



# Analyse exploratoire de données

Fabrice Rossi

Télécom ParisTech

## Introduction

- Exploration

- Modélisation

- Modèle des données

## Analyses univariées

- Variables numériques

  - Histogramme

  - Boxplot et statistiques

- Variables nominales

## Analyses multivariées

- Diagramme de dispersion

- Matrice de corrélation

- Diagramme mosaïque

- Coordonnées parallèles

- Interaction

## Introduction

Exploration

Modélisation

Modèle des données

## Analyses univariées

Variables numériques

Histogramme

Boxplot et statistiques

Variables nominales

## Analyses multivariées

Diagramme de dispersion

Matrice de corrélation

Diagramme mosaïque

Coordonnées parallèles

Interaction

# Exploiter des données

- Que faire d'un paquet de données ?
- Comment exploiter le contenu d'un entrepôt de données ?

# Exploiter des données

- Que faire d'un paquet de données ?
- Comment exploiter le contenu d'un entrepôt de données ?
- recensement
- 32561 personnes
- 15 attributs par personne

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	Age	Workclass	Privet	Education	Educat	Marital-status	Occupation	Relationship	Ethnicity	Gender	Capita	Capit	Hour	perthive country	Salary
2	40	State-gov	77511	Bachelors	13	Never married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<50K
3	23	Self-emp-inc	53311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<50K
4	39	Private	170644	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<50K
5	51	Private	214071	11th	7	Married-civ-spouse	Handers-cleaners	Husband	Black	Male	0	0	40	United-States	<50K
6	20	Private	338408	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<50K
7	37	Private	364280	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<50K
8	46	Private	160167	9th	5	Married-spouse-absent	Other-service	Not-in-family	Black	Female	0	0	16	Jamaica	<50K
9	52	Self-emp-inc	309642	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	45	United-States	>50K
10	31	Private	437071	Masters	14	Never married	Prof-specialty	Not-in-family	White	Female	14804	0	50	United-States	>50K
11	42	Private	159449	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	5170	0	40	United-States	>50K
12	37	Private	280444	Some-college	10	Married-civ-spouse	Exec-managerial	Husband	Black	Male	0	0	60	United-States	>50K
13	30	State-gov	141260	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband	Asian-Pac-Islander	Male	0	0	40	Other	>50K
14	23	Private	122072	Bachelors	13	Never married	Adm-clerical	Own-child	White	Female	0	0	30	United-States	<50K
15	30	Private	205016	Assoc-acadm	12	Never married	Sales	Not-in-family	Black	Male	0	0	50	United-States	<50K
16	40	Private	121772	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband	Asian-Pac-Islander	Male	0	0	40	?	>50K
17	34	Private	345487	7th-8th	6	Married-civ-spouse	Transport-moving	Husband	Asian-Indian-Sakimo	Male	0	0	45	Mexico	<50K
18	25	Self-emp-inc	110780	HS-grad	9	Never married	Farming-fishing	Own-child	White	Male	0	0	35	United-States	<50K
19	32	Private	180024	HS-grad	9	Never married	Machine-op-inspct	Unmarried	White	Male	0	0	40	United-States	<50K
20	38	Private	23887	11th	7	Married-civ-spouse	Sales	Husband	White	Male	0	0	50	United-States	<50K
21	43	Self-emp-inc	262178	Masters	14	Divorced	Exec-managerial	Unmarried	White	Female	0	0	45	United-States	>50K
22	40	Private	193024	Doctorate	16	Married-civ-spouse	Prof-specialty	Husband	White	Male	0	0	60	United-States	>50K
23	54	Private	302148	HS-grad	6	Separated	Other-service	Unmarried	Black	Female	0	0	20	United-States	<50K
24	35	Federal-gov	76465	9th	5	Married-civ-spouse	Farming-fishing	Husband	Black	Male	0	0	40	United-States	<50K
25	43	Private	117027	11th	7	Married-civ-spouse	Transport-moving	Husband	White	Male	0	2062	40	United-States	<50K
26	56	Private	100015	HS-grad	9	Divorced	Technician	Unmarried	White	Female	0	0	40	United-States	<50K
27	36	Local-gov	216051	Bachelors	13	Married-civ-spouse	Technician	Husband	White	Male	0	0	40	United-States	>50K
28	19	Private	160024	HS-grad	9	Never married	Craft-repair	Own-child	White	Male	0	0	40	United-States	<50K
29	54	?	190211	Some-college	10	Married-civ-spouse	?	Husband	Asian-Pac-Islander	Male	0	0	60	South	>50K
30	39	Private	367280	HS-grad	9	Divorced	Exec-managerial	Not-in-family	White	Male	0	0	60	United-States	<50K
31	49	Private	193366	HS-grad	9	Married-civ-spouse	Craft-repair	Husband	White	Male	0	0	40	United-States	<50K
32	23	Local-gov	190078	Assoc-acadm	12	Never married	Protective-serv	Not-in-family	White	Male	0	0	52	United-States	<50K
33	30	Private	260015	Some-college	10	Never married	Sales	Own-child	Black	Male	0	0	44	United-States	<50K
34	45	Private	389449	Bachelors	13	Divorced	Exec-managerial	Own-child	White	Male	0	1400	40	United-States	<50K

# Exploiter des données

- Que faire d'un paquet de données ?
- Comment exploiter le contenu d'un entrepôt de données ?

■ recensement

■ 32561

personnes

■ 15 attributs  
par personne

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
1	Age	Workclass	Private	Education	Education	Marital-status	Occupation	Relationship	Ethnicity	Gender	Capital	Capital	Hours	pernative	country	Salary
2	39	State-gov	77018	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K	
3	50	Self-emp-inc	53311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K	
4	39	Private	17048	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K	
5	51	Private	21427	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K	
6	20	Private	33840	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K	
7	37	Private	26250	Masters	16	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	>=50K	
8	49	Private	90167	9th	5	Married-spouse-absent	Other-service	Not-in-family	Black	Female	0	0	16	Jamaica	>=50K	
9	52	Self-emp-inc	30942	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	45	United-States	>=50K	
10	31	Private	43751	Masters	14	Never-married	Prof-specialty	Not-in-family	White	Female	14504	0	50	United-States	>=50K	
11	42	Private	15944	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	5170	0	40	United-States	>=50K	
12	37	Private	29444	Some-college	10	Married-civ-spouse	Exec-managerial	Husband	Black	Male	0	0	60	United-States	>=50K	
13	30	State-gov	14120	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband	Asian-Pac-Islander	Male	0	0	40	Other	>=50K	
14	23	Private	12272	Bachelors	13	Never-married	Adm-clerical	Own-child	White	Female	0	0	30	United-States	<=50K	
15	30	Private	20501	Assoc-acadm	12	Never-married	Staten	Not-in-family	Black	Male	0	0	50	United-States	<=50K	
16	40	Private	12172	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband	Asian-Pac-Islander	Male	0	0	40	?	>=50K	
17	34	Private	34542	7th-8th	4	Married-civ-spouse	Transport-moving	Husband	Asian-Indian-Sikimo	Male	0	0	45	Mexico	>=50K	
18	29	Self-emp-inc	11078	HS-grad	9	Never-married	Farming-fishing	Own-child	White	Male	0	0	35	United-States	>=50K	
19	32	Private	19504	HS-grad	9	Never-married	Machine-op-inspct	Unmarried	White	Male	0	0	40	United-States	>=50K	
20	38	Private	23857	11th	7	Married-civ-spouse	Sales	Husband	White	Male	0	0	50	United-States	>=50K	
21	43	Self-emp-inc	26217	Masters	14	Divorced	Exec-managerial	Unmarried	White	Female	0	0	45	United-States	>=50K	
22	40	Private	19324	Doctorate	16	Married-civ-spouse	Prof-specialty	Husband	White	Male	0	0	60	United-States	>=50K	
23	54	Private	30144	HS-grad	9	Separated	Other-service	Unmarried	Black	Female	0	0	20	United-States	<=50K	
24	35	Federal-gov	76645	9th	5	Married-civ-spouse	Farming-fishing	Husband	Black	Male	0	0	40	United-States	<=50K	
25	43	Private	117027	11th	7	Married-civ-spouse	Transport-moving	Husband	White	Male	0	2062	40	United-States	<=50K	
26	56	Private	150015	HS-grad	9	Divorced	Technician	Unmarried	White	Female	0	0	40	United-States	>=50K	
27	36	Local-gov	21605	Bachelors	13	Married-civ-spouse	Technician	Husband	White	Male	0	0	40	United-States	>=50K	
28	19	Private	19504	HS-grad	9	Never-married	Craft-repair	Own-child	White	Male	0	0	40	United-States	>=50K	
29	54	?	19021	Some-college	10	Married-civ-spouse	?	Husband	Asian-Pac-Islander	Male	0	0	60	South	>=50K	
30	39	Private	36720	HS-grad	9	Divorced	Exec-managerial	Not-in-family	White	Male	0	0	60	United-States	<=50K	
31	49	Private	19346	HS-grad	9	Married-civ-spouse	Craft-repair	Husband	White	Male	0	0	40	United-States	<=50K	
32	23	Local-gov	19078	Assoc-acadm	12	Never-married	Protective-serv	Not-in-family	White	Male	0	0	52	United-States	<=50K	
33	30	Private	26017	Some-college	10	Never-married	Sales	Own-child	Black	Male	0	0	44	United-States	<=50K	
34	45	Private	39042	Bachelors	13	Divorced	Exec-managerial	Own-child	White	Male	0	1400	40	United-States	>=50K	

