

# Statistique : Statistiques descriptives

## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon



# Statistique : Statistiques descriptives

Joseph Salmon

Septembre 2014



## Statistique : Statistiques descriptives

### Introduction générale

Notion de statistique  
Résumés basiques d'un jeu de données  
Corrélation

### Introduction générale

Notion de statistique  
Résumés basiques d'un jeu de données  
Corrélation

Joseph Salmon

## Statistique : Statistiques descriptives

### Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Statistique

- ▶ On observe des réalisations  $(y_1, \dots, y_n)$  de variables aléatoires inconnues (éventuellement vectorielles)
- ▶ On suppose ici que les variables sont indépendantes et identiquement distribuées (*i.i.d.*) selon une loi  $\mathbb{P}_Y$

## But de l'estimation

Comment apprendre certaines caractéristiques de  $\mathbb{P}_Y$  à partir de  $(y_1, \dots, y_n)$  ?

Souvent : on se prépare à observer  $y_{n+1}$ .

## Cas de la prédiction

Que peut-on attendre de  $y_{n+1}$  ? (en moyenne, ou avec une certaine probabilité ?)

## Statistique : Statistiques descriptives

### Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

- ▶ Observations  $\mathbf{y} = y_{1:n} = (y_1, \dots, y_n)$  : **échantillon** de **taille**  $n$ .
- ▶ Grandeurs **théoriques** : dépendant de la loi  $\mathbb{P}_Y$  **inconnue**  
**Exemple**: l'espérance de la variable  $y$  sous la loi  $\mathbb{P}_Y$ .
- ▶ Grandeurs **empiriques** : calculées à partir des observations  $y_i$ .  
**Exemple**:  $\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$  est la moyenne empirique
- ▶ Objectif général : apprendre les caractéristiques théoriques de  $\mathbb{P}_Y$  à partir de résumés empiriques.

## Statistique : Statistiques descriptives

### Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

- ▶ Première analyse sans hypothèse sur la loi  $\mathbb{P}_Y$ .
- ▶ Analyse qualitative du jeu de données /échantillon

### Définition : Statistique

Une **statistique** est une fonction des observations  $(y_1, \dots, y_n)$ .

## Statistique : Statistiques descriptives

### Introduction générale

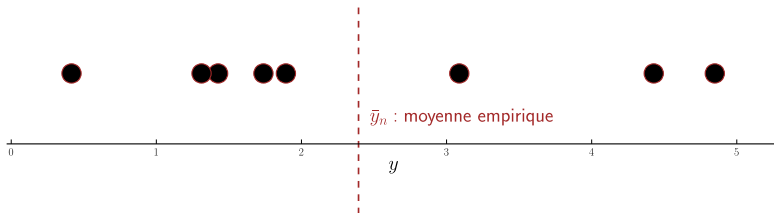
Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Moyenne



## Définition : Moyenne

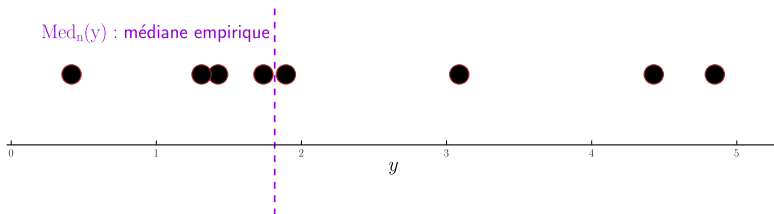
$$\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$$

Notons  $\mathbf{1}_n$  le vecteur  $(1, \dots, 1) \in \mathbb{R}^n$ . La moyenne est (à facteur  $1/n$  près) un produit scalaire dans  $\mathbb{R}^n$  :

$$\bar{y}_n = \langle \mathbf{y}, \mathbf{1}_n/n \rangle$$

cf. McKinney (2012) pour les statistiques avec python

# Médiane empirique

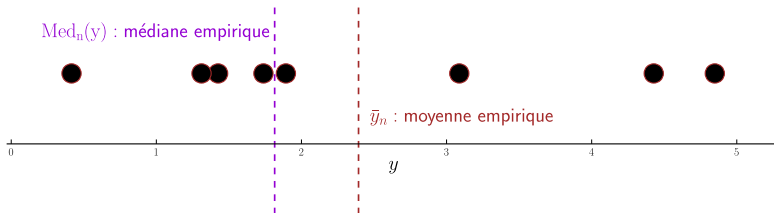


On ordonne les  $y_i$  :  $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$

Définition : Médiane (NON-UNIQUE!)

