

# Relation entre variables

## Analyse de régression

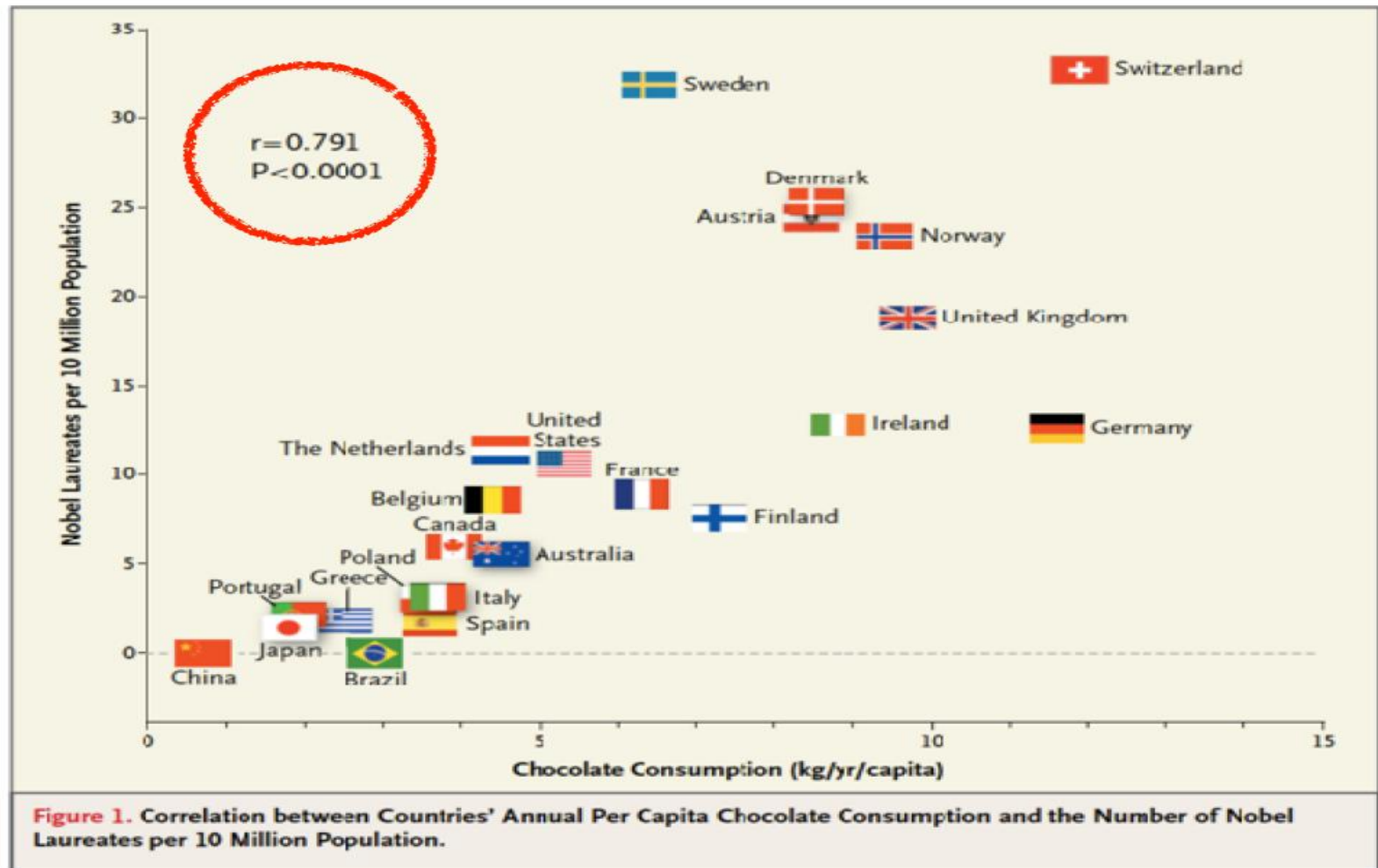
variables avec lien de causalité

ex. variable indépendante (x) et dépendante (y)