■ Volume classique : millions à millions de lignes, dizaine à centaines de colonnes

■ Exploration systématique impossible (même pour de petits paquets de données)

- Support informatique et mathématique :
  - outils d'exploitation des données
  - but : diminuer la charge cognitive pour l'analyste
- Deux grandes classes d'outils :
  1. exploration
  2. modélisation

- Support informatique et mathématique :
  - outils d'exploitation des données
  - but : diminuer la charge cognitive pour l'analyste
- Deux grandes classes d'outils :
  1. exploration
    - pas d'idée *a priori* sur les données
    - recherche de régularité (dépendances, groupes homogènes, etc.)
  2. modélisation



- Support informatique et mathématique :
  - outils d'exploitation des données
  - but : diminuer la charge cognitive pour l'analyste
- Deux grandes classes d'outils :
  1. exploration
    - pas d'idée *a priori* sur les données
    - recherche de régularité (dépendances, groupes homogènes, etc.)
  2. modélisation
    - idée précise sur les données
    - construction de modèles prédictifs

- Support informatique et mathématique :
  - outils d'exploitation des données
  - but : diminuer la charge cognitive pour l'analyste
- Deux grandes classes d'outils :
  1. exploration
    - pas d'idée *a priori* sur les données
    - recherche de régularité (dépendances, groupes homogènes, etc.)
  2. modélisation
    - idée précise sur les données
    - construction de modèles prédictifs
- outil utilisé : R (<http://R-project.org/>)

## ■ Objectifs :

- obtenir une vision globale d'un jeu de données
- découvrir des formes de régularité

## ■ Moyens :

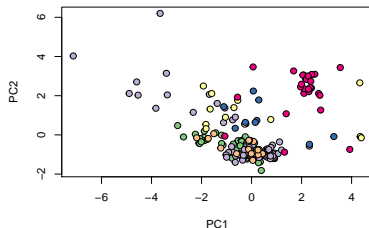
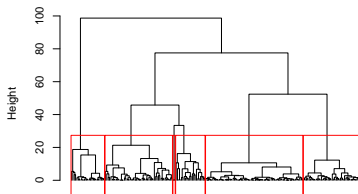
- représentations visuelles (et interactives) des données
- recherche automatique de régularités :
  - corrélation et dépendance entre variables
  - groupes homogènes (classification)
  - schémas fréquents

## ■ Objectifs :

- obtenir une vision globale d'un jeu de données
- découvrir des formes de régularité

## ■ Moyens :

- représentations visuelles (et interactives) des données
- recherche automatique de régularités :
  - corrélation et dépendance entre variables
  - groupes homogènes (classification)
  - schémas fréquents



## ■ Objectifs :

- inférer des informations inconnues
- prédire l'évolution des données

## ■ Moyens :

- données d'apprentissage :
  - connaître l'évolution d'une grandeur dans le passé pour prédire son évolution future (données historiques)
  - connaître une propriété de certains objets (par exemple le salaire de certains clients) pour inférer sa valeur pour les autres objets
- méthodes d'apprentissage : construire un modèle à partir des données d'apprentissage

## ■ Objectifs :

- inférer des informations inconnues
- prédire l'évolution des données

## ■ Moyens :

- données d'apprentissage :
  - connaître l'évolution d'une grandeur dans le passé pour prédire son évolution future (données historiques)
  - connaître une propriété de certains objets (par exemple le salaire de certains clients) pour inférer sa valeur pour les autres objets
- méthodes d'apprentissage : construire un modèle à partir des données d'apprentissage

## ■ Stratégie :

- analyse exploratoire
- formulation d'hypothèses
- construction de modèles pour valider les hypothèses

- On a  $N$  observations, les  $z_i \in \mathcal{Z}$
- Modèle statistique/probabiliste
  - il existe une distribution  $P_Z$  sur  $\mathcal{Z}$  inconnue
  - les  $z_i$  sont des réalisations de variables aléatoires avec cette distribution
  - les variables aléatoires sont indépendantes (en général)

- On a  $N$  observations, les  $z_i \in \mathcal{Z}$
- Modèle statistique/probabiliste
  - il existe une distribution  $P_Z$  sur  $\mathcal{Z}$  inconnue
  - les  $z_i$  sont des réalisations de variables aléatoires avec cette distribution
  - les variables aléatoires sont indépendantes (en général)
- En général
  - $\mathcal{Z} = \prod_{p=1}^P \mathcal{Z}_p$  :  $P$  variables pour décrire chaque objet
  - quand  $\mathcal{Z}_p \subset \mathbb{R}$  : variable numérique (ou ordonnée)
  - quand  $\mathcal{Z}_p = \{a, b, \dots\}$  : variable nominale (un nombre fini de valeurs possibles non ordonnées)



## Introduction

Exploration

Modélisation

Modèle des données

## Analyses univariées

Variables numériques

Histogramme

Boxplot et statistiques

Variables nominales

## Analyses multivariées

Diagramme de dispersion

Matrice de corrélation

Diagramme mosaïque

Coordonnées parallèles

Interaction

- Première étape d'une analyse exploratoire
  - travailler variable par variable
  - numériquement et graphiquement
- Variable numérique
  - à valeurs dans  $\mathbb{R}$
  - statistiques classiques : moyenne, variance, médiane, etc.
  - représentations associées : histogramme, *boxplot*

# Analyses élémentaires

## ■ Première étape d'une analyse exploratoire

- travailler variable par variable
- numériquement et graphiquement

## ■ Variable numérique

- à valeurs dans  $\mathbb{R}$
- statistiques classiques : moyenne, variance, médiane, etc.
- représentations associées : histogramme, *boxplot*

Age	Workclass	Education	Education-Num	Married-Status	Occupation	Relationship	Married	Sex	Gender	Capital-Gain	Hours-Per-Week	Native-country	Salary
1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9
10	10	10	10	10	10	10	10	10	10	10	10	10	10
11	11	11	11	11	11	11	11	11	11	11	11	11	11
12	12	12	12	12	12	12	12	12	12	12	12	12	12
13	13	13	13	13	13	13	13	13	13	13	13	13	13
14	14	14	14	14	14	14	14	14	14	14	14	14	14
15	15	15	15	15	15	15	15	15	15	15	15	15	15
16	16	16	16	16	16	16	16	16	16	16	16	16	16
17	17	17	17	17	17	17	17	17	17	17	17	17	17
18	18	18	18	18	18	18	18	18	18	18	18	18	18
19	19	19	19	19	19	19	19	19	19	19	19	19	19
20	20	20	20	20	20	20	20	20	20	20	20	20	20
21	21	21	21	21	21	21	21	21	21	21	21	21	21
22	22	22	22	22	22	22	22	22	22	22	22	22	22
23	23	23	23	23	23	23	23	23	23	23	23	23	23
24	24	24	24	24	24	24	24	24	24	24	24	24	24
25	25	25	25	25	25	25	25	25	25	25	25	25	25
26	26	26	26	26	26	26	26	26	26	26	26	26	26
27	27	27	27	27	27	27	27	27	27	27	27	27	27
28	28	28	28	28	28	28	28	28	28	28	28	28	28
29	29	29	29	29	29	29	29	29	29	29	29	29	29
30	30	30	30	30	30	30	30	30	30	30	30	30	30
31	31	31	31	31	31	31	31	31	31	31	31	31	31
32	32	32	32	32	32	32	32	32	32	32	32	32	32
33	33	33	33	33	33	33	33	33	33	33	33	33	33
34	34	34	34	34	34	34	34	34	34	34	34	34	34
35	35	35	35	35	35	35	35	35	35	35	35	35	35
36	36	36	36	36	36	36	36	36	36	36	36	36	36
37	37	37	37	37	37	37	37	37	37	37	37	37	37
38	38	38	38	38	38	38	38	38	38	38	38	38	38
39	39	39	39	39	39	39	39	39	39	39	39	39	39
40	40	40	40	40	40	40	40	40	40	40	40	40	40
41	41	41	41	41	41	41	41	41	41	41	41	41	41
42	42	42	42	42	42	42	42	42	42	42	42	42	42
43	43	43	43	43	43	43	43	43	43	43	43	43	43
44	44	44	44	44	44	44	44	44	44	44	44	44	44
45	45	45	45	45	45	45	45	45	45	45	45	45	45
46	46	46	46	46	46	46	46	46	46	46	46	46	46
47	47	47	47	47	47	47	47	47	47	47	47	47	47
48	48	48	48	48	48	48	48	48	48	48	48	48	48
49	49	49	49	49	49	49	49	49	49	49	49	49	49
50	50	50	50	50	50	50	50	50	50	50	50	50	50