$$\text{Med}_n(y) = \begin{cases} \frac{y_{(\lfloor \frac{n}{2} \rfloor)} + y_{(\lfloor \frac{n}{2} \rfloor + 1)}}{2} & \text{Si } n \text{ est pair} \\ y_{(\frac{n+1}{2})} & \text{Si } n \text{ est impair} \end{cases}$$

# Moyenne vs médiane



- Les deux statistiques ne coïncident pas
- Une médiane est plus robuste aux points atypiques (en anglais : *outliers*)



## Statistique : Statistiques descriptives

### Introduction générale

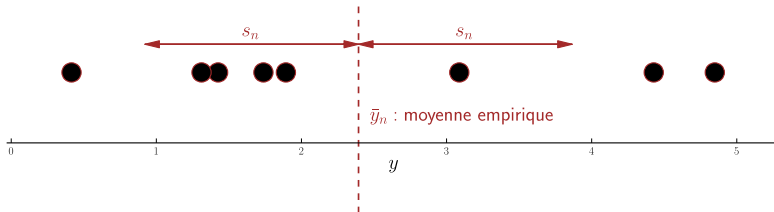
Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Dispersion



## Variance empirique

$$\text{var}_n(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2 = \frac{1}{n} \|\mathbf{y} - \bar{y}_n \mathbf{1}_n\|^2$$

## Écart-type empirique

$$s_n(\mathbf{y}) = \sqrt{\text{var}_n(\mathbf{y})} \quad \left( = \frac{1}{\sqrt{n}} \|\mathbf{y} - \bar{y}_n \mathbf{1}_n\| \right)$$

## Statistique : Statistiques descriptives

### Introduction générale

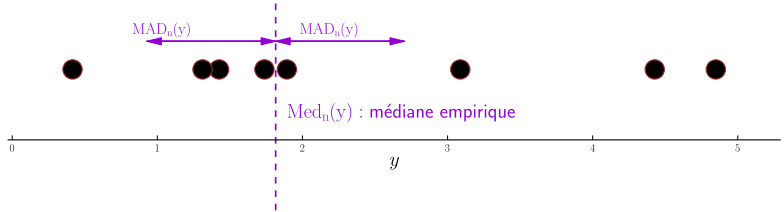
Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

## Dispersion



### Mean Absolute deviation

Déviations médiane absolue :

$$MAD_n(y) = \text{Med} (|\text{Med}(y) - y|),$$

# Histogramme

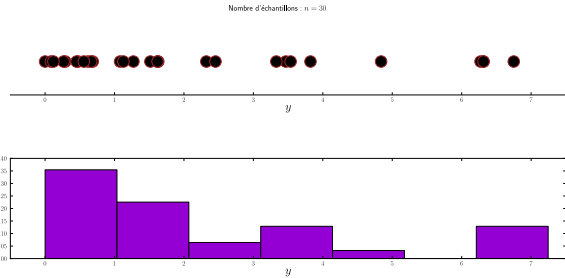
Statistique : Statistiques  
descriptives

## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation



Répartition des données dans des « cases »

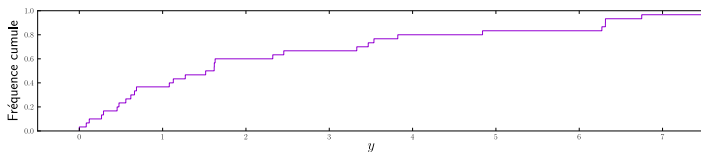
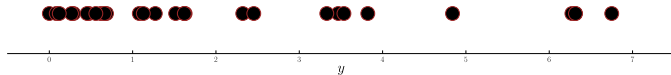
L'**aire** de chaque case est proportionnelle à la fraction des données qui « tombent » dans la case.

L'histogramme est une approximation de la **densité** de  $y$

Joseph Salmon

# Fonction de répartition empirique

Nombre d'échantillons :  $n = 30$



- ▶ *Rappel* : Fonction de répartition :  $F(u) = \mathbb{P}_Y(-\infty, u]$
- ▶ Version empirique : proportion des données en-dessous de  $u$

$$F_n(u) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i \leq u\}}$$

Statistique : Statistiques  
descriptives

Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Quantiles empiriques

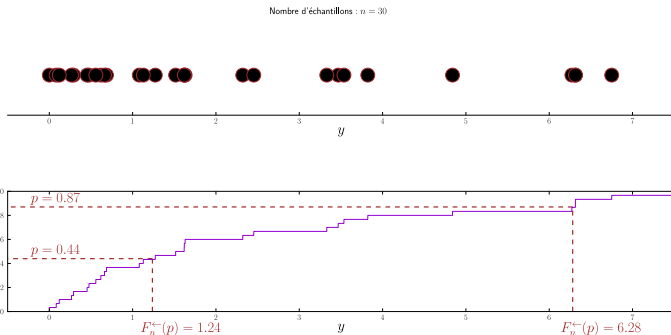
Statistique : Statistiques  
descriptives

## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation



- Inverse de la fonction de répartition empirique.
- Soit  $\lceil u \rceil$  le nombre entier tel que  $\lceil u \rceil - 1 < u \leq \lceil u \rceil$ .