## Analyse de corrélation

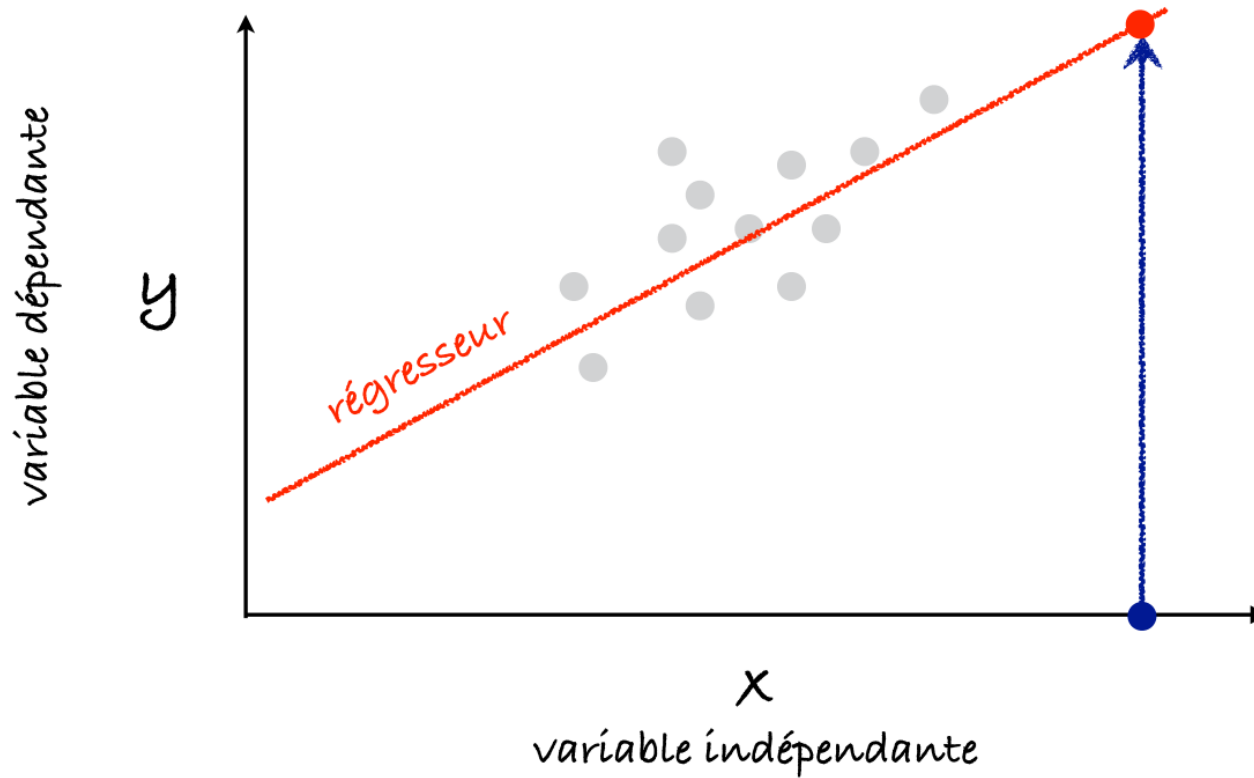
variable interdépendantes, sans hiérarchie