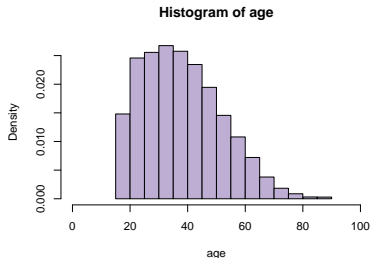
Variable âge : numérique

- travailler variable par variable
- numériquement et graphiquement

- à valeurs dans  $\mathbb{R}$
- statistiques classiques : moyenne, variance, médiane, etc.
- représentations associées : histogramme, *boxplot*

[illegible]

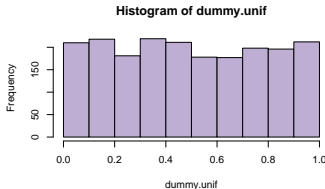
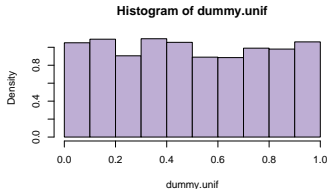
Variable âge : numérique



- Un histogramme représente une estimation de la **distribution** d'une variable
- Principe de construction :
  - division de l'intervalle  $[\min, \max]$  en  $K$  sous-intervalles (diverses règles pour  $K$ , par exemple  $\sim \log N$ )
  - dénombrement des objets pour lesquels la valeur de la variable tombe dans chacun des intervalles
  - représentation par des barres de **surfaces** proportionnelles aux décomptes

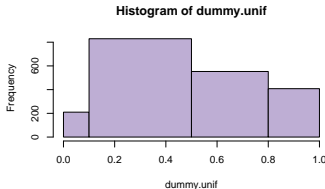
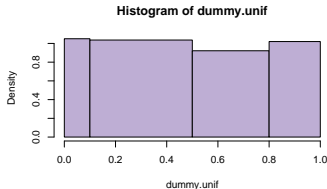
# Histogramme

- Un histogramme représente une estimation de la **distribution** d'une variable
- Principe de construction :
  - division de l'intervalle  $[\min, \max]$  en  $K$  sous-intervalles (diverses règles pour  $K$ , par exemple  $\sim \log N$ )
  - dénombrement des objets pour lesquels la valeur de la variable tombe dans chacun des intervalles
  - représentation par des barres de **surfaces** proportionnelles aux décomptes
- Attention aux intervalles de longueurs différentes



# Histogramme

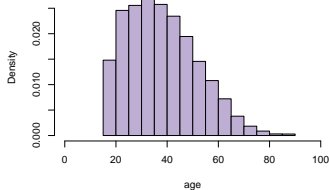
- Un histogramme représente une estimation de la **distribution** d'une variable
- Principe de construction :
  - division de l'intervalle  $[\min, \max]$  en  $K$  sous-intervalles (diverses règles pour  $K$ , par exemple  $\sim \log N$ )
  - dénombrement des objets pour lesquels la valeur de la variable tombe dans chacun des intervalles
  - représentation par des barres de **surfaces** proportionnelles aux décomptes
- Attention aux intervalles de longueurs différentes





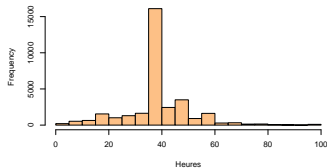
# Intérêts

Histogram of age



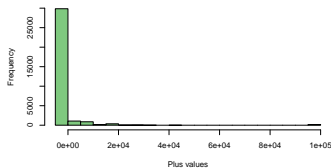
Âge

Histogramme des heures travaillées par semaine



Temps de travail

Histogramme des plus values



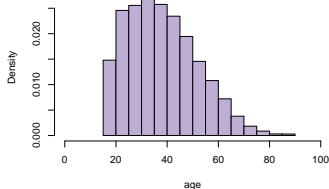
Plus values





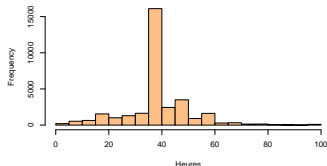
# Intérêts

Histogram of age



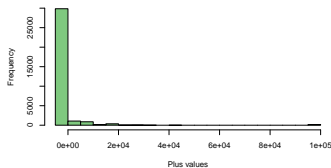
Âge

Histogramme des heures travaillées par semaine



Temps de travail

Histogramme des plus values



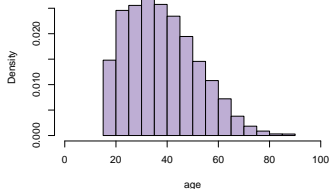
Plus values

■ Idée générale de la distribution



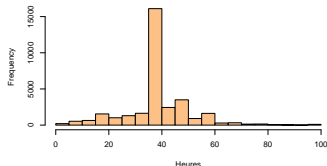
# Intérêts

Histogram of age



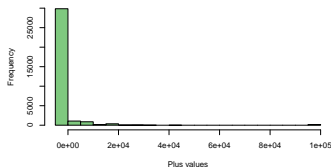
Âge

Histogramme des heures travaillées par semaine



Temps de travail

Histogramme des plus values



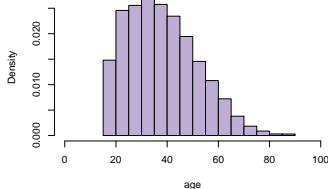
Plus values

- Idée générale de la distribution
- “irrégularités”



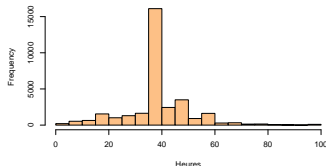


Histogram of age



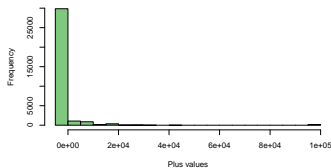
Âge

Histogramme des heures travaillées par semaine



Temps de travail

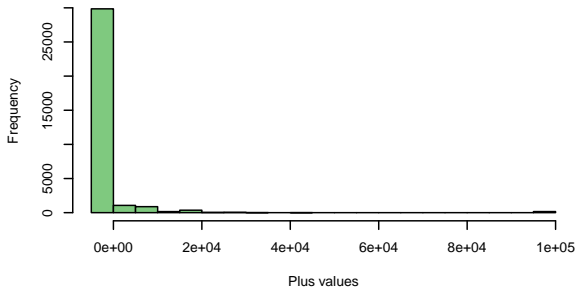
Histogramme des plus values



Plus values

- Idée générale de la distribution
- “irrégularités”
- distribution complètement atypique

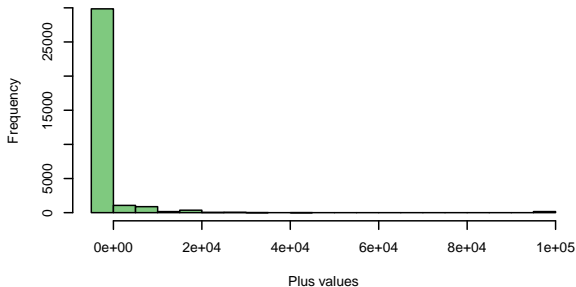
Histogramme des plus values



■ presque aucune information :

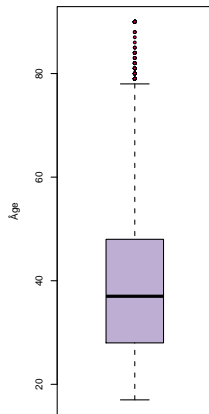
- presque toutes les valeurs sont négatives
- quelques valeurs très grandes

