

# Formation à l'analyse statistique

## Module 1 : Statistique élémentaire

- A** Indicateurs statistiques
- B** Représentations graphiques
- C** Notions de probabilités
- D** Estimation et test statistique
- E** Modélisation

# **A** Indicateurs statistiques

- Indicateurs pour 1 série de données
  - Position / Tendances centrale
  - Dispersion
- Indicateurs pour 2 séries de données

# Objectif

Résumer une série de chiffres  
par un ou plusieurs chiffres (mais  
pas beaucoup)

# Indicateurs de tendance centrale

**Moyenne** : somme des observations divisées par le nombre d'observations  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

**Médiane** : Valeur qui sépare l'échantillon en 2 sous-ensembles de tailles égales.

**Mode** : Valeur la plus fréquente dans un échantillon.

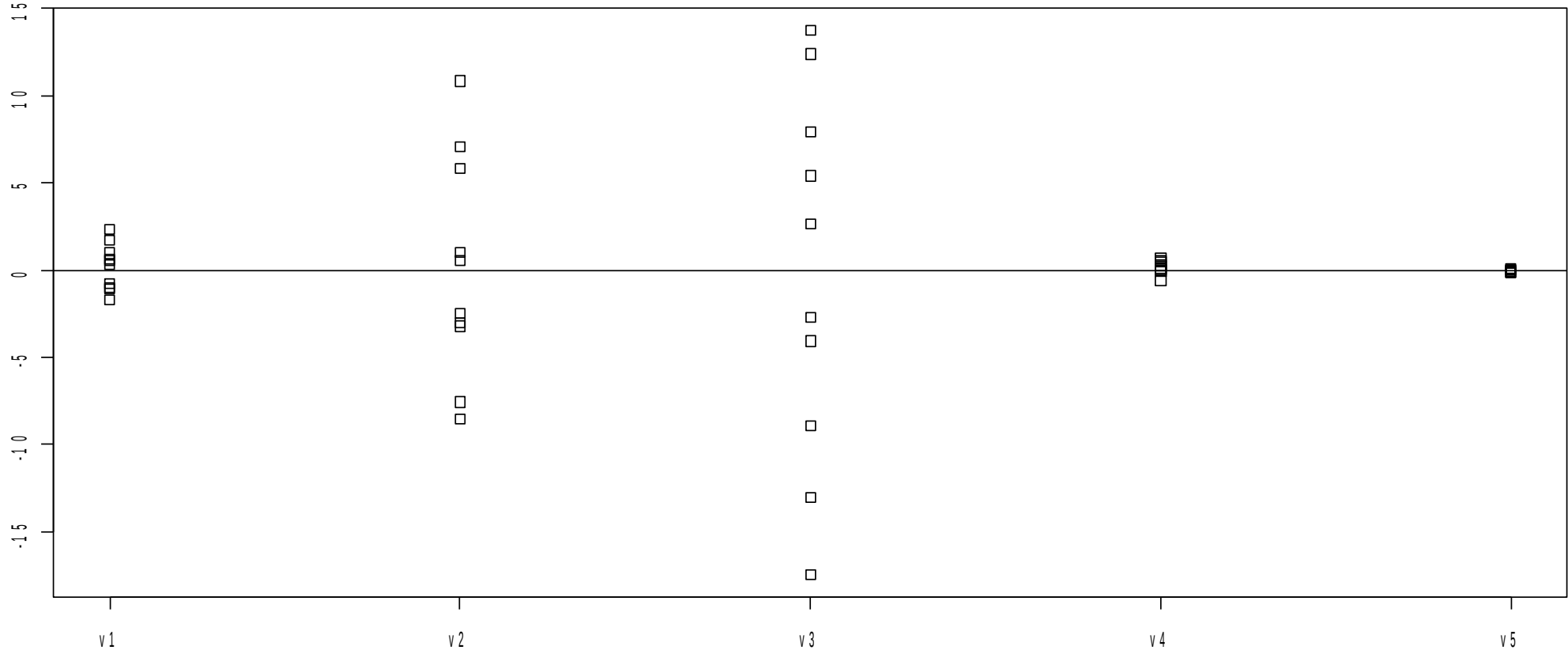
**Quartiles** : **3** valeurs qui sépare l'échantillon en **4** sous-ensembles de tailles égales.

**Déciles** : **9** valeurs qui sépare l'échantillon en **10** sous-ensembles...

**Percentiles** : **99** valeurs qui sépare l'échantillon en **100** sous-ensembles...

**Quantiles** : généralise les précédents

# Limites des indicateurs de position



Les 5 séries de 10 observations représentées ici ont, à peu près, les mêmes moyennes et médianes mais on voit bien qu'elles ne se « ressemblent » pas. Les valeurs de chaque série sont plus ou moins dispersées.

# Indicateurs de dispersion

## Variance

$$\text{var}(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

## Écart-type

$$\sigma(X) = \sqrt{\text{var}(X)}$$

.....

**Étendue** : différence entre la plus grande et la plus petite valeur d'un échantillon

**Espace inter-quartile** : différence entre le 1er et le 3ème quartile, correspond à l'étendue de l'échantillon privé de la moitié de ces observations (le  $\frac{1}{4}$  le plus élevé et le  $\frac{1}{4}$  le plus faible)

# Variance et écart-type

$$\text{var}(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Moyenne des carrés des écarts à la moyenne

$$\sigma(X) = \sqrt{\text{var}(X)}$$

Racine carrée de la variance

Quelques propriétés de l'écart-type :

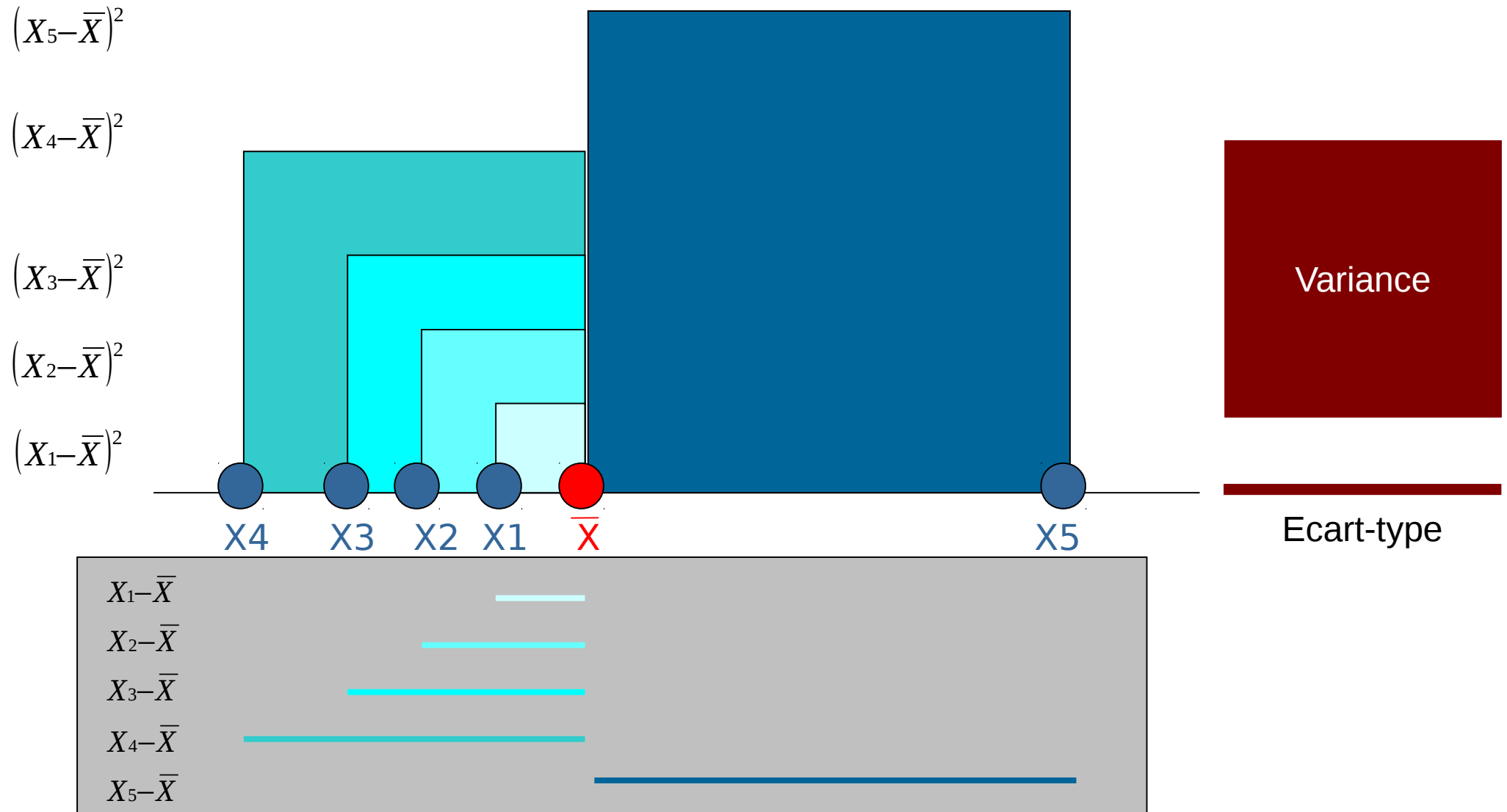
- Positif (nul si la série est constante)
- Invariant par translation
- Sensible aux valeurs extrêmes
- **De la même unité que la donnée** (et que la moyenne) :

*Si l'échantillon est constitué de mesures en  $m$  alors l'écart-type s'exprime également en  $m$  (tout comme la moyenne) ; ce qui n'est pas le cas de la variance  $m^2$  !*

On peut ainsi additionner moyenne et écart-type (*mais pas moyenne et variance*), ce qui est fondamental pour la construction d'intervalle de confiance.

# Variance et écart-type

Racine carrée de la moyenne des carrés des écarts à la moyenne

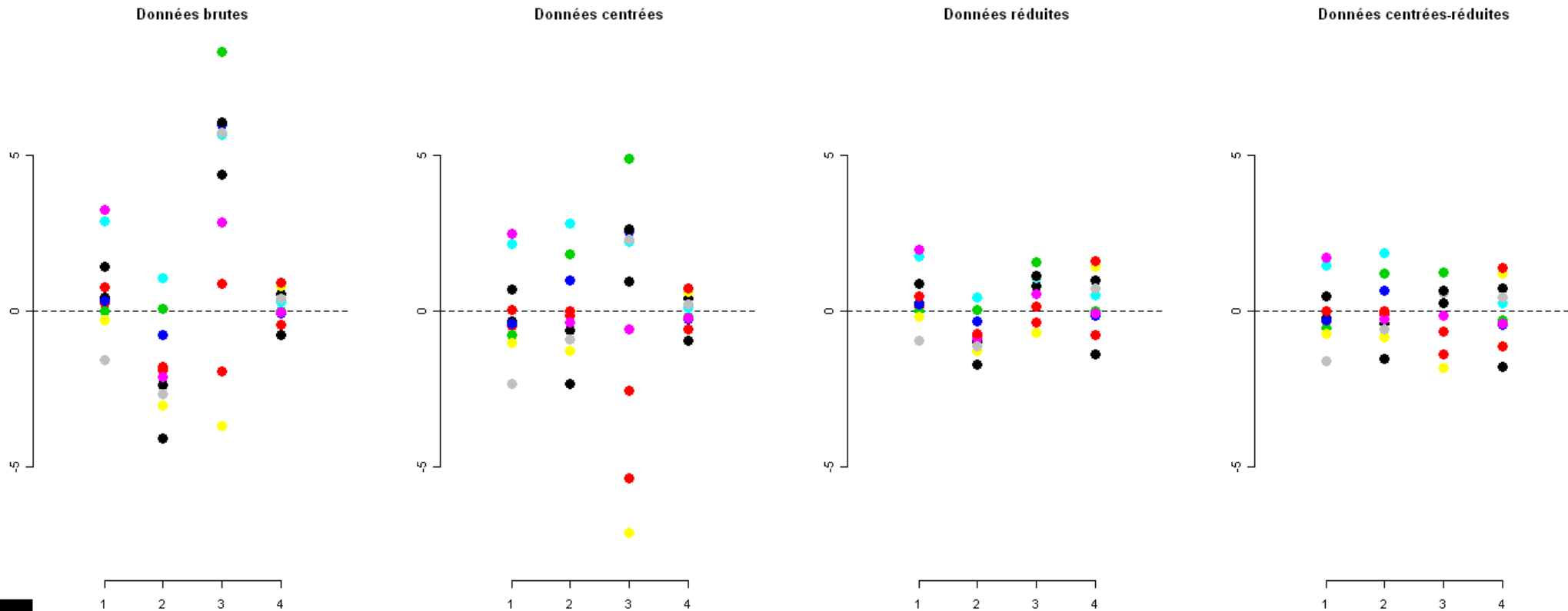




# Centrage-Réduction

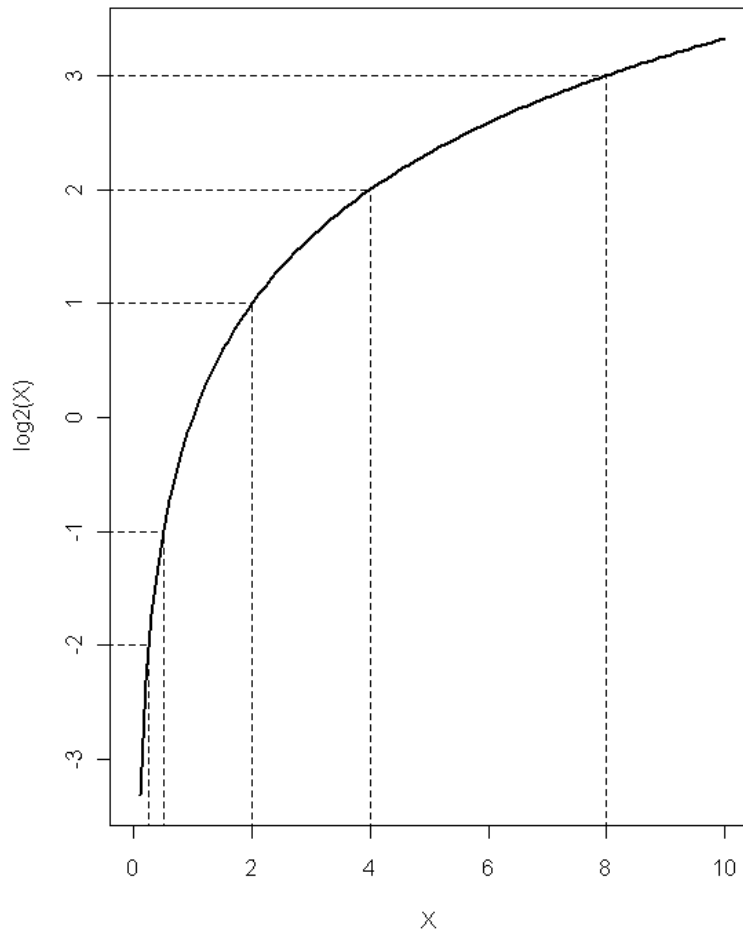
$$Z_i = \frac{X_i - \bar{X}}{\sigma_X}$$

- **Centrer** : retrancher la moyenne
- **Réduire** : diviser par l'écart-type
- Permet d'exprimer des variables différentes sur une échelle commune, en les débarrassant de leurs unités physiques : les observations s'expriment en **nombre d'écart-type par rapport à la moyenne**.
- Après centrage-réduction, la moyenne des observations est nulle et l'écart-type vaut 1 (ainsi que la variance).
- Appelé parfois « z-transformation » ou « z-score »



# Conversion « log »

<b>X</b>	0,125 = $2^{-3}$	0,25 = $2^{-2}$	0,5 = $2^{-1}$	1 = $2^0$	2 = $2^1$	4 = $2^2$	8 = $2^3$
<b><math>\log_2(X)</math></b>	-3	-2	-1	0	1	2	3



- Utile pour la conversion de ratio
- Exemple : Une sur-expression double et une sous-expression de moitié sont rendues symétriques par rapport à 0
- Pour les p-values, on aura plus intérêt à utiliser  $\log_{10}$ .
- La fonction réciproque de « log » est la fonction puissance :

$$Y = \log_2(X) \leftrightarrow X = 2^Y$$

$$Y = \log_{10}(X) \leftrightarrow X = 10^Y$$

$$Y = \ln(X) \leftrightarrow X = e^Y = \exp(Y)$$

# Indicateurs statistiques 2D

Covariance

$$\text{cov}(X,Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

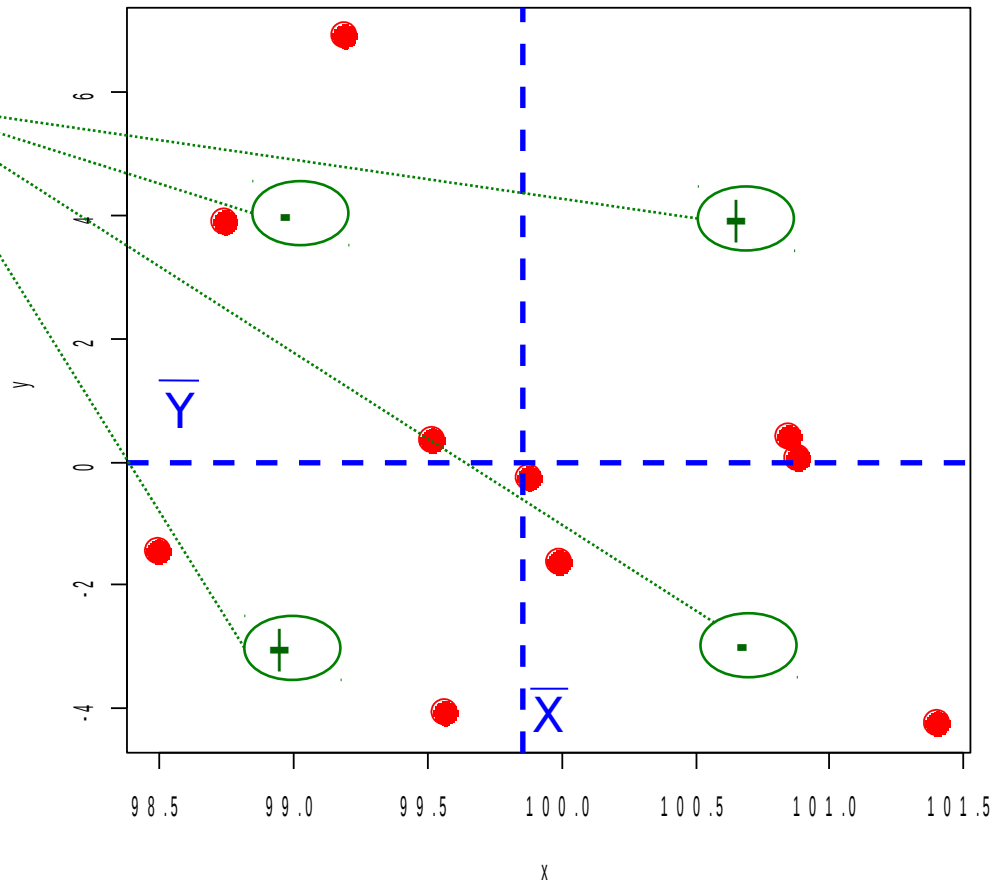
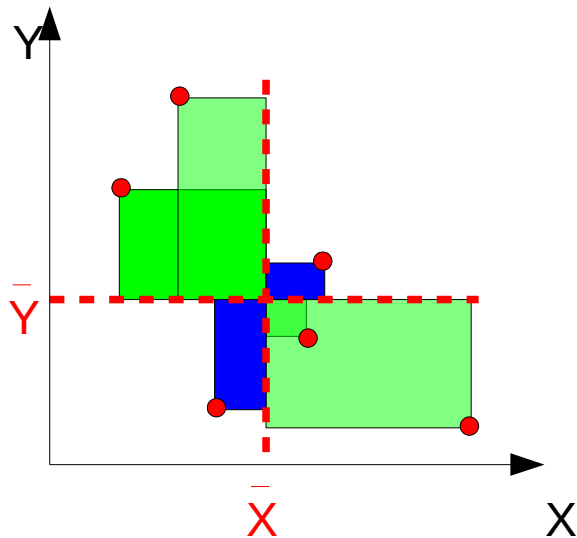
$$\text{cov}(X,X) = \text{var}(X)$$

Signe du produit  $(X_i - \bar{X})(Y_i - \bar{Y})$

Intuitivement :

- Si les + l'emportent  
→ liaison linéaire positive
- Si les - l'emportent  
→ liaison linéaire négative

Sur cet exemple :  $\text{cov}(X,Y) = -1.36$



La covariance dépend des unités de mesure →  
coefficient de corrélation

# Coefficient(s) de corrélation

- Coefficient de corrélation **linéaire** de Pearson

$$\rho(X, Y) = \text{cov}(X, Y) / (\sigma_X \sigma_Y)$$

- Coefficient de corrélation de Spearman  
 → Robustesse due au travail sur les rangs

$$\rho_s(X, Y) = \rho(RX, RY) = 1 - 6 \frac{\sum_{i=1}^n d^2}{n(n^2 - 1)}$$

**(1)** calcul des rangs

**(2)** différence des rangs

**(3)** carrés des différences des rangs

**(4)** somme des carré des différences des rangs

			<b>(1)</b>	<b>(2)</b>	<b>(3)</b>	
X	Y		RX	RY	d = RX-RY	d <sup>2</sup>
20,6	20,7		6	7	-1	1
21,6	21,8		2	2	0	0
18,8	20,4		9	8	1	1
20,8	21,1		3	4	-1	1
17,5	18,3		10	10	0	0
19,5	18,9		7	9	-2	4
20,8	21,1		4	4	0	0
20,6	21,2		5	3	2	4
19,2	20,9		8	6	2	4
22,2	22,9		1	1	0	0
					<b>Somme</b>	<b>15</b>

**(4)**

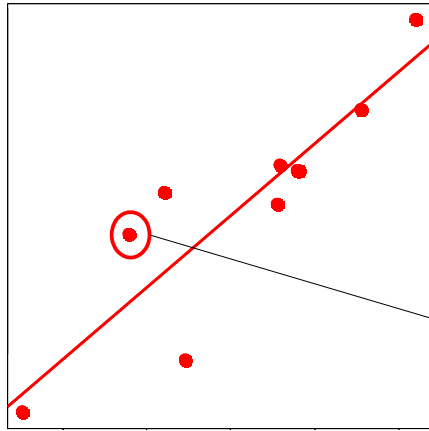
# Coefficient(s) de corrélation

Quelques propriétés des coefficients de corrélation :

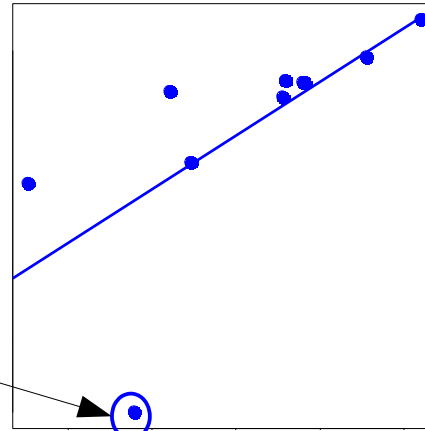
- Coefficient de corrélation de Pearson : **relation linéaire**
- Coefficient de corrélation de Spearman : considère les rangs, relation "**monotone**"
- Compris entre  $-1$  et  $1$ .
- Les valeurs extrêmes  $-1$  et  $1$  indique des corrélations parfaites entre les 2 variables.
- Si le coefficient est **positif** : quand une variable est élevée, l'autre l'est également. Quand une variable est faible, l'autre l'est également.
- Si le coefficient est **négatif** : quand une variable est élevée (resp. faible), l'autre est faible (resp. élevée).

# Corrélation : exemples

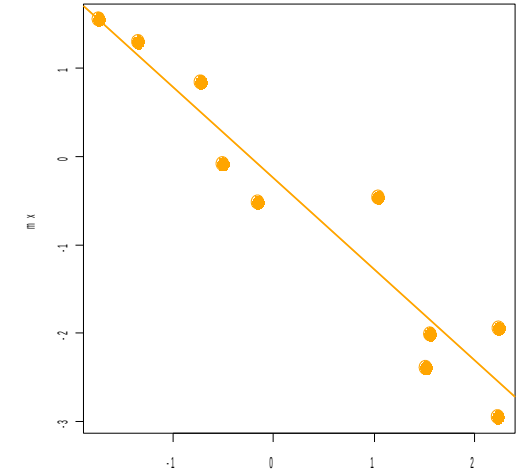
$\rho$  : Pearson -  $\rho_s$  : Spearman



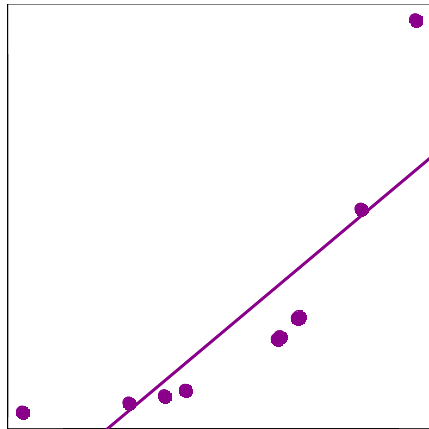
$$\rho = 0.884 - \rho_s = 0.9$$



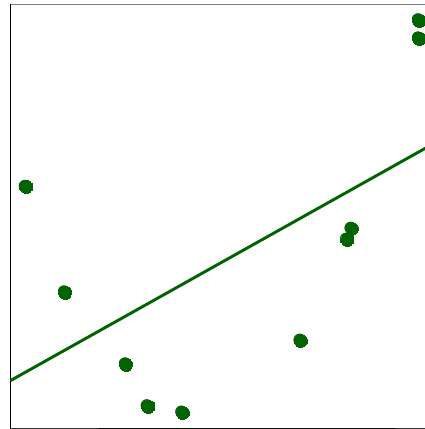
$$\rho = 0.676 - \rho_s = 0.912$$



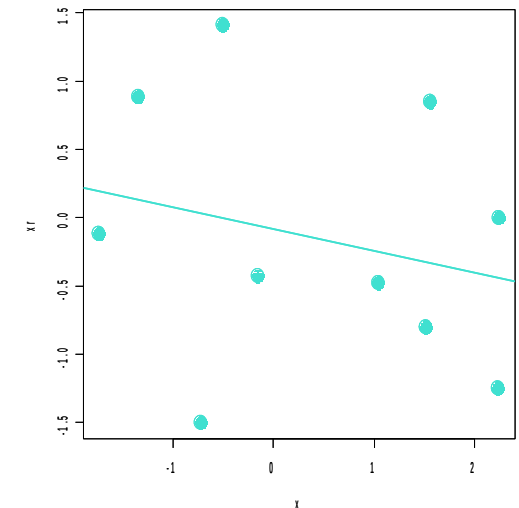
$$\rho = -0.954 - \rho_s = -0.903$$



$$\rho = 0.822 - \rho_s = 1$$



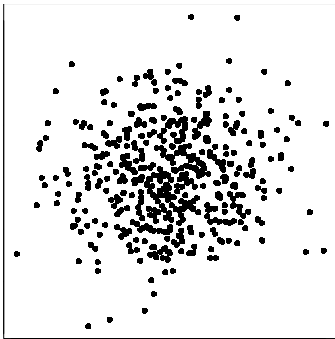
$$\rho = 0.584 - \rho_s = 0.491$$



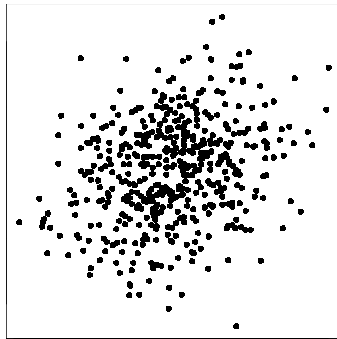
$$\rho = -0.248 - \rho_s = -0.164$$

# Corrélation : exemples

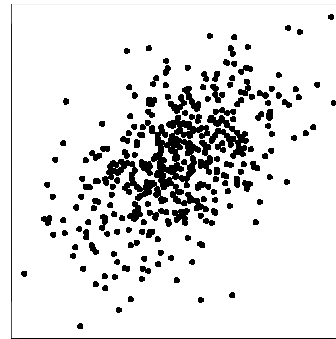
0.05



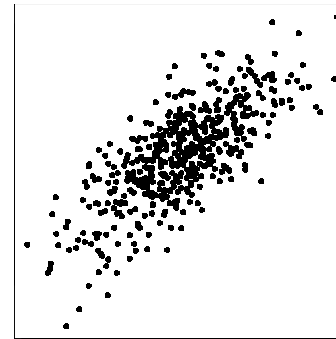
0.26



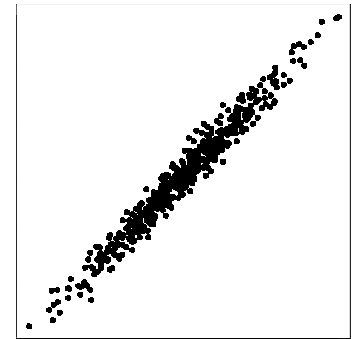
0.49



0.75



0.98



Valeurs  
significatives du  
coefficient de  
corrélacion

$v / \alpha$	0.10	0.05	0.02	$v / \alpha$	0.10	0.05	0.02
<b>1</b>	0.9877	0.9969	0.9995	<b>16</b>	0.4000	0.4683	0.5425
<b>2</b>	0.9000	0.9500	0.980	<b>17</b>	0.3887	0.4555	0.5285
<b>3</b>	0.8054	0.8783	0.9343	<b>18</b>	0.3783	0.4438	0.5155
<b>4</b>	0.7293	0.8114	0.8822	<b>19</b>	0.3687	0.4329	0.5034
<b>5</b>	0.6694	0.7545	0.8329	<b>20</b>	0.3598	0.4227	0.4921
<b>6</b>	0.6215	0.7067	0.7887	<b>25</b>	0.3233	0.3809	0.4451
<b>7</b>	0.5822	0.6664	0.7498	<b>30</b>	0.2960	0.3494	0.4093
<b>8</b>	0.5494	0.6319	0.7155	<b>35</b>	0.2746	0.3246	0.3810
<b>9</b>	0.5214	0.6021	0.6851	<b>40</b>	0.2573	0.3044	0.3578
<b>10</b>	0.4973	0.5750	0.6581	<b>45</b>	0.2428	0.2875	0.3384
<b>11</b>	0.4762	0.5529	0.6339	<b>50</b>	0.2306	0.2732	0.3218
<b>12</b>	0.4575	0.5324	0.6120	<b>60</b>	0.2108	0.2500	0.2948
<b>13</b>	0.4409	0.5139	0.5923	<b>70</b>	0.1954	0.2319	0.2737
<b>14</b>	0.4259	0.4973	0.5742	<b>80</b>	0.1829	0.2172	0.2565
<b>15</b>	0.4124	0.4821	0.5577	<b>90</b>	0.1726	0.2050	0.2422
				<b>100</b>	0.1638	0.1946	0.2301

# Corrélation $\neq$ causalité

Quelques exemples d'événements corrélés n'impliquant pas forcément une causalité :

- Le fait de dormir avec ses chaussures est fortement corrélé avec le fait de se réveiller avec la « gueule de bois ». Doit-on en conclure que dormir avec ses chaussures donne la « gueule de bois » ? Ou un troisième facteur est-il impliqué ?
- La fréquence des attaques de requins est fortement corrélée à la vente de glaces sur les plages ! Manger des glaces rend-il plus appétissant pour les requins ?
- Les personnes qui meurent ont très fréquemment vu un médecin dans les jours qui ont précédé. Est-il si dangereux de rencontrer un médecin ?
- Dans la plupart des villes, on constate une forte corrélation positive entre le nombre de nids de cigognes et la natalité. Nous cacheraient-on des choses ?
- ...