# « Correlation is not causation »



# Régression

En deux mots



# Modèles linéaires

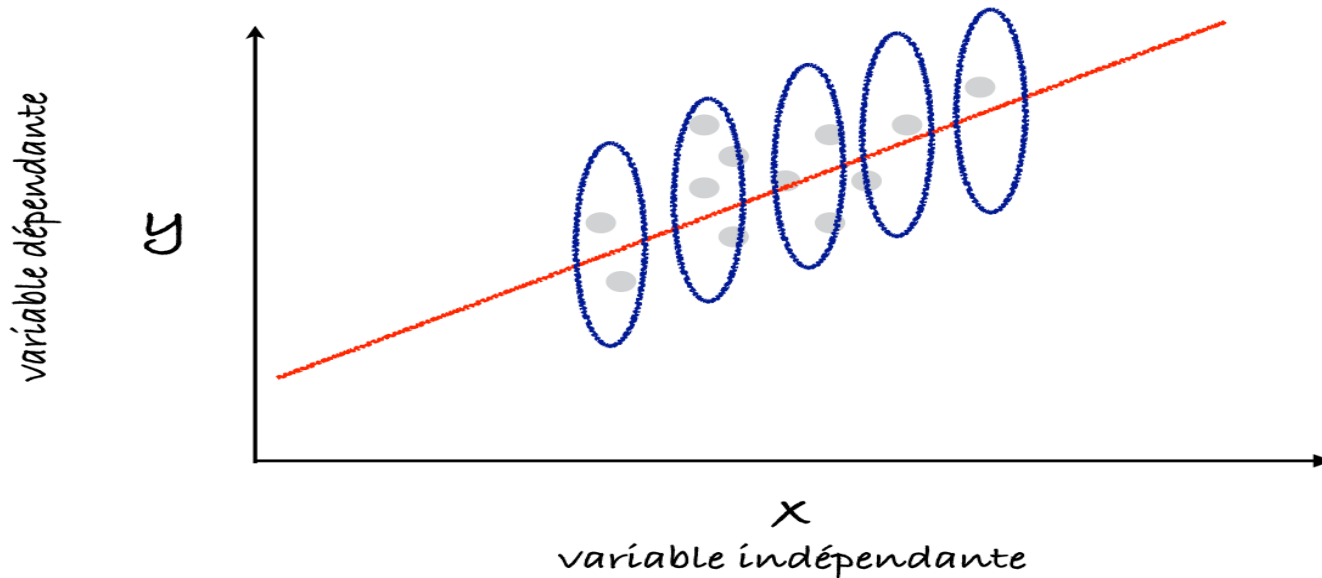
## Régresseur

ligne théorique  $y = a + b_1x_1 + \cdots + b_nx_n$

## Modèle

observations  $y_i = a + b_1x_{1,i} + \cdots + b_nx_{n,i} + \epsilon_i$