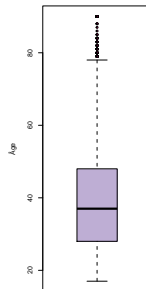
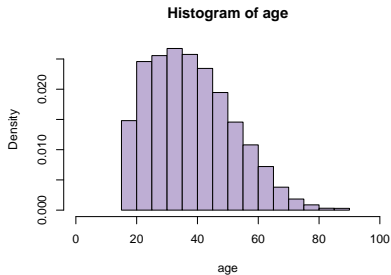
Histogramme des plus values

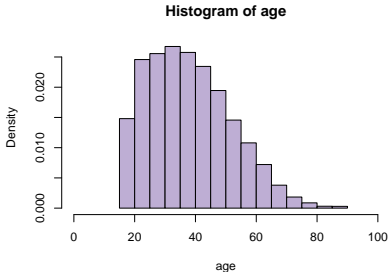


- presque aucune information :
  - presque toutes les valeurs sont négatives
  - quelques valeurs très grandes
- comparaisons difficiles (cf la suite)

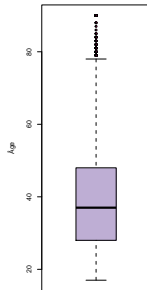
- a.k.a. boîte à moustaches ou boîte à pattes
- Représentation compacte d'une distribution
  - ligne centrale : médiane
  - ligne basse : premier quartile
  - ligne haute : troisième quartile
  - moustaches :
    - le max du min et de la médiane - 1.5 l'intervalle interquartile
    - le min du max et de la médiane + 1.5 l'intervalle interquartile
  - points atypiques (*outliers*) : au delà des moustaches





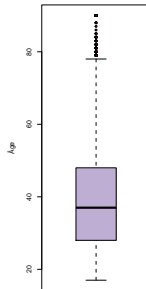
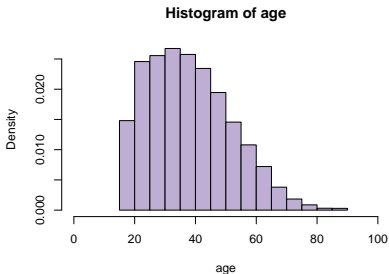


■ plus d'information



- plus dépouillé





- plus d'information
- inférence moins précise

- plus dépouillé
- quelques informations très précises

## ■ Indicateurs classiques :

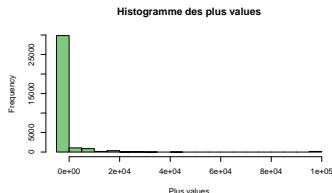
- tendance : moyenne et médiane
- dispersion : écart-type, intervalle interquartile

## ■ Indicateurs classiques :

- tendance : moyenne et médiane
- dispersion : écart-type, intervalle interquartile

## ■ Interprétation parfois délicate :

- moyenne = 990
- médiane = 0
- écart-type = 7410
- intervalle interquartile = 0

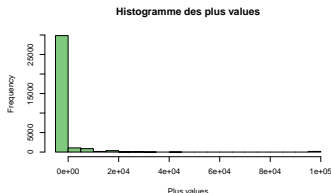


## ■ Indicateurs classiques :

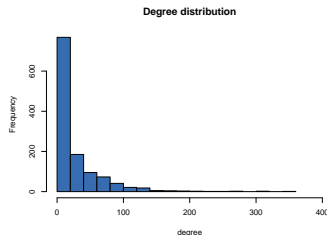
- tendance : moyenne et médiane
- dispersion : écart-type, intervalle interquartile

## ■ Interprétation parfois délicate :

- moyenne = 990
- médiane = 0
- écart-type = 7410
- intervalle interquartile = 0
- meilleurs choix ici :
  - 87 % des personnes ont une plus value nulle, 8.3 % positive et 4.7 % négative
  - puis statistiques sur les deux groupes (par ex., perte médiane 1887)



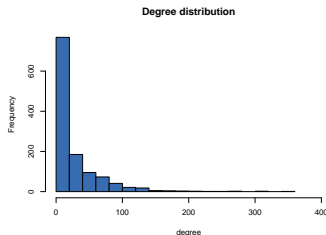
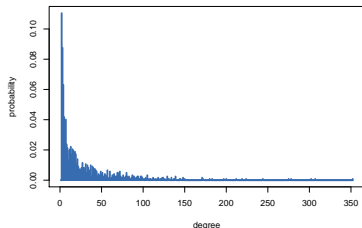
- La pertinence de la statistique dépend de la distribution
- Exemple :
  - blogs politiques
  - graphe des liens entre les blogs (blogroll)
  - distribution des degrés des noeuds



■  $\mu = 27.36, \sigma = 38.42$

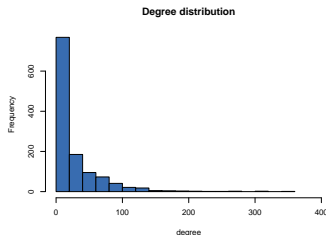
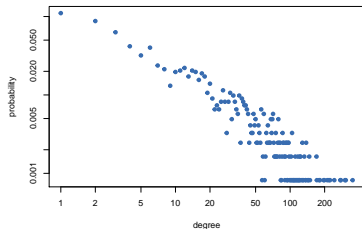
■  $m = 13, \delta = 33$

- La pertinence de la statistique dépend de la distribution
- Exemple :
  - blogs politiques
  - graphe des liens entre les blogs (blogroll)
  - distribution des degrés des noeuds



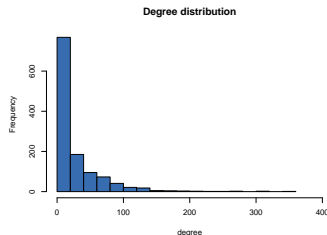
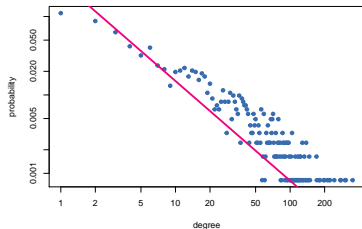
- $\mu = 27.36, \sigma = 38.42$
- $m = 13, \delta = 33$
- loi puissance :  $P(x) \simeq x^{-\alpha}$

- La pertinence de la statistique dépend de la distribution
- Exemple :
  - blogs politiques
  - graphe des liens entre les blogs (blogroll)
  - distribution des degrés des noeuds



- $\mu = 27.36, \sigma = 38.42$
- $m = 13, \delta = 33$
- loi puissance :  $P(x) \simeq x^{-\alpha}$
- sans échelle : la moyenne informe peu

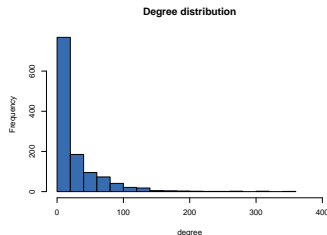
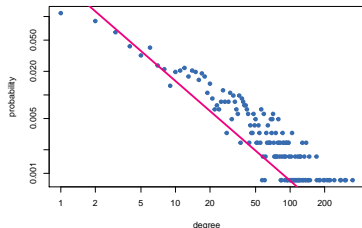
- La pertinence de la statistique dépend de la distribution
- Exemple :
  - blogs politiques
  - graphe des liens entre les blogs (blogroll)
  - distribution des degrés des noeuds



- $\mu = 27.36, \sigma = 38.42$
- $m = 13, \delta = 33$
- loi puissance :  $P(x) \simeq x^{-\alpha}$
- sans échelle : la moyenne informe peu
- ici  $\alpha \simeq 1.27$



- La pertinence de la statistique dépend de la distribution
- Exemple :
  - blogs politiques
  - graphe des liens entre les blogs (blogroll)
  - distribution des degrés des noeuds

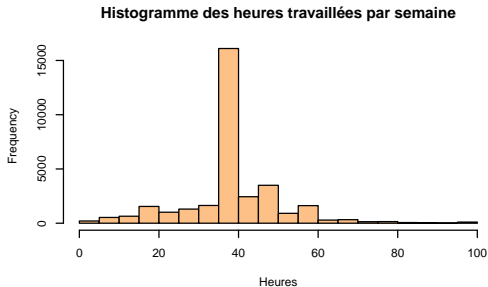


- $\mu = 27.36, \sigma = 38.42$
- $m = 13, \delta = 33$
- loi puissance :  $P(x) \simeq x^{-\alpha}$
- sans échelle : la moyenne informe peu
- ici  $\alpha \simeq 1.27$