# Un peu de théorie

Un graphique devrait :

- Montrer les données
- Inciter celui qui regarde à penser
- Éviter de distordre ce que les données ont à dire
- Présenter beaucoup de données sur une petite surface
- Encourager l'œil à comparer différents morceaux des données
- Révéler les données à des niveaux différents : d'un aperçu global à des structures plus fines
- Servir un objectif clair et raisonnable
- Être étroitement intégré à une description statistique du jeu de données

**Edward Tufte**, *The Visual Display of Quantitative Information*, Cheshire, CT, Graphics Press, 2001, 2e éd. (1re éd. 1983)

# Représentations graphiques

Données de type « effectif »

A	30
B	15
C	30
D	20
E	25

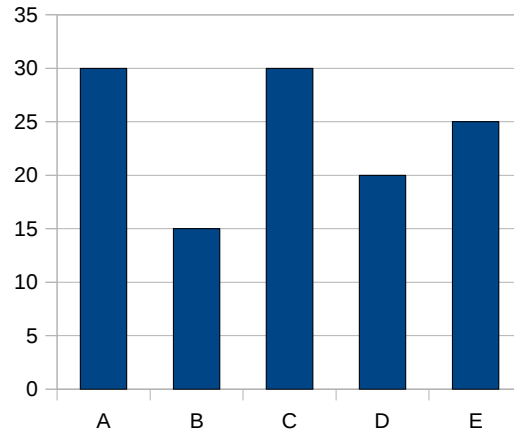
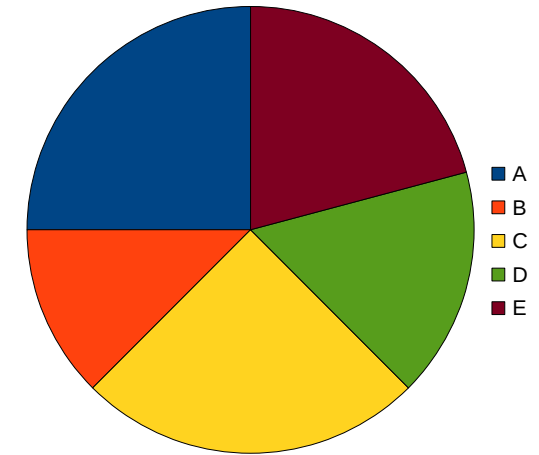
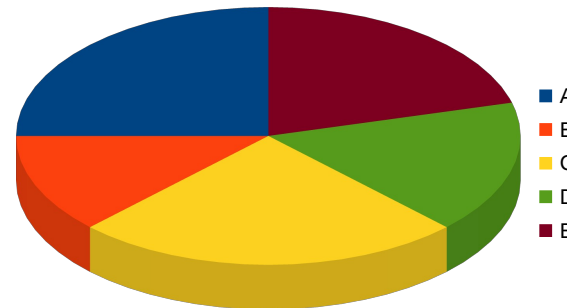
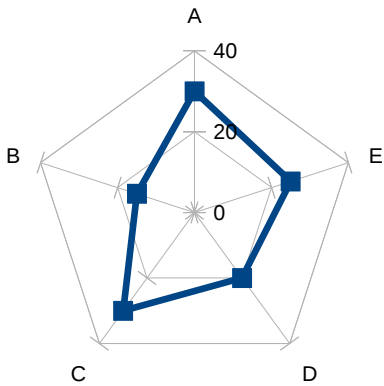


Diagramme en bâtons

Diagramme circulaire



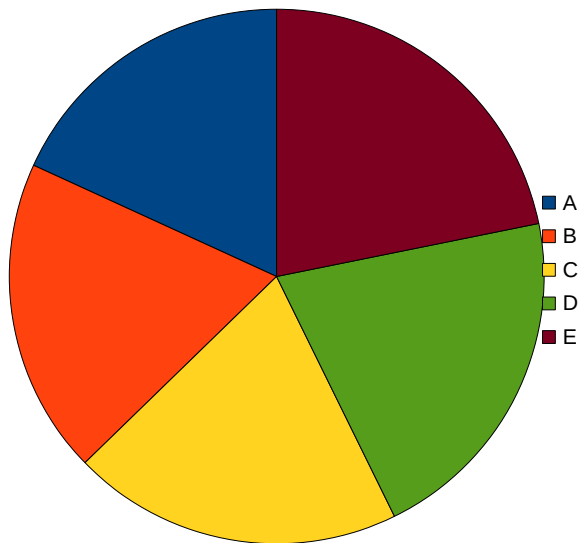
Spiderman plot :-)



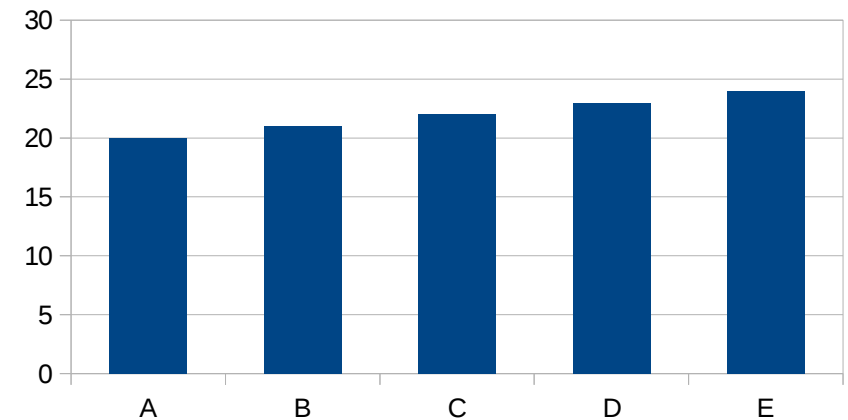
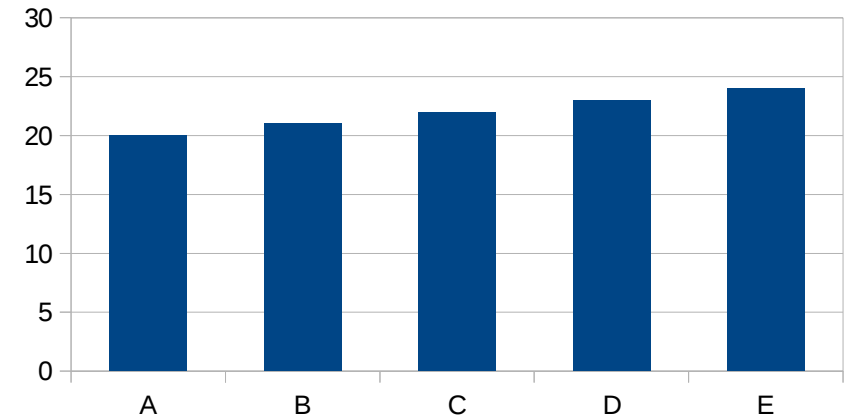
Attention aux représentations 3D. Est-ce si évident que A et C ont la même valeur ?

# Représentations graphiques

## Données de type « effectif »



A	20
B	21
C	22
D	23
E	24



# Représentations graphiques

## Données de type « effectif »

	X	Y
A	30	25
B	15	15
C	30	20
D	20	30
E	25	35

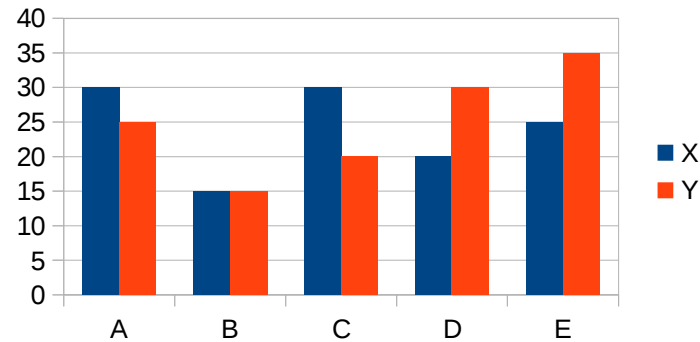
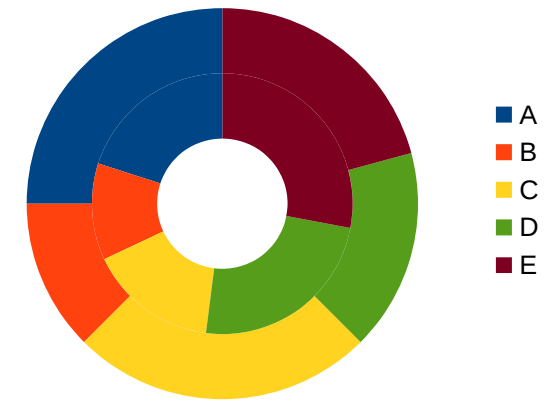
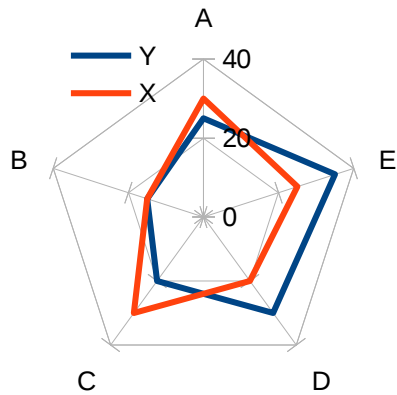
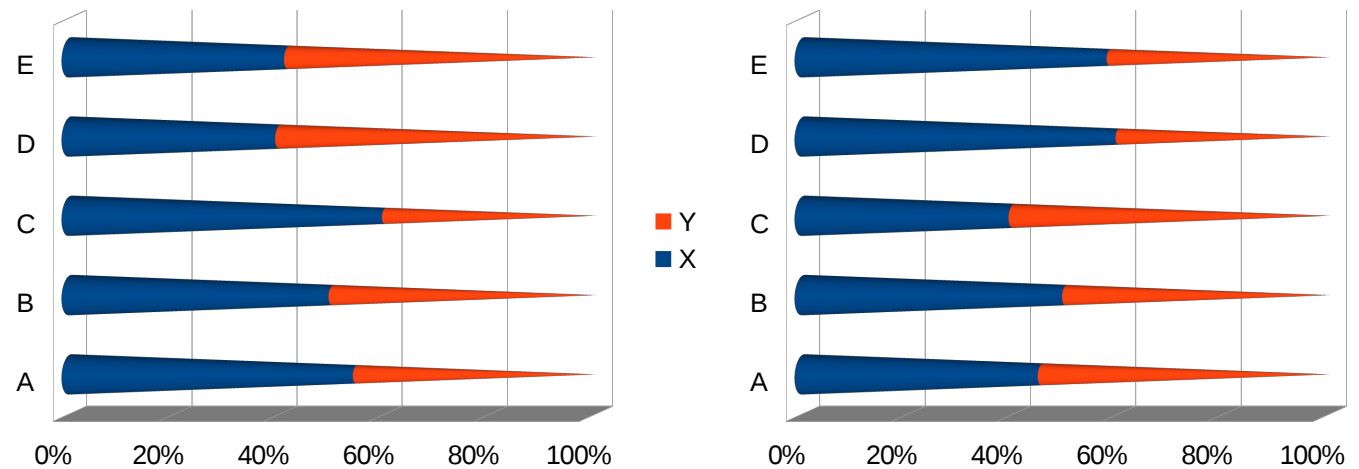


Diagramme en bâtons

Diagramme « tore »



Cylindres 3D (!?)



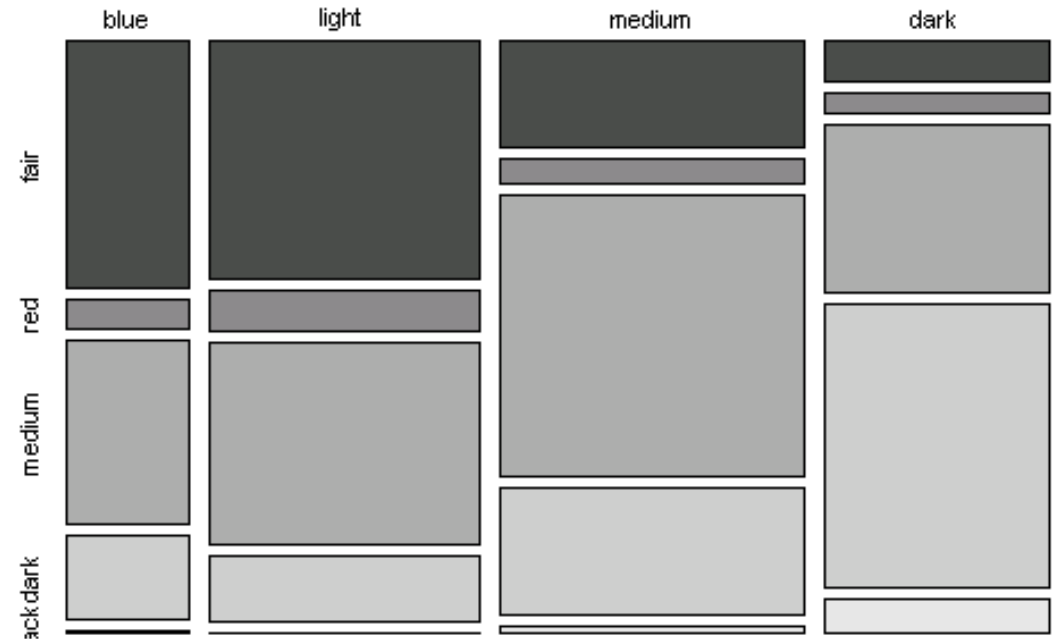
# Représentations graphiques

## Données de type « effectif » - Table de contingence

(données *caith*, R)

	fair	red	medium	dark	black
blue	326	38	241	110	3
light	688	116	584	188	4
medium	343	84	909	412	26
dark	98	48	403	681	85

Diagramme mosaïque



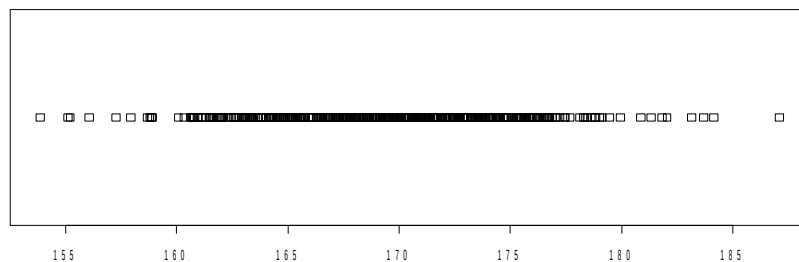
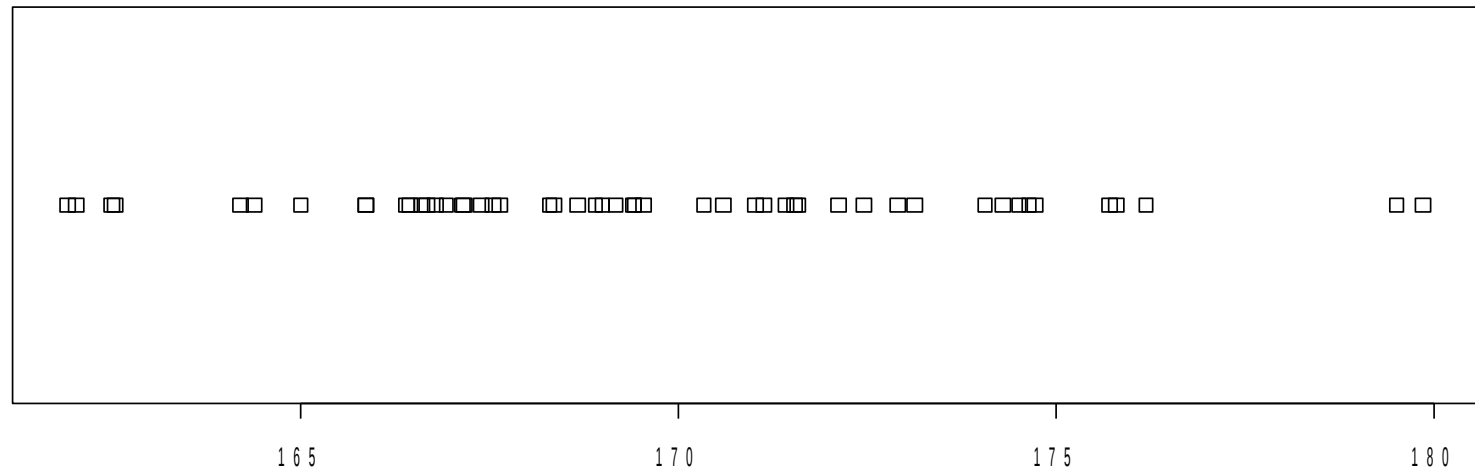
# Représentations graphiques

## Données quantitatives

50 observations : taille en cm de 50 individus

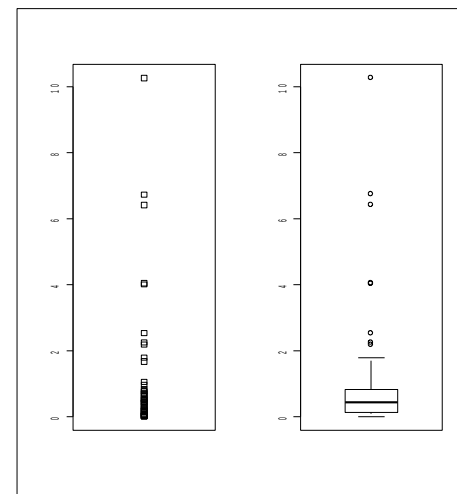
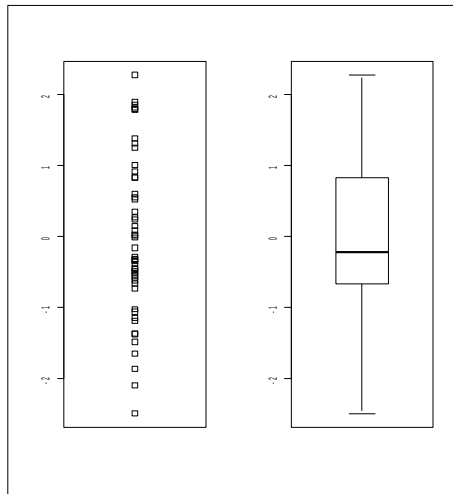
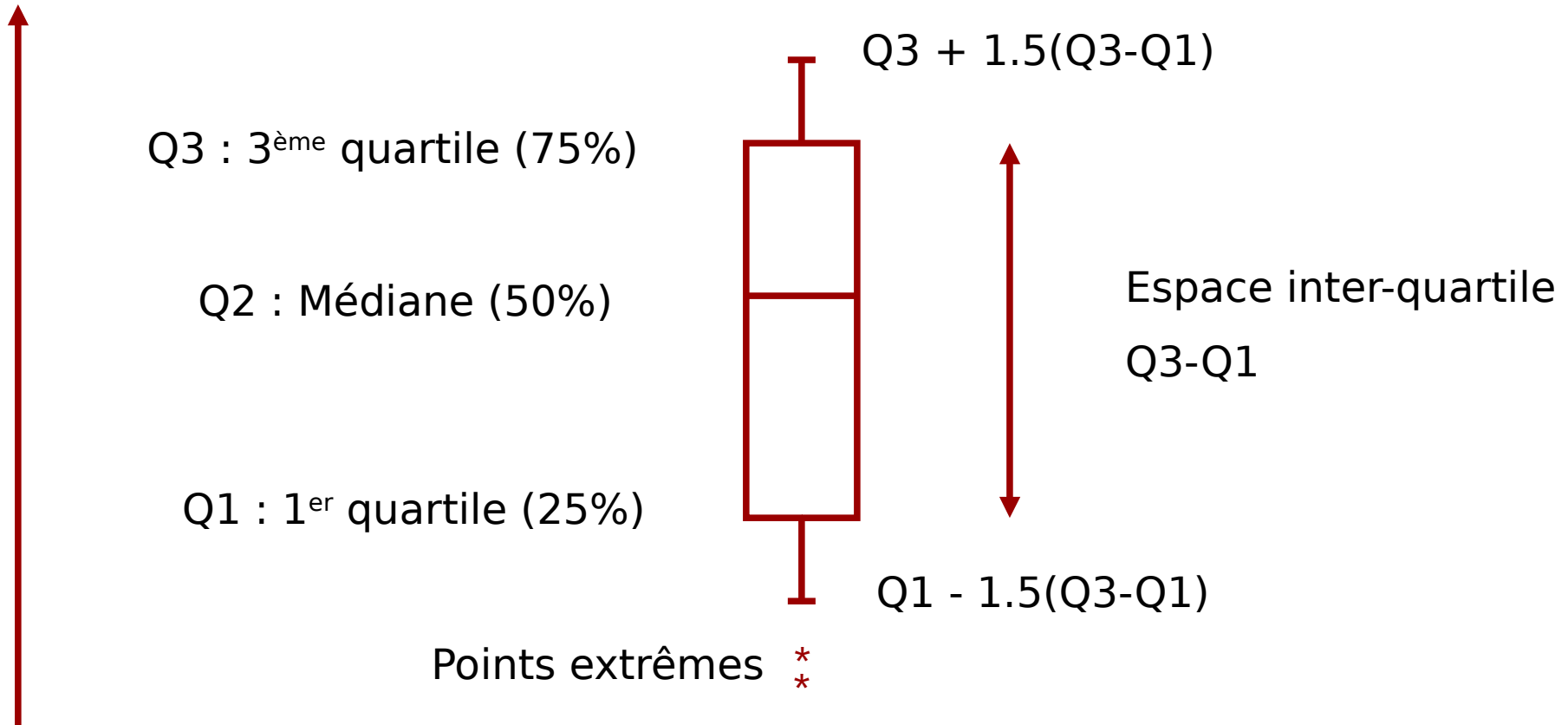
```
171.4 175.8 168.9 166.8 174.3 169.2 165.9 171.6 167.1 179.8
174.5 169.0 170.6 173.1 169.4 167.4 162.5 172.9 174.7 164.2
175.7 167.2 169.5 174.1 162.5 172.4 166.4 162.0 167.5 168.7
166.9 171.0 176.2 172.1 166.4 168.3 170.3 164.4 167.6 168.3
165.9 174.6 171.5 169.4 166.6 179.5 166.7 171.1 161.9 165.0
```

Nuage de points 1D  
(*stripchart*)



Avec 500 observations, la lisibilité devient délicate.

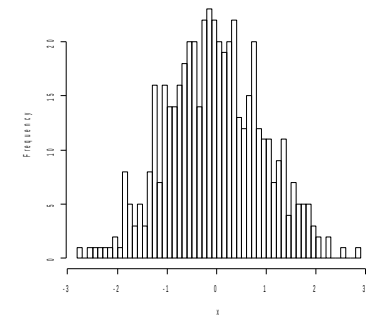
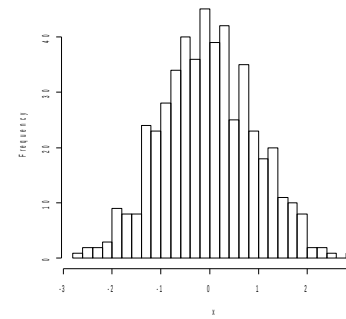
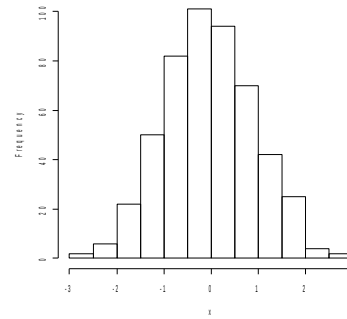
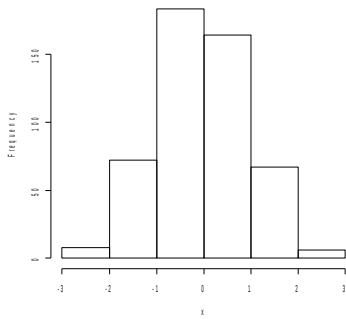
# Boxplot



# Histogramme

≠

diagramme en bâtons ☠



Histogramme

~

Densité de  
probabilité

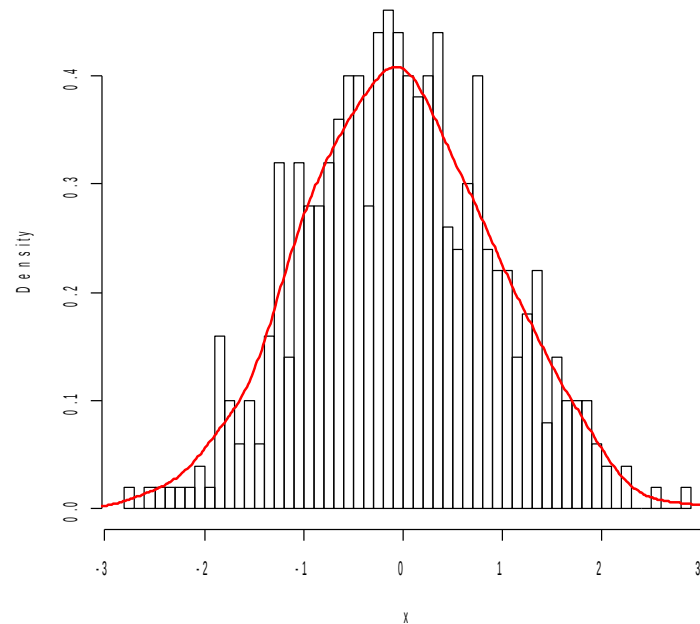
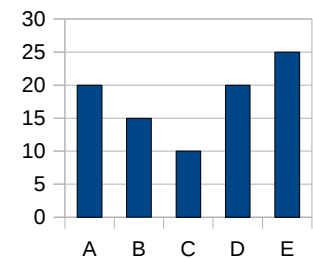
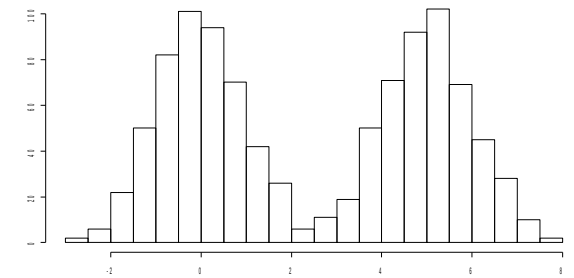
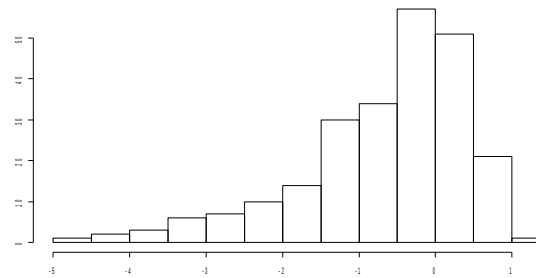
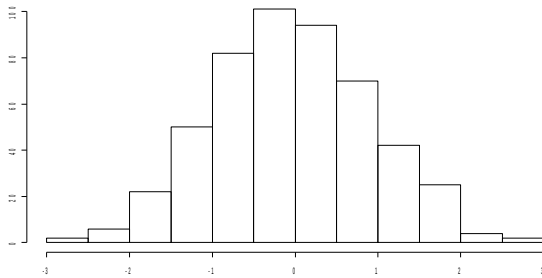
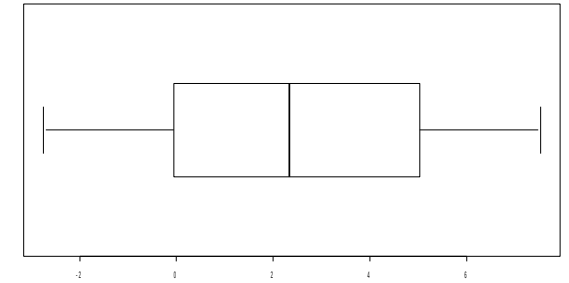
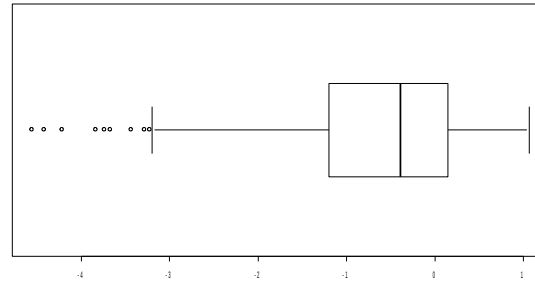
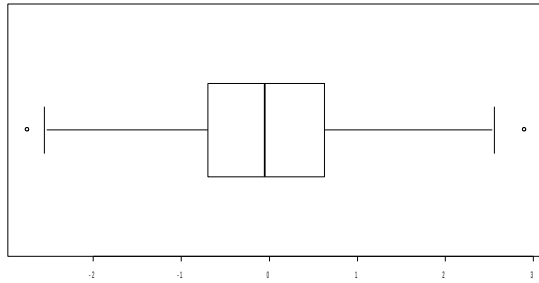


Diagramme  
en bâtons





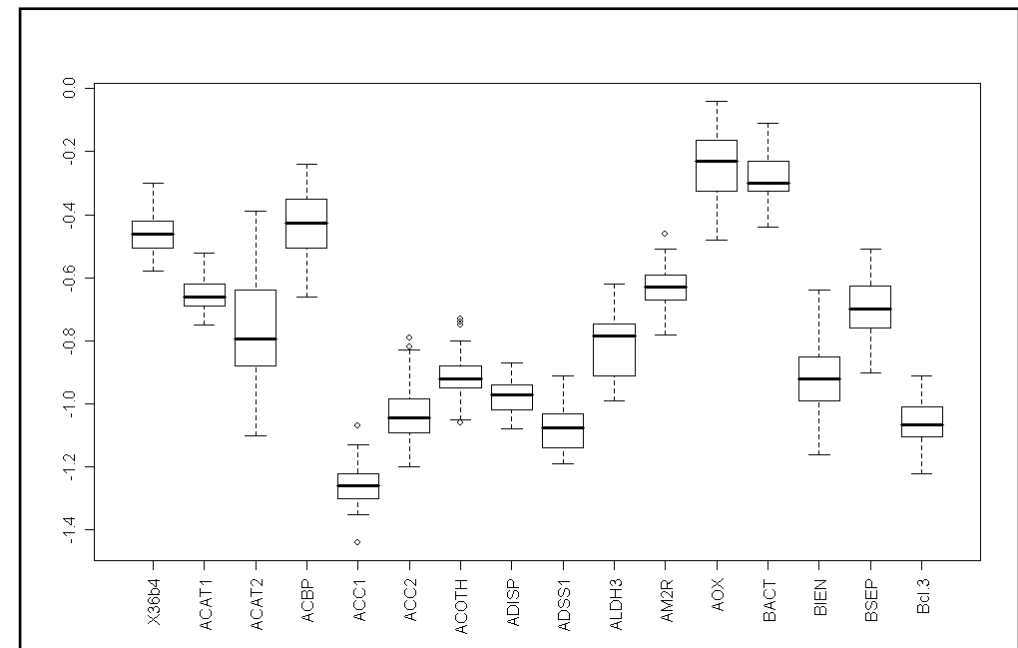
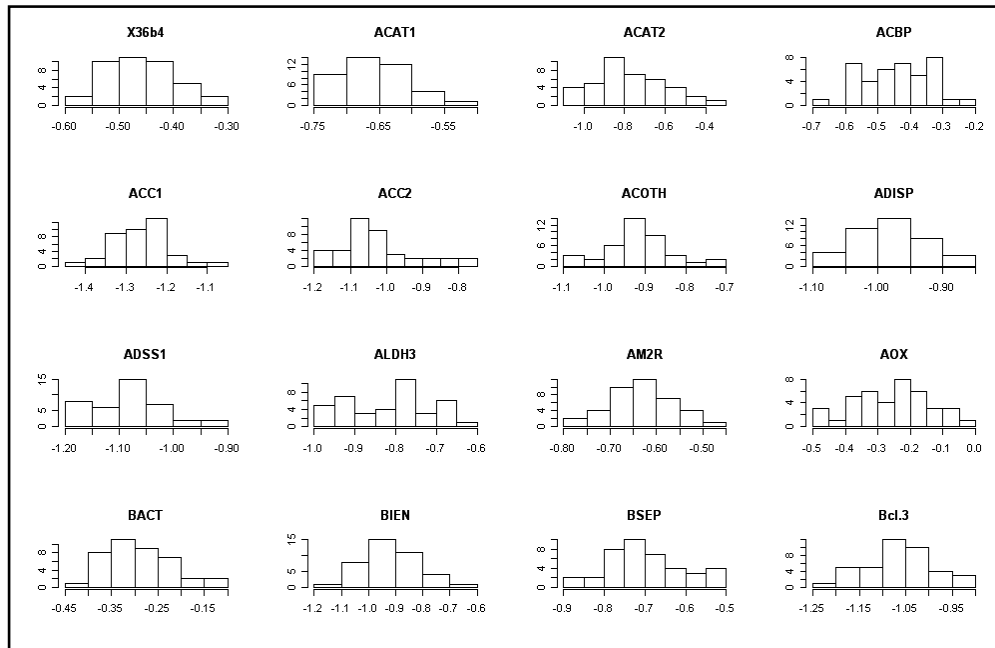
# Boxplot et histogramme



Individuellement, ces graphiques apportent la même information sur la position et la répartition des données. Attention cependant aux cas de multi-modalités que le boxplot ne peut pas capter.

