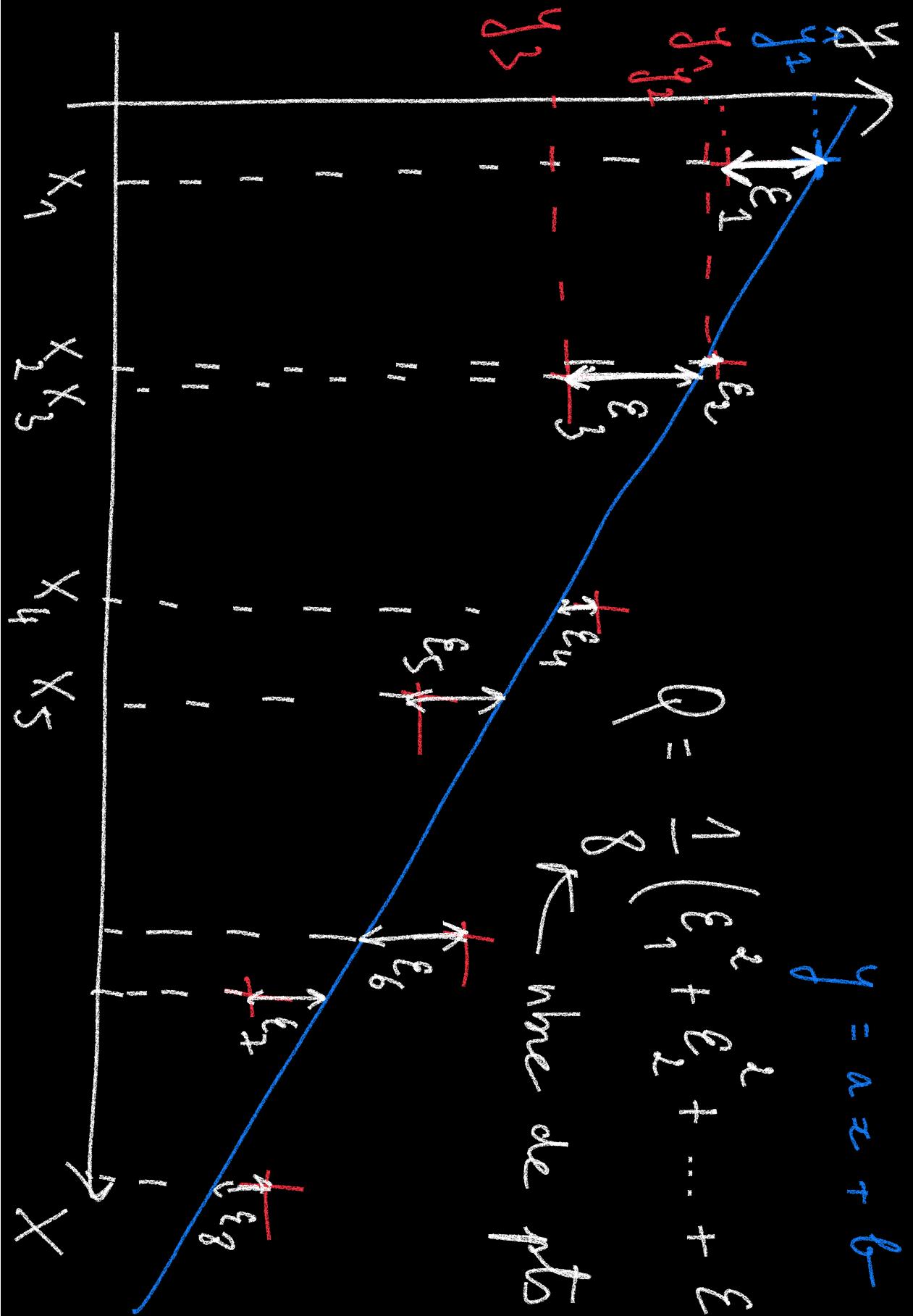
$$\frac{1}{N} \sum_{i=1}^N \text{erreurs}^2$$