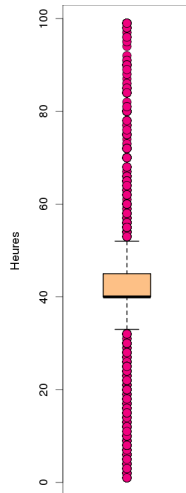
Adapter les statistiques  
aux données



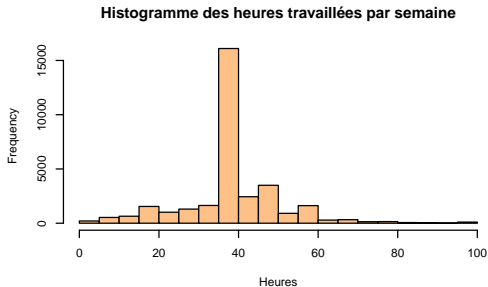
# Trois points de vue



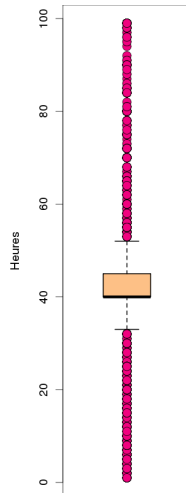
- Moyenne : 40.44, Écart-type : 12.35
- Médiane : 40, Interquartile : 5



# Trois points de vue



- Moyenne : 40.44, Écart-type : 12.35
- Médiane : 40, Interquartile : 5
- Compléments :
  - 47 % = 40 heures
  - 29 % > 40 heures
  - 24 % < 40 heures



## Variables nominales

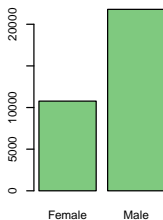
- variable **nominale** (ou **qualitative**) : variable à valeurs dans un ensemble fini quelconque (les **modalités**)
- quand les modalités sont ordonnées : variable **ordinaire**

## Variables nominales

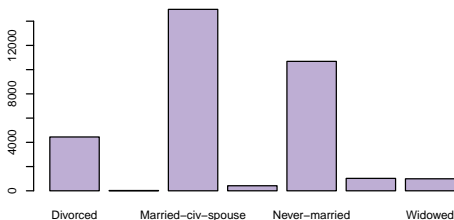
- variable **nominale** (ou **qualitative**) : variable à valeurs dans un ensemble fini quelconque (les **modalités**)
- quand les modalités sont ordonnées : variable **ordinaire**
- représentation par un **diagramme à bâtons** :
  - un bâton par modalité
  - hauteur proportionnelle à la fréquence de la modalité
  - ordre arbitraire sauf dans la cas ordinal

# Variables nominales

- variable **nominale** (ou **qualitative**) : variable à valeurs dans un ensemble fini quelconque (les **modalités**)
- quand les modalités sont ordonnées : variable **ordinaire**
- représentation par un **diagramme à bâtons** :
  - un bâton par modalité
  - hauteur proportionnelle à la fréquence de la modalité
  - ordre arbitraire sauf dans la cas ordinal

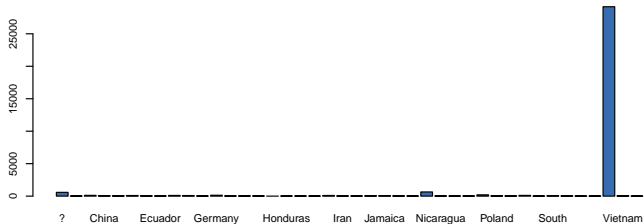


Genre

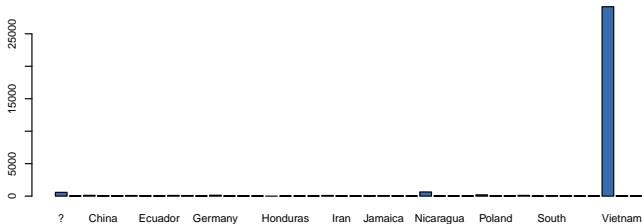


Statut marital

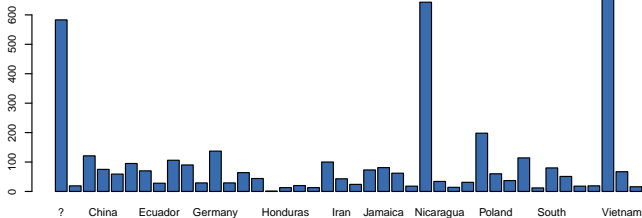
## ■ Déséquilibre



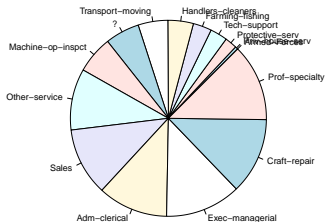
## ■ Déséquilibre



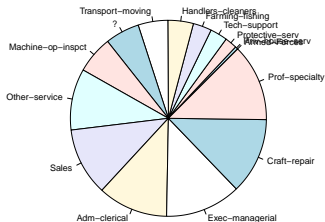
## ■ Grand nombre de modalités



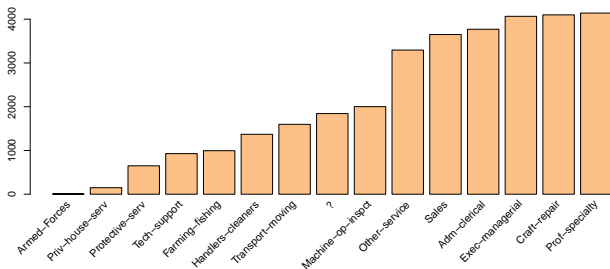




- représentation très classique
- versions “créatives” (3D...)
- mauvaise solution : lecture des surfaces et des angles difficiles



- représentation très classique
- versions “créatives” (3D...)
- mauvaise solution : lecture des surfaces et des angles difficiles



## Introduction

Exploration

Modélisation

Modèle des données

## Analyses univariées

Variables numériques

Histogramme

Boxplot et statistiques

Variables nominales

## Analyses multivariées

Diagramme de dispersion

Matrice de corrélation

Diagramme mosaïque

Coordonnées parallèles

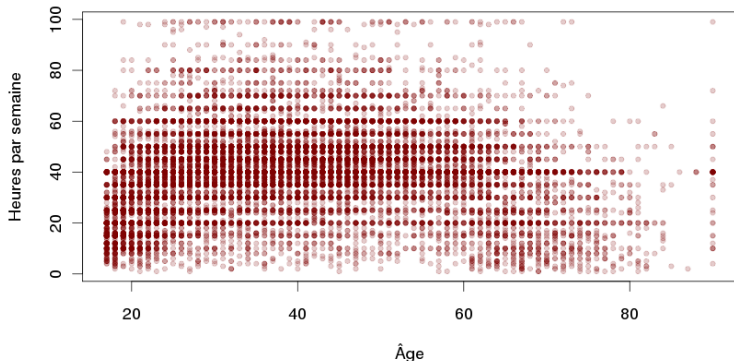
Interaction

- Relativement peu d'information dans chaque variable
- Analyse croisée nécessaire
- Difficultés :
  - vision humaine limitée (2D ou 3D, formes et couleurs)
  - beaucoup de combinaisons possibles
  - variables incompatibles
- Solutions :
  - outils de la visualisation de l'information (interaction)
  - outils de l'apprentissage automatique (automatisation)



# Diagramme de dispersion

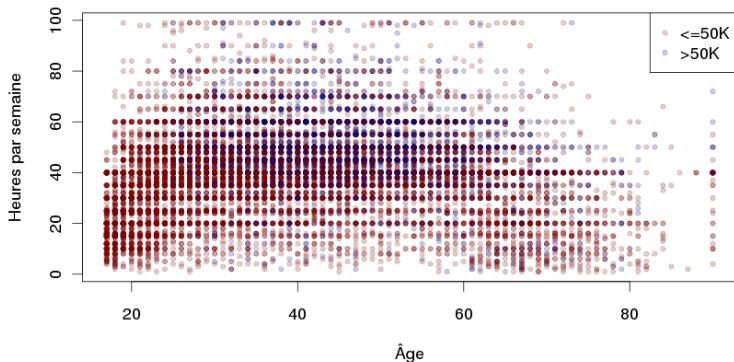
- Deux variables numériques : l'une en fonction de l'autre
- *scatter plot*