⇒ *avantage histogramme ?*

# Boxplot et histogramme



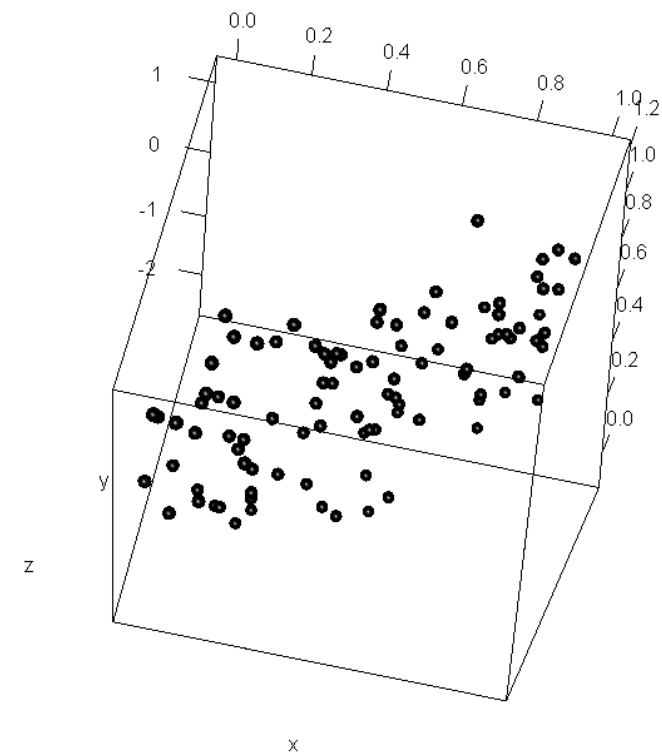
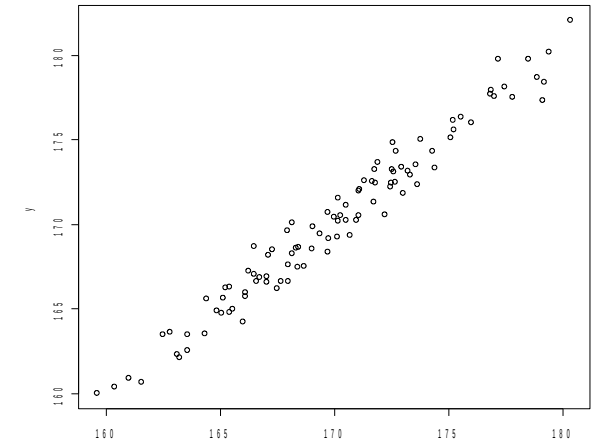
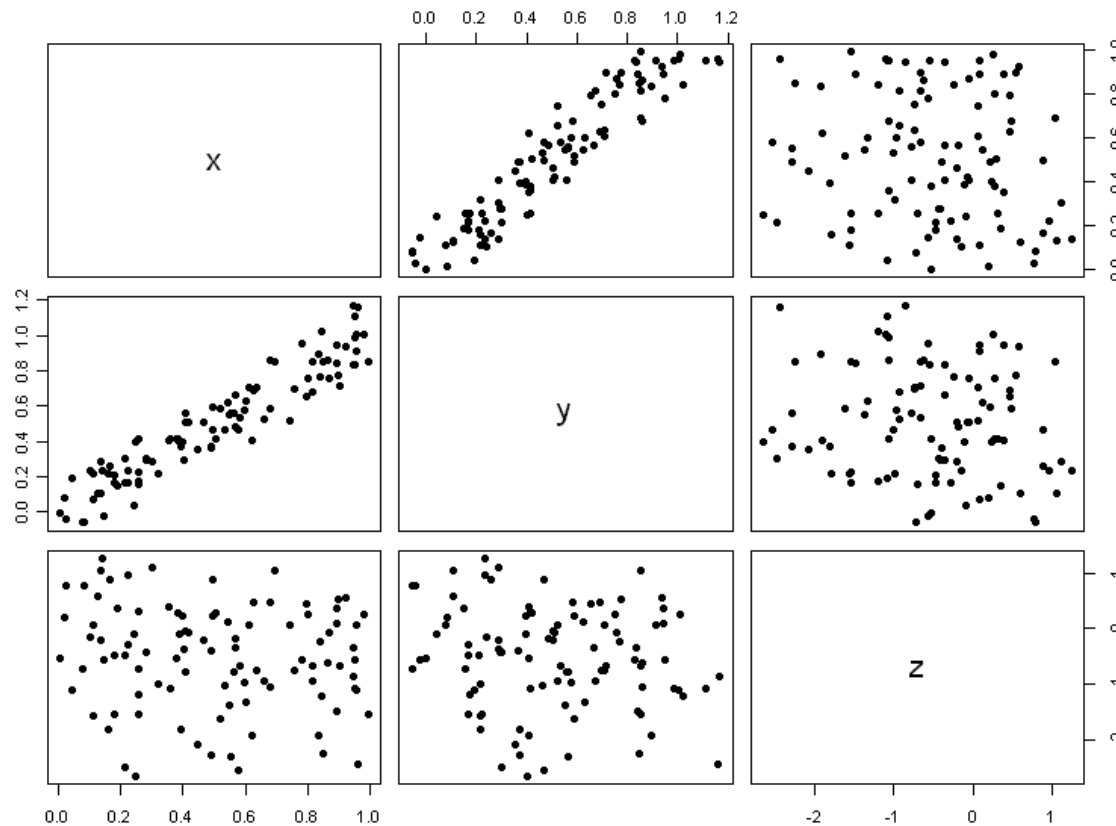
⇒ égalité ?



# Représentations graphiques

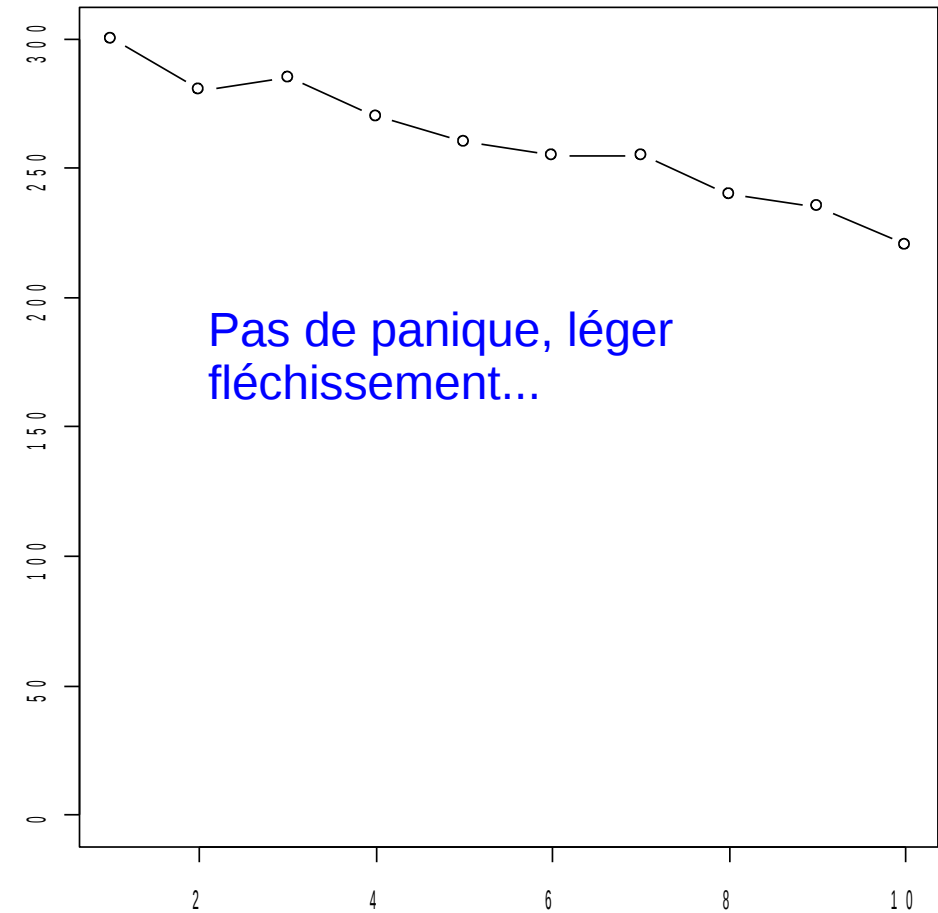
## 2 variables ou plus

Nuage de points, diagramme de dispersion (*scatterplot*)



# Échelles des axes

Temps	1	2	3	4	5	6	7	8	9	10
Chômage, CAC40, chiffre d'affaires...	300	280	285	270	260	255	255	240	235	220



# Notions de probabilité

*Probable : qui a une apparence de vérité ; dont la vérité a plus de raisons d'être confirmée que d'être infirmée. Qu'il est raisonnable de supposer, de conjecturer, de prévoir; qui a beaucoup de chances de se produire.*

- « Quantifier » des événements dont la réalisation est incertaine
- « Formule » : nombre de cas favorables / nombre de cas possibles
- Exemple classique : jeu de « pile ou face »
  - Nombre de cas possibles : 2 (pile et face)
  - Nombre de cas favorables : 1 (le côté sur lequel on a misé)
  - Probabilité de gagner (ou de perdre) :  $1/2$



# Espérance

*Espérance : disposition de l'âme qui porte l'homme à considérer dans l'avenir un bien important qu'il désire et qu'il croit pouvoir se réaliser.*

- Espérance mathématique : Valeur moyenne théorique d'une variable aléatoire

$$E(x) = x_1.P(X=x_1) + x_2.P(X=x_2) + \dots$$

Exemples :

- Gain de 10€ si le dé indique une valeur paire 0€ sinon :

Espérance de gain :  $10 \cdot \frac{3}{6} + 0 \cdot \frac{3}{6} = 5\text{€}$

- Gain de 20€ si le dé tombe sur 6 et 10€ sur 5, 0 sinon

Espérance de gain :  $20 \cdot \frac{1}{6} + 10 \cdot \frac{1}{6} + 0 \cdot \frac{4}{6} = 5\text{€}$



*Gagner peu souvent  $\Leftrightarrow$  gagner beaucoup rarement.*

# Variable aléatoire

Une **variable aléatoire** est utilisée pour modéliser le résultat d'une expérience non déterministe qui génère un résultat aléatoire.

Exemples :

- lancement d'une pièce, d'un dé (v.a. **discrète**)
- taille d'un individu pris au hasard dans une population (v.a. **continue**)

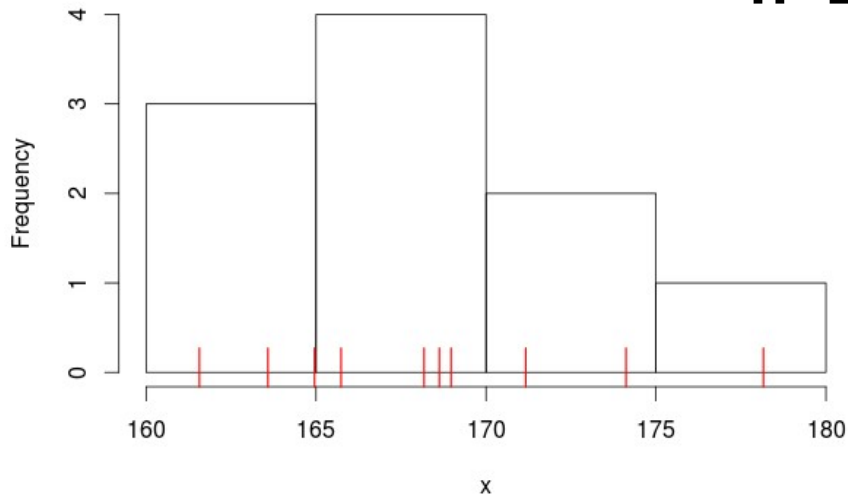
La répartition des valeurs prises par une v.a. conduit à la notion de **loi de probabilité**.



# Exemple 1 : taille d'individus

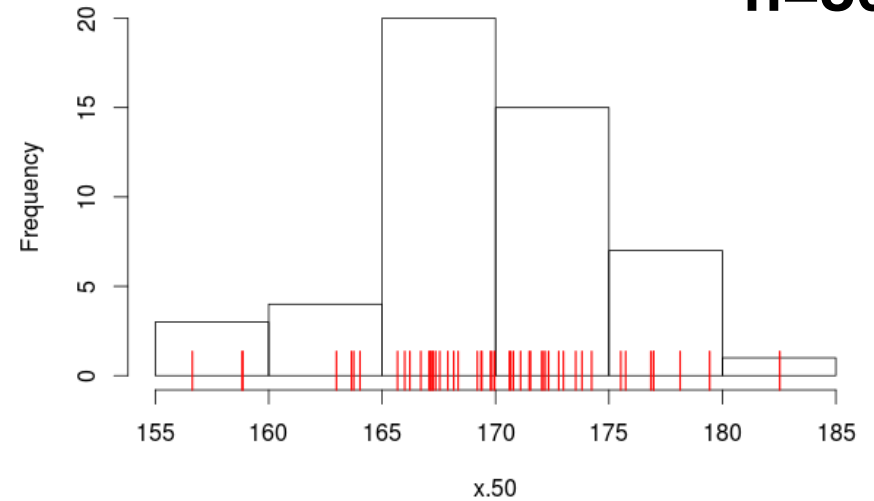
Histogram of x

**n=10**



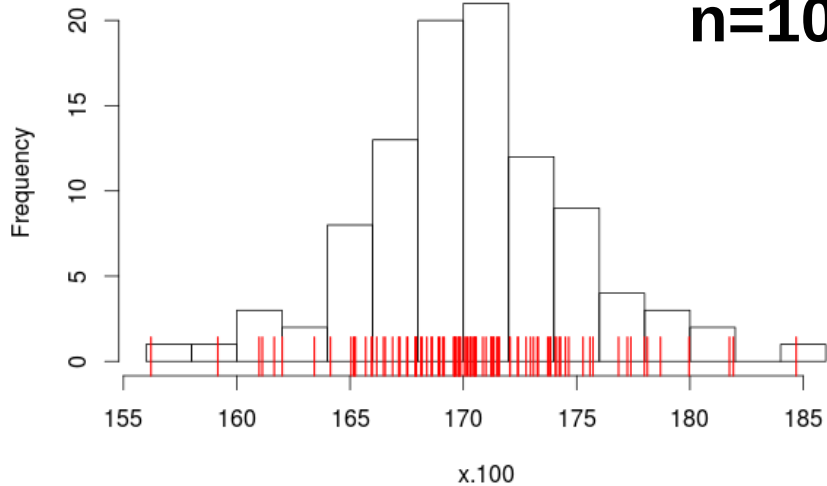
Histogram of x.50

**n=50**



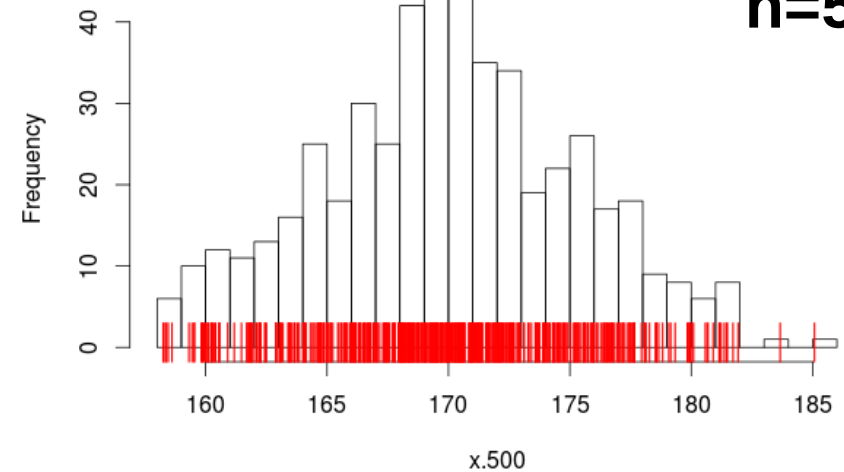
Histogram of x.100

**n=100**

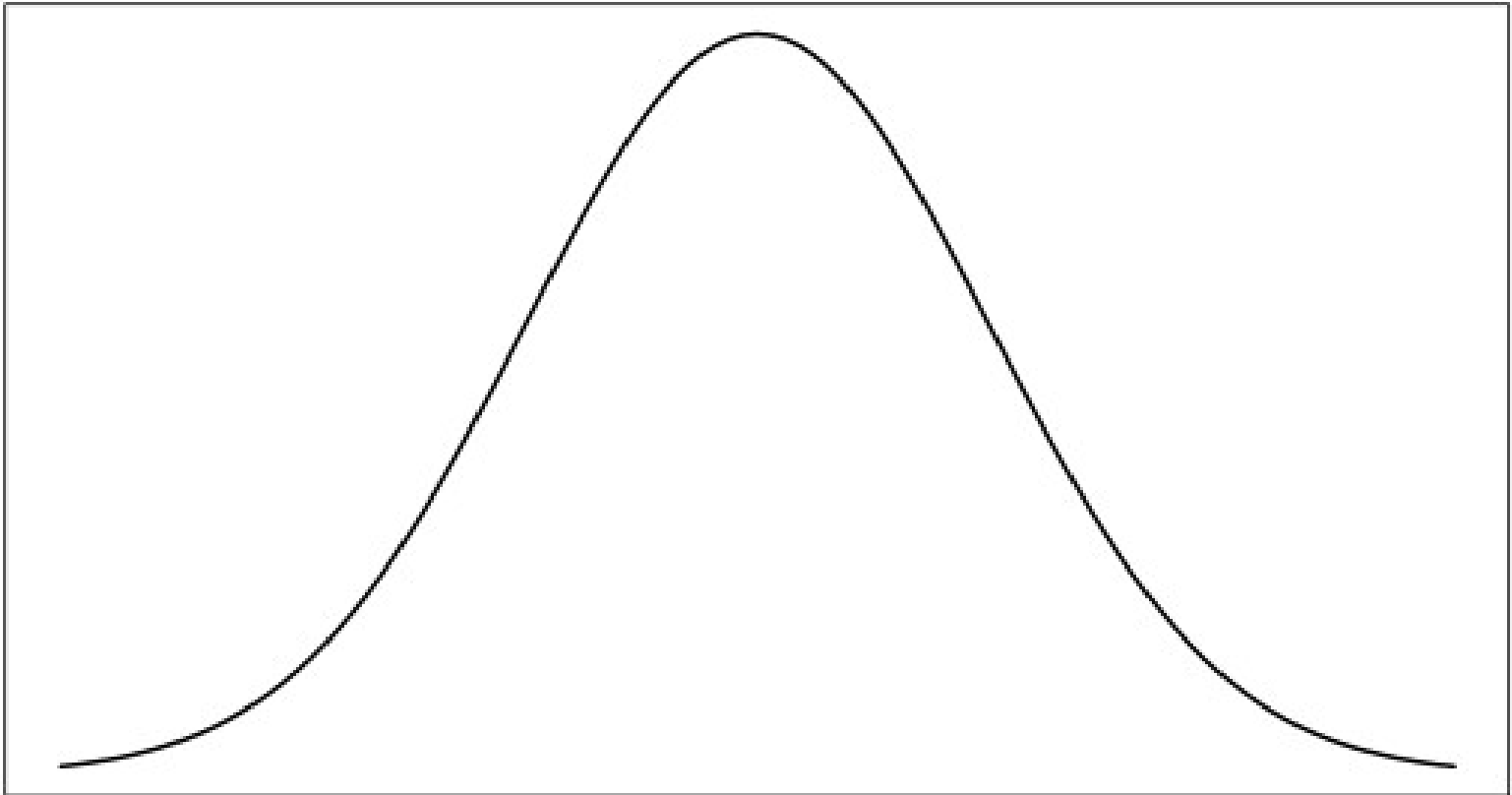


Histogram of x.500

**n=500**



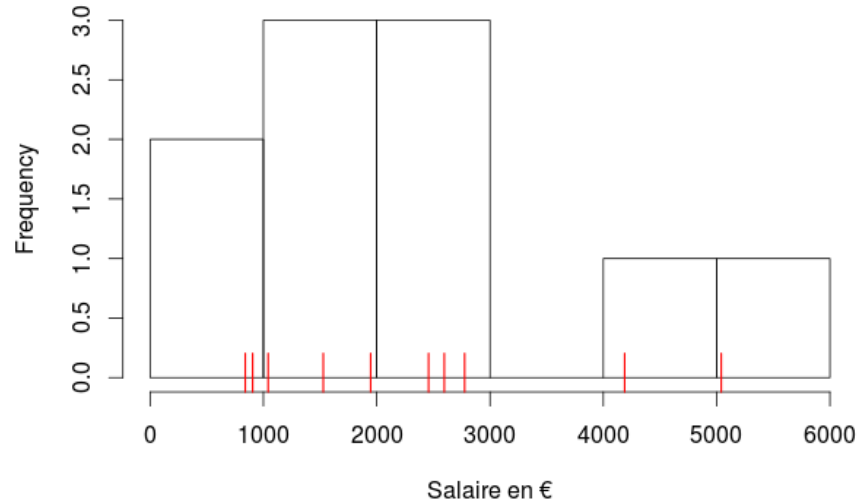
# Loi normale, loi de Student



# Exemple 2 : salaires

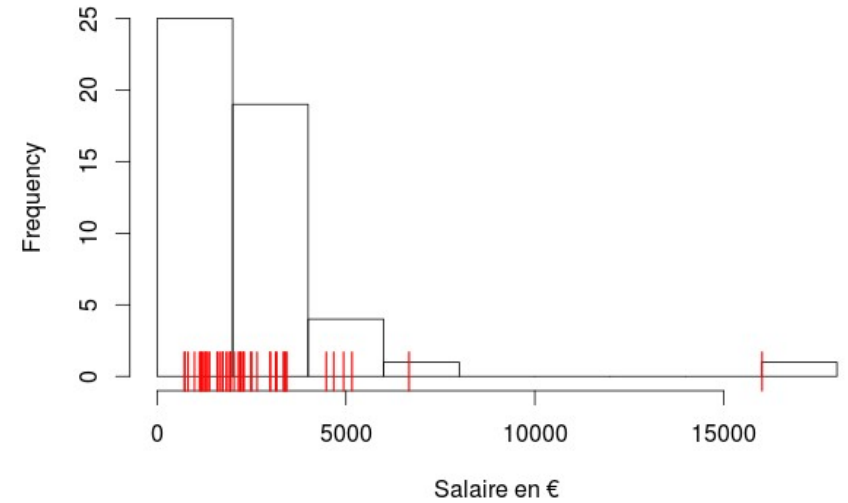
Histogram of  $y_{.10}$

**n=10**



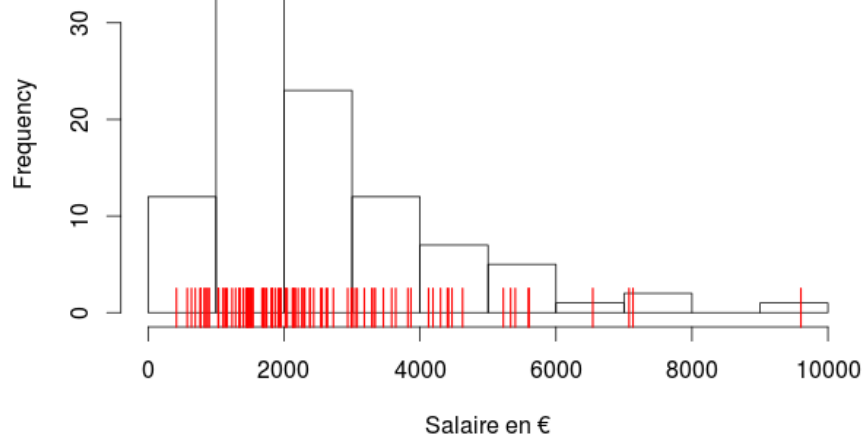
Histogram of  $y_{.50}$

**n=50**



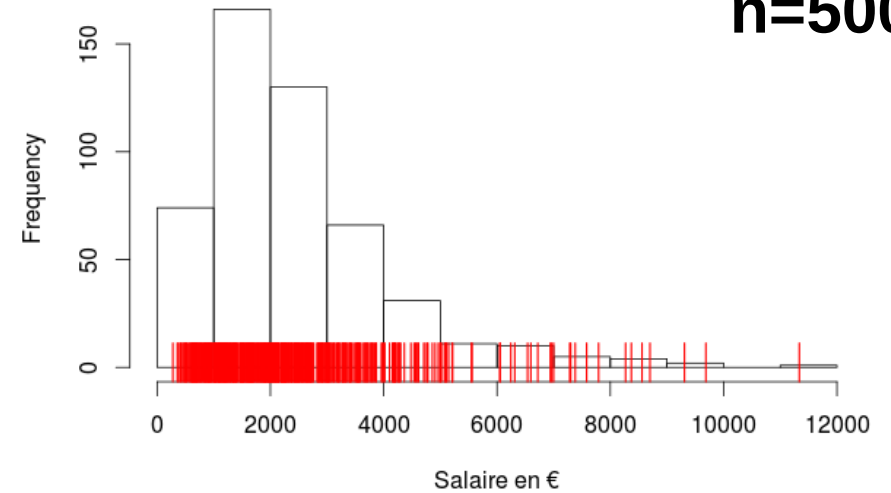
Histogram of  $y_{.100}$

**n=100**

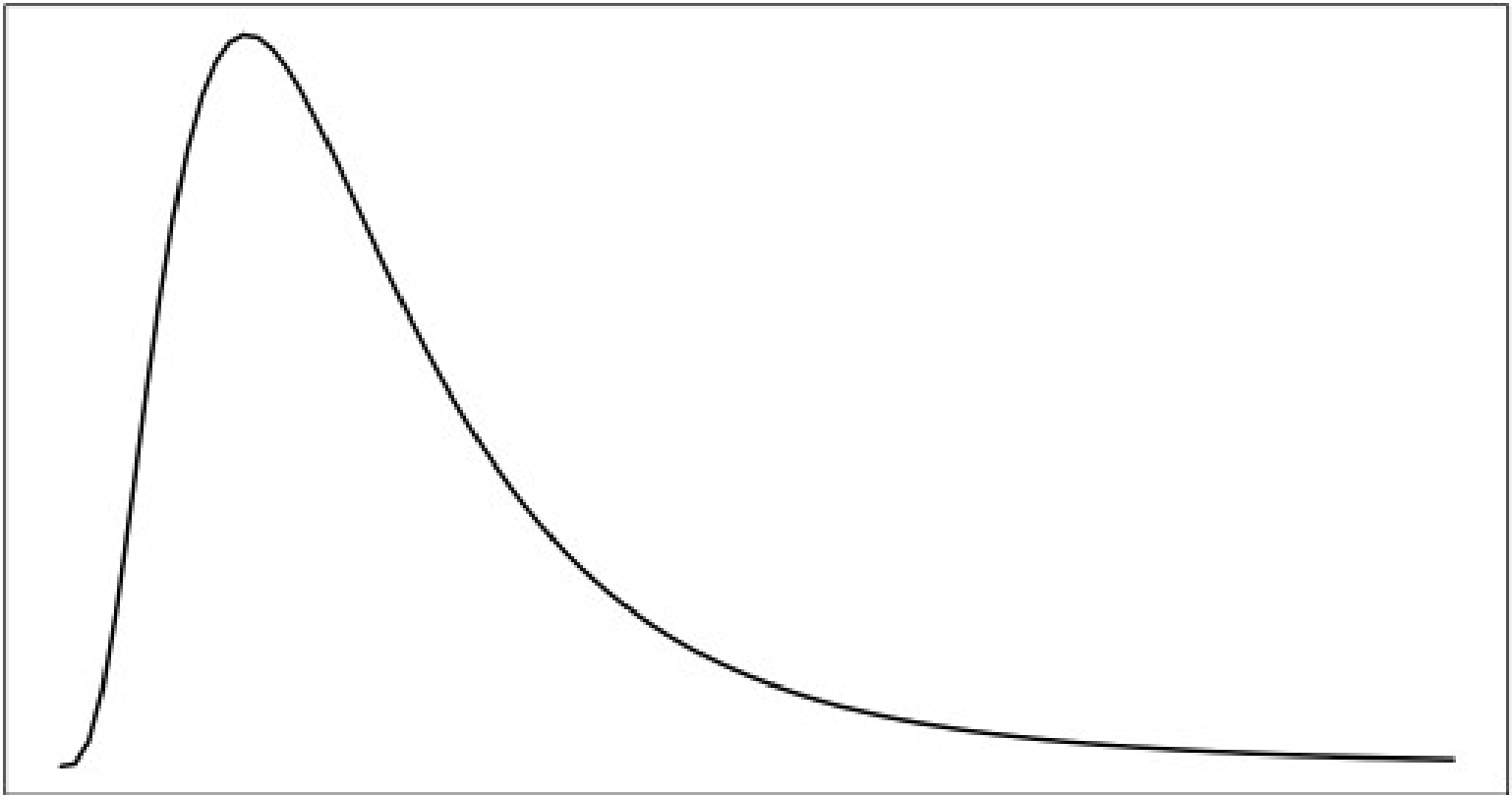


Histogram of  $y_{.500}$

**n=500**



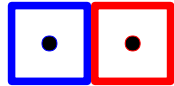
# Loi de $\chi^2$ , loi de Fisher



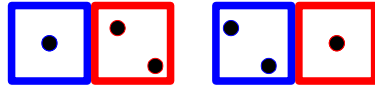
# Exemple 3 : lancer de 2 dés

- $P(X < 0) = P(X = 0) = P(X = 1) = 0$

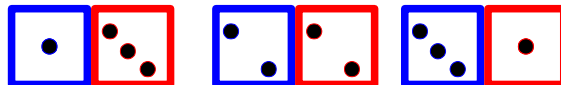
- $P(X = 2) = 1/36$



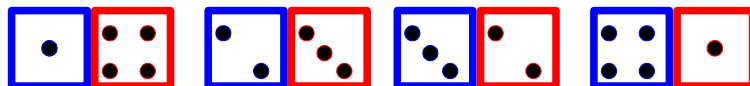
- $P(X = 3) = 2/36$



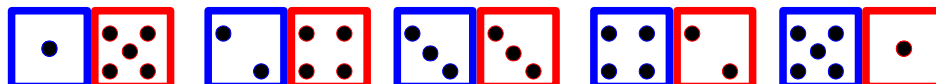
- $P(X = 4) = 3/36$



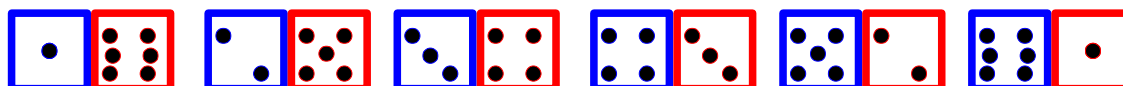
- $P(X = 5) = 4/36$



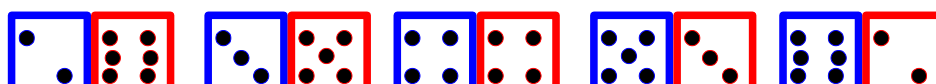
- $P(X = 6) = 5/36$



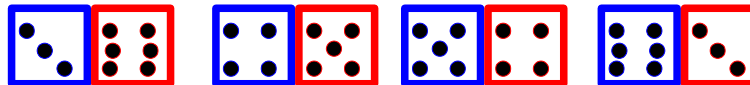
- $P(X = 7) = 6/36$



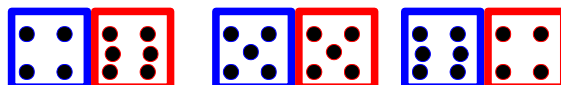
- $P(X = 8) = 5/36$



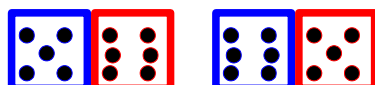
- $P(X = 9) = 4/36$



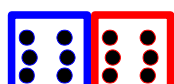
- $P(X = 10) = 3/36$



- $P(X = 11) = 2/36$

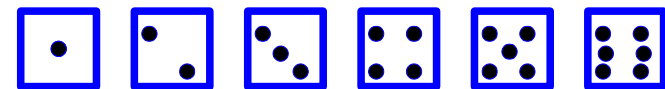
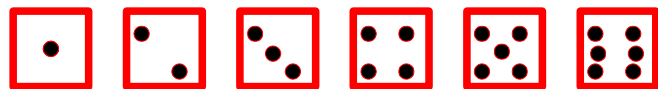


- $P(X = 12) = 1/36$

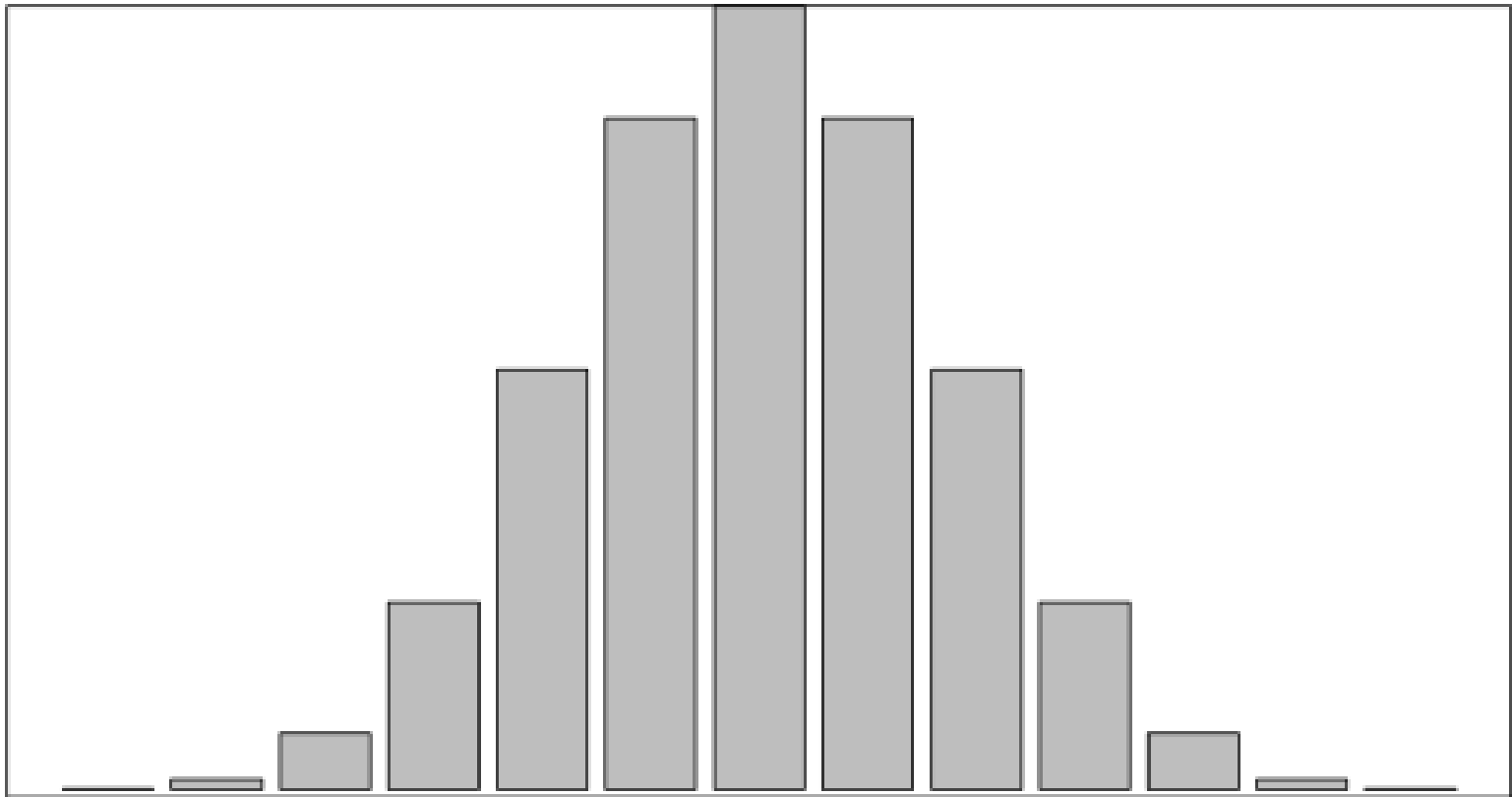


- $P(X > 12) = 0$

Soit  $X$  la variable aléatoire traduisant la somme du lancer de 2 dés à 6 faces. Établissons la loi de probabilité de  $X$ . Il y a 36 combinaisons possibles, il suffit de décompter les cas favorables pour calculer les probabilités de chaque événement.