nbre de pts

$$\text{erreur}[i] = \underbrace{y[i] - (\hat{y}[i])}_{\text{erreur}}$$

où $\hat{y}[i] = (a X[i] + b)$

prediction



$$Q = \frac{1}{8} (e_1^2 + e_2^2 + \dots + e_8^2)$$

where pts

Mean

Square

=

$$\frac{1}{N} \sum_{i=1}^N$$

$$(y_{[i]} - \hat{y}_{[i]})^2 \geq 0$$

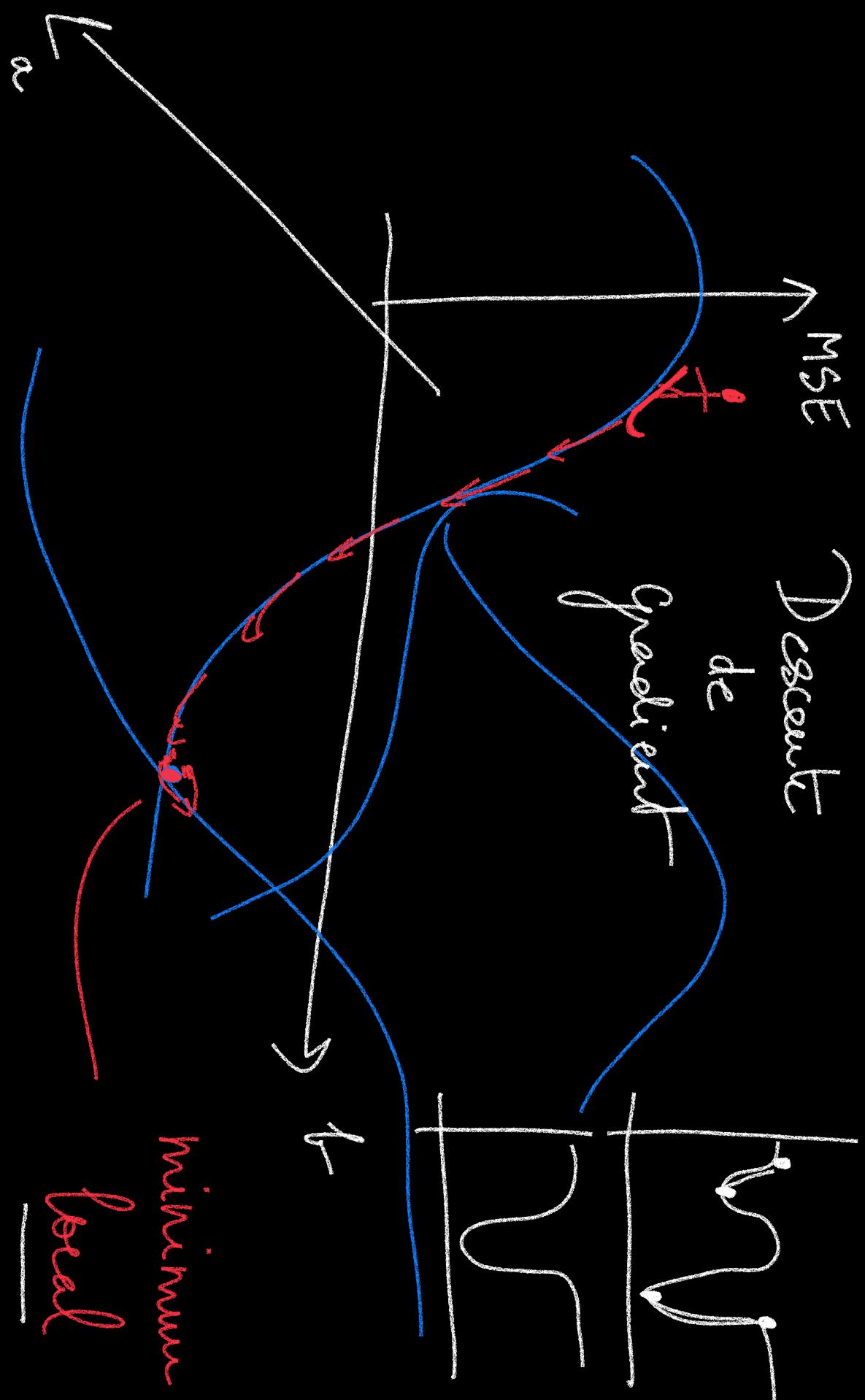
on

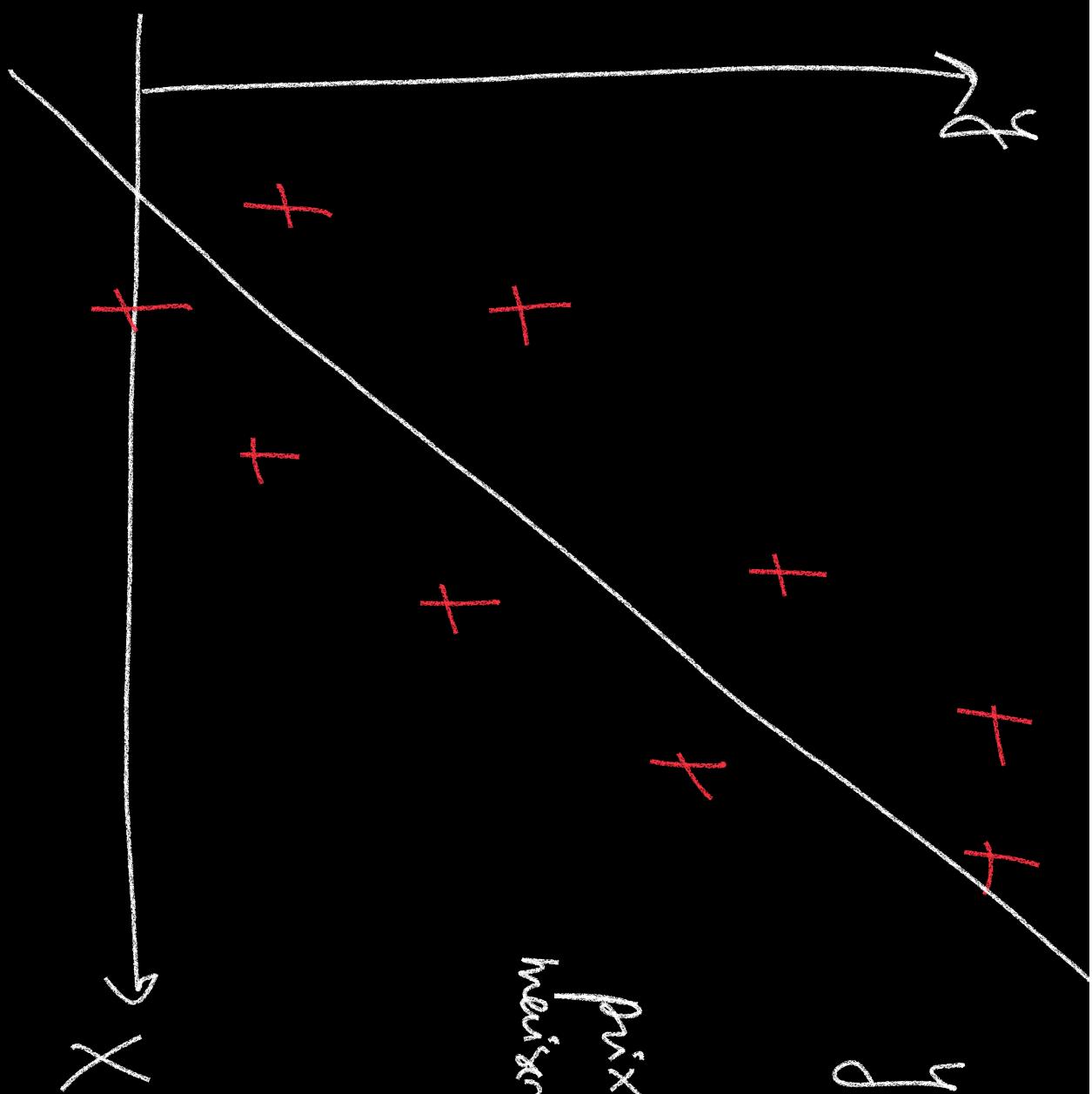
\hat{y} = prediction

y = real

$$\text{Reg Lin w/o Noise term: } \underset{a, b}{\operatorname{arg\,min}} \frac{1}{N} \sum_{i=1}^N (y_{[i]} - (ax_{[i]} + b))^2$$

Reg Lin w/ Noise term:





$y = ax + b$

fix = $13,7 \cdot$ surface
mechanism

$\uparrow m^2 \Rightarrow \uparrow 13,7$
ht

$$y = d_1 x_1 + \dots + d_p x_p + d_{p+1}$$

mix = 13.7. surface

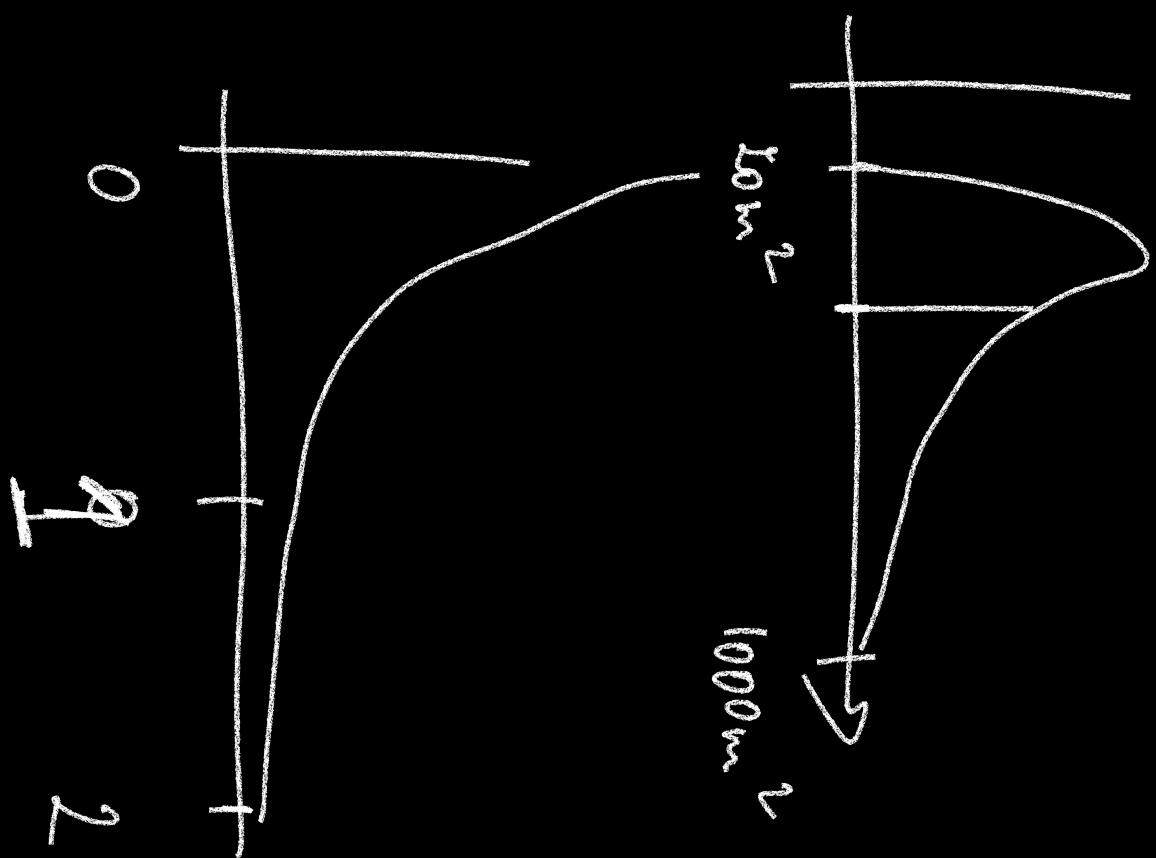
+ 3.8. # picks

- 2.8. # desibels

+ ...

Surfaces

surfaces



'cent - type

$\approx 200 \text{ m}^2$

'ext - type

$\approx 0,15$

l

l

l

10000

l

10000000

l

l

prix
maxim

3700 E

→ a chose
with

37

+

(d)

cm

→ want

2400

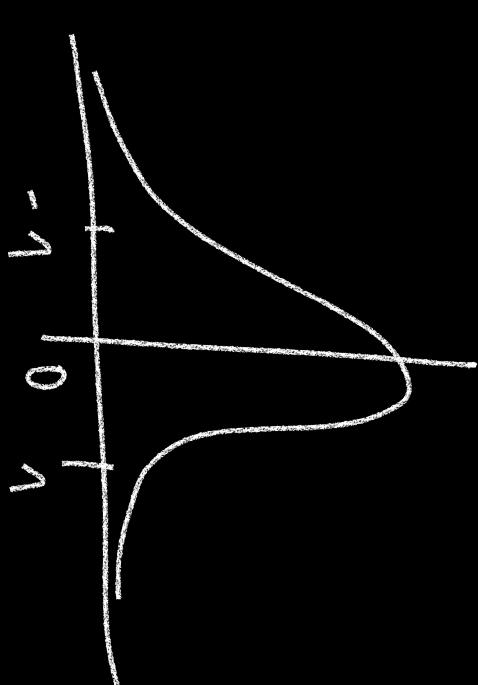
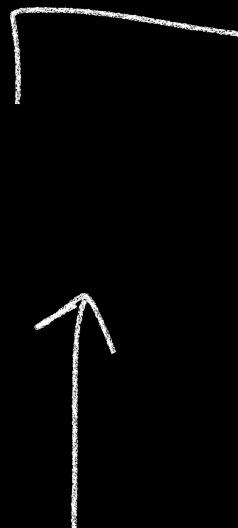
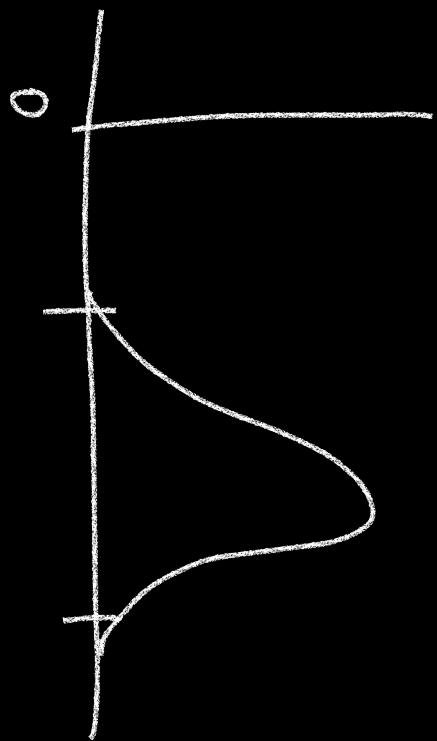
(L)

m

+

...