- Superposition : *alpha blending*

## ■ Compléments du diagramme :

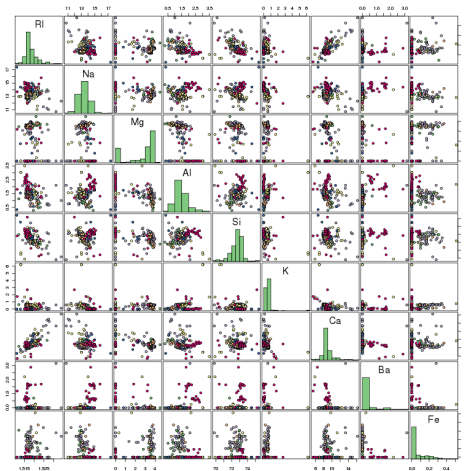
- couleur en fonction d'une autre variable
- symbole en fonction d'une autre variable



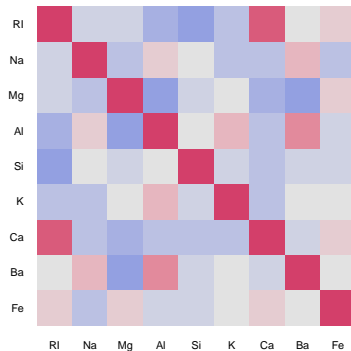
## ■ Assez limité

# Matrice de diagrammes

- matrice de diagrammes de dispersion
- tous les couples de variables numériques
- limités à quelques variables (croissance quadratique)
- décorations possibles
- ici : 7 types de verre décrits par 9 variables

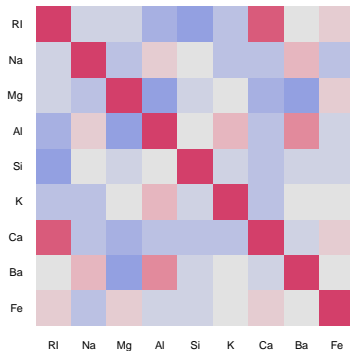


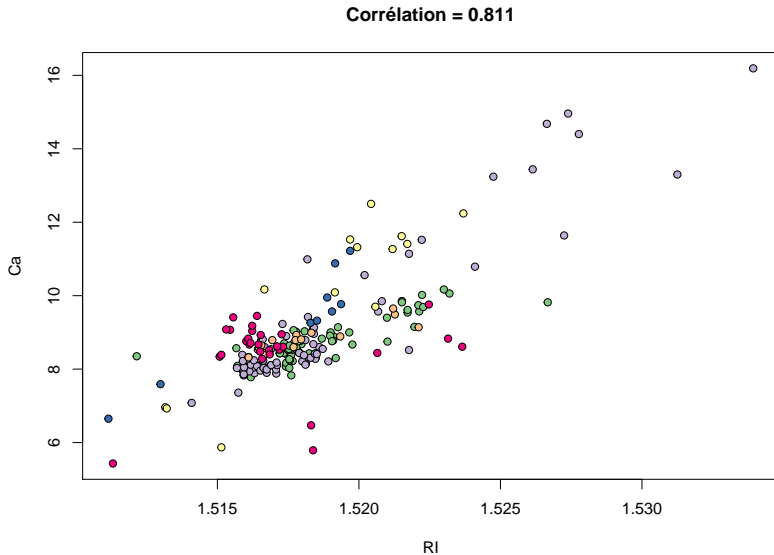
- Recherche de corrélations
- Représentation graphique de la matrice de corrélation :
  - rouge : forte corrélation positive
  - bleu : forte corrélation négative





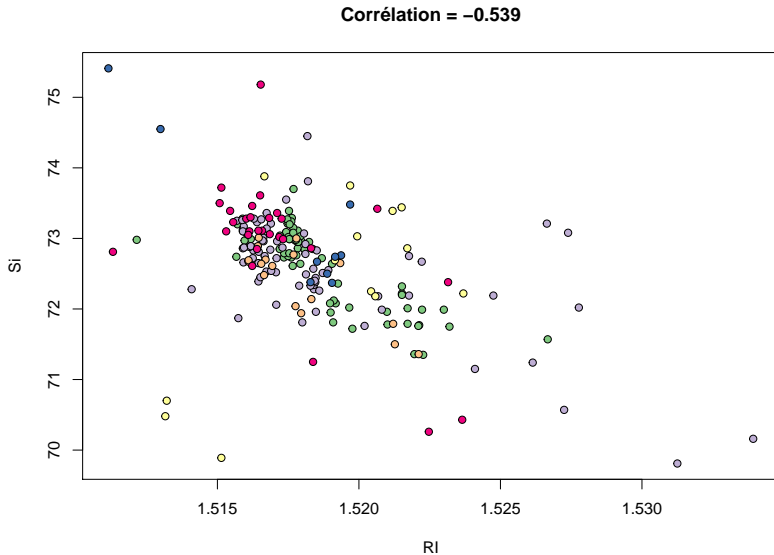
- Recherche de corrélations
- Représentation graphique de la matrice de corrélation :
  - rouge : forte corrélation positive
  - bleu : forte corrélation négative
- Ici :
  - RI corrélé avec Ca
  - Mg anti-corrélé avec Al
  - RI anti-corrélé avec Si
  - Aucun lien entre Al et Si