# Méthode des moindres carrés

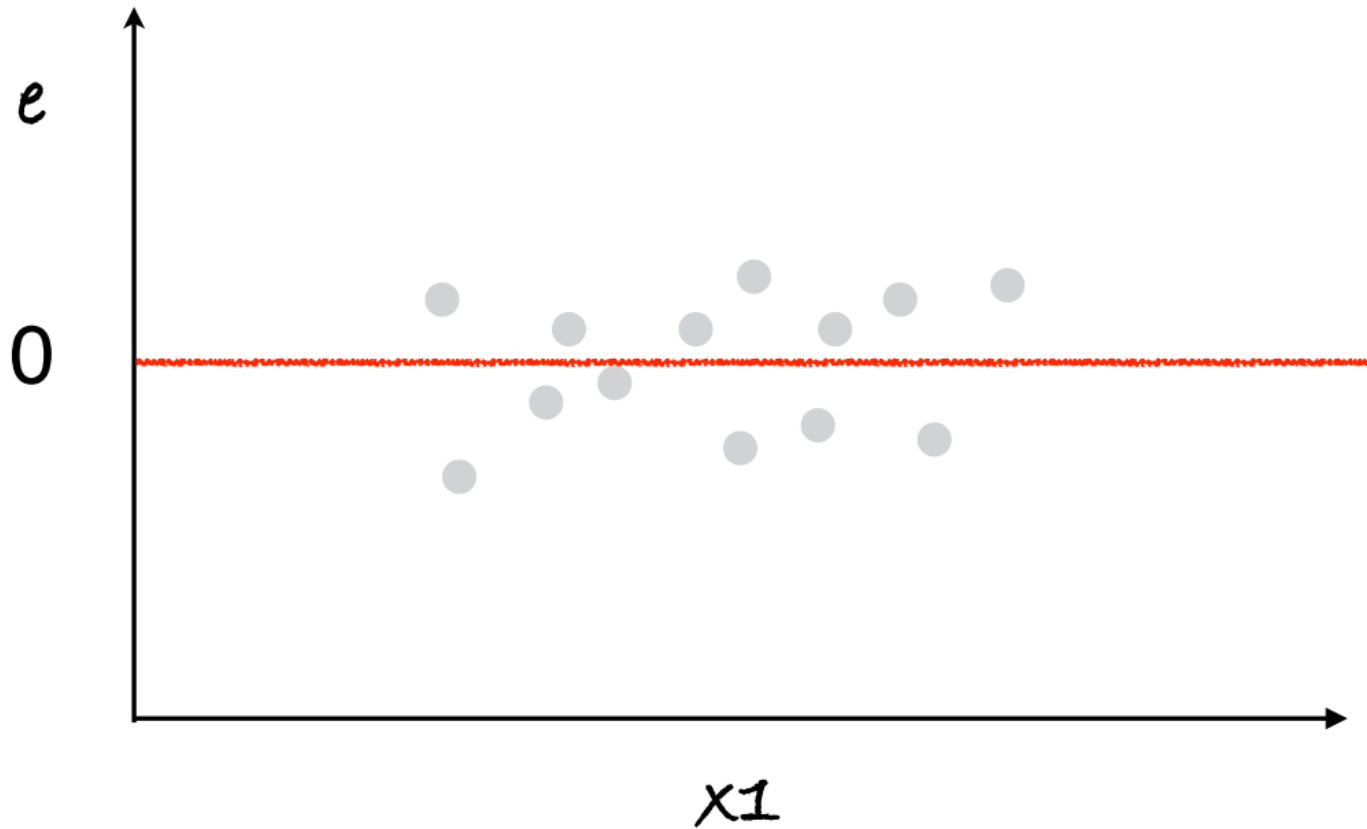


minimiser la somme quadratique des déviations  
des observations aux prédictions

$$\sum_i (y_i - y)^2$$

# Résidus

$$e_i = y_i - \hat{y}$$