## Quantiles empiriques

quantile d'ordre  $p = y_{(\lceil np \rceil)} = F_n^{←}(p) \quad (p \in [0, 1])$

Joseph Salmon

# Covariance et corrélation empirique

## Covariance empirique

Pour deux échantillons  $x_{1:n}$  et  $y_{1:n}$  de moyennes et variances empiriques  $\mathbf{x} = \bar{x}_n$ ,  $\mathbf{y} = \bar{y}_n$  et  $\text{var}_n(\mathbf{x})$ ,  $\text{var}_n(\mathbf{y})$  :

$$\text{cov}_n(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n) \quad \text{c'est-à-dire}$$

$$\text{cov}_n(x, y) = \frac{1}{n} \langle x_{1:n} - \bar{x}_n \mathbf{1}_n, y_{1:n} - \bar{y}_n \mathbf{1}_n \rangle$$

## Corrélation empirique

$$\rho = \text{corr}_n(x, y) = \frac{\text{cov}_n(x, y)}{\sqrt{\text{var}_n(\mathbf{x})} \sqrt{\text{var}_n(\mathbf{y})}}, \quad \text{c'est-à-dire}$$

$$\rho = \frac{\langle x_{1:n} - \bar{x}_n \mathbf{1}_n, y_{1:n} - \bar{y}_n \mathbf{1}_n \rangle}{\|x - \bar{x}_n\| \|y - \bar{y}_n\|} = \cos(x_{1:n} - \bar{x}_n \mathbf{1}_n, y_{1:n} - \bar{y}_n \mathbf{1}_n)$$

Statistique : Statistiques  
descriptives

Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Interprétation pour $n = 3$ et $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$

Statistique : Statistiques  
descriptives

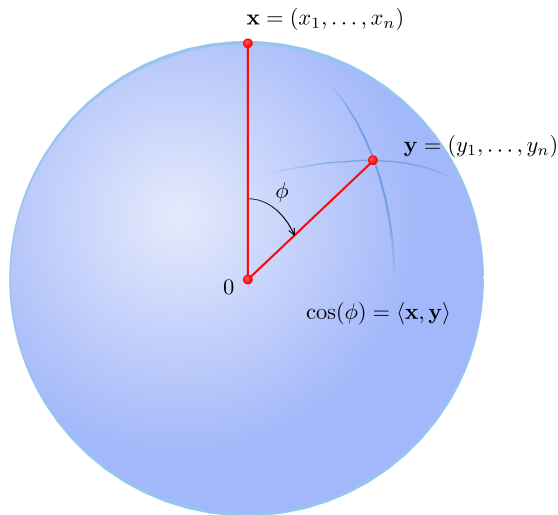
## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon



# Exemples de corrélations

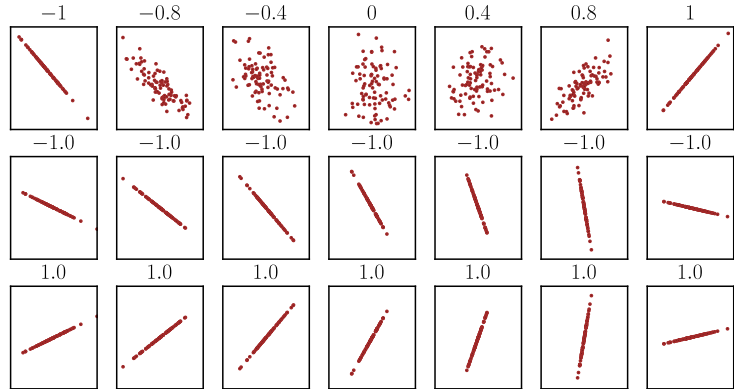
Statistique : Statistiques  
descriptives

## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

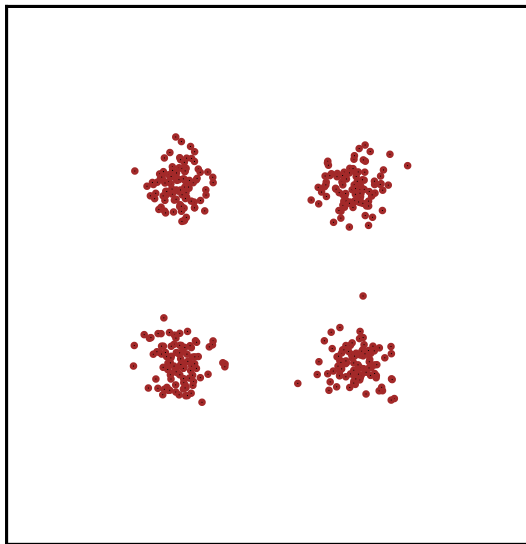


Joseph Salmon



# Exemples de corrélations proches de zéros

Corrélation =  $-0.021$



Statistique : Statistiques  
descriptives

## Introduction générale

Notion de statistique

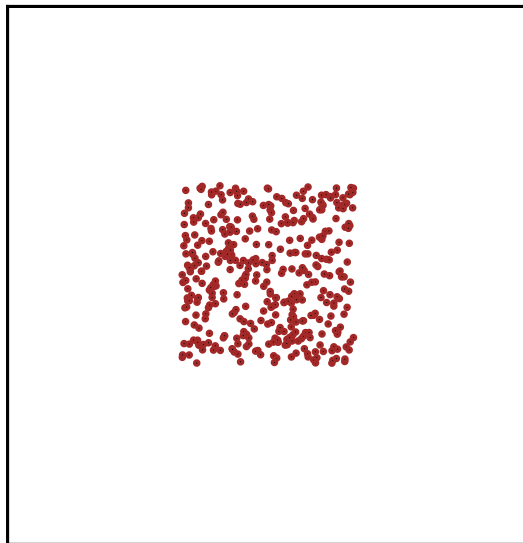
Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Exemples de corrélations proches de zéros

Corrélation = 0.007



Statistique : Statistiques  
descriptives

Introduction générale

Notion de statistique

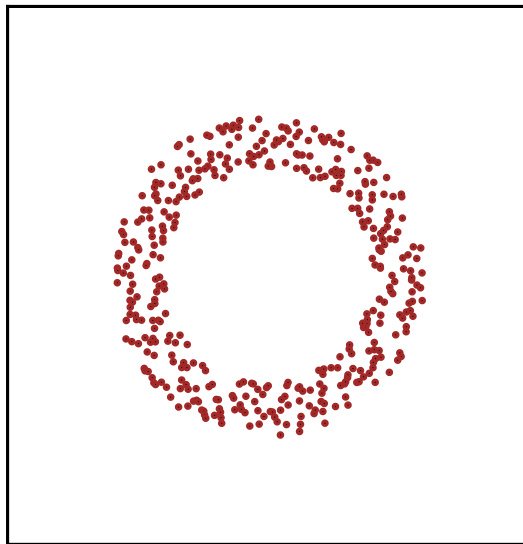
Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

# Exemples de corrélations proches de zéros

Corrélation = 0.011



Statistique : Statistiques  
descriptives

Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

## Statistique : Statistiques descriptives

### Introduction générale

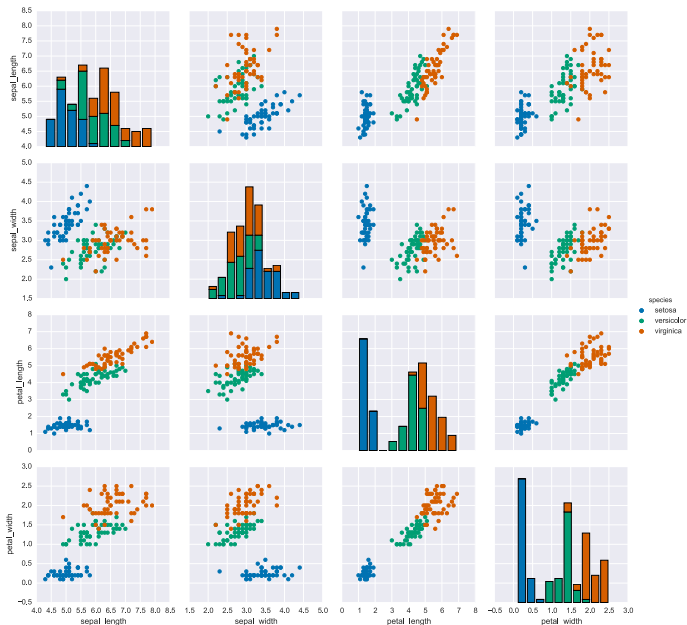
Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon

## Exemples de visualisation



## Statistique : Statistiques descriptives

## Introduction générale

Notion de statistique

Résumés basiques d'un jeu de données

Corrélation

Joseph Salmon



W. McKinney.

*Python for Data Analysis : Data Wrangling with Pandas, NumPy, and IPython.*

O'Reilly Media, 2012.