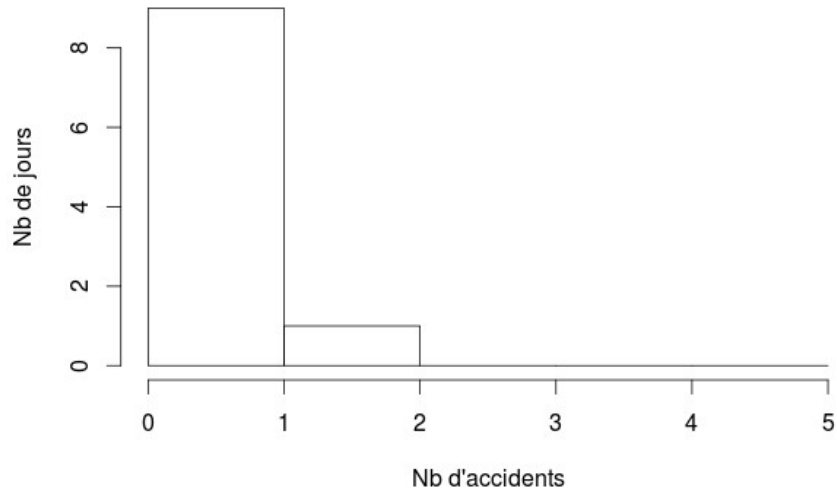
# Loi binomiale



# Exemple 4 : accidents

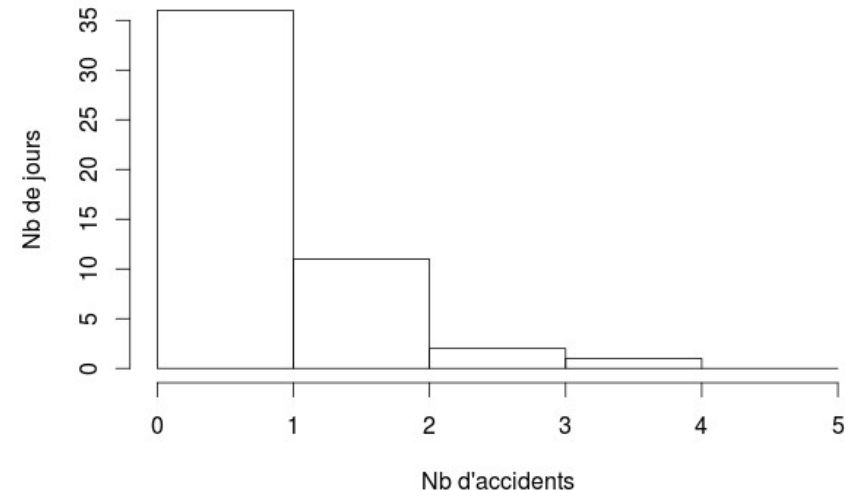
Histogram of z.10

**n=10**



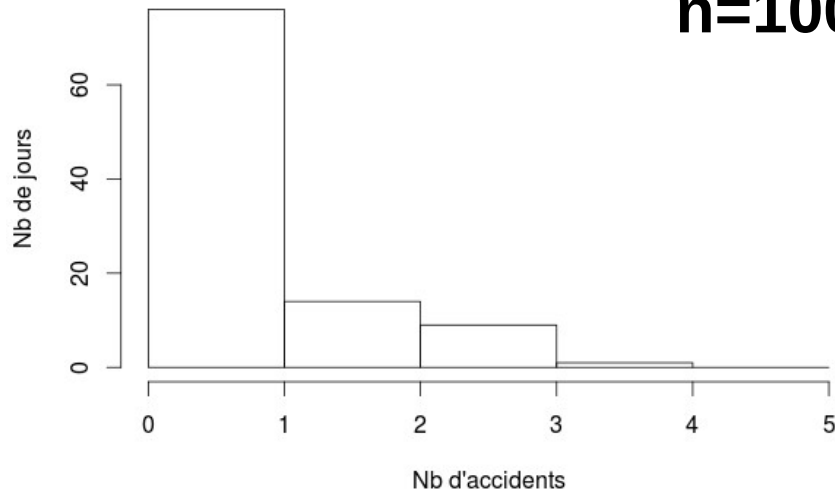
Histogram of z.50

**n=50**



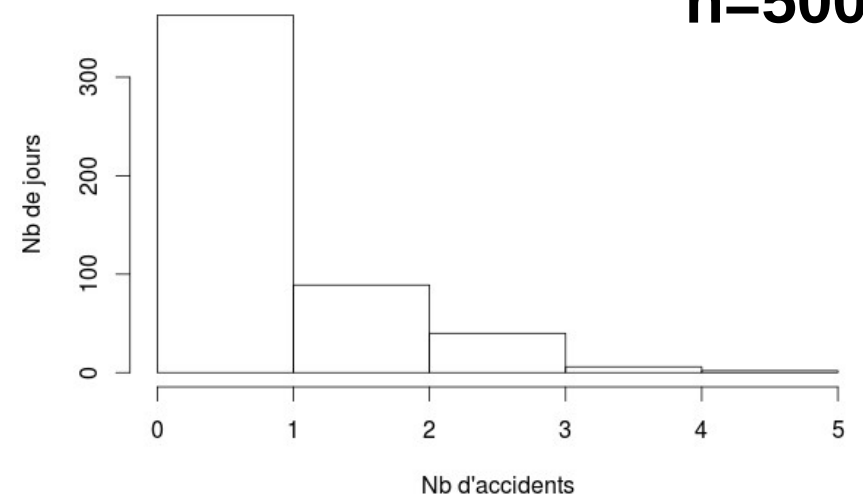
Histogram of z.100

**n=100**

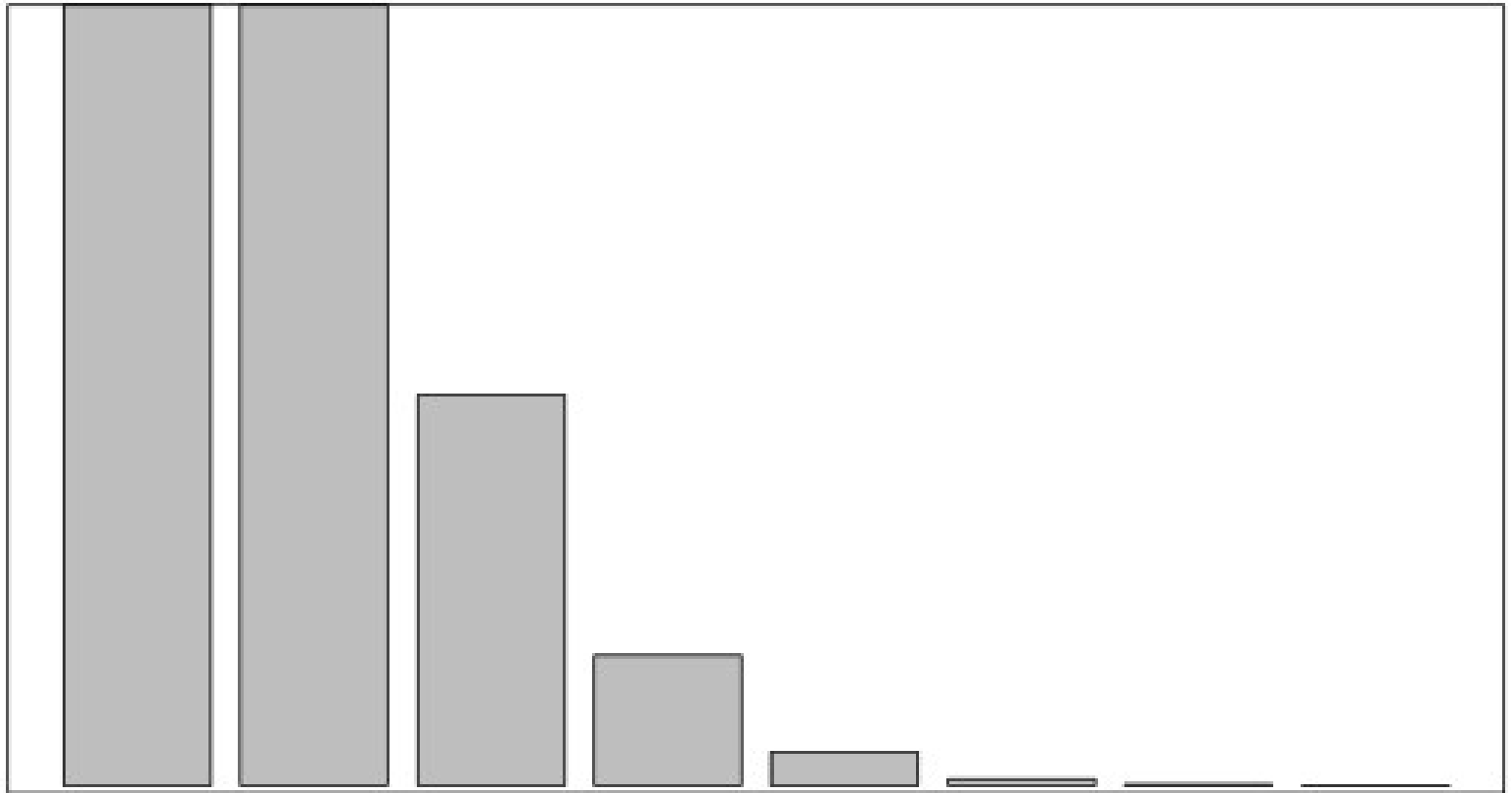


Histogram of z.500

**n=500**



# Loi de Poisson, loi binomiale négative

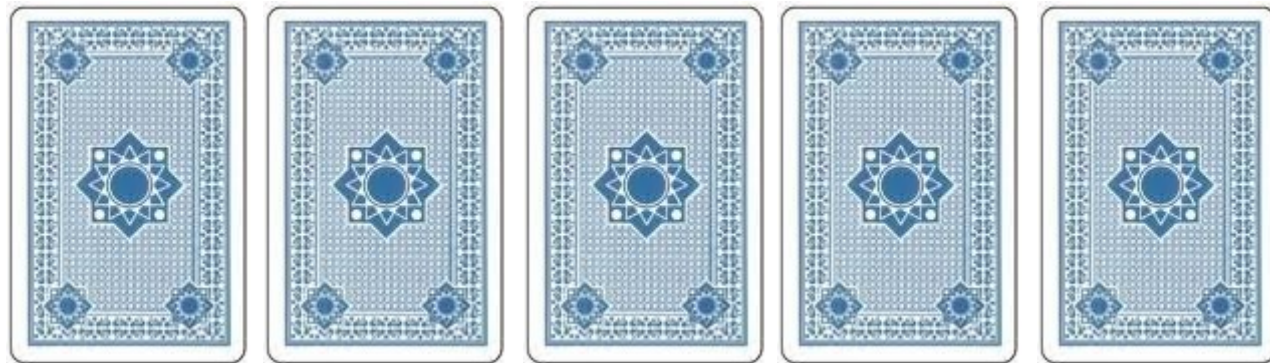




# Exemple 5 : poker (simplifié)

Jeu de 52 cartes

Chaque joueur reçoit 5 cartes



Quelle est la probabilité  
d'obtenir un carré d'as ?



# Loi hypergéométrique

Probabilité d'avoir 0 As dans une main de 5 cartes alors qu'il y en a 4 parmi 52 (48+4).

```
> dhyper(0, 4, 48, 5) 0.6588
```

Probabilité d'avoir 1 As dans une main de 5 cartes alors qu'il y en a 4 parmi 52 (48+4).

```
> dhyper(1, 4, 48, 5) 0.2995
```

Probabilité d'avoir 2 As dans une main de 5 cartes alors qu'il y en a 4 parmi 52 (48+4).

```
> dhyper(2, 4, 48, 5) 0.0399
```

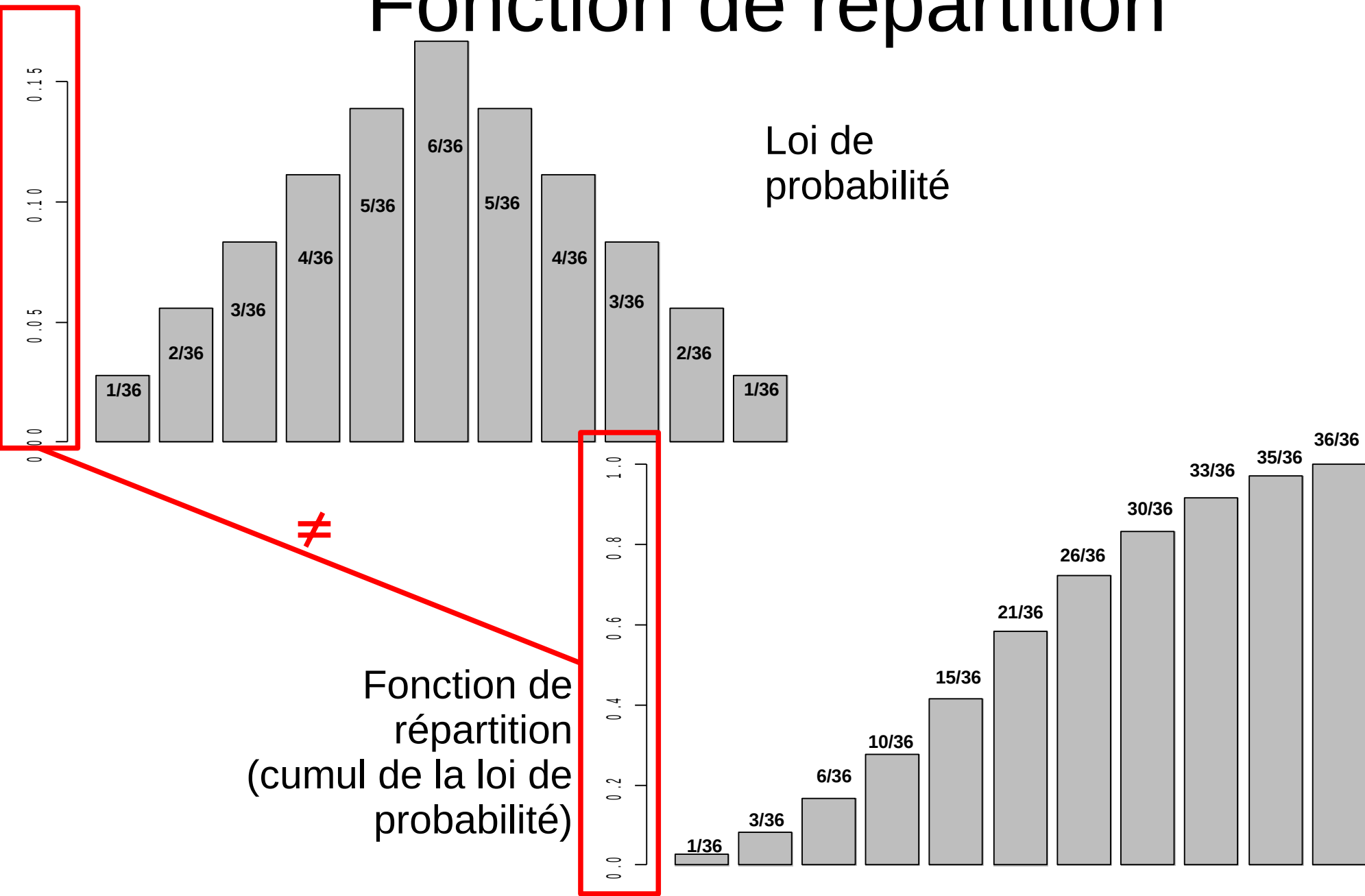
Probabilité d'avoir 3 As dans une main de 5 cartes alors qu'il y en a 4 parmi 52 (48+4).

```
> dhyper(3, 4, 48, 5) 0.0017
```

Probabilité d'avoir 4 As dans une main de 5 cartes alors qu'il y en a 4 parmi 52 (48+4).

```
> dhyper(4, 4, 48, 5) 0.00002 2 chances sur 100 000
```

# Distribution de probabilité – Fonction de répartition



# Distribution de probabilité – Fonction de répartition

Loi de probabilité

- dérivée de la fonction de répartition
- à valeur positive ou nulle

Fonction de répartition

- croissante sur  $]-\infty; +\infty[$
- tend vers 0 en  $-\infty$  et 1 en  $+\infty$
- continue à droite en tout point

Cas discret

$$P_X(k) = P(X = k)$$

$$\sum_{k \in \mathbb{Z}} P(X = k) = 1$$

$$F_X(n) = \sum_{k=-\infty}^n P(X = k)$$

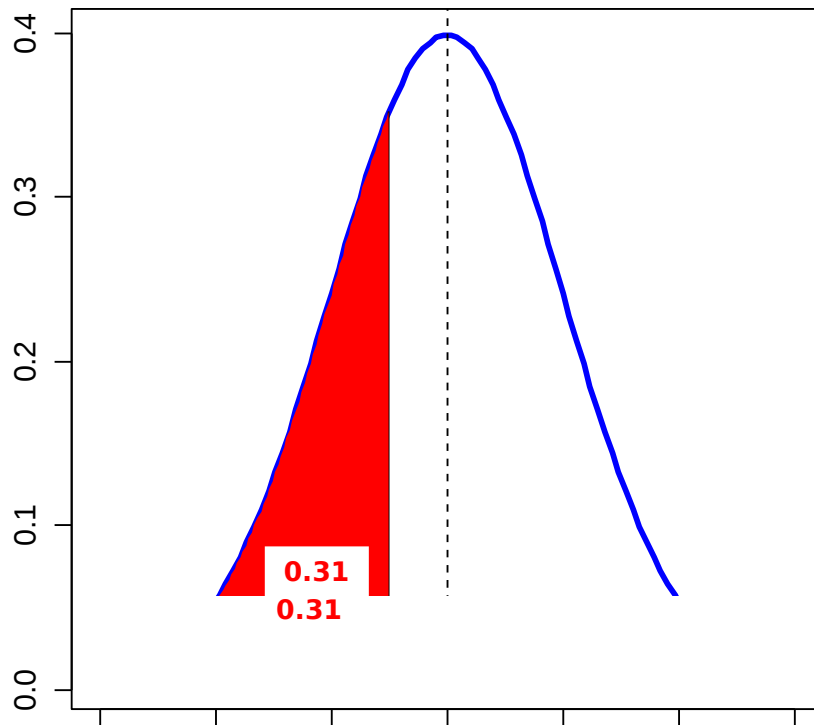
Cas continu (densité  $f$ )

$$P(a < X < b) = \int_a^b f(x) dx$$

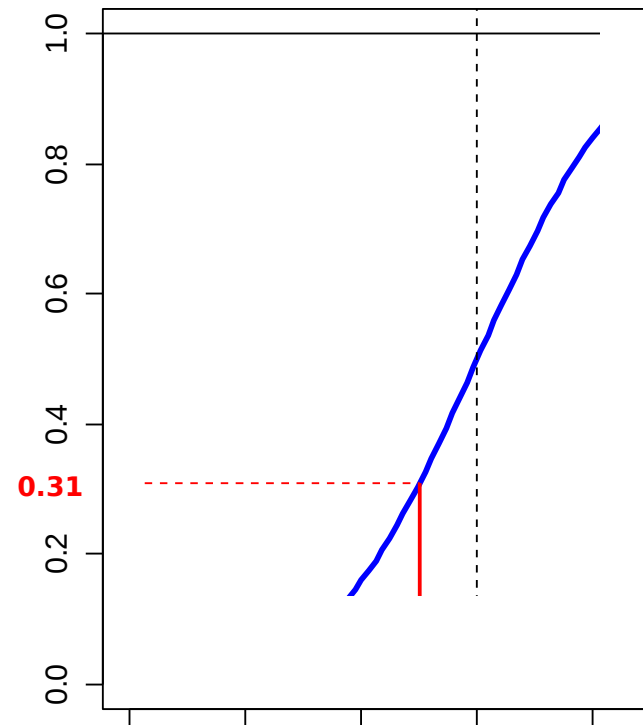
$$\int_{-\infty}^{+\infty} f(x) dx = 1$$

$$F_X(x) = \int_{-\infty}^x f(t) dt$$

# Distribution de probabilité – Fonction de répartition



Loi de  
probabilité

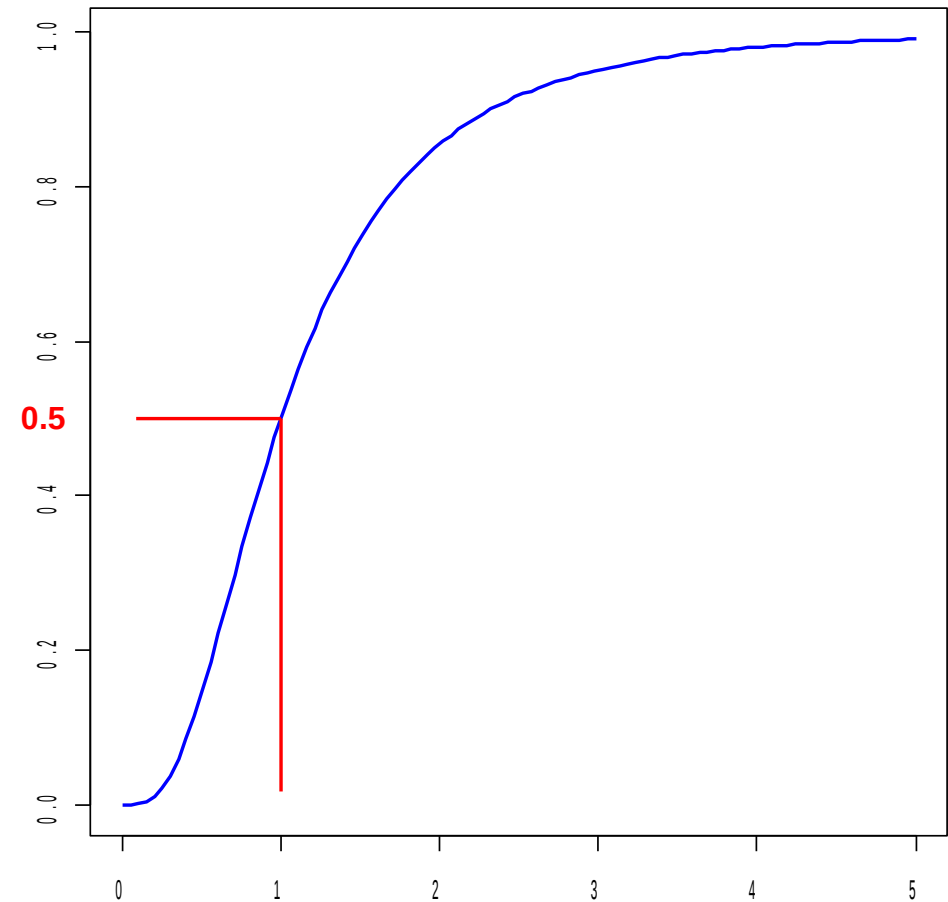
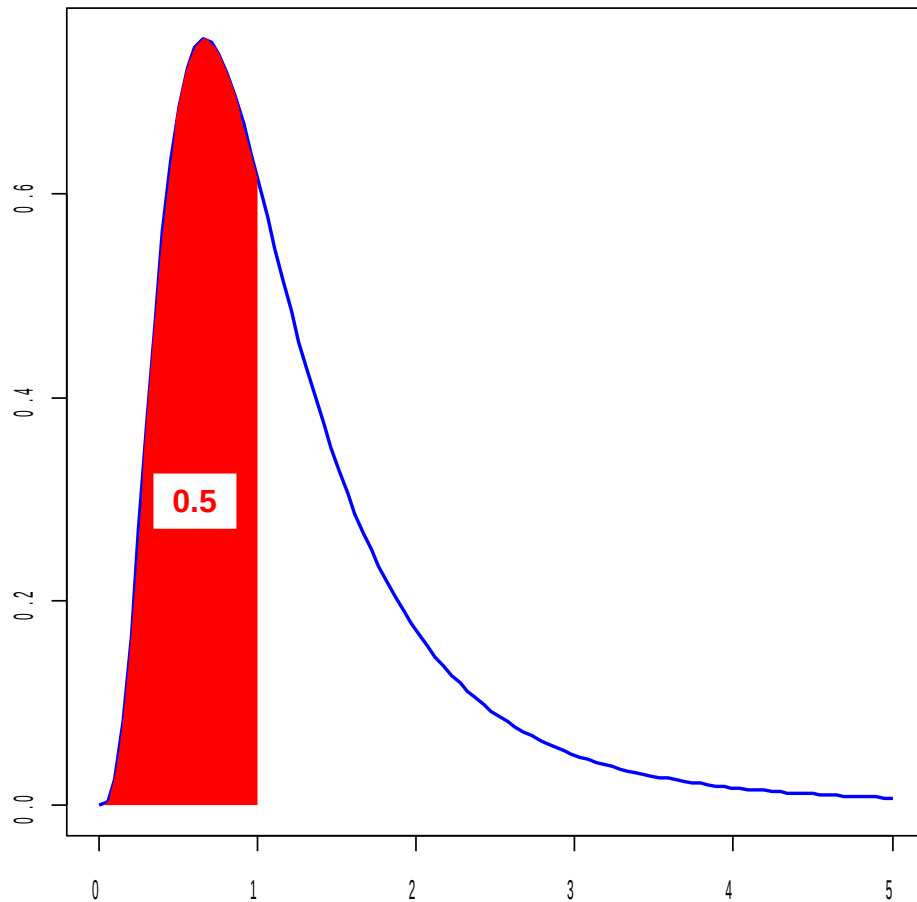


Fonction de répartition

Aire sous la courbe = probabilité

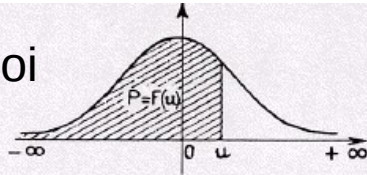
# Distribution de probabilité – Fonction de répartition

Exemple avec une loi asymétrique :  $F(10,10)$

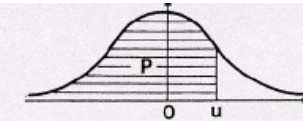


# Histoire : les tables statistiques

Fonction de répartition de la loi normale  $N(0,1)$



Fractiles de la loi normale  $N(0,1)$



u	0,00	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09
0,0	0,5000	0,5040	0,5080	0,5120	0,5160	0,5199	0,5239	0,5279	0,5319	0,5359
0,1	0,5398	0,5438	0,5478	0,5517	0,5557	0,5596	0,5636	0,5675	0,5714	0,5753
0,2	0,5793	0,5832	0,5871	0,5910	0,5948	0,5987	0,6026	0,6064	0,6103	0,6141
0,3	0,6179	0,6217	0,6255	0,6293	0,6331	0,6368	0,6406	0,6443	0,6480	0,6517
0,4	0,6554	0,6591	0,6628	0,6664	0,6700	0,6736	0,6772	0,6808	0,6844	0,6879
0,5	0,6915	0,6950	0,6985	0,7019	0,7054	0,7088	0,7123	0,7157	0,7190	0,7224
0,6	0,7257	0,7290	0,7324	0,7357	0,7389	0,7422	0,7454	0,7486	0,7517	0,7549
0,7	0,7580	0,7611	0,7642	0,7673	0,7704	0,7734	0,7764	0,7794	0,7823	0,7852
0,8	0,7881	0,7910	0,7939	0,7967	0,7995	0,8023	0,8051	0,8078	0,8106	0,8133
0,9	0,8159	0,8186	0,8212	0,8238	0,8264	0,8289	0,8315	0,8340	0,8365	0,8389
1,0	0,8413	0,8438	0,8461	0,8485	0,8508	0,8531	0,8554	0,8577	0,8599	0,8621
1,1	0,8643	0,8665	0,8686	0,8708	0,8729	0,8749	0,8770	0,8790	0,8810	0,8830
1,2	0,8849	0,8869	0,8888	0,8907	0,8925	0,8944	0,8962	0,8980	0,8997	0,9015
1,3	0,9032	0,9049	0,9066	0,9082	0,9099	0,9115	0,9131	0,9147	0,9162	0,9177
1,4	0,9192	0,9207	0,9222	0,9236	0,9251	0,9265	0,9279	0,9292	0,9306	0,9319
1,5	0,9332	0,9345	0,9357	0,9370	0,9382	0,9394	0,9406	0,9418	0,9429	0,9441
1,6	0,9452	0,9463	0,9474	0,9484	0,9495	0,9505	0,9515	0,9525	0,9535	0,9545
1,7	0,9554	0,9564	0,9573	0,9582	0,9591	0,9599	0,9608	0,9616	0,9625	0,9633
1,8	0,9641	0,9649	0,9656	0,9664	0,9671	0,9678	0,9686	0,9693	0,9699	0,9706
1,9	0,9713	0,9719	0,9726	0,9732	0,9738	0,9744	0,9750	0,9756	0,9761	0,9767
2,0	0,9772	0,9779	0,9783	0,9788	0,9793	0,9798	0,9803	0,9808	0,9812	0,9817
2,1	0,9821	0,9826	0,9830	0,9834	0,9838	0,9842	0,9846	0,9850	0,9854	0,9857
2,2	0,9861	0,9864	0,9868	0,9871	0,9875	0,9878	0,9881	0,9884	0,9887	0,9890
2,3	0,9893	0,9896	0,9898	0,9901	0,9904	0,9906	0,9909	0,9911	0,9913	0,9916
2,4	0,9918	0,9920	0,9922	0,9925	0,9927	0,9929	0,9931	0,9932	0,9934	0,9936
2,5	0,9938	0,9940	0,9941	0,9943	0,9945	0,9946	0,9948	0,9949	0,9951	0,9952
2,6	0,9953	0,9955	0,9956	0,9957	0,9959	0,9960	0,9961	0,9962	0,9963	0,9964
2,7	0,9965	0,9966	0,9967	0,9968	0,9969	0,9970	0,9971	0,9972	0,9973	0,9974
2,8	0,9974	0,9975	0,9976	0,9977	0,9977	0,9978	0,9979	0,9979	0,9980	0,9981
2,9	0,9981	0,9982	0,9982	0,9983	0,9984	0,9984	0,9985	0,9985	0,9986	0,9986