Standardisation



$$L = \mu$$

moyenne 0
écart-type 1

l (cm)

L (m)

$$\text{pix} = 37 \cdot l \\ + 1400 \cdot L$$

3000

42

$$L \uparrow 1 \Rightarrow \text{pix} \uparrow 37$$

10000

53

$$\text{pix} \uparrow 2400$$

100000

eff $L >$ eff l

$\frac{l}{L}$	9
42500	17
39153	14
34	

57000

24

30000 - 46500
33150

100000 - 42500
57500

~~75 - 300 000~~

57000 - 000t5

$$\frac{42}{45} = \frac{14}{15}$$

$$\begin{array}{r} 63 \\ \underline{-} 34 \\ 34 \end{array}$$

$$\frac{m_1}{m_2} = \frac{17}{34}$$

25 - 35

-

g k
o t
o t

10

27

四百四十一

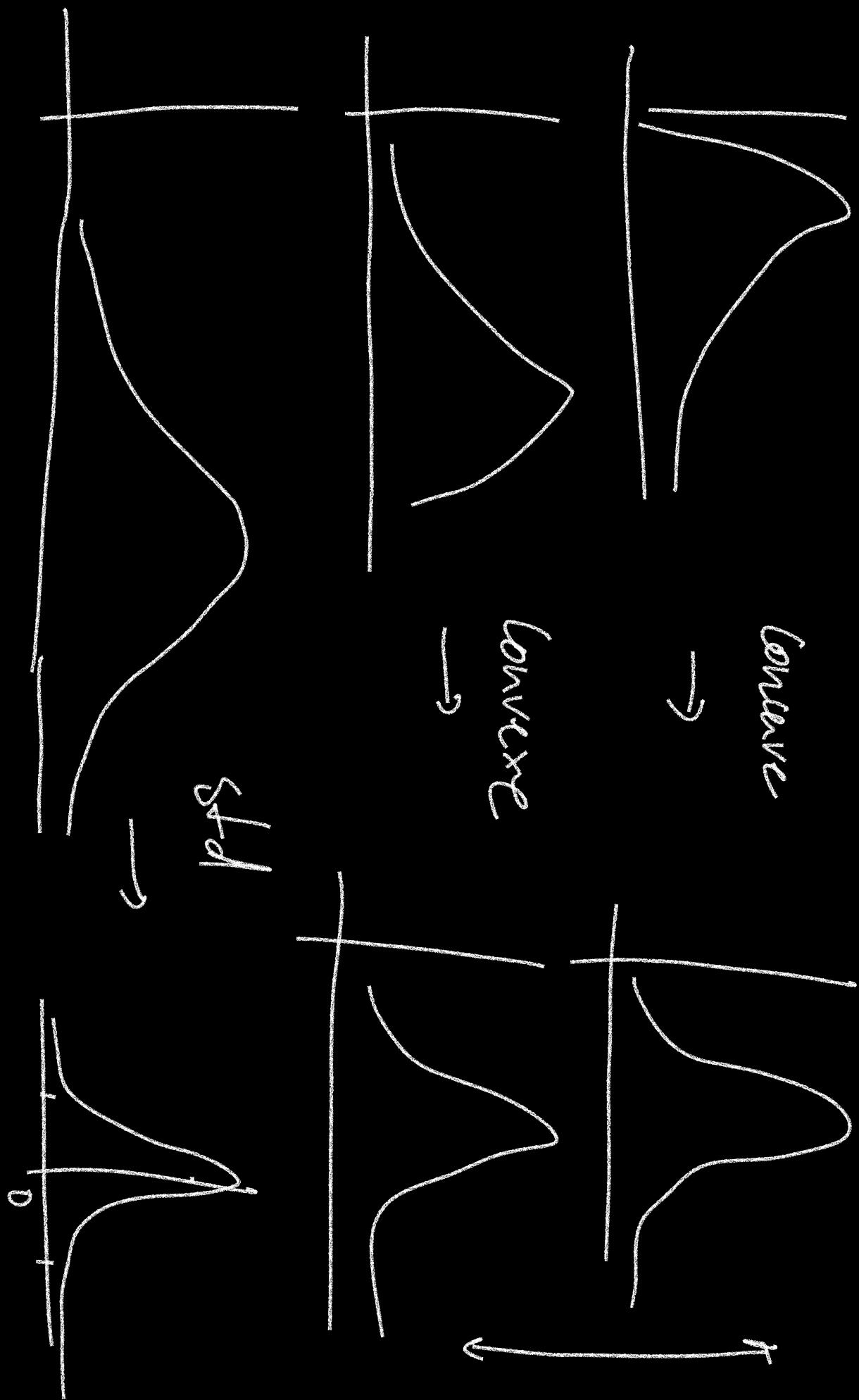
۱۰۵

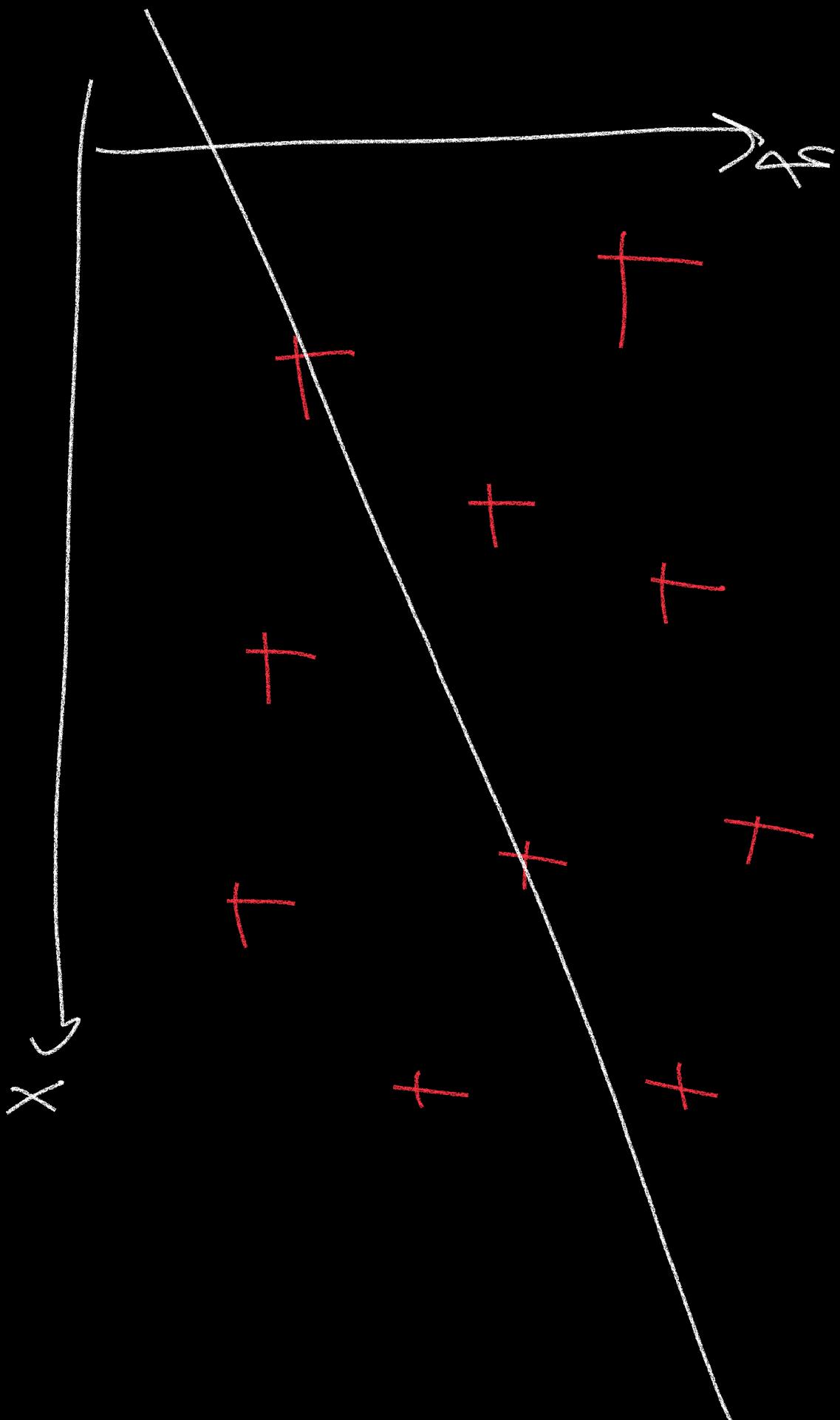
18

$$p_{\text{mix}} = 3,7 \cdot t + 2,4 \cdot L$$

effet L > effet T

régression linéaire appliquée
aux colonnes standardisées





$$\text{Mean Square Error} = \frac{1}{n} \sum \text{errors}^2$$