# Résidus

$$e_i = y_i - \hat{y}$$

$y, y_i$

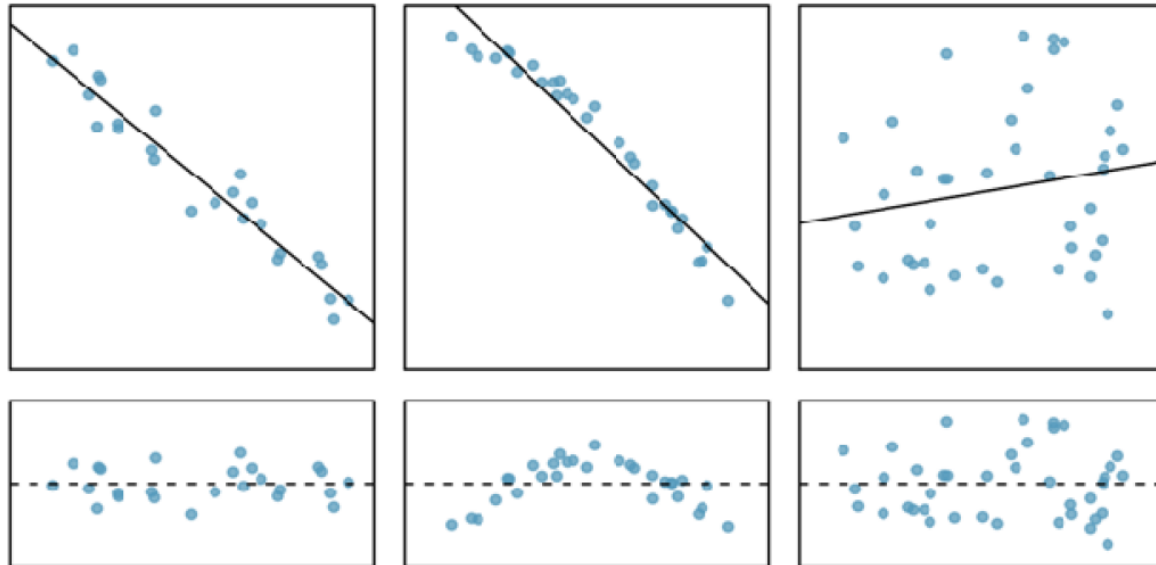


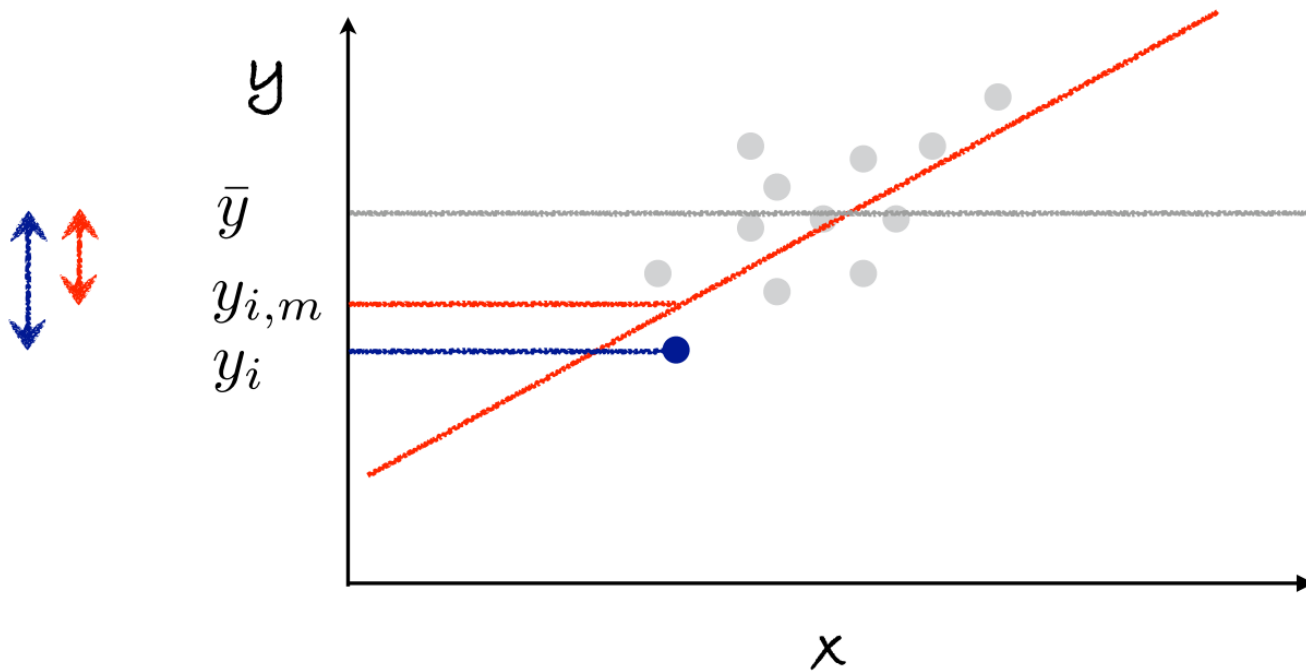
Figure 7.9: Sample data with their best fitting lines (top row) and their corresponding residual plots (bottom row).