Table pour les grandes valeurs de u

u	3,0	3,1	3,2	3,3	3,4	3,5	3,6	3,8	4,0	4,5
F(u)	0,99865	0,99804	0,99931	0,99952	0,99966	0,99976	0,999841	0,999928	0,999968	0,999997

P	0,000	0,001	0,002	0,003	0,004	0,005	0,006	0,007	0,008	0,009	0,010
0,00	∞	3,0902	3,5782	2,7478	2,6521	2,3758	2,5121	2,4573	2,4089	2,3658	2,3263
0,01	2,3263	2,2904	2,2571	2,2262	2,1973	2,1701	2,1444	2,1201	2,0969	2,0749	2,0537
0,02	2,0537	2,0335	2,0141	1,9954	1,9774	1,9600	1,9431	1,9268	1,9110	1,8957	1,8808
0,03	1,8808	1,8663	1,8522	1,8384	1,8250	1,8119	1,7991	1,7866	1,7744	1,7624	1,7507
0,04	1,7507	1,7392	1,7278	1,7169	1,7060	1,6954	1,6849	1,6747	1,6646	1,6546	1,6449
0,05	1,6449	1,6352	1,6258	1,6164	1,6072	1,5982	1,5893	1,5805	1,5718	1,5632	1,5548
0,06	1,5548	1,5464	1,5382	1,5301	1,5220	1,5141	1,5063	1,4985	1,4909	1,4833	1,4758
0,07	1,4758	1,4684	1,4611	1,4538	1,4466	1,4395	1,4325	1,4255	1,4187	1,4118	1,4051
0,08	1,4051	1,3984	1,3917	1,3852	1,3787	1,3722	1,3658	1,3595	1,3532	1,3469	1,3408
0,09	1,3408	1,3346	1,3285	1,3225	1,3165	1,3106	1,3047	1,2988	1,2930	1,2873	1,2816
0,10	1,2816	1,2759	1,2702	1,2646	1,2591	1,2536	1,2481	1,2426	1,2372	1,2319	1,2265
0,11	1,2265	1,2212	1,2160	1,2107	1,2055	1,2004	1,1952	1,1901	1,1850	1,1800	1,1750
0,12	1,1750	1,1700	1,1650	1,1601	1,1552	1,1503	1,1455	1,1407	1,1359	1,1311	1,1264
0,13	1,1264	1,1217	1,1170	1,1123	1,1077	1,1031	1,0985	1,0939	1,0893	1,0848	1,0803
0,14	1,0803	1,0758	1,0714	1,0669	1,0625	1,0581	1,0537	1,0494	1,0450	1,0407	1,0364
0,15	1,0364	1,0322	1,0279	1,0237	1,0194	1,0152	1,0110	1,0069	1,0027	0,9986	0,9945
0,16	0,9945	0,9904	0,9863	0,9822	0,9782	0,9741	0,9701	0,9661	0,9621	0,9581	0,9542
0,17	0,9542	0,9502	0,9463	0,9424	0,9385	0,9346	0,9307	0,9269	0,9230	0,9192	0,9154
0,18	0,9154	0,9116	0,9078	0,9040	0,9002	0,8965	0,8927	0,8890	0,8853	0,8816	0,8779
0,19	0,8779	0,8742	0,8705	0,8669	0,8633	0,8598	0,8562	0,8526	0,8491	0,8456	0,8421
0,20	0,8421	0,8386	0,8351	0,8316	0,8281	0,8246	0,8211	0,8176	0,8141	0,8106	0,8071
0,21	0,8071	0,8036	0,7999	0,7964	0,7929	0,7894	0,7859	0,7824	0,7789	0,7754	0,7720
0,22	0,7720	0,7686	0,7651	0,7617	0,7582	0,7548	0,7514	0,7480	0,7446	0,7412	0,7378
0,23	0,7378	0,7345	0,7312	0,7279	0,7246	0,7213	0,7180	0,7147	0,7114	0,7081	0,7048
0,24	0,7048	0,7015	0,6982	0,6949	0,6916	0,6883	0,6850	0,6817	0,6784	0,6751	0,6718
0,25	0,6718	0,6685	0,6652	0,6619	0,6586	0,6553	0,6520	0,6487	0,6454	0,6421	0,6388
0,26	0,6388	0,6355	0,6322	0,6289	0,6256	0,6223	0,6190	0,6157	0,6124	0,6091	0,6058
0,27	0,6058	0,6025	0,5992	0,5959	0,5926	0,5893	0,5860	0,5827	0,5794	0,5761	0,5728
0,28	0,5728	0,5695	0,5662	0,5629	0,5596	0,5563	0,5530	0,5497	0,5464	0,5431	0,5398
0,29	0,5398	0,5365	0,5332	0,5299	0,5266	0,5233	0,5200	0,5167	0,5134	0,5101	0,5068
0,30	0,5068	0,5035	0,5002	0,4969	0,4936	0,4903	0,4870	0,4837	0,4804	0,4771	0,4738
0,31	0,4738	0,4705	0,4672	0,4639	0,4606	0,4573	0,4540	0,4507	0,4474	0,4441	0,4408
0,32	0,4408	0,4375	0,4342	0,4309	0,4276	0,4243	0,4210	0,4177	0,4144	0,4111	0,4078
0,33	0,4078	0,4045	0,4012	0,3979	0,3946	0,3913	0,3880	0,3847	0,3814	0,3781	0,3748
0,34	0,3748	0,3715	0,3682	0,3649	0,3616	0,3583	0,3550	0,3517	0,3484	0,3451	0,3418
0,35	0,3418	0,3385	0,3352	0,3319	0,3286	0,3253	0,3220	0,3187	0,3154	0,3121	0,3088
0,36	0,3088	0,3055	0,3022	0,2989	0,2956	0,2923	0,2890	0,2857	0,2824	0,2791	0,2758
0,37	0,2758	0,2725	0,2692	0,2659	0,2626	0,2593	0,2560	0,2527	0,2494	0,2461	0,2428
0,38	0,2428	0,2395	0,2362	0,2329	0,2296	0,2263	0,2230	0,2197	0,2164	0,2131	0,2098
0,39	0,2098	0,2065	0,2032	0,1999	0,1966	0,1933	0,1900	0,1867	0,1834	0,1801	0,1768
0,40	0,1768	0,1735	0,1702	0,1669	0,1636	0,1603	0,1570	0,1537	0,1504	0,1471	0,1438
0,41	0,1438	0,1405	0,1372	0,1339	0,1306	0,1273	0,1240	0,1207	0,1174	0,1141	0,1108
0,42	0,1108	0,1075	0,1042	0,1009	0,9976	0,9943	0,9910	0,9877	0,9844	0,9811	0,9778
0,43	0,9778	0,9745	0,9712	0,9679	0,9646	0,9613	0,9580	0,9547	0,9514	0,9481	0,9448
0,44	0,9448	0,9415	0,9382	0,9349	0,9316	0,9283	0,9250	0,9217	0,9184	0,9151	0,9118
0,45	0,9118	0,9085	0,9052	0,9019	0,8986	0,8953	0,8920	0,8887	0,8854	0,8821	0,8788
0,46	0,8788	0,8755	0,8722	0,8689	0,8656	0,8623	0,8590	0,8557	0,8524	0,8491	0,8458
0,47	0,8458	0,8425	0,8392	0,8359	0,8326	0,8293	0,8260	0,8227	0,8194	0,8161	0,8128
0,48	0,8128	0,8095	0,8062	0,8029	0,7996	0,7963	0,7930	0,7897	0,7864	0,7831	0,7798
0,49	0,7798	0,7765	0,7732	0,7699	0,7666	0,7633	0,7600	0,7567	0,7534	0,7501	0,7468

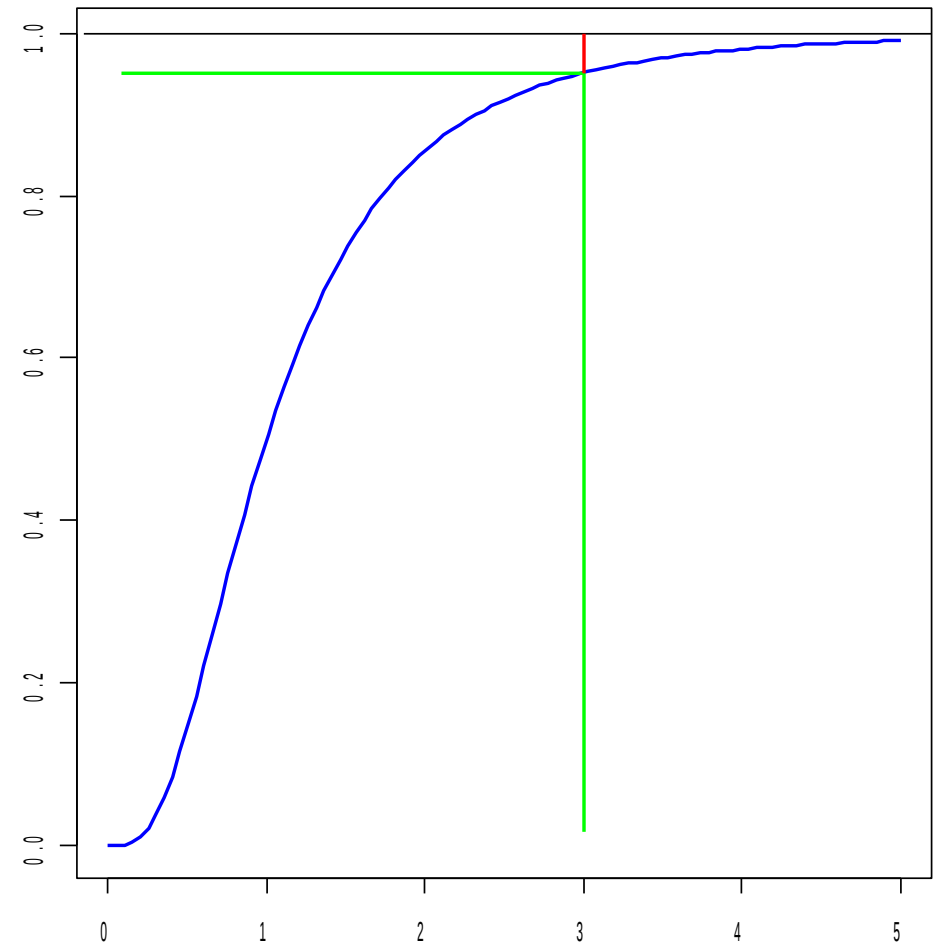
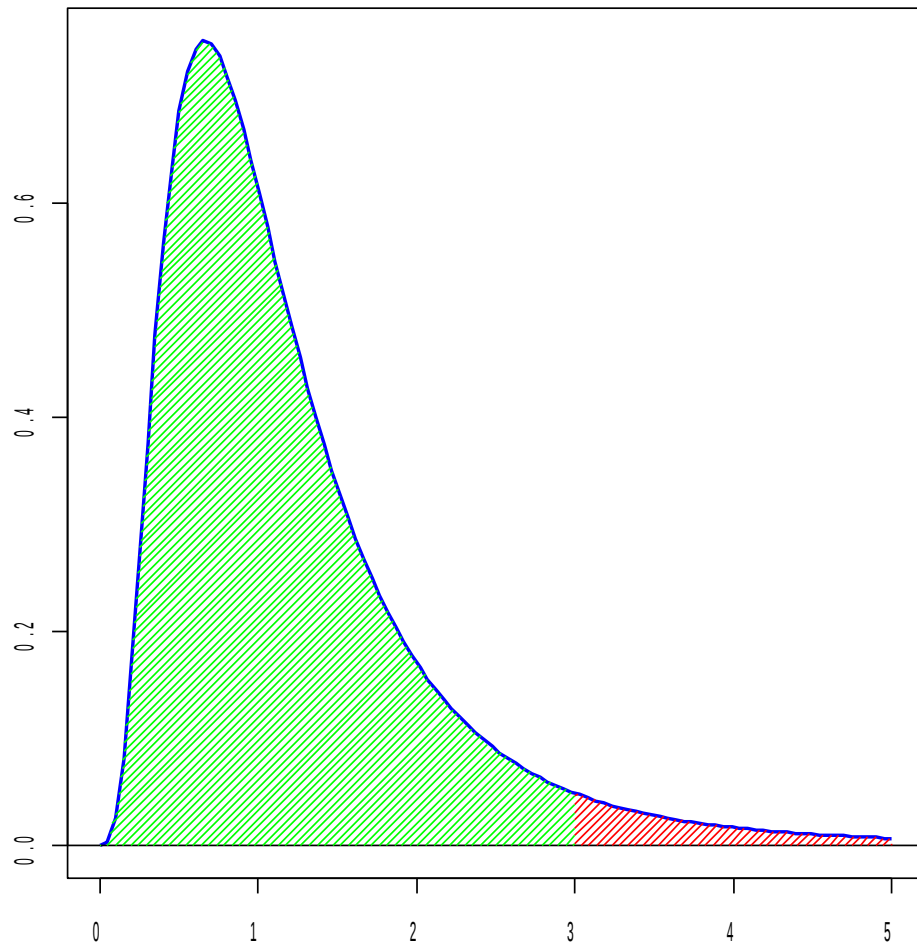
Grandes valeurs de u

P	0,9999	0,99999	0,999999	0,9999999	0,99999999	0,999999999
u	3,7190	4,2649	4,7534	5,1993	5,6120	5,9978

N.B. Si  $P < 0,5$ , u est négatif.

# Distribution de probabilité – Fonction de répartition

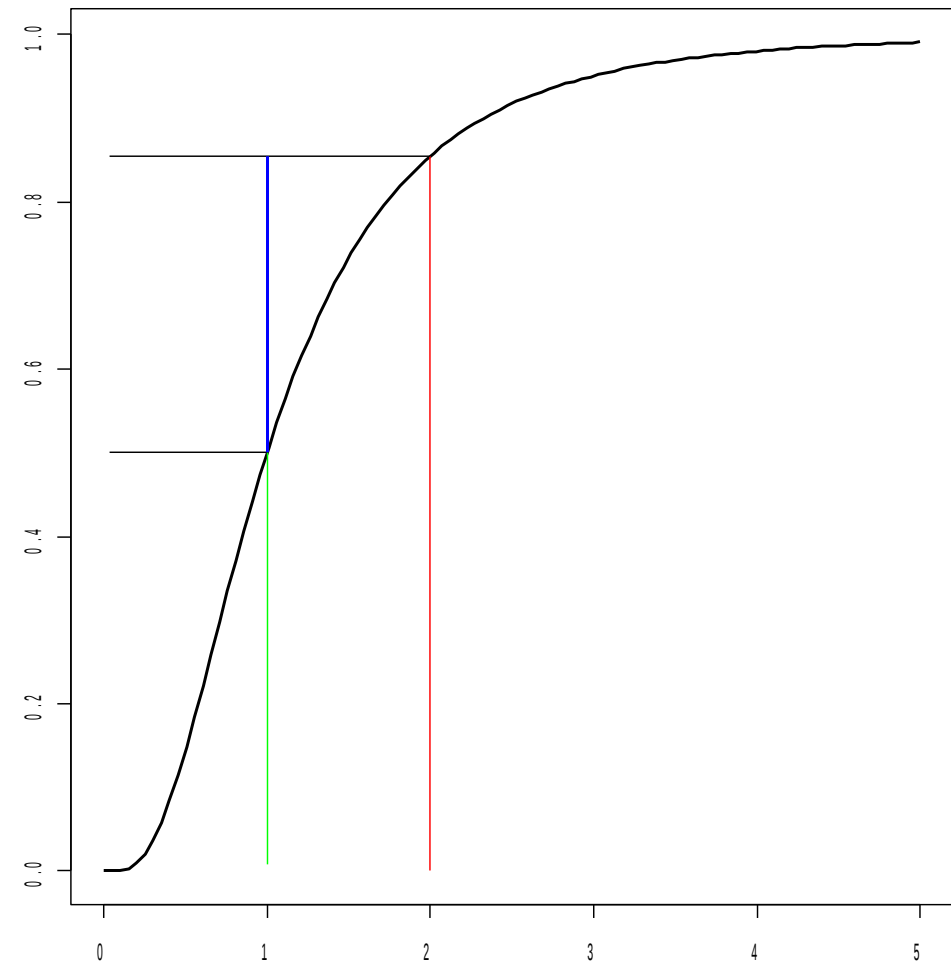
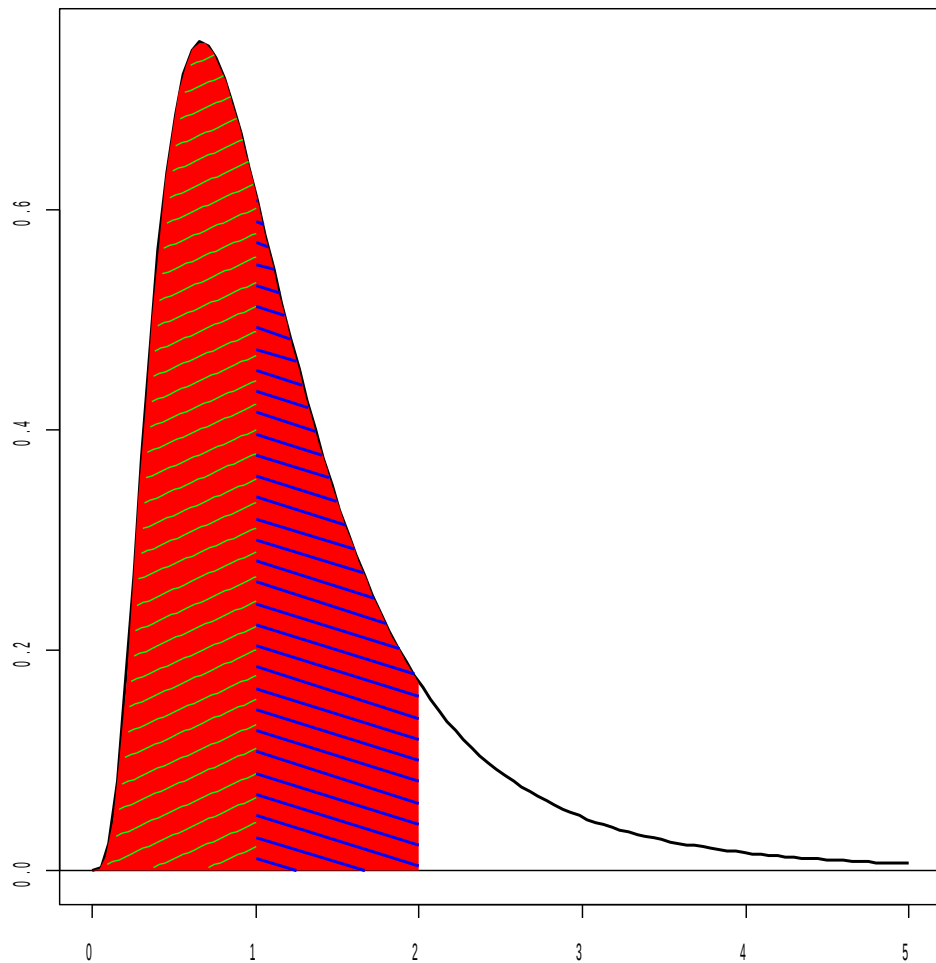
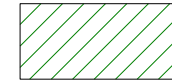
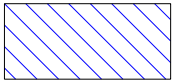
« aire à droite » =  $1 -$  « aire à gauche »





# Distribution de probabilité – Fonction de répartition

« aire centrale » = « grande aire à gauche » - « petite aire à gauche »



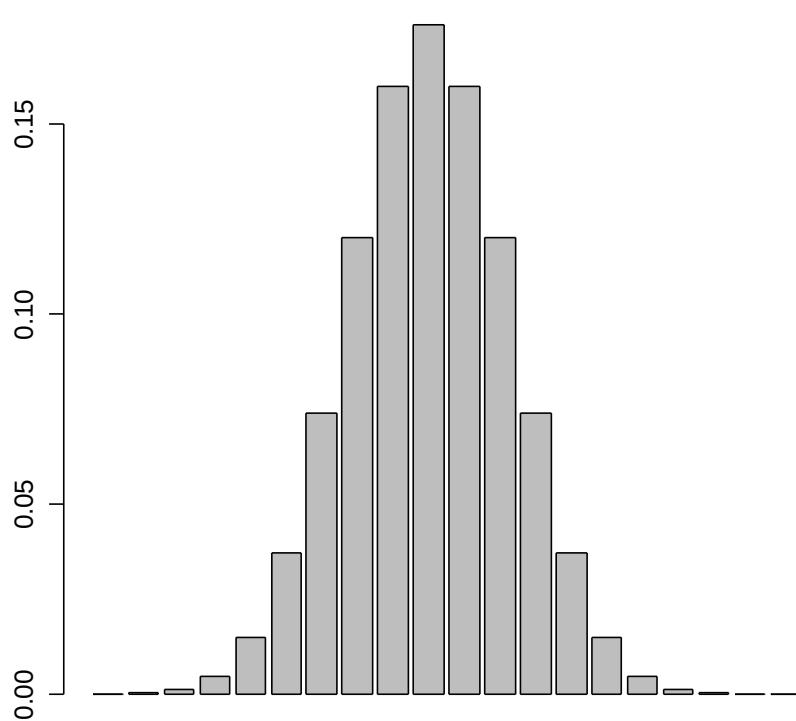
# Quelques lois usuelles

Loi binomiale  $B(n,p)$

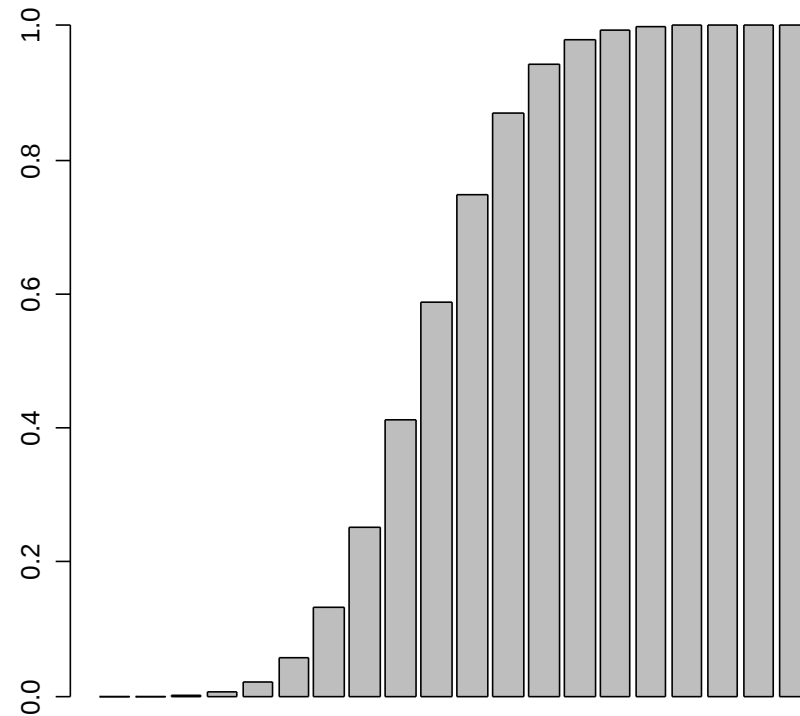
$$P(X=k)=C_n^k p^k(1-p)^{n-k}$$

Ex :  $n=20, p=0.5$

$$C_n^k = \frac{n!}{k!(n-k)!}$$



Loi de probabilité



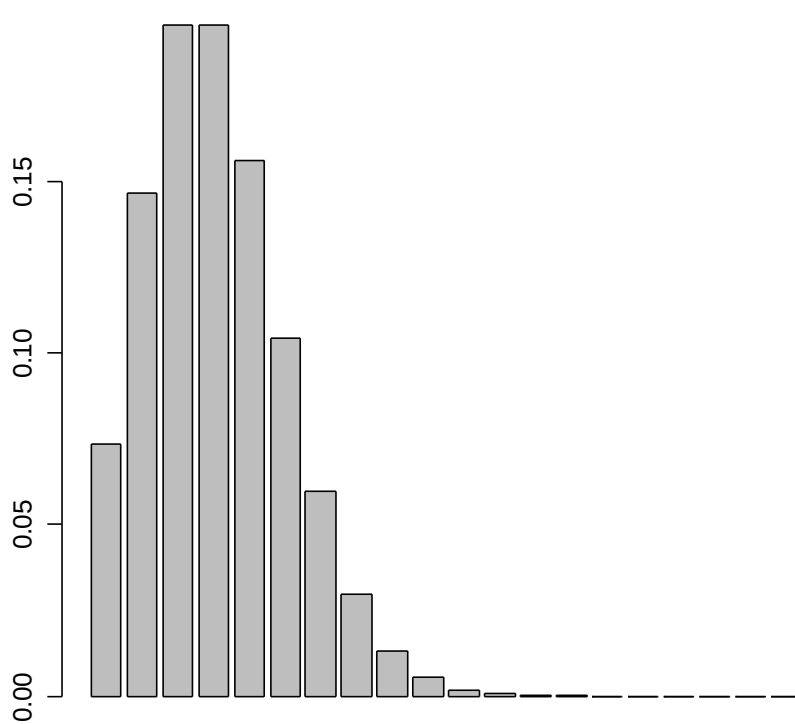
Fonction de répartition

# Quelques lois usuelles

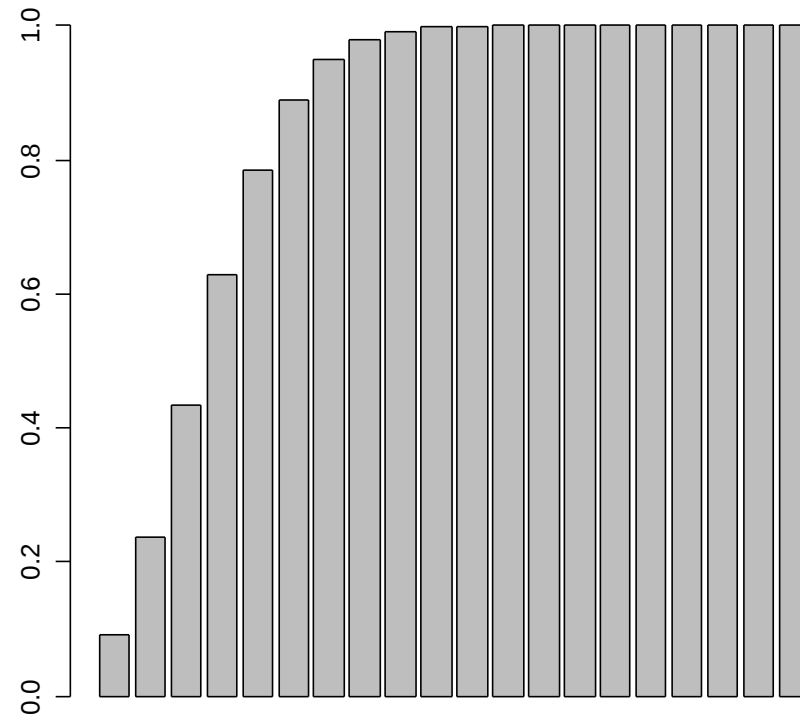
Loi de Poisson  $P(\lambda)$

$$P(X=k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

Ex :  $\lambda=4$



Loi de probabilité



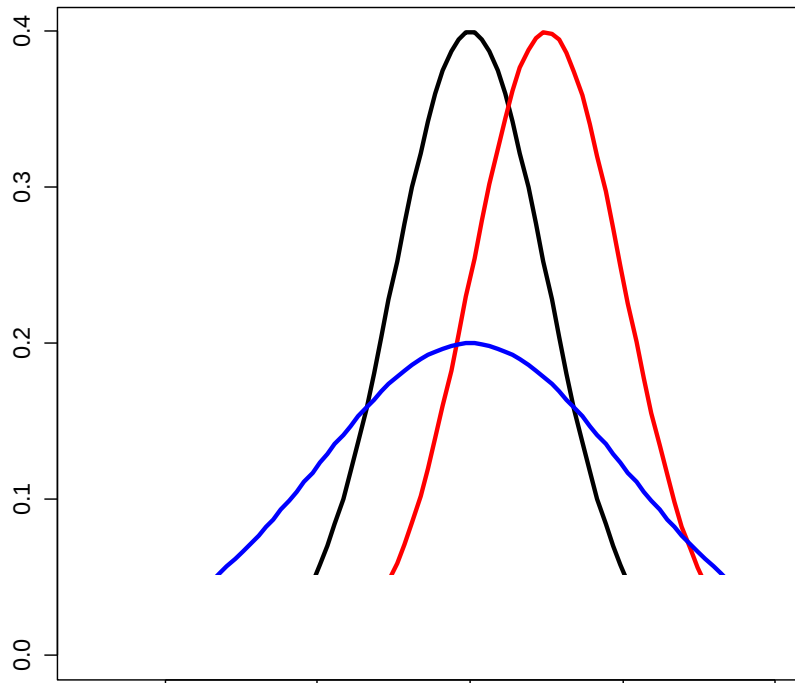
Fonction de répartition

# Quelques lois usuelles

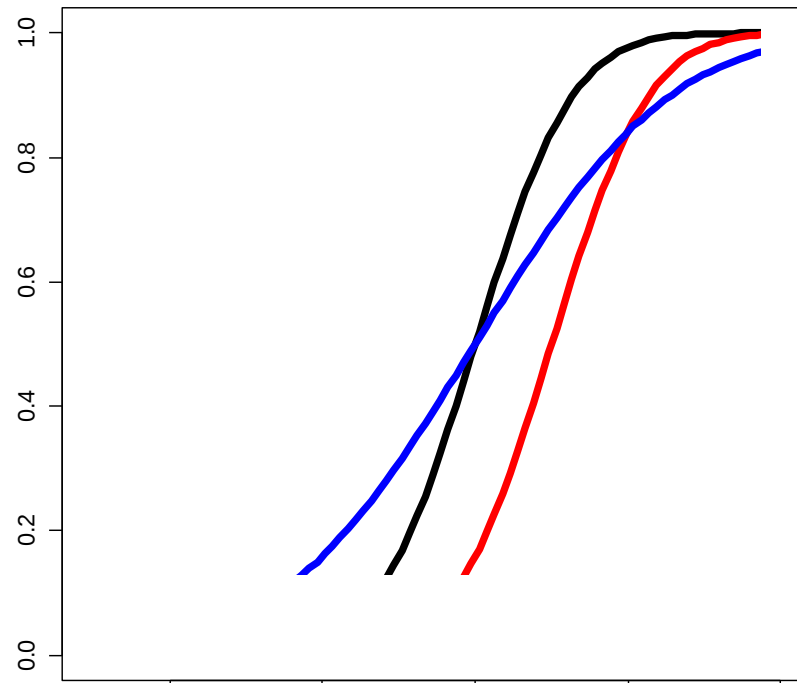
Loi normale  $\mathbf{N}(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Ex : **N**(0,1)   **N**(1,1)   **N**(0,2)