$$(\text{MS E}) = \bar{\epsilon}^2$$

$$\text{Root Mean Square Error} = \sqrt{\text{MS E}}$$

$$(\text{RMSE}) = \bar{\epsilon}$$

RMSE

$$\text{RMSE} = 1500000 \in \epsilon [0, 100]$$

RSE \Rightarrow mix

Coefficient de

= R^2

Determination

$$= \frac{1}{N}$$

$$= \sum_{i=1}^N \text{erreurs}^2$$

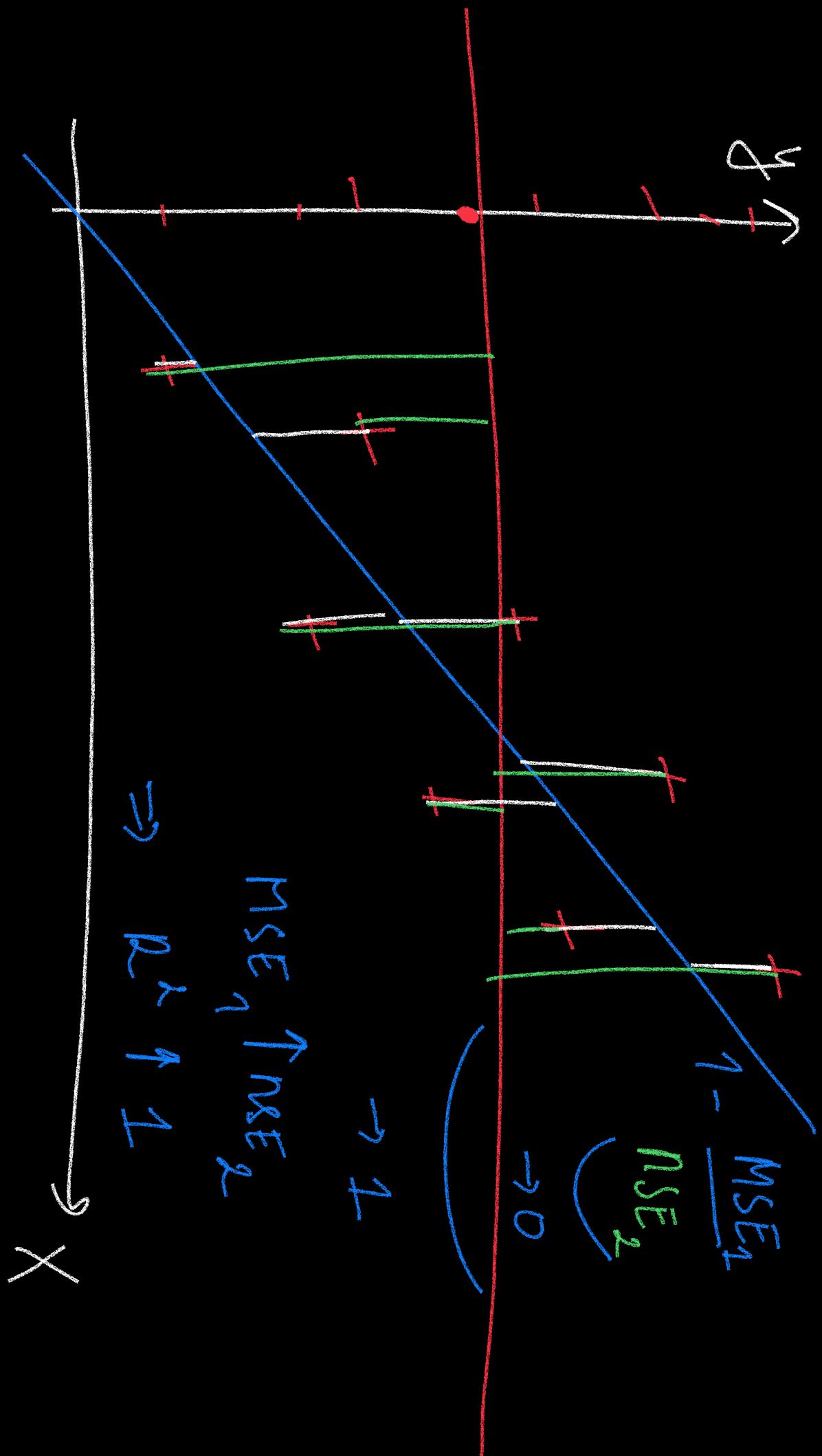
$$\frac{\sum_{i=1}^N (y^{[i]} - \bar{y}^{[i]})^2}{\sum_{i=1}^N (y^{[i]})^2}$$