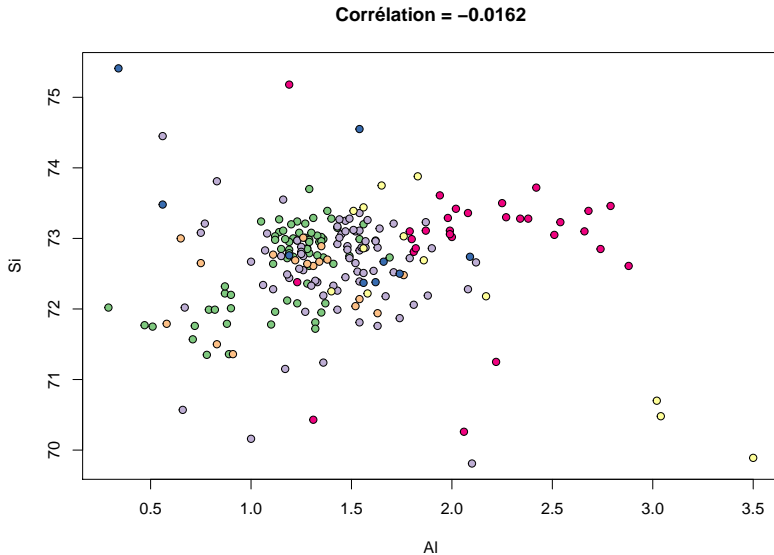




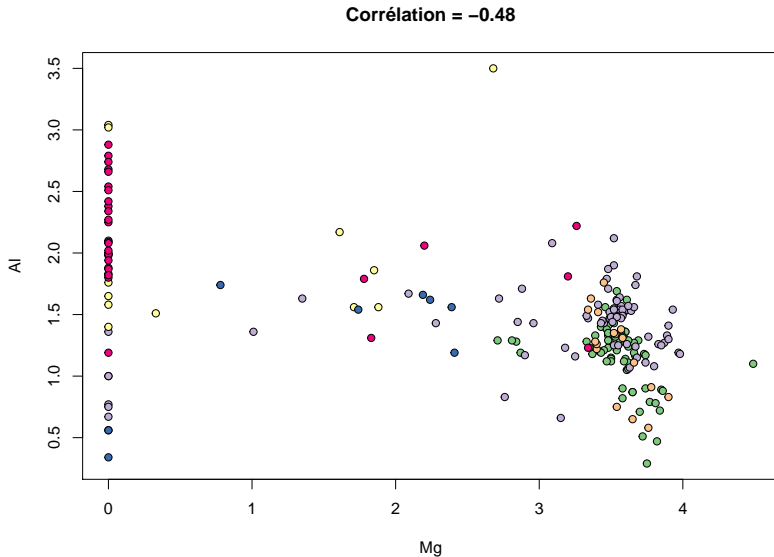


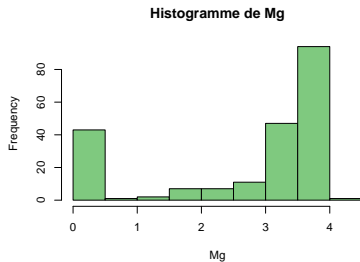
# Corrélation RI et Si



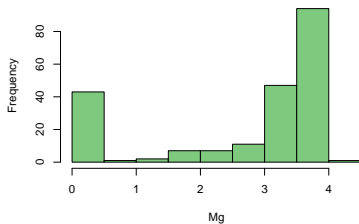


# Corrélation Mg et Al

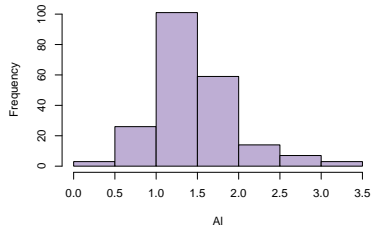




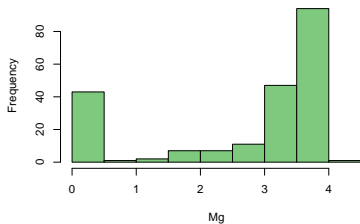
Histogramme de Mg



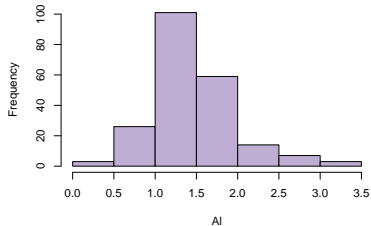
Histogramme de Al



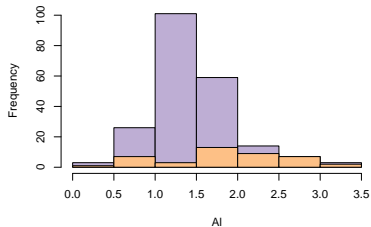
Histogramme de Mg



Histogramme de Al

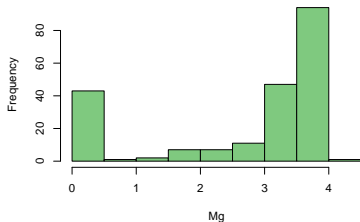


Histogramme de Al

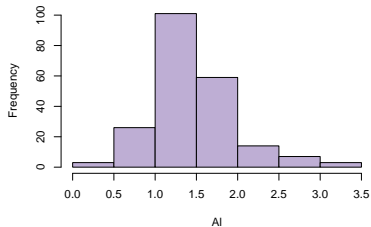




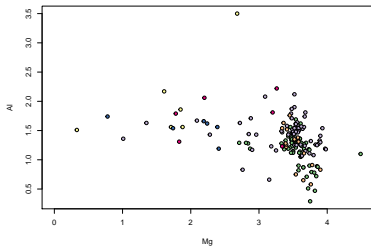
Histogramme de Mg



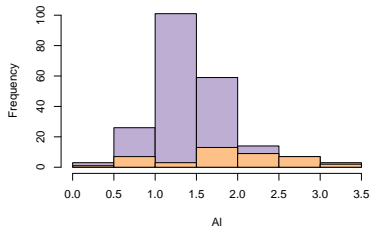
Histogramme de Al



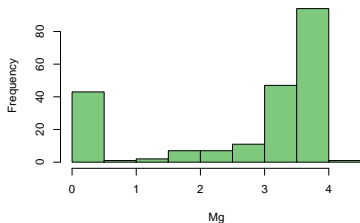
Corrélation = -0.367



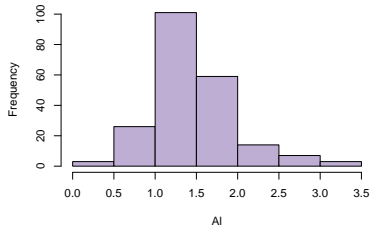
Histogramme de Al



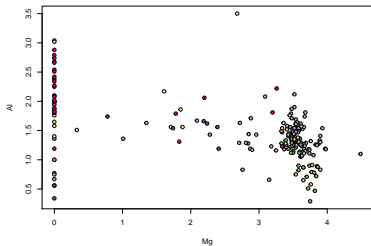
Histogramme de Mg



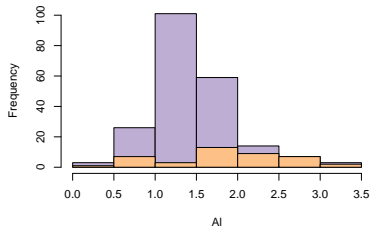
Histogramme de Al



Corrélation = -0.48

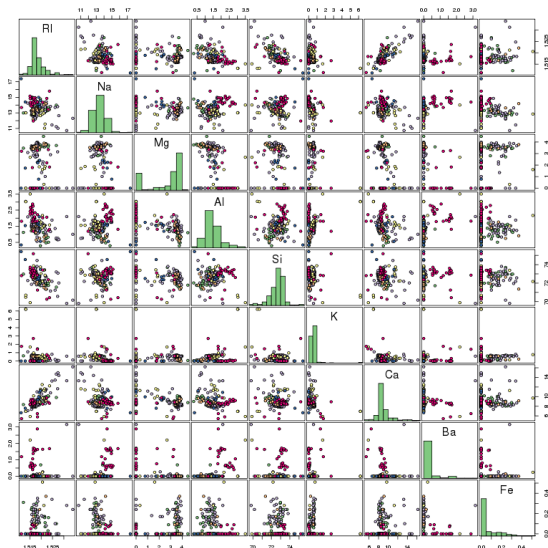


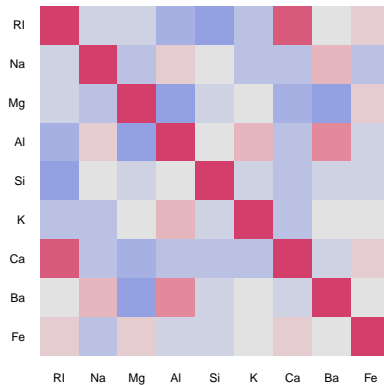
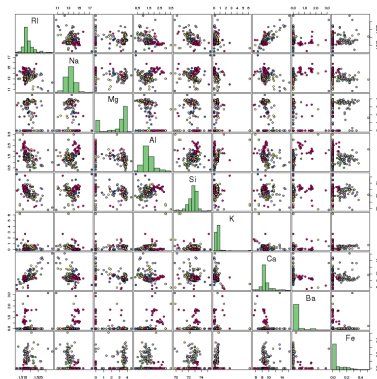
Histogramme de Al





# Vision globale

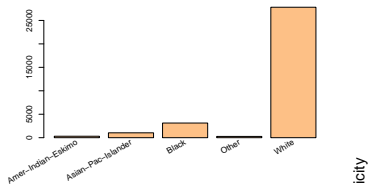




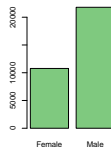
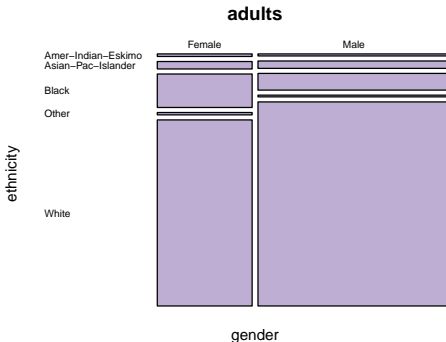


## Mosaic plot

Équivalent du *scatter plot* pour les variables qualitatives



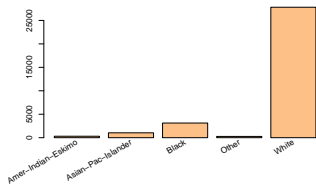
- découpage récursif
- surface proportionnelle à la fréquence



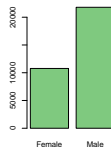
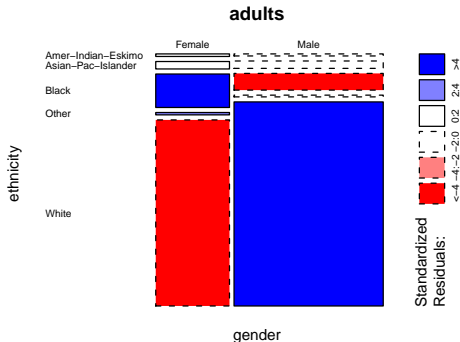


## Mosaic plot

Équivalent du *scatter plot* pour les variables qualitatives



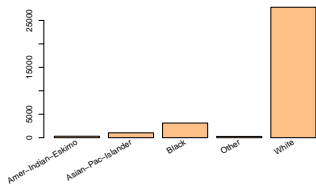
- découpage récursif
- surface proportionnelle à la fréquence
- significativité