Loi de  
probabilité



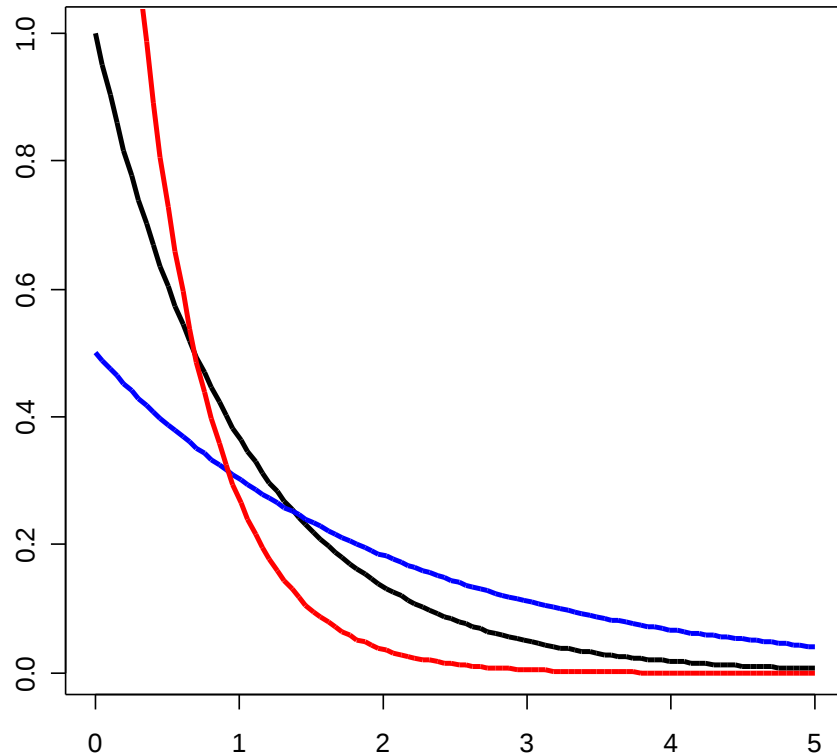
Fonction de répartition

# Quelques lois usuelles

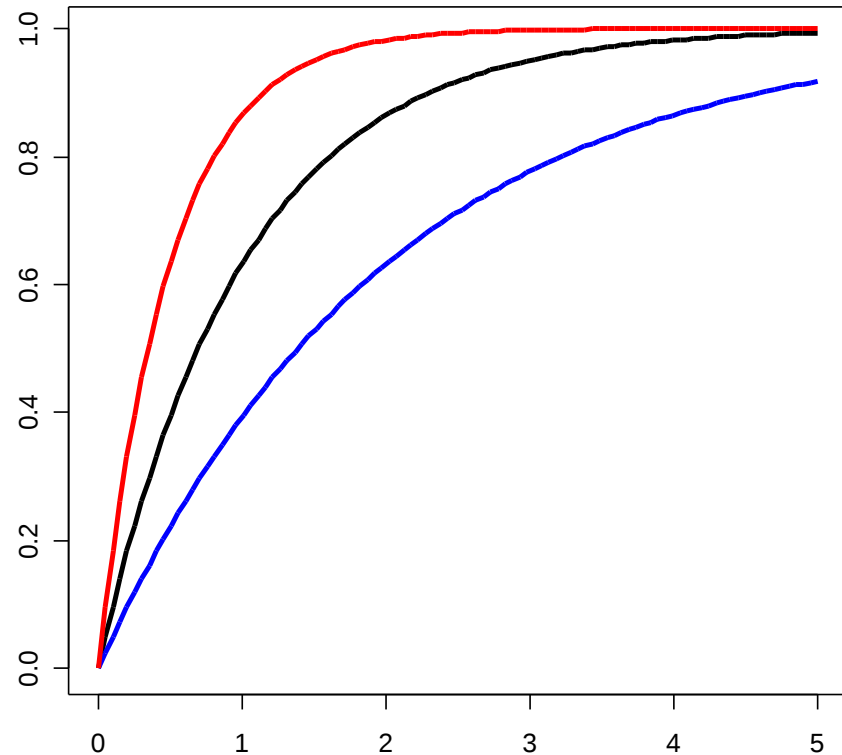
Loi exponentielle  $E(\lambda)$

$$f(x) = \lambda \exp(-\lambda x)$$

Ex :  $\lambda = 1 - 0.5 - 2$



Loi de probabilité



Fonction de répartition

# Quelques lois usuelles

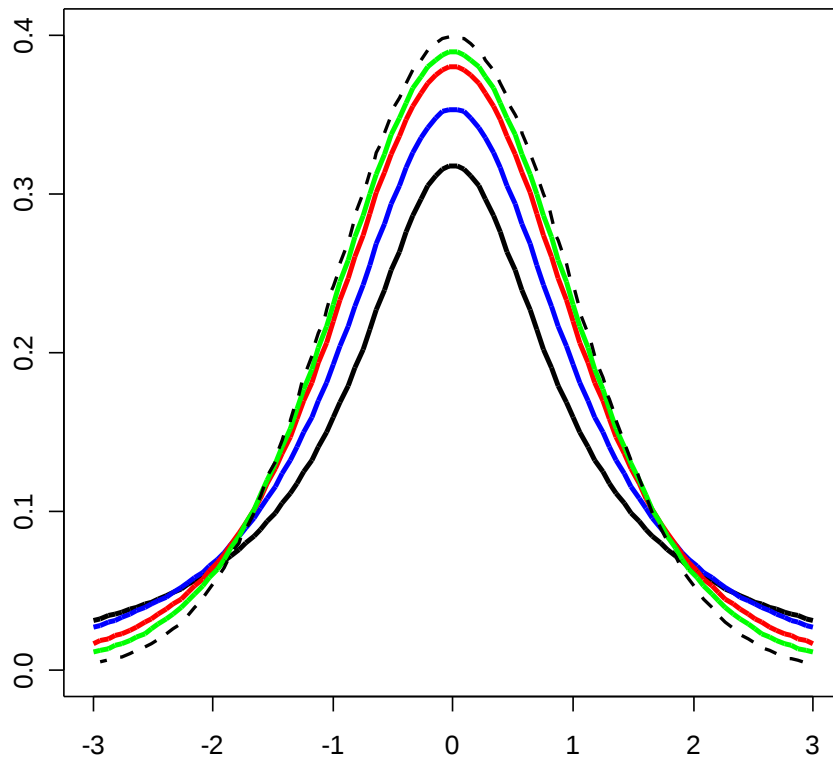
Loi de Student  $\text{St}_{(k)}$

$$f(x) = \frac{\Gamma\left(\frac{k+1}{2}\right)}{\sqrt{k\pi} \Gamma\left(\frac{k}{2}\right) \left(1 + \frac{t^2}{k}\right)^{\frac{k+1}{2}}}$$

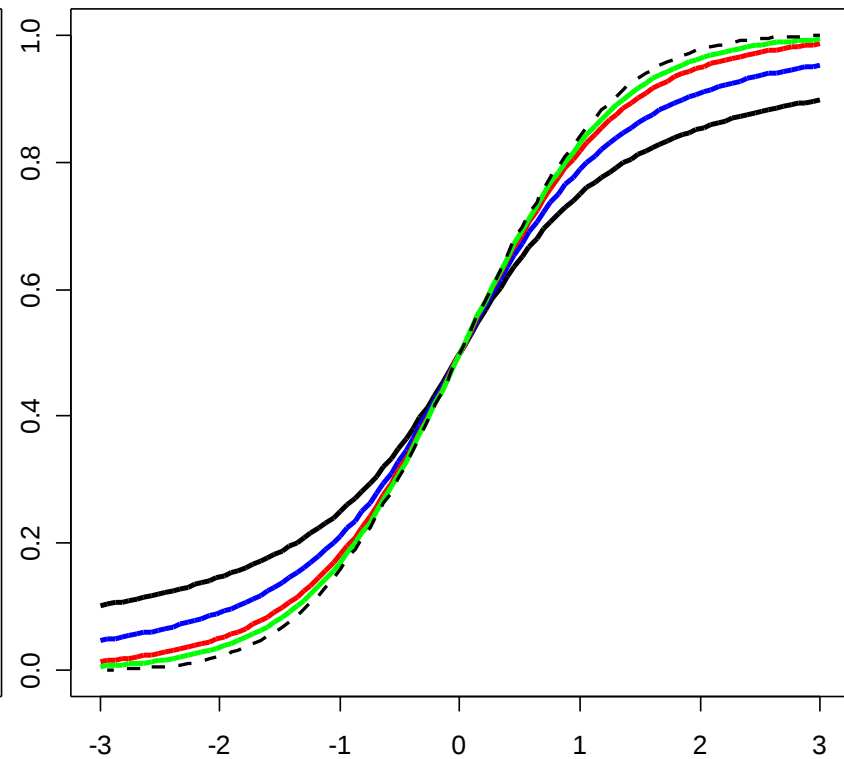
$\Gamma$  fonction  
Gamma d'Euler

Ex :  $k=1$  - 2 - 5 - 10

Loi normale  $\mathbf{N}(0,1)$  - - -



Loi de  
probabilité



Fonction de répartition

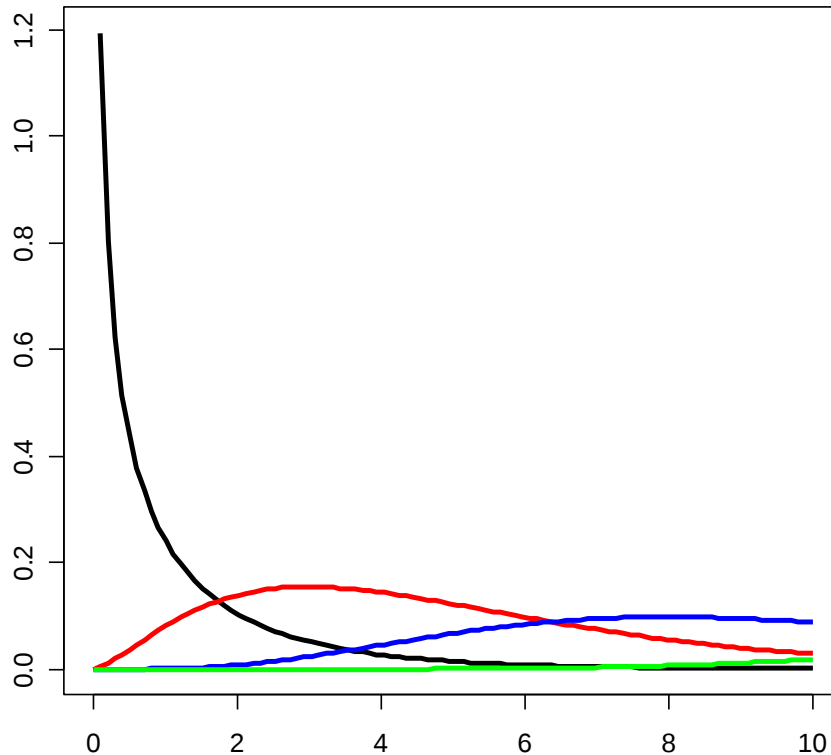
# Quelques lois usuelles

Loi du khi-deux  $\chi^2(k)$

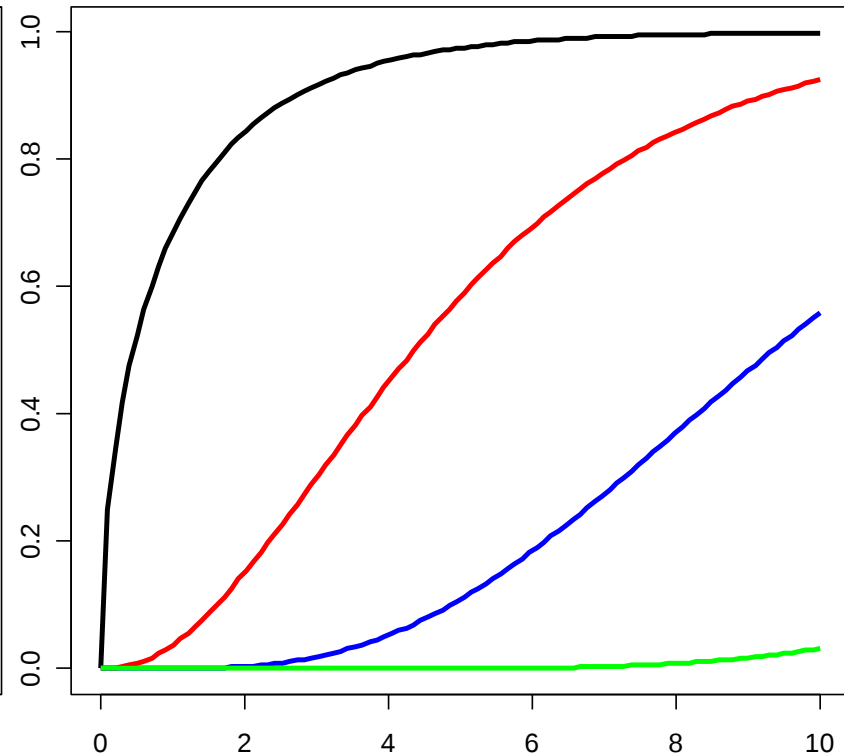
$$f(x) = \frac{1}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)} x^{\frac{k}{2}-1} e^{-\frac{x}{2}}$$

$\Gamma$  fonction  
Gamma d'Euler

Ex :       $k=1$  - 5 - 10 - 20



Loi de  
probabilité



Fonction de répartition

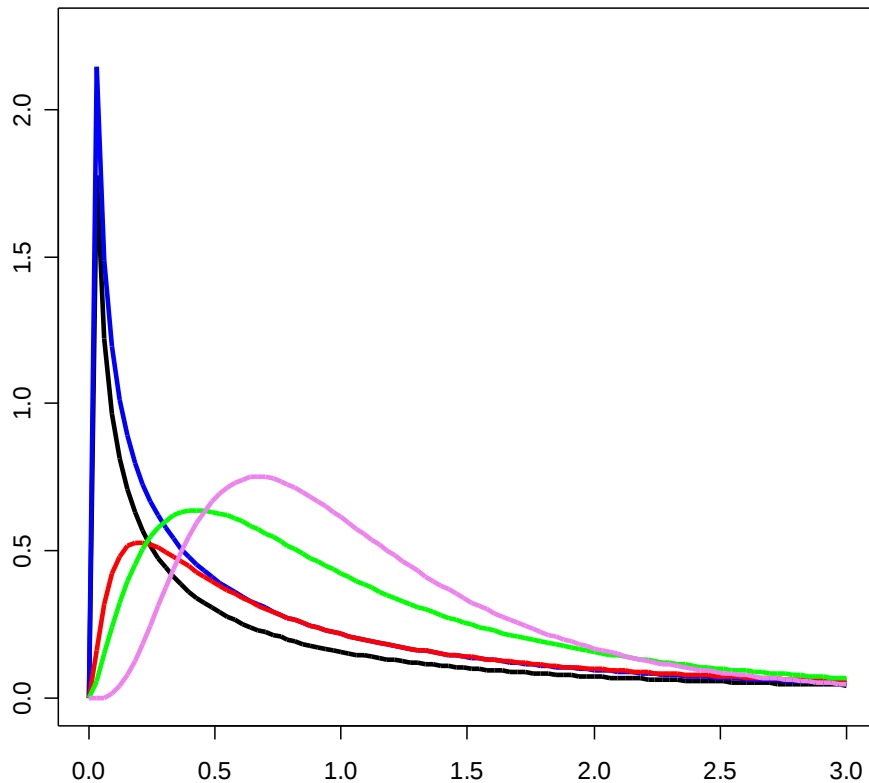
# Quelques lois usuelles

Loi de Fisher  $F_{(n_1, n_2)}$

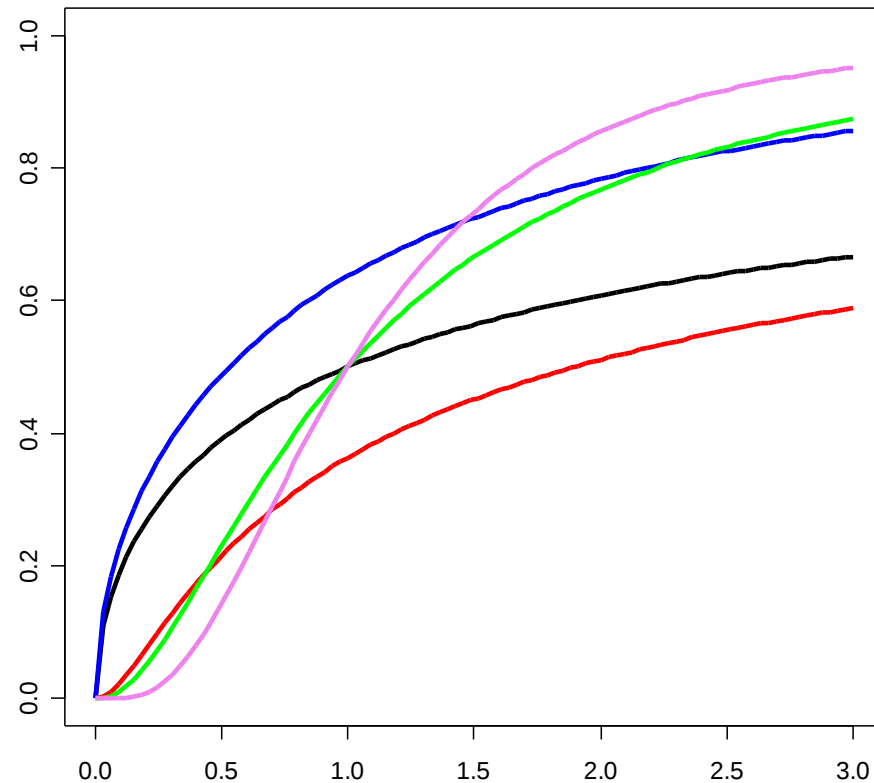
$$f(x) = n_1^{\frac{n_1}{2}} n_2^{\frac{n_2}{2}} \frac{\Gamma\left(\frac{n_1+n_2}{2}\right) x^{\frac{n_1}{2}-1}}{\Gamma\left(\frac{n_1}{2}\right) \Gamma\left(\frac{n_2}{2}\right) (n_1 x + n_2)^{\frac{n_1+n_2}{2}}}$$

$\Gamma$  fonction  
Gamma d'Euler

Ex : **F**(1,1)   **F**(1,5)   **F**(5,1)   **F**(5,5)   **F**(10,10)



Loi de  
probabilité



Fonction de répartition

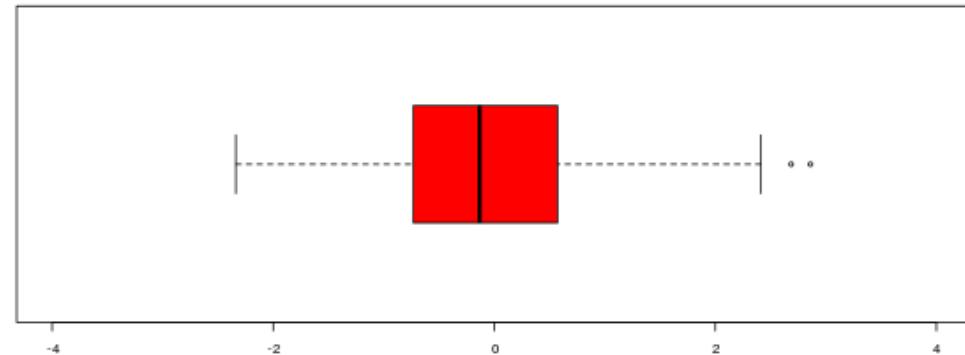
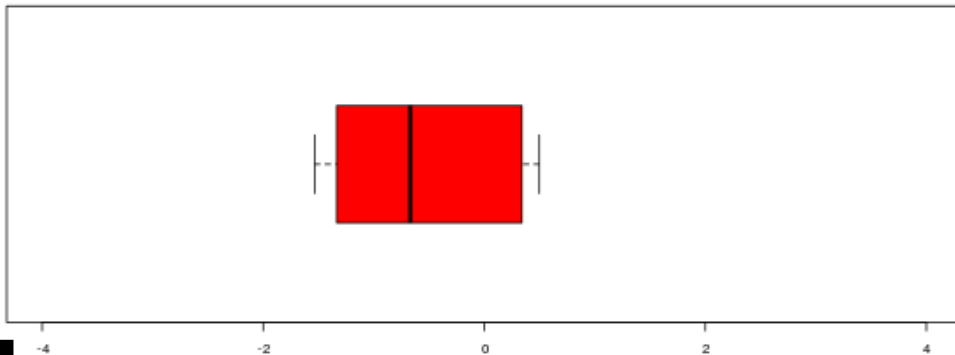
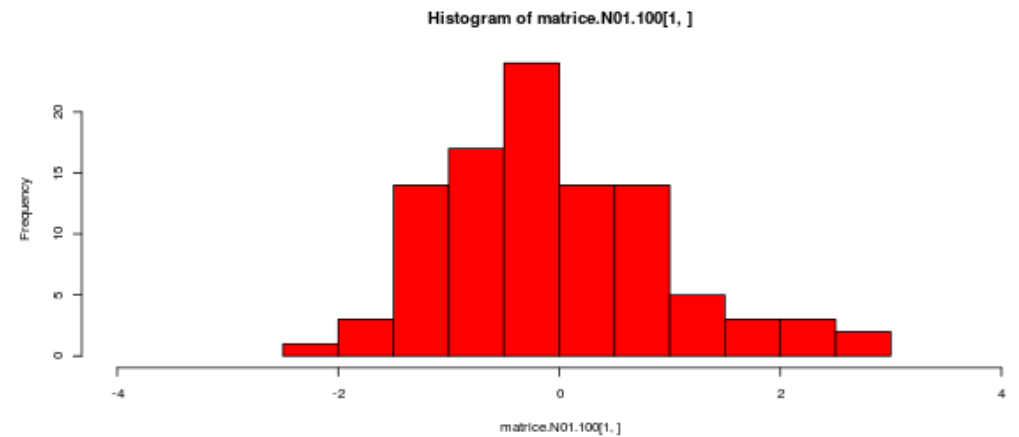
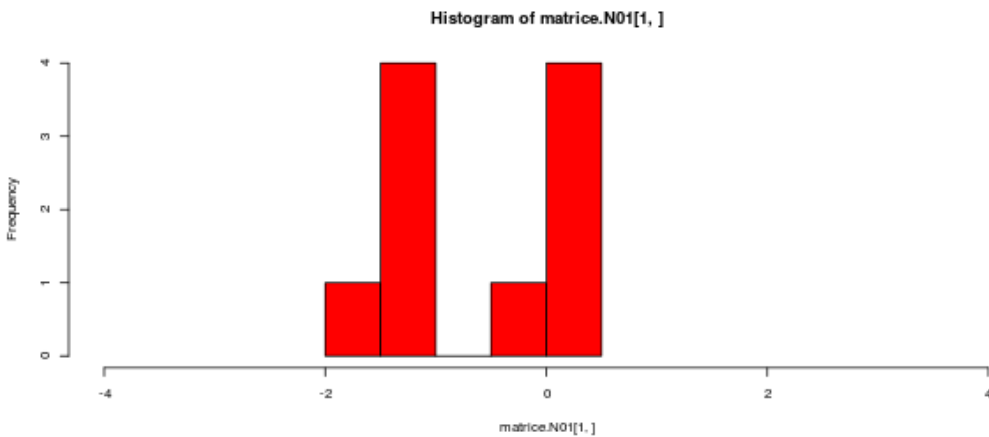
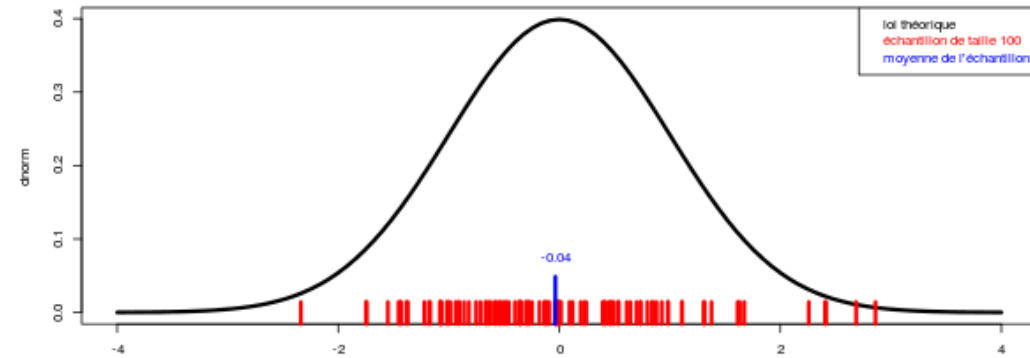
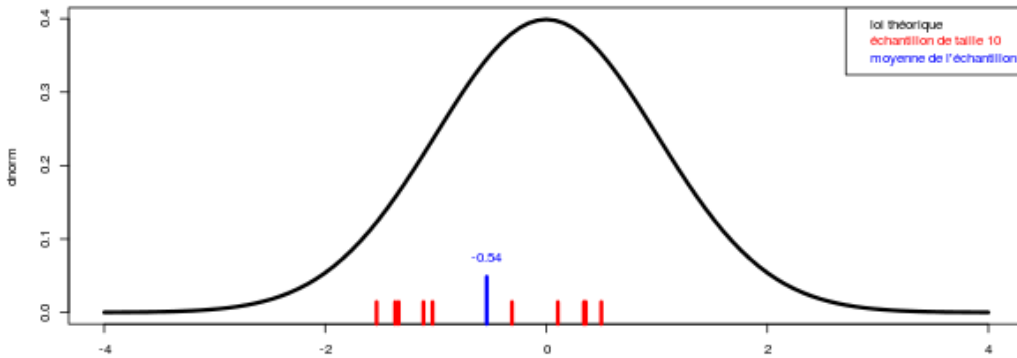


# Des simulations

Échantillon  
de longueur  
**10**

Tirages aléatoires selon une loi normale  $N(0,1)$

Échantillon  
de longueur  
**100**

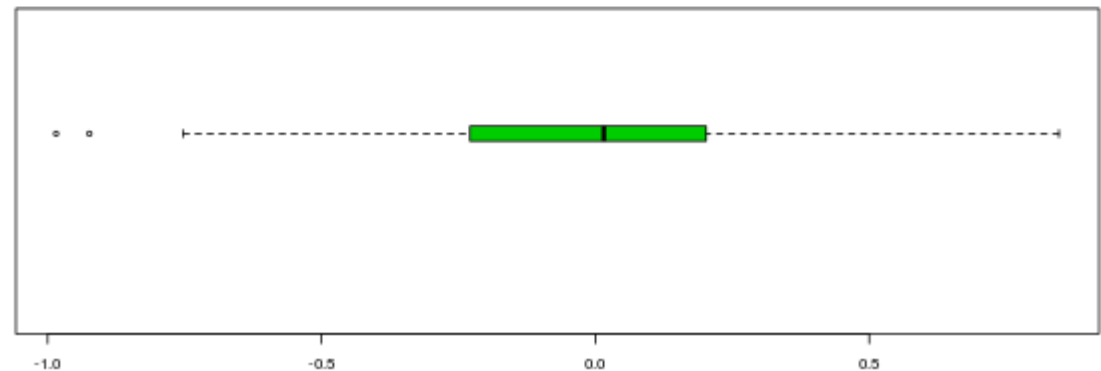
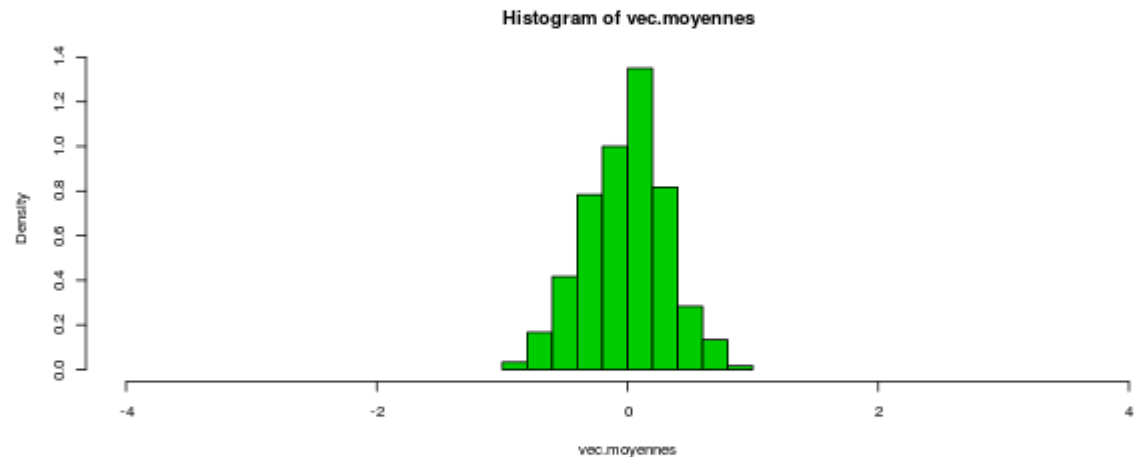
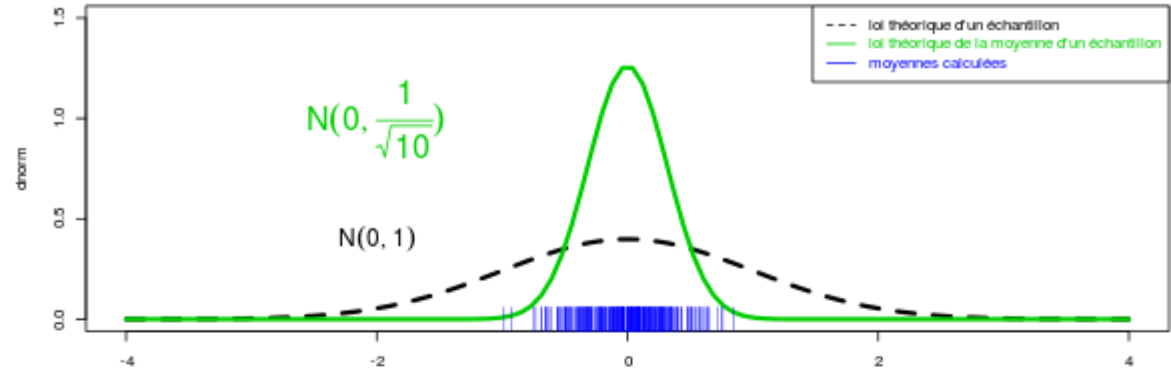


# Des simulations

Supposons maintenant que nous disposons de **300** échantillons de longueur **10** tous issus d'une loi normale  $N(0,1)$ . On peut calculer la moyenne de chaque échantillon et avoir ainsi un vecteur de moyennes de longueur 300.