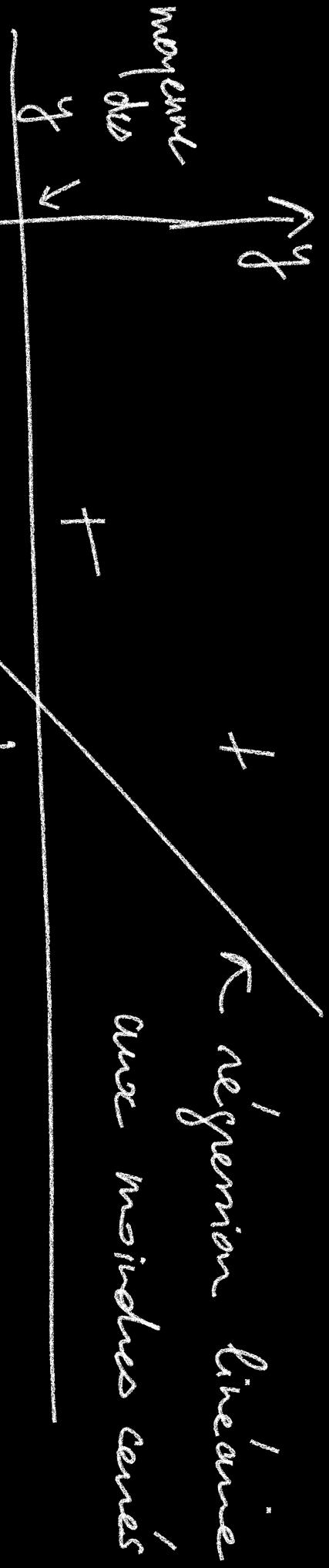
Variance
de \hat{y}
de X

y

$$R^2 \in [-\infty, 1]$$

$R^2 \uparrow \Rightarrow$ meilleure
régression
linéaire





$$\begin{aligned}
 & R^2 \rightarrow 1 \\
 & \text{NEU} \rightarrow 0 \\
 & \text{NEN} \rightarrow 0
 \end{aligned}$$

R predict
 bsp mince
 que la moyenne

Coefficient

$$= \left(\text{Corr} \text{ Pearson} \right)^2$$

Dcf

"

le taux de variance de la variable cible
expliquée par

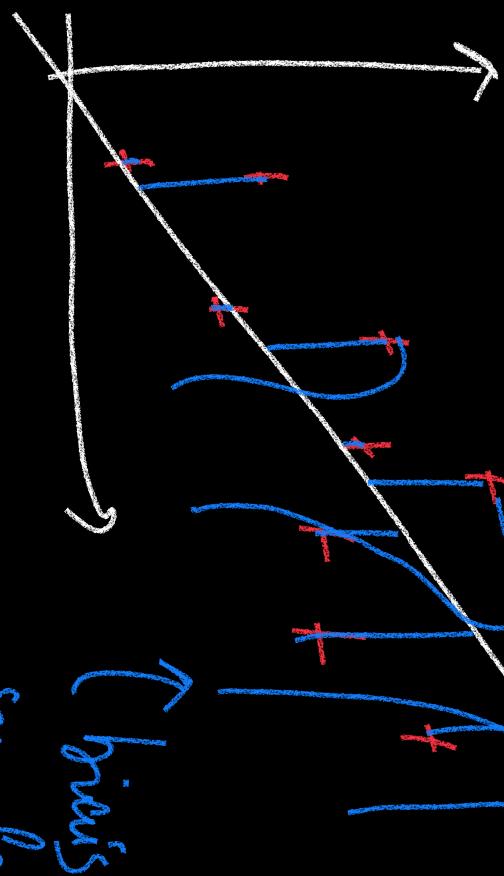
le taux de variance de la variable
explicatrice .

①

ϵ

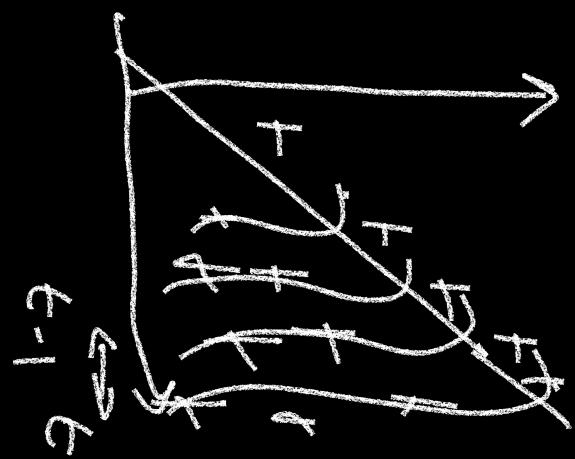
ϵ

$$y = \omega x + b + \epsilon \sim \text{gaussian}$$



bias
sin la
droite

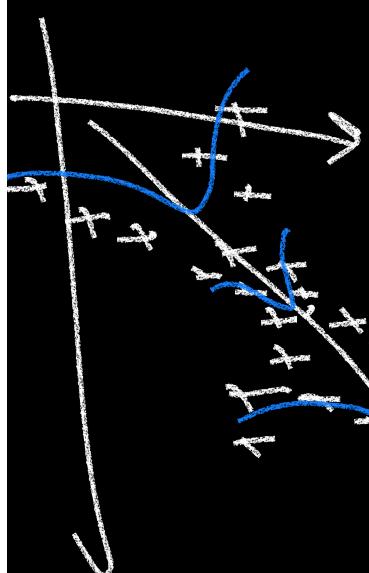
\rightarrow bmit non gaussian



②

bmit independant

③ Homoscedastic



$y = x_1, \dots, x_p$
Variables collinear

x^c - On vert des

Variables explicatives

les redondantes

x^o - possibles.