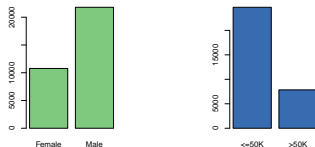
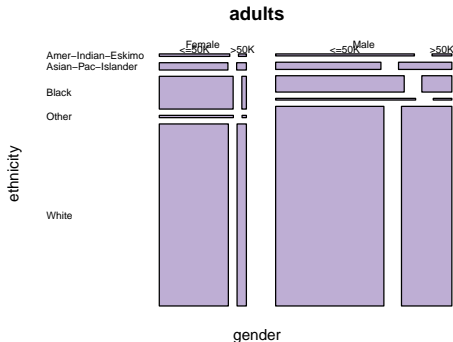


# Mosaic plot

Équivalent du *scatter plot* pour les variables qualitatives



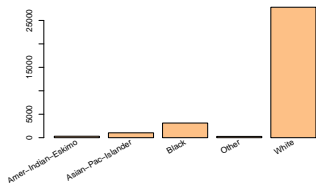
- découpage récursif
- surface proportionnelle à la fréquence
- significativité
- plus de 2 variables



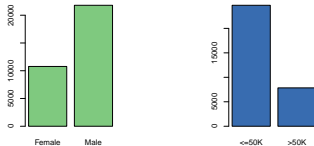
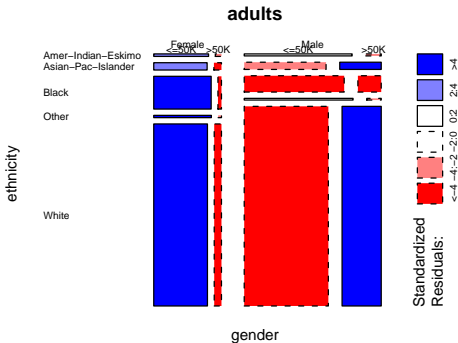


# Mosaic plot

Équivalent du *scatter plot* pour les variables qualitatives



- découpage récursif
- surface proportionnelle à la fréquence
- significativité
- plus de 2 variables

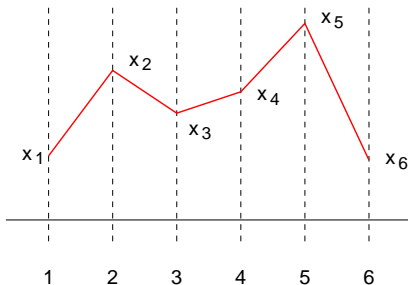


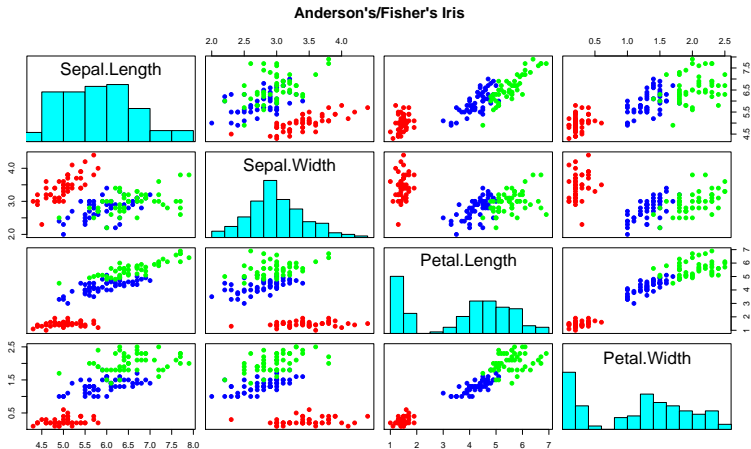


# Coordonnées parallèles

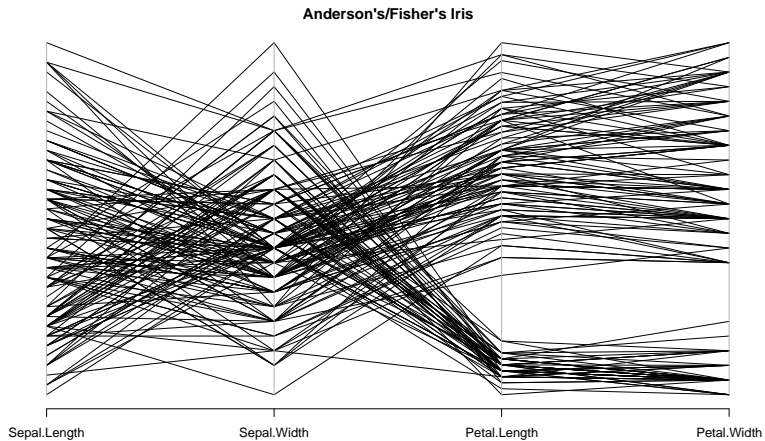
Méthode proposée en 1985 par A. Inselberg

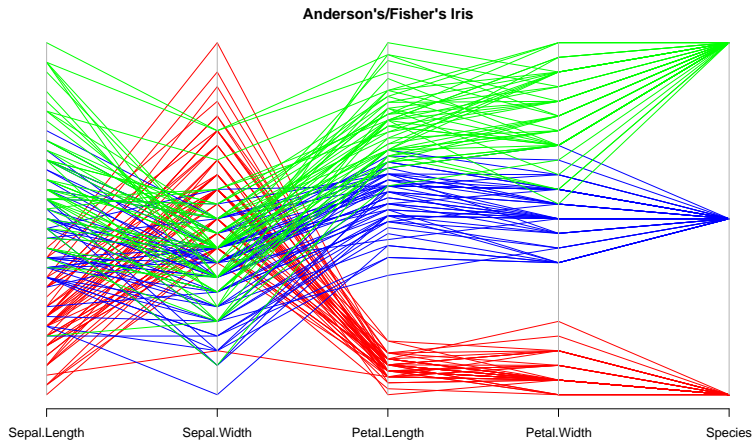
- un axe vertical par variable
- un objet devient une ligne brisée
- $(x_1, \dots, x_p)$  est représenté par la ligne brisée passant par  $(1, x_1), (2, x_2), \dots, (p, x_p)$

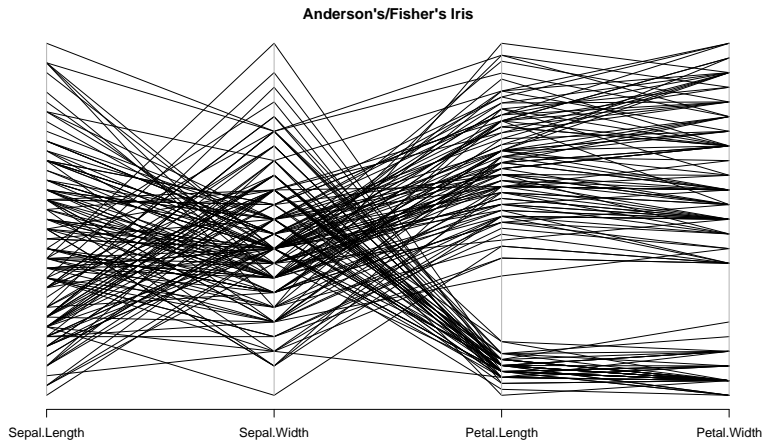




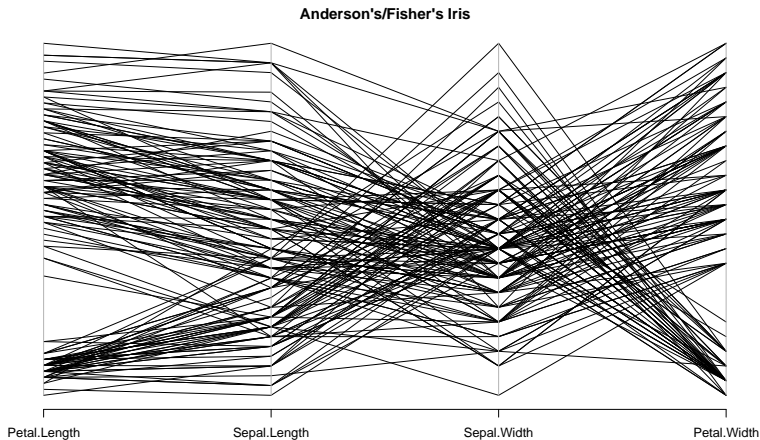
4+1 variables, 150 objets







Les variables Petal sont elles corrélées ?



Les variables Petal sont elles corrélées ?

## ■ problèmes :

- surcharge de l'écran
- surcharge cognitive

## ■ solution par interaction :

- zoom
- vues multiples
- sélection et lien :
  - sélection d'une zone (*brushing*)
  - affichage des résultats sur toutes les vues (*linking*)

## ■ en R

- iplots
- ggobi et rggobi