- La moyenne des 300 moyennes vaut (sur l'illustration) :  
**-0.017**
- L'écart-type des 300 moyennes vaut (sur l'illustration) :  
**0.303**

ce qui proche de la valeur théorique donnée par  $\frac{1}{\sqrt{10}} \approx 0.316$

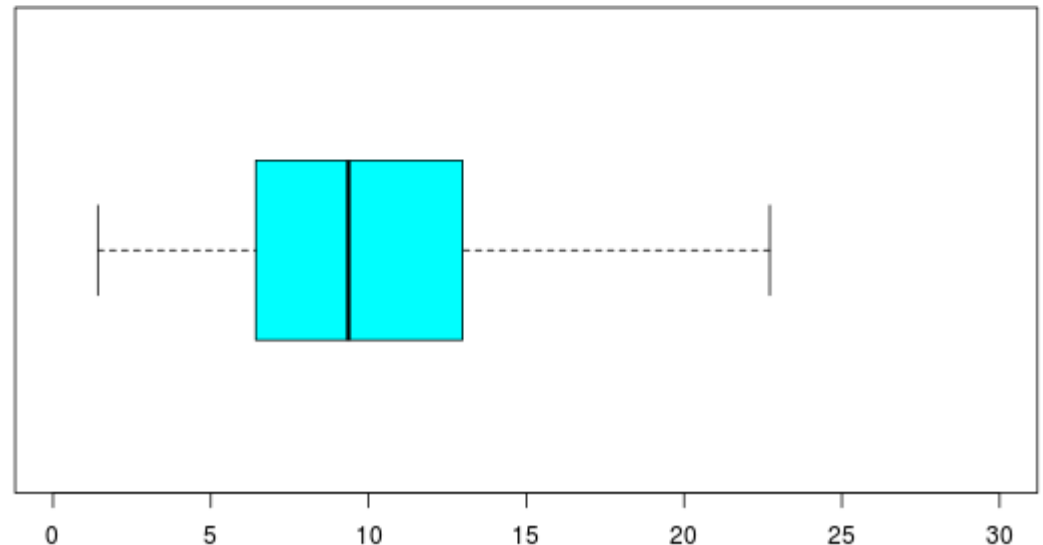
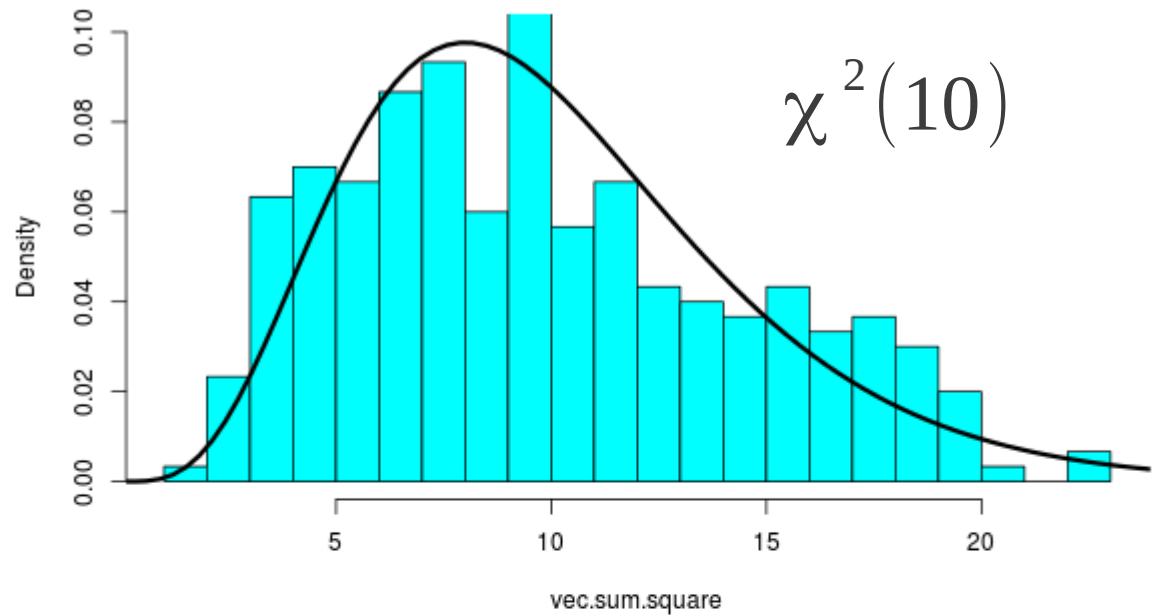


## *Théorie*

$$\begin{cases} X_i \sim N(\mu, \sigma^2), i=1, \dots, n \\ \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \end{cases}$$

# Des simulations

Intéressons-nous maintenant à la somme des carrés des **10** valeurs de chaque échantillon. Nous disposons d'un vecteur de longueur **300** dont l'histogramme est représenté ci-contre. La courbe représente la loi théorique de cette distribution c'est à dire une loi de  $\chi^2$  à 10 degrés de liberté.



## ***Théorie***

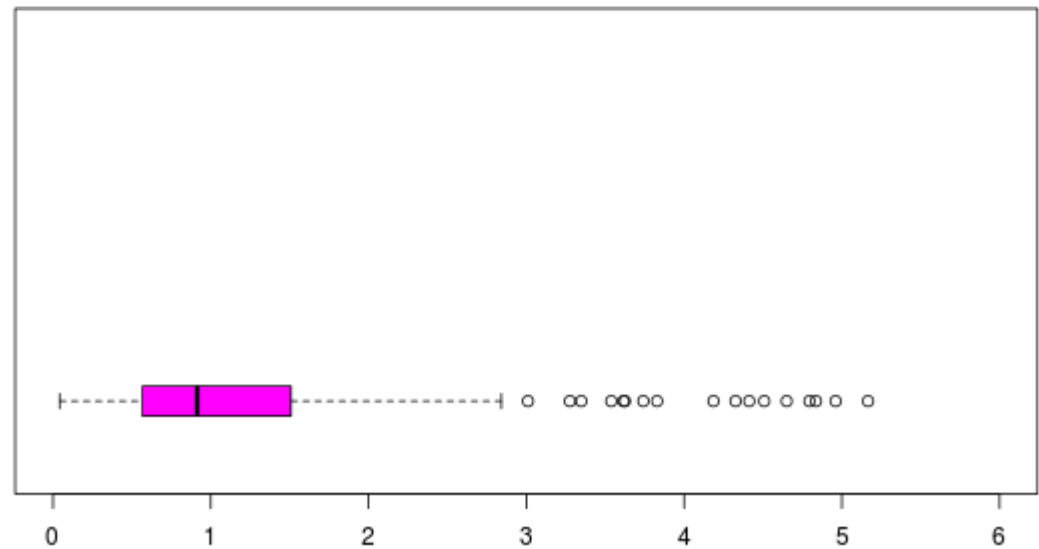
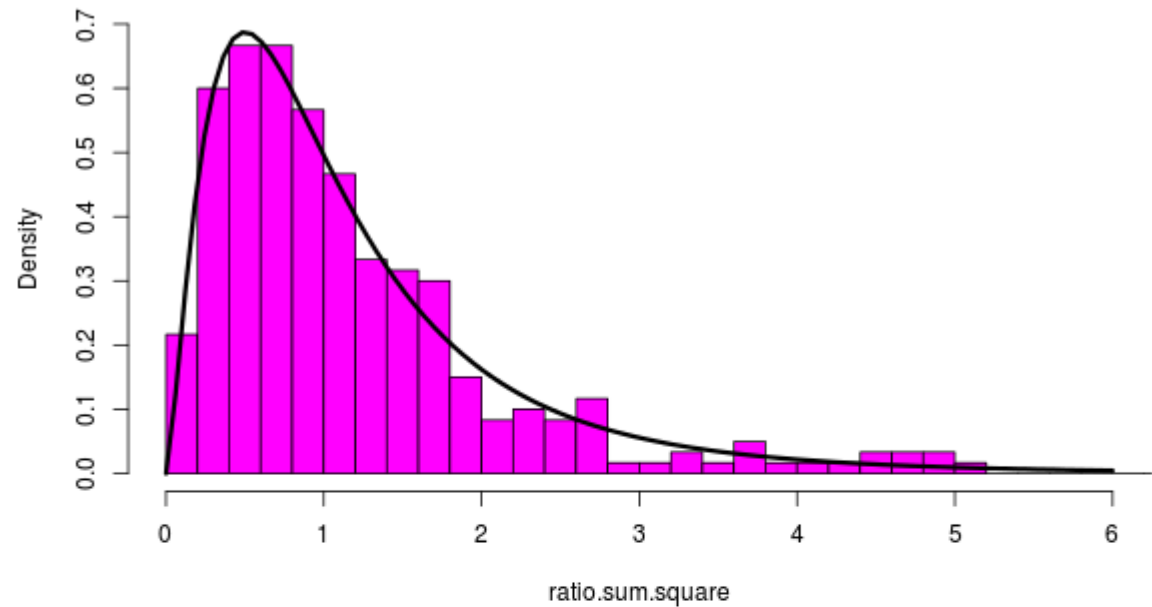
$$X_i \sim N(\mu, \sigma^2), i=1, \dots, n$$

$$\sum_{i=1}^n X_i^2 \sim \chi^2(n)$$

# Des simulations

Supposons maintenant que nous disposons de **300** échantillons de longueur **10** et de **300** échantillons de longueur **5** tous issus d'une loi normale  $N(0,1)$ . On sait que les sommes de carrés des éléments des vecteurs de longueur 10 et 5 suivent des lois de  $\chi^2$  à 5 et 10 degrés de liberté.

Le rapport des 2 variables « sommes de carrés » divisées par la taille des échantillons suit alors une loi de Fisher à 5 et 10 degrés de liberté.



## ***Théorie***

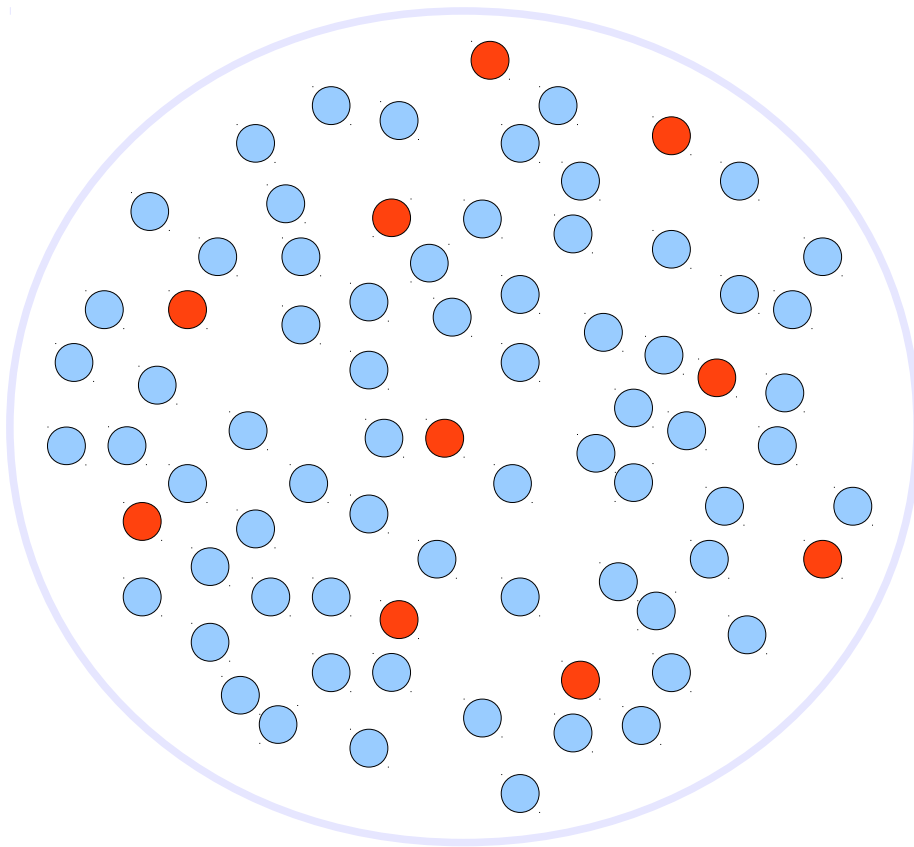
$$C_1 \sim \chi^2(n_1); C_2 \sim \chi^2(n_2)$$

$$\frac{C_1/n_1}{C_2/n_2} \sim F(n_1, n_2)$$

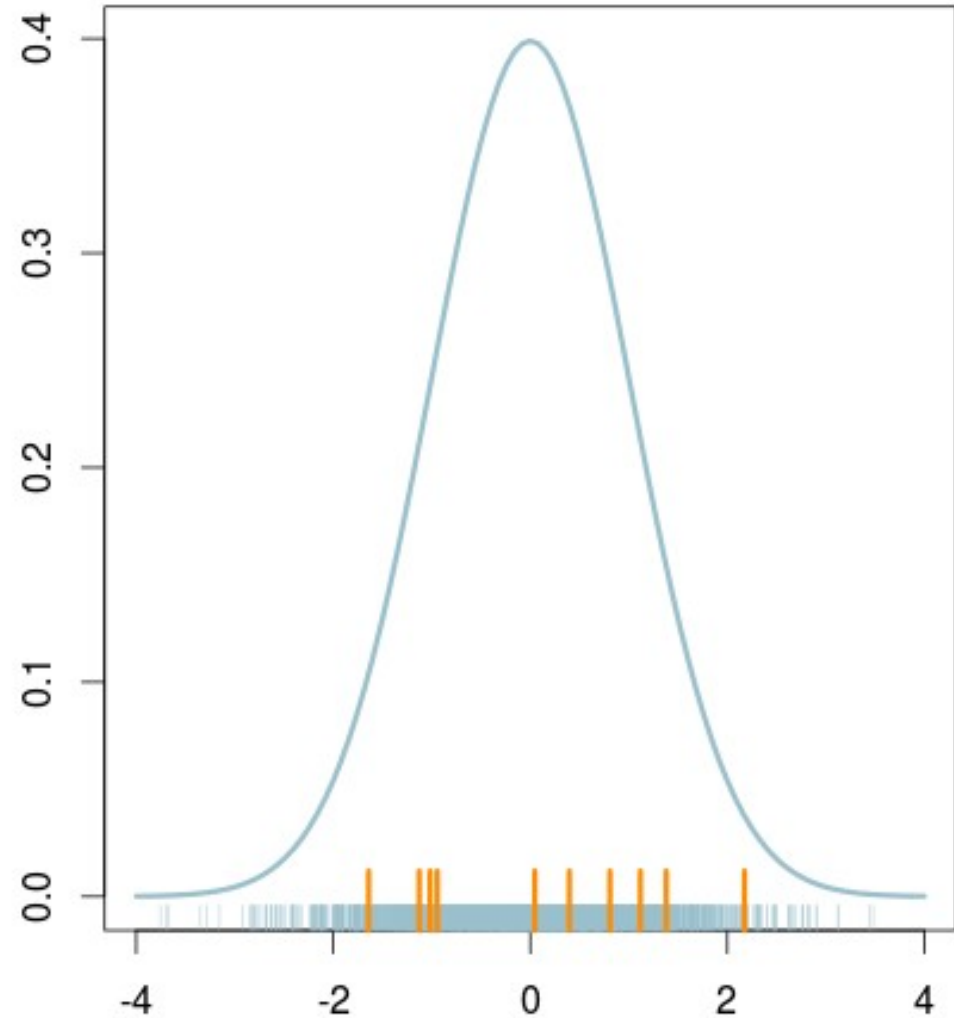
# Statistique inférentielle

- Tirer des conclusions à l'échelle d'une **population** à partir d'informations recueillies sur un **échantillon**.
- Sondage, recensement, échantillon représentatif...
- Lorsque l'on avance des informations quantitatives à l'échelle de la population, on ne parle plus de mesure mais d'**estimation**.
- Les mesures effectuées sur l'échantillon sont des **observations** de la variable aléatoire traduisant le phénomène à l'échelle de la population.

# Statistique inférentielle

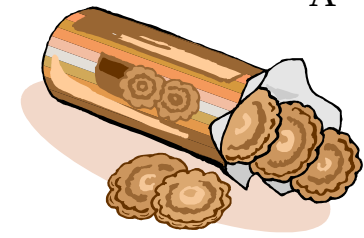


● Population  
● Échantillon



Test statistique : ce qui est observé sur un échantillon permet-il d'invalider une hypothèse faite sur la population ?

# Estimation



Dans une fabrique de biscuits, le procédé mis en œuvre pour vérifier l'aspect moelleux du produit fini consiste à plier le biscuit et à mesurer l'angle d'inclinaison nécessaire pour le casser (un tel test est dit destructif). La règle étant qu'un bon biscuit doit avoir un angle de rupture de  $50^\circ$  (valeur fictive) : si l'angle est inférieur, le biscuit est trop sec, s'il est supérieur, le biscuit est trop moelleux. Tout lot de biscuits doit être validé avant d'être commercialisé.

Il va de soi qu'un biscuit cassé n'est pas commercialisable ainsi qu'un biscuit n'étant pas convenablement moelleux (angle de rupture  $\neq 50^\circ$ ).

Dans de telles conditions, **il est impossible de tester l'ensemble des biscuits** (test destructif). Il est donc nécessaire d'effectuer les mesures sur un **échantillon représentatif** de la population des biscuits (éviter par exemple de prendre les  $n$  premiers ou les  $n$  derniers biscuits fabriqués dans une journée ou, sur une même ligne de production si plusieurs fonctionnent en parallèle). L'angle moyen de rupture calculé sur l'échantillon est un **estimateur** de cet angle chez les biscuits du même lot (aux conditions de fabrication analogues).

# Test statistique

**Exemple** : Fabrication industrielle de biscuits dont l'angle de rupture doit être de  $50^\circ$ . Des facteurs incontrôlés font que cet angle est aléatoire.

**Question** : comment décider qu'un lot est conforme ?

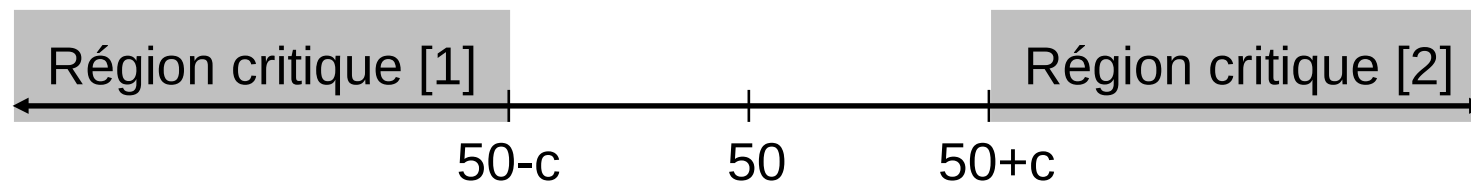
**Hypothèses** :

**H0** : le lot est conforme ( $\mu=50$ )

**H1** : le lot n'est pas conforme ( $\mu \neq 50$ )

Pour trancher entre les 2 hypothèses, on tire au hasard un échantillon de  $n$  biscuits et on en mesure l'angle de rupture  $(X_i)_{i=1,\dots,n}$ . Chaque  $X_i$  suit une loi  $N(\mu, \sigma^2)$ .

**Règle de décision** (principe): Rejet de  $H_0 \Leftrightarrow \bar{X} \notin [50-c ; 50+c]$





# Test statistique : risques d'erreur

Rappel :  
 H0 : lot conforme  
 H1 : lot non conforme

		Décision	
		H1 (rejet de H0)	H0 (accept. H0)
Réalité	H0	$\alpha$	Bonne décision
	H1	Bonne décision	$\beta$

Interprétation des risques (en termes de biscuit) :

- $\alpha$  : rejeter le lot de biscuits alors qu'il est conforme (gaspillage !)  
 → *Le patron ne va pas être content.*
- $\beta$  : déclarer conforme (et donc vendre) des biscuits « défectueux »  
 → *Dans ce cas, c'est le client qui n'est pas content.*

# Région critique et risque $\alpha$

**Règle de décision** : Rejet de  $H_0 \Leftrightarrow \bar{X} \notin [50-c ; 50+c]$

$$\alpha = P[\text{Rejeter } H_0 // H_0 \text{ vraie}] = P[\bar{X} \notin [50-c ; 50+c] // \mu = 50]$$

Sous  $H_0$  ( $\mu = 50$ )

$$X_i \sim N(\mu, \sigma^2) \rightarrow \bar{X} \sim N(\mu, \sigma^2/n)$$

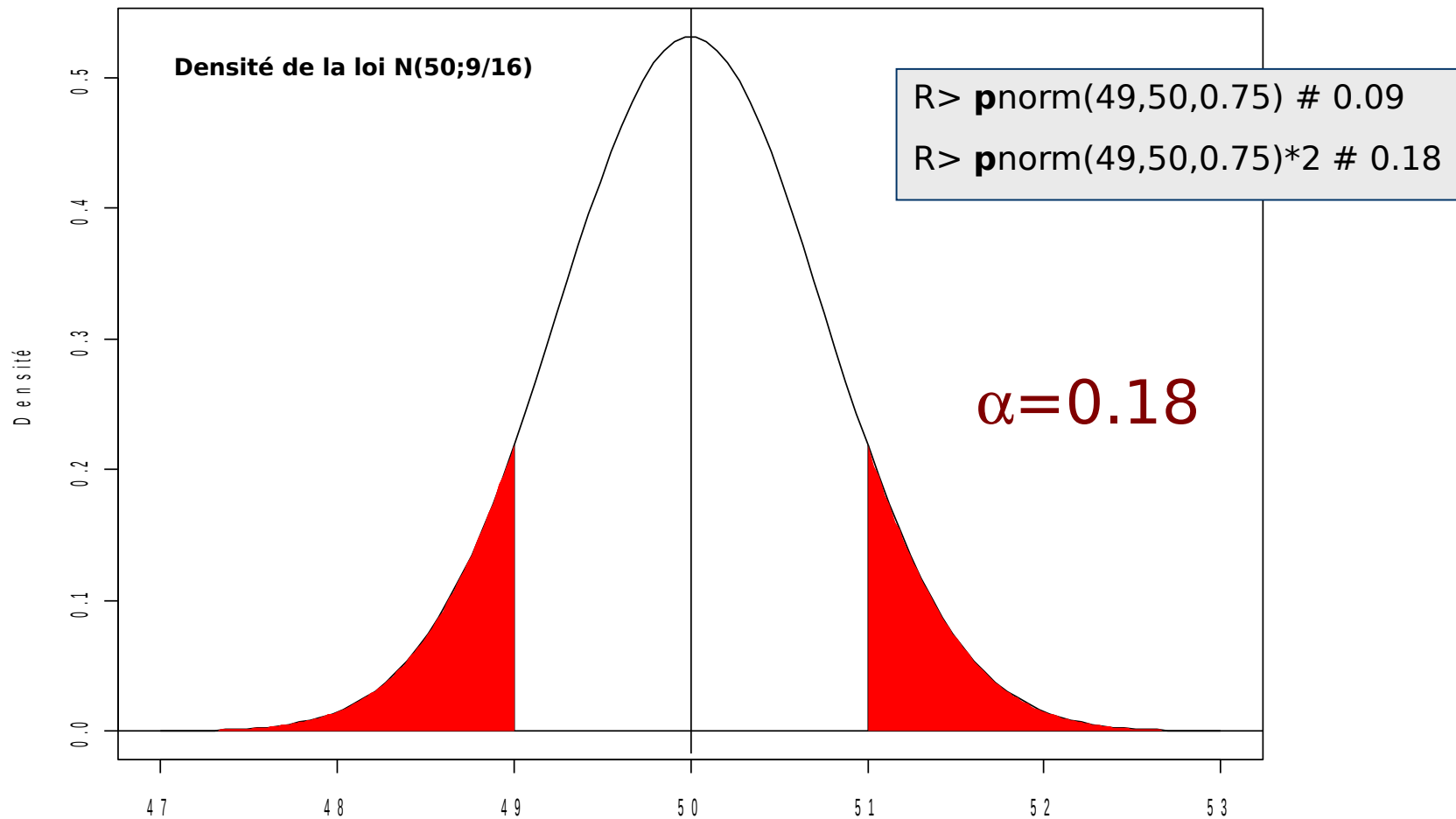
Pour l'application numérique :  $n=16$  et  $\sigma^2=9$

$$\alpha = P[\ll N(50, 9/16) \gg \notin [50-c ; 50+c] ]$$

Le risque  $\alpha$  est la probabilité qu'une variable aléatoire suivant une loi normale de moyenne 50 et de variance 9/16 n'appartienne pas à l'intervalle  $[50-c ; 50+c]$ .

# Région critique et risque $\alpha$

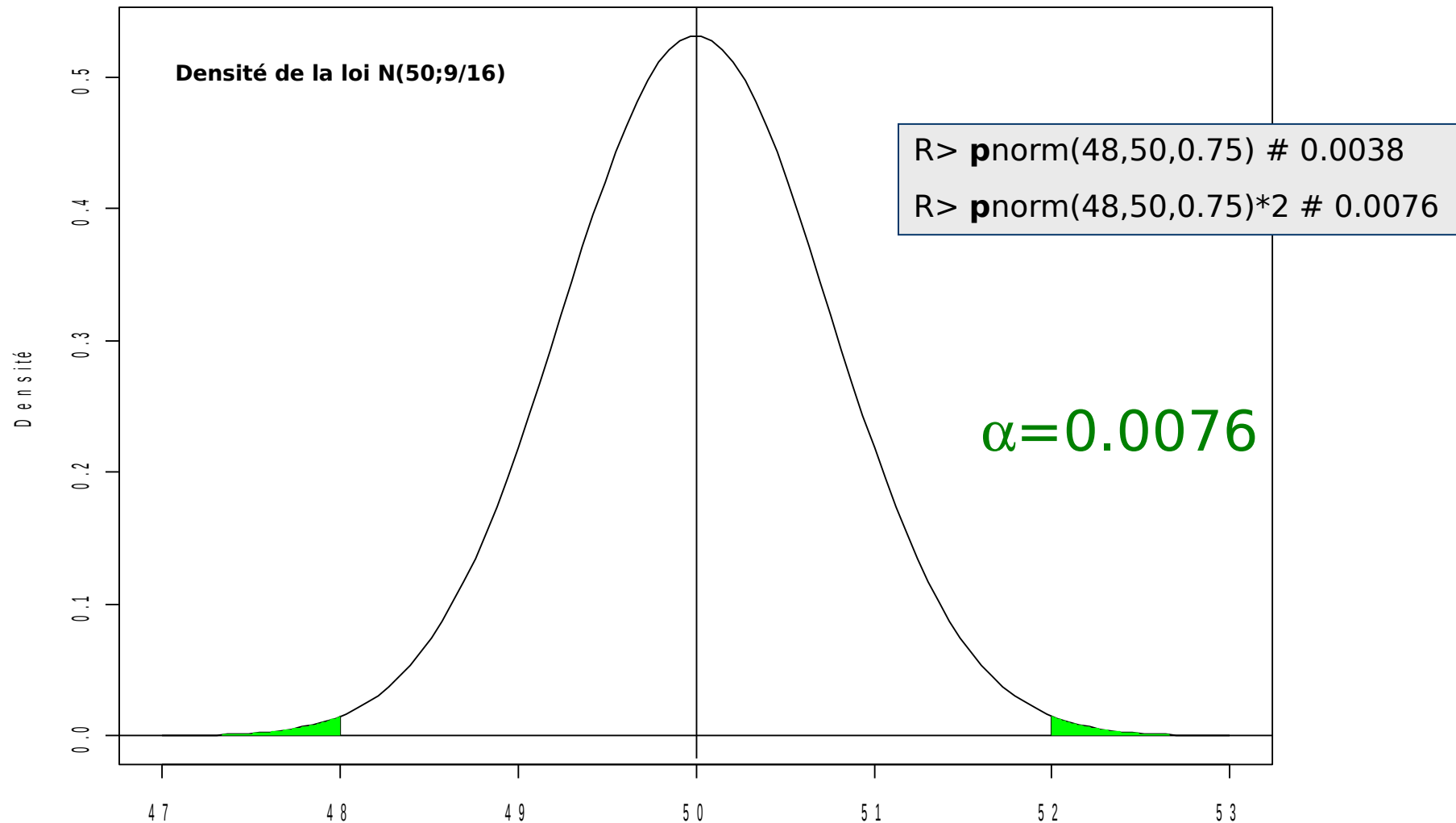
→ Exemple : prenons  $c=1$ , la région critique est  $! [49 ; 51]$ . Calculons le risque  $\alpha$  associé.  
 $\alpha = P[\ll N(50, 9/16) \gg \notin [49 ; 51 ]$



$! [49 ; 51] = ]-\infty ; 49] \cup [51 ; +\infty [$

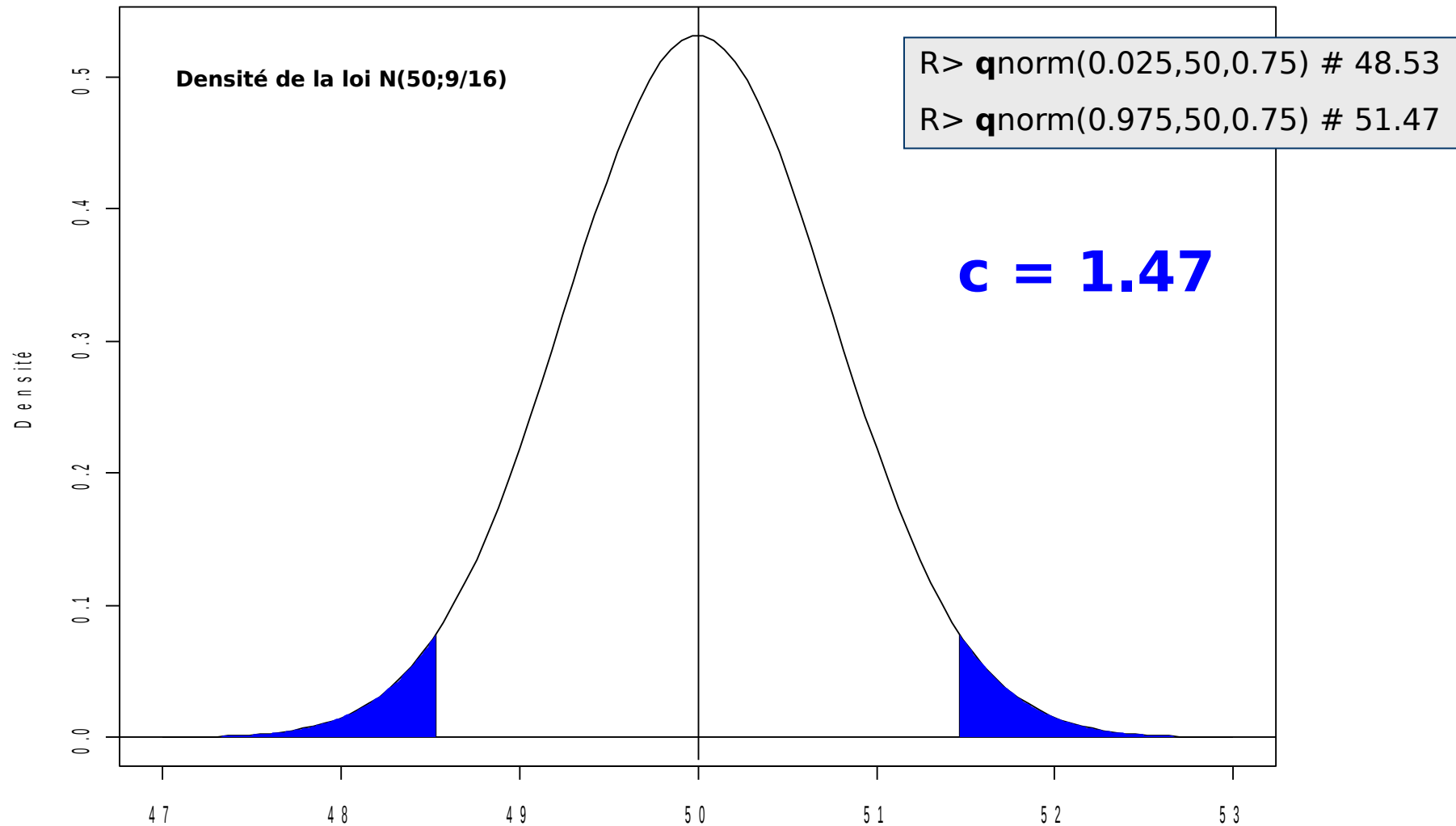
# Région critique et risque $\alpha$

→ Exemple : prenons  $c=2$ , la région critique est  $! [48 ; 52]$ . Calculons le risque  $\alpha$  associé.  $\alpha = P[\ll N(50, 9/16) \gg \notin [48 ; 52 ]$



# Région critique et risque $\alpha$

→ Trouver  $c$  tel que :  $\alpha = P[\ll N(50, 9/16) \gg \notin [50-c ; 50+c]] = \underline{\underline{0.05}}$



# P-value

- « Degré de significativité »
- C'est la plus petite des valeurs de  $\alpha$  pour lesquelles la valeur observée de la statistique de test conduit au rejet de  $H_0$ .
- C'est donc la probabilité d'obtenir, sous  $H_0$ , la valeur observée de la statistique de test ou une valeur plus extrême.
- Il est équivalent de :
  - Comparer la statistique de test avec la valeur limite de la région critique (calculée à partir de  $\alpha$ ).
  - Comparer la p-value et  $\alpha$  : rejet de  $H_0 \Leftrightarrow \text{p-value} < \alpha$

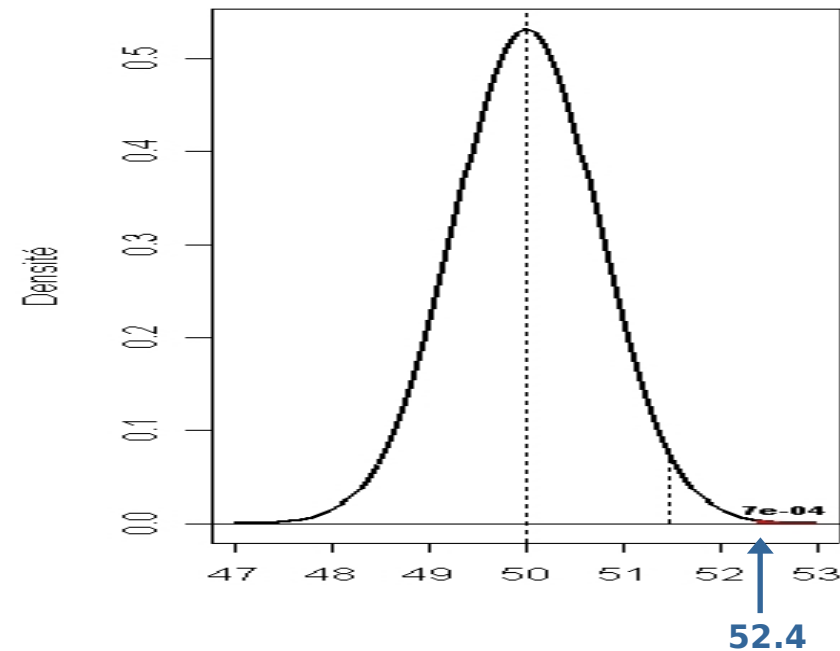
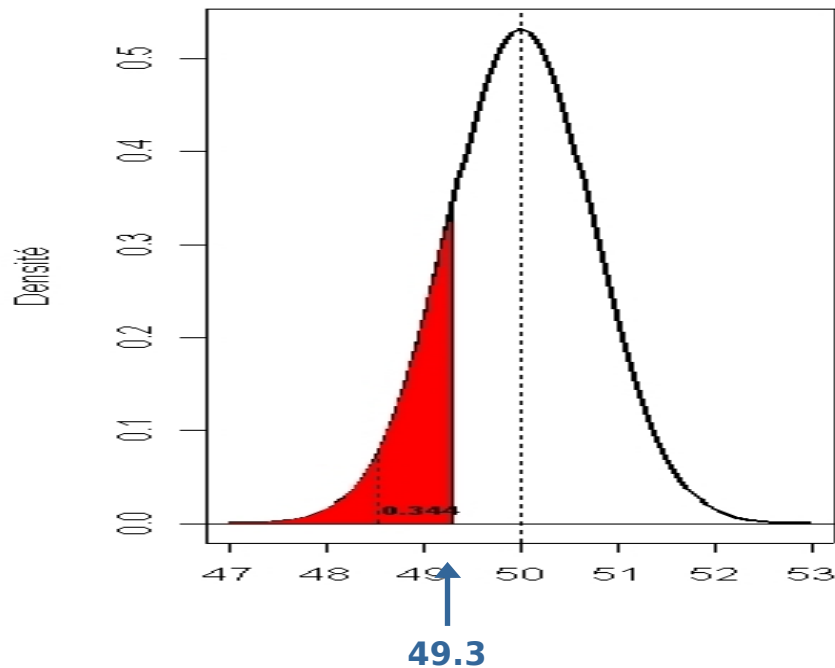
# P-value (exemple)

**Cas 1)** : Jour J1, 16 biscuits tirés au hasard, angle moyen : **49.3**

Cette valeur **n'est pas dans la région critique** ( $[48.53 ; 51.47]$ ), on ne peut pas rejeter  $H_0$ , la production du jour est probablement conforme. La p-value associée à la valeur 49.3 est environ 0.17 ce qui est supérieur au seuil de 5%.