x1

# Coefficient de corrélation (R)

$$R^2 = \frac{\sum_i (y_{i,m} - \bar{y})^2}{\sum_i (y_i - \bar{y})^2}$$

*fraction de la variation de y  
expliquée par la régression*



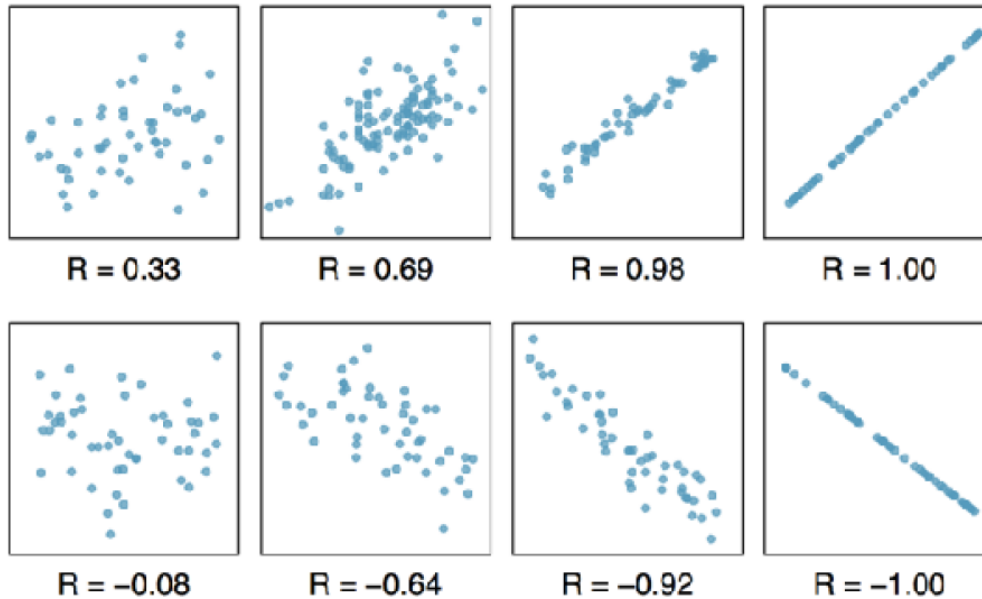


# Coefficient de corrélation (R)

$$R^2 = \frac{\sum_i (y_{i,m} - \bar{y})^2}{\sum_i (y_i - \bar{y})^2}$$

*fraction de la variation de y  
expliquée par la régression*

$$0 \leq R^2 \leq 1$$

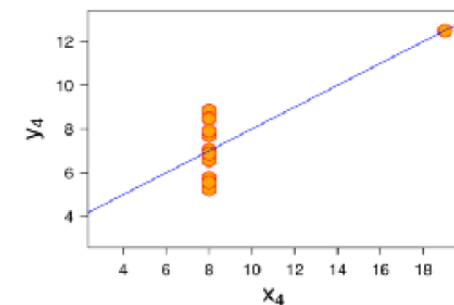
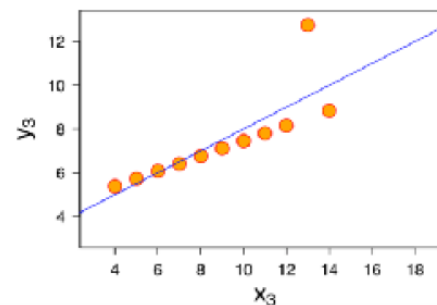
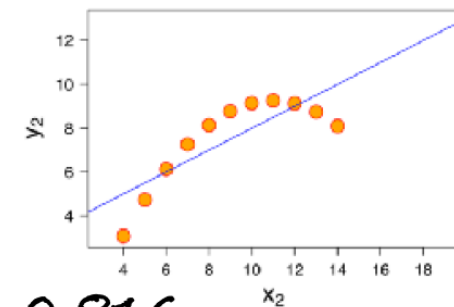
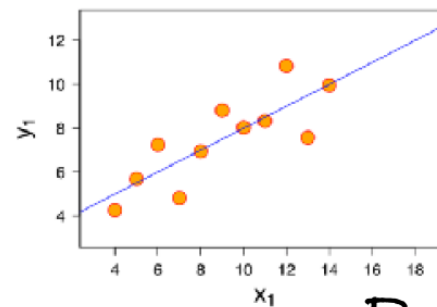


ATTENTION

Regardez les  
données avant  
d'interpréter R

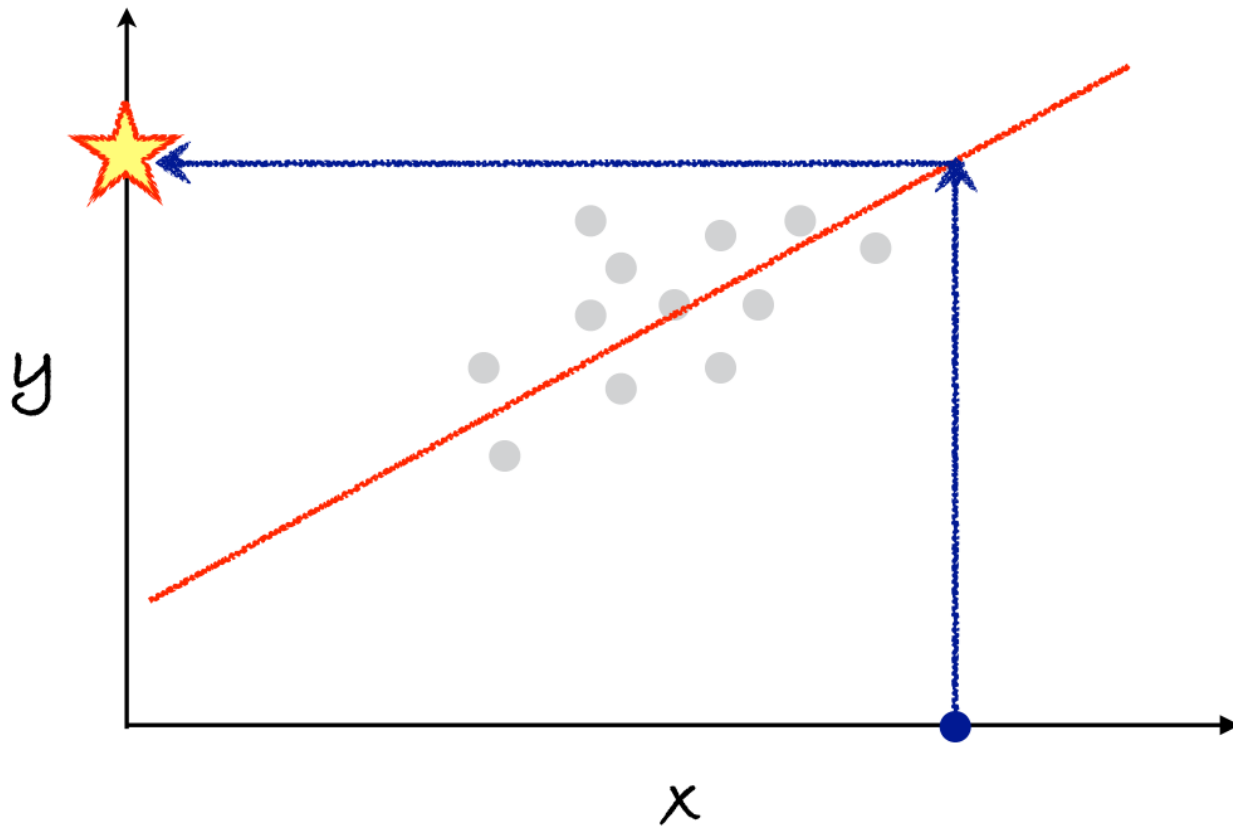
## Quartet d'Anscombe

4 datasets très différents, mais avec  
même propriétés statistiques de base

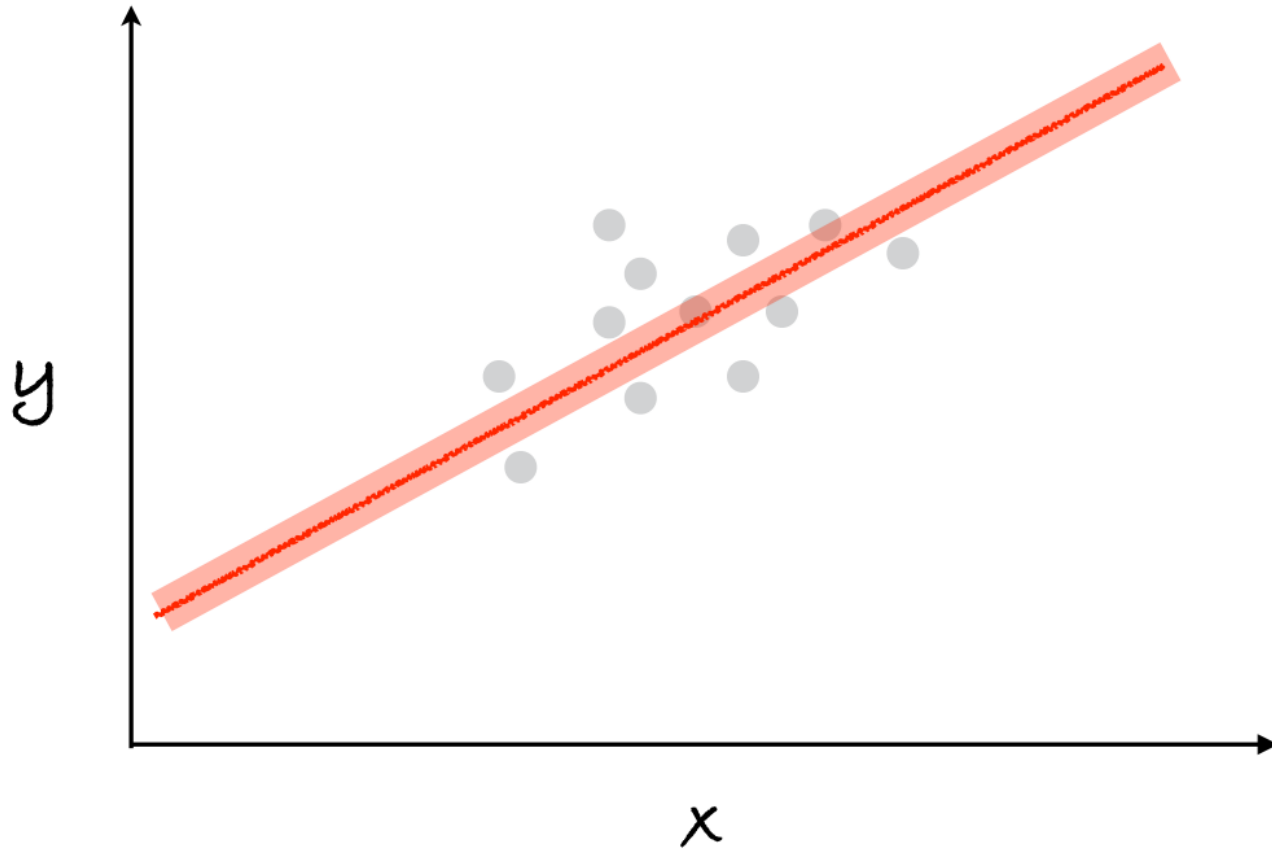


$$R = 0.816$$

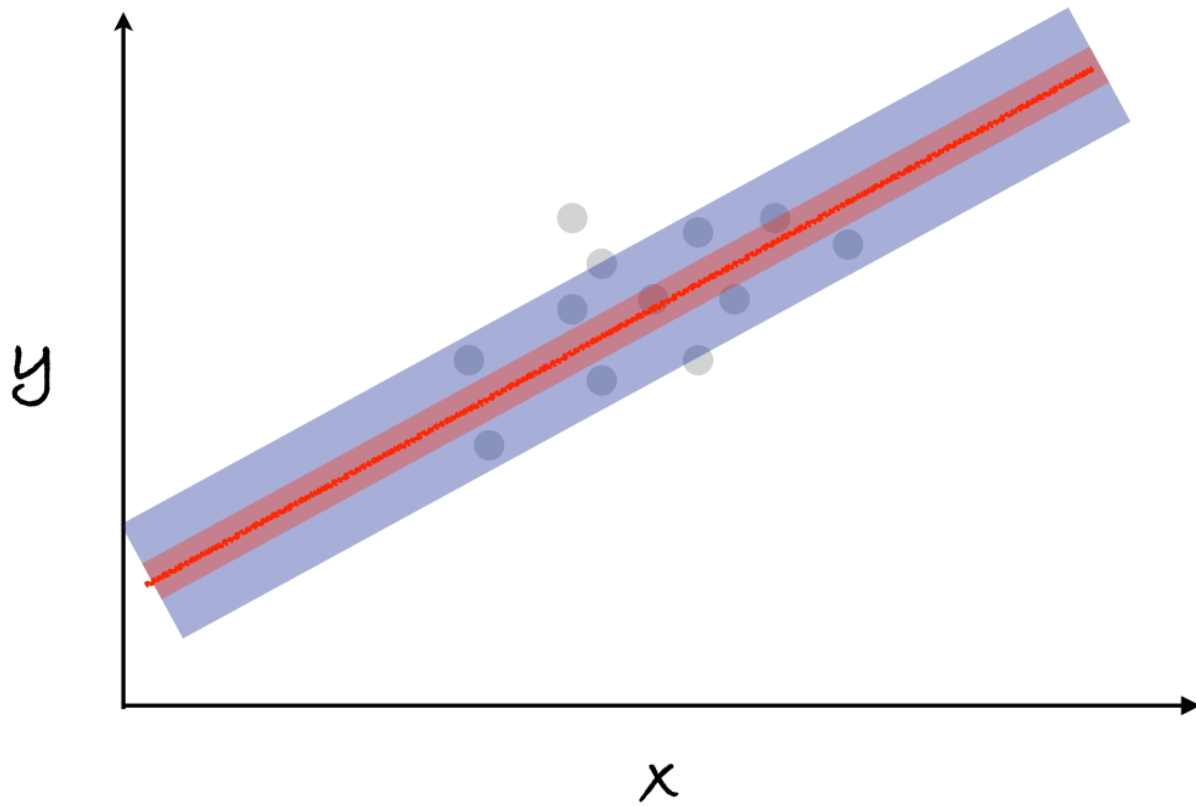
# Prédictions



# Erreur d'échantillonnage



Erreur de prédiction



# PREDICTION

*Erreur de prédiction*