- **Cas 2)** : Jour J2, 16 biscuits tirés au hasard, angle moyen : **52.4**

Cette valeur **est dans la région critique**, on rejette  $H_0$ , la production du jour n'est pas conforme (au seuil de 5%). La p-value associée à la valeur 52.4 est de l'ordre de  $7 \times 10^{-4}$  ce qui est inférieur au seuil de 5%.



# Et le risque $\beta$ ?

$$\beta = P[\text{Accepter } H_0 // H_1 \text{ vraie}] = P[\bar{X} \in [50-c ; 50+c] // \mu = \text{???}]$$

Le calcul explicite du risque  $\beta$  nécessite des valeurs de  $\mu$ .

$$\beta(\mu) = P[\langle N(\mu, 9/16) \rangle \in [50-c; 50+c]]$$

Pour  $c=1$  : RC=[49 ; 51] -  $\alpha=0.18$

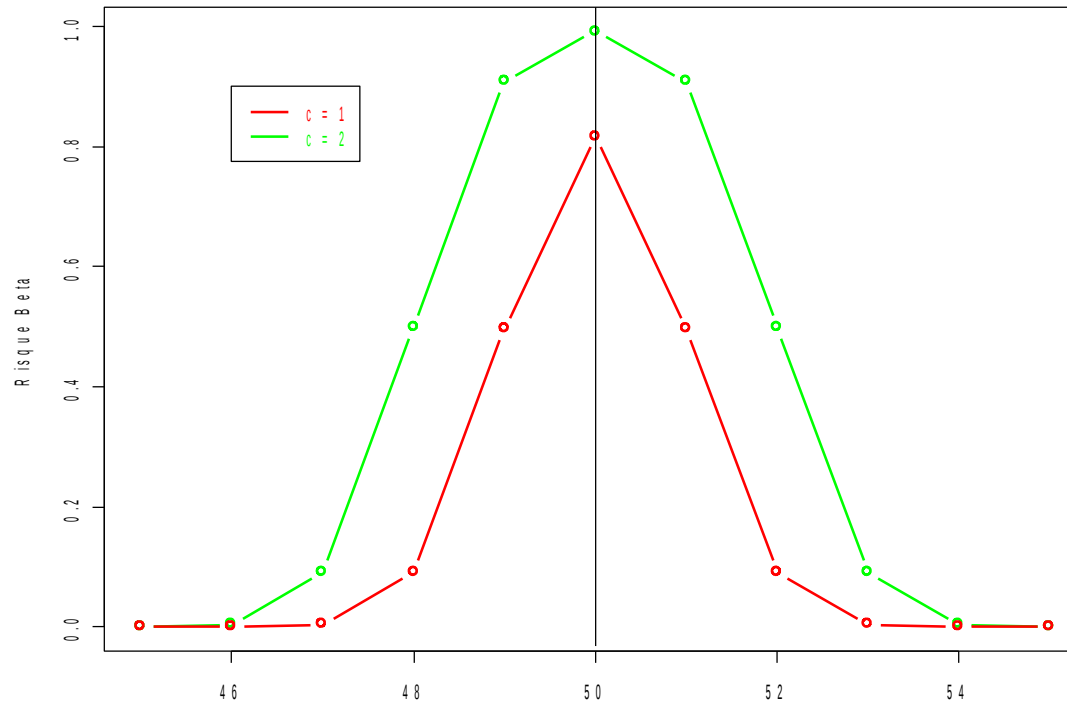
```
R> pnorm(51,45,0.75)-pnorm(49,45,0.75)
```

$\beta(45)=4.10^{-8}$   
 $\beta(46)=3.10^{-5}$   
 $\beta(47)=0.0038$   
 $\beta(48)=0.091$   
 $\beta(49)=0.496$   
 $\beta(50)=0.818$   
 $\beta(51)=0.496$   
 $\beta(52)=0.091$   
 $\beta(53)=0.0038$   
 $\beta(54)=3.10^{-5}$   
 $\beta(55)=4.10^{-8}$

Pour  $c=2$  : RC=[48 ; 52] -  $\alpha=0.0076$

```
R> pnorm(52,45,0.75)-pnorm(48,45,0.75)
```

$\beta(45)=3.10^{-5}$   
 $\beta(46)=0.0038$   
 $\beta(47)=0.091$   
 $\beta(48)=0.5$   
 $\beta(49)=0.91$   
 $\beta(50)=0.99$   
 $\beta(51)=0.91$   
 $\beta(52)=0.5$   
 $\beta(53)=0.091$   
 $\beta(54)=0.0038$   
 $\beta(55)=3.10^{-5}$





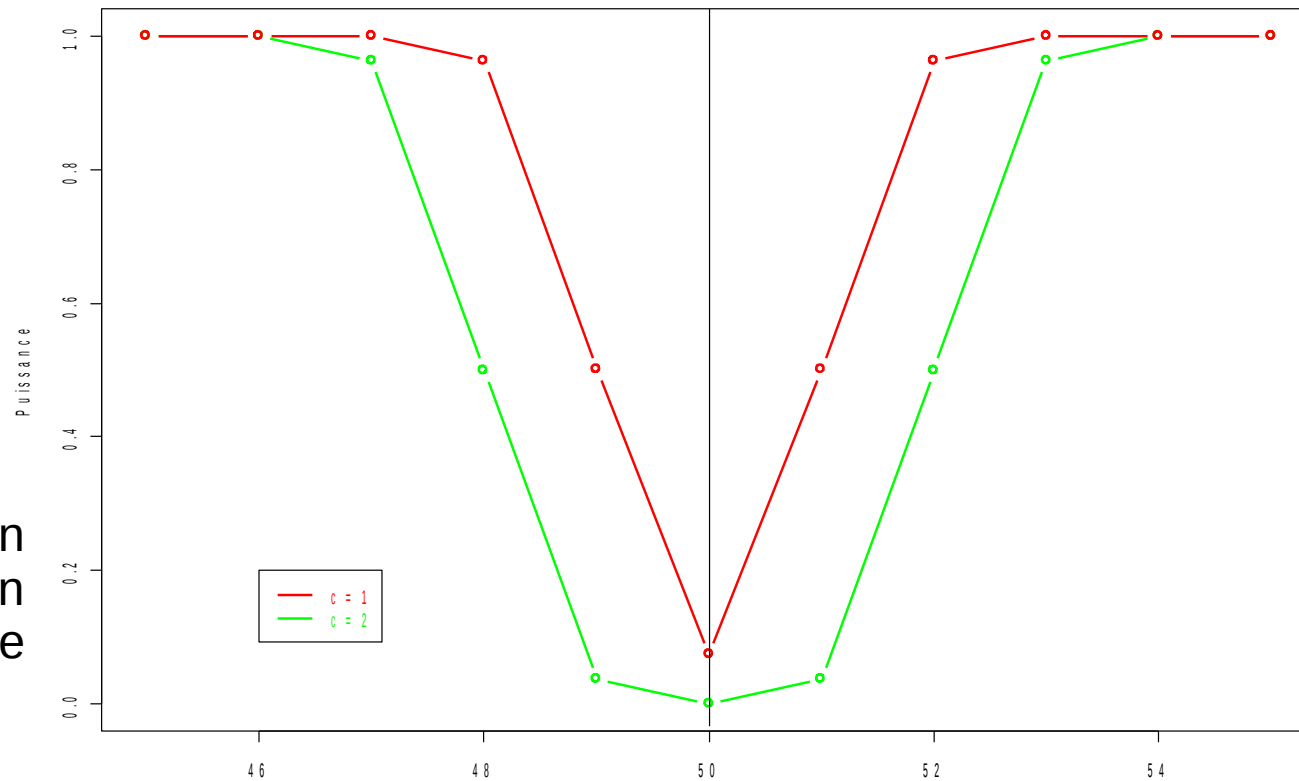
# Puissance d'un test

La puissance d'un test est la probabilité de détecter une différence (rejeter  $H_0$ ) lorsqu'elle existe.

$$\mathbf{P} = 1 - \beta$$

$$= 1 - P[\text{Accepter } H_0 // H_1 \text{ vraie}]$$

$$= P[\text{Rejeter } H_0 // H_1 \text{ vraie}]$$



Représentation  
de la fonction  
puissance

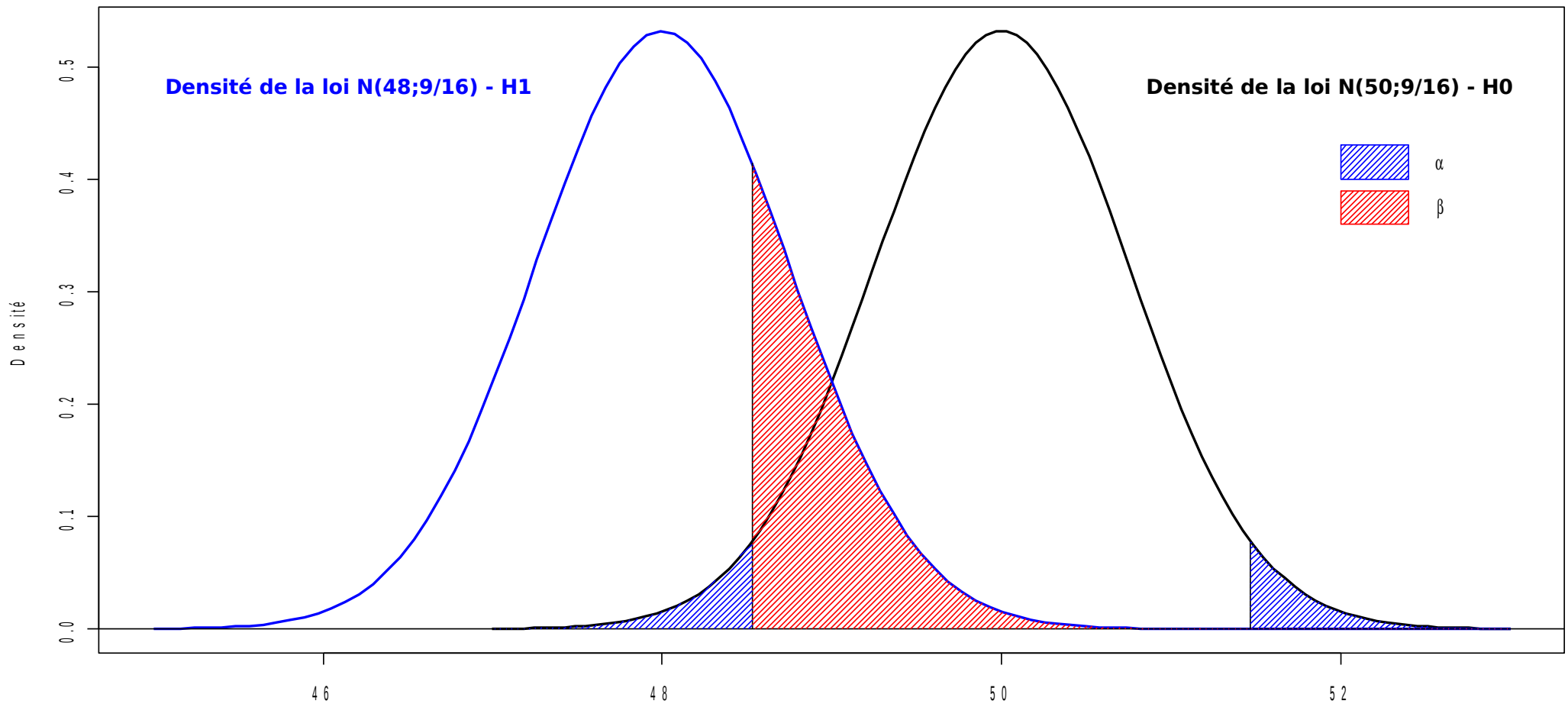
# Représentation graphique de $\alpha$ et $\beta$

H0:  $\mu=50$

H1:  $\mu=48$

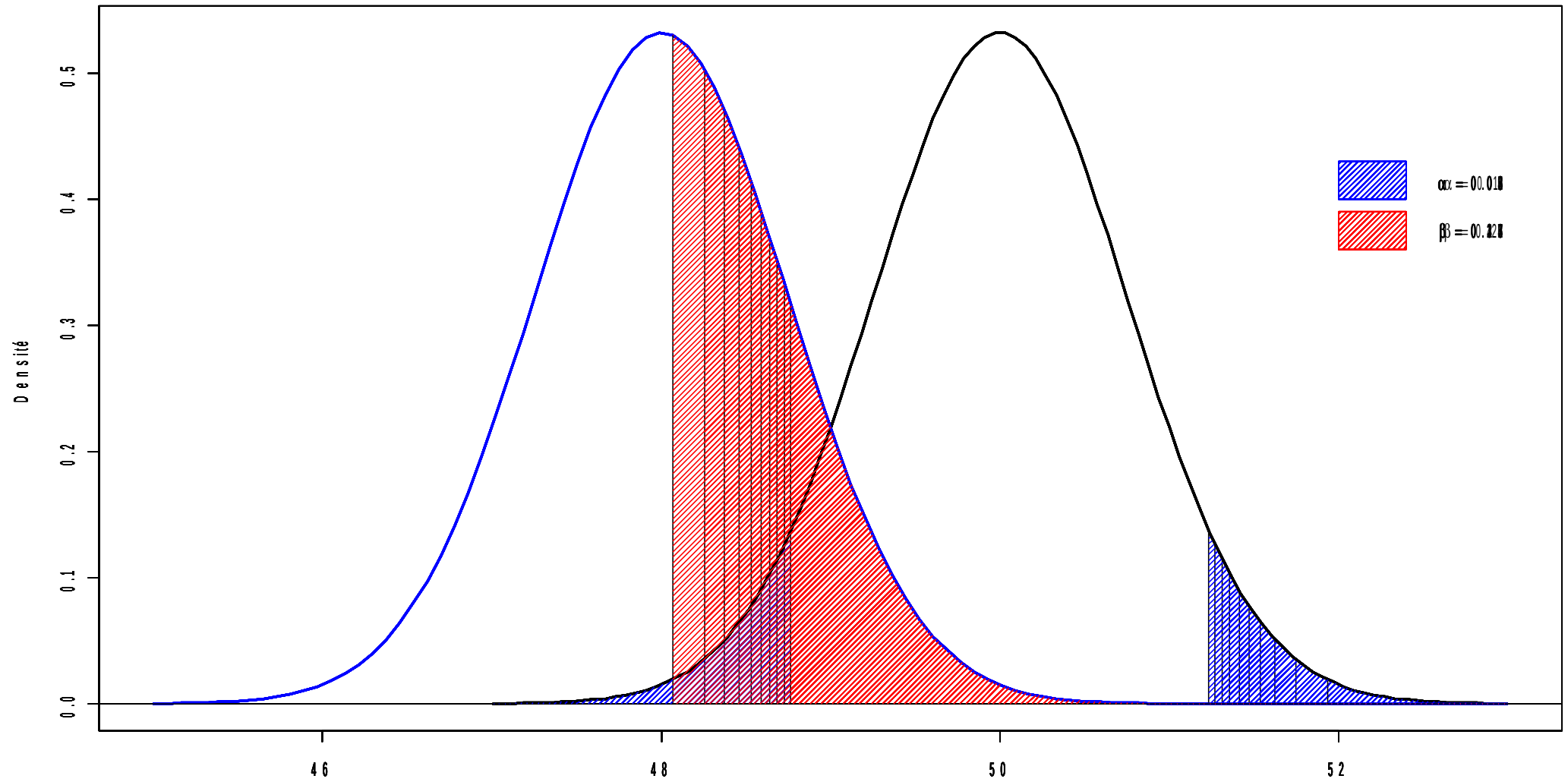
Rappel : pour  $\alpha=5\%$ , la région critique est  $]-\infty ; 48.53] \cup [51.47 ; +\infty [$

Dans ces conditions,  $\beta = 0.24$  `R> 1- pnorm(48.53,48,0.75)`



# Représentation graphique de $\alpha$ et $\beta$

Variations de  $\alpha$  de 0.01 à 0.1

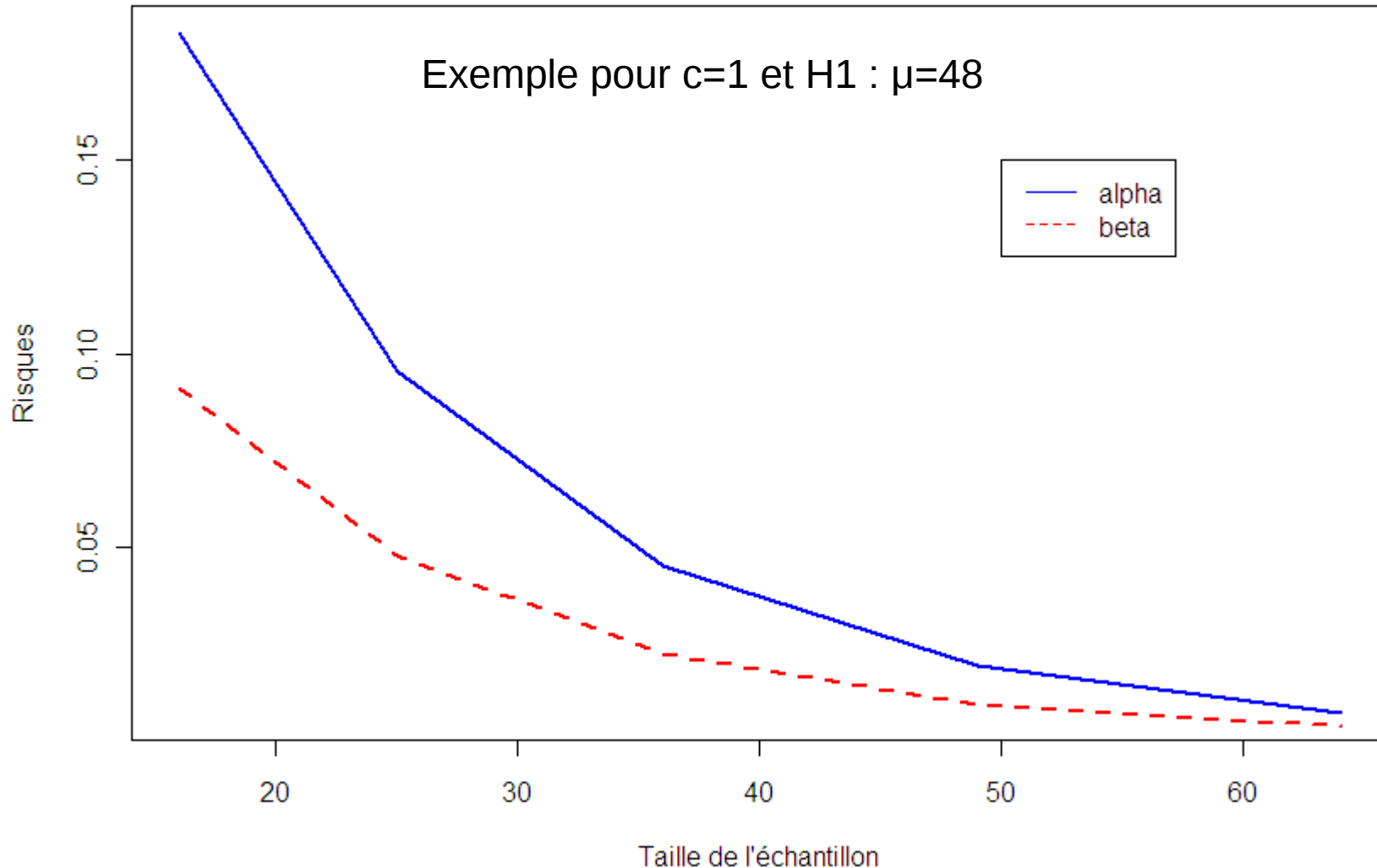


# Test « significatif »

- Si le test conduit à rejeter  $H_0$ , le risque de se tromper ( $\alpha$ ) est faible. La conclusion en faveur de  $H_1$  est solide. Le test est dit **significatif**.
- Si le test conduit à accepter  $H_0$ , le risque de se tromper ( $\beta$ ) peut être grand (selon l'hypothèse alternative). Cette conclusion est moins solide. Dans ce cas, le test est dit **non significatif**. D'où l'habitude d'affirmer « on ne peut pas rejeter  $H_0$  » plutôt que « on accepte  $H_0$  ».
- D'où la nécessaire réflexion du choix des hypothèses  $H_0$  et  $H_1$  ;  $H_1$  étant celle que l'on souhaite voir satisfaite avec un faible risque de se tromper.

# Diminuer $\alpha$ et $\beta$

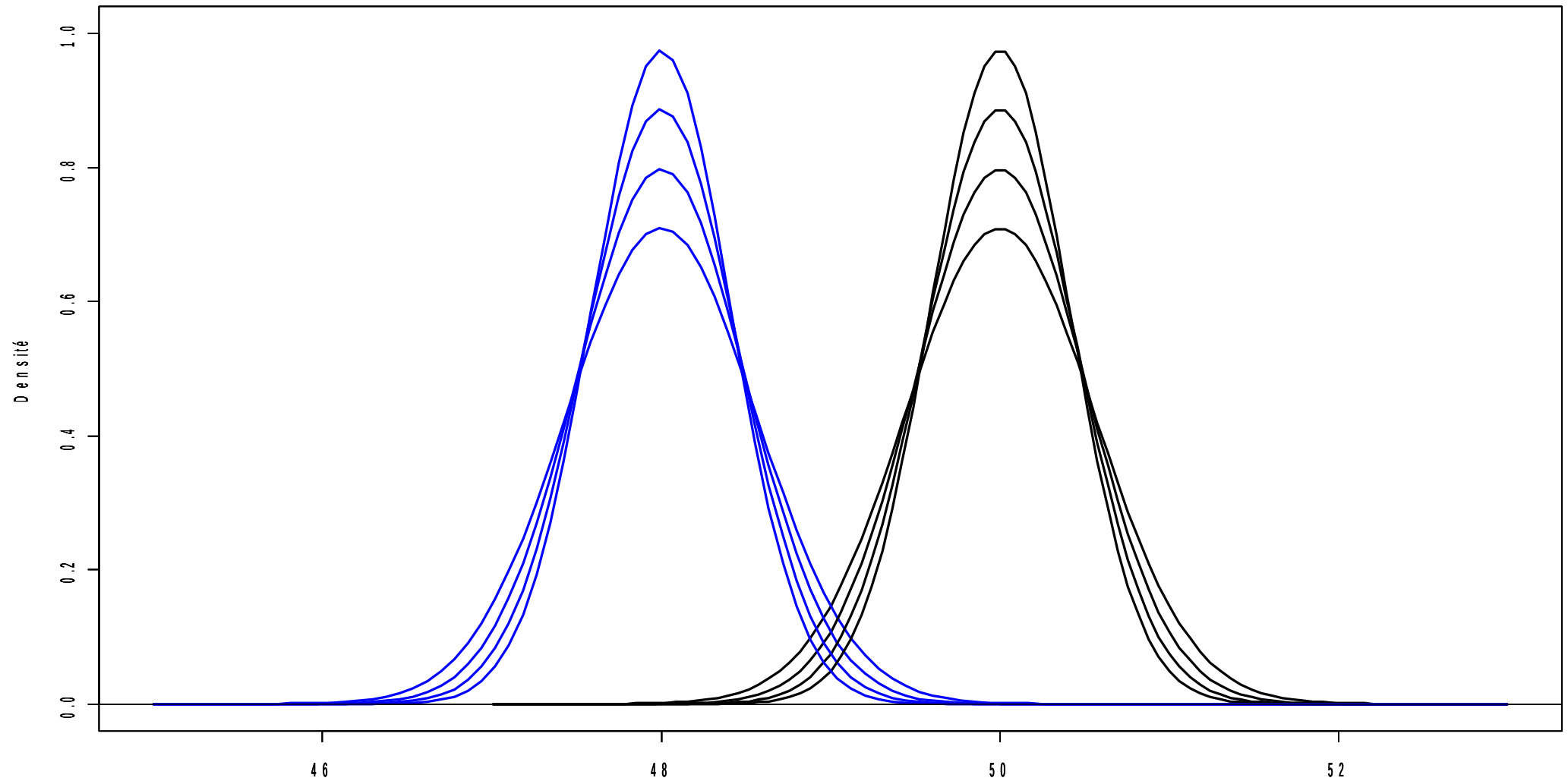
Le seul moyen de diminuer simultanément les risques  $\alpha$  et  $\beta$  consiste à augmenter la taille de l'échantillon (ce qui implique une diminution de la variance de  $\bar{X}$  et donc diminue le recouvrement des 2 courbes).



# Représentation graphique de $\alpha$ et $\beta$

Variations de la taille de l'échantillon

$n = 20$



# Calcul d'effectif : principe

La puissance diminue (= le risque  $\beta$  augmente) quand :

- Le risque  $\alpha$  diminue et / ou
- La taille de l'échantillon diminue et / ou
- La taille de l'effet recherché diminue. En d'autres termes, l'hypothèse alternative  $H_1$  « se rapproche » de l'hypothèse nulle  $H_0$ .

# Calcul d'effectif : principe

Une fois fixés :

- Le risque  $\alpha$
- Le risque  $\beta$  ou la puissance du test ( $1-\beta$ )
- La taille de l'effet à mettre en évidence

seule la taille de l'échantillon reste comme inconnue.

On peut donc la définir **a priori** compte tenu des autres informations.



# Calcul d'effectif : exemple

Extrait de l'aide en ligne de la fonction `power.t.test()` de R

## Power calculations for one and two sample t tests

### Description

Compute power of test, or determine parameters to obtain target power.

### Usage

```
power.t.test(n = NULL, delta = NULL, sd = 1, sig.level = 0.05,
             power = NULL,
             type = c("two.sample", "one.sample", "paired"),
             alternative = c("two.sided", "one.sided"),
             strict = FALSE)
```

### Arguments

<code>n</code>	Number of observations (per group)
<code>delta</code>	True difference in means
<code>sd</code>	Standard deviation
<code>sig.level</code>	Significance level (Type I error probability)
<code>power</code>	Power of test (1 minus Type II error probability)
<code>type</code>	Type of t test alternative One- or two-sided test
<code>strict</code>	Use strict interpretation in two-sided case

### Details

Exactly one of the parameters `n`, `delta`, `power`, `sd`, and `sig.level` must be passed as `NULL`, and that parameter is determined from the others.

...

# Calcul d'effectif : exemple

```
R> power.t.test(n = 25, delta = 1, sd = 1, sig.level = 0.05, power =
NULL)
```

```
Two-sample t test power calculation
```

```
      n = 25
  delta = 1
     sd = 1
sig.level = 0.05
  power = 0.9337076
alternative = two.sided
```

NOTE: n is number in \*each\* group

```
R> power.t.test(n = 30, delta = 1, sig.level = NULL ,power= 0.9)
```

```
      n = 30
  delta = 1
     sd = 1
sig.level = 0.01286591
  power = 0.90
```

```
R> power.t.test(n = NULL, delta = 1, sd = 1, sig.level = 0.05, power = 0,9)
```

```
Two-sample t test power calculation
```

```
      n = 22.02110
  delta = 1
     sd = 1
sig.level = 0.05
  power = 0.9
alternative = two.sided
```

NOTE: n is number in \*each\* group

```
R> power.t.test(n = 30, delta = NULL, sig = 0.05 ,power= 0.9)
```

```
      n = 30
  delta = 0.8511743
     sd = 1
sig.level = 0.05
  power = 0.9
```

# Panorama de quelques tests statistiques

## Problèmes à 1 échantillon

Type de test	Test paramétrique	Test non paramétrique
Conformité à 1 standard	Test de comparaison de moyenne (Student), d'écart-type, d'une proportion à une valeur de référence	
Adéquation à une loi		<ul style="list-style-type: none"> <li>• Kolmogorov-Smirnov</li> <li>• <math>\chi^2</math> d'adéquation</li> <li>• Shapiro-Wilk</li> </ul>

## Association entre variables

Type de test	Test paramétrique	Test non paramétrique
2 variables quantitatives	Coefficient de corrélation de Pearson	
2 variables qualitatives		$\chi^2$ d'indépendance

# Panorama de quelques tests statistiques

## Problèmes à K échantillons : comparaison de population

Type de test	Test paramétrique	Test non paramétrique
Comparaison de populations, les fonctions de répartition sont les mêmes dans les groupes		<ul style="list-style-type: none"> <li>• <b>Kolmogorov-Smirnov</b></li> <li>• Cramer – von Mises</li> </ul>
Tests de comparaison de K échantillons indépendants (différenciation selon les caractéristiques de tendance centrale)	<ul style="list-style-type: none"> <li>• <b>Test de comparaison de moyennes</b> (K=2)</li> <li>• <b>ANOVA</b> (analyse de variance) à 1 facteur</li> </ul>	<ul style="list-style-type: none"> <li>• <b>somme des rangs de Wilcoxon</b> (K=2)</li> <li>• <b>Mann - Whitney</b> (K=2)</li> <li>• Kruskal - Wallis</li> <li>• Test des médianes</li> </ul>
Tests de comparaison de K échantillons indépendants (différenciation selon les caractéristiques de dispersion)	<ul style="list-style-type: none"> <li>• <b>Fisher</b> (K=2)</li> <li>• Bartlett</li> <li>• Cochran</li> <li>• F-max de Hartley</li> </ul>	<ul style="list-style-type: none"> <li>• Ansari - Bradley</li> <li>• Siegel-Tukey</li> <li>• Test des différences extrêmes de Moses</li> </ul>
Tests pour K échantillons appariés (mesures répétées ou blocs aléatoires complets)	<ul style="list-style-type: none"> <li>• <b>Test de Student de comparaison de moyennes pour échantillons appariés</b> (K=2)</li> <li>• Test de comparaison de variances pour échantillons appariés (K=2)</li> <li>• ANOVA pour blocs aléatoires complets</li> </ul>	<ul style="list-style-type: none"> <li>• Test des signes (K=2)</li> <li>• <b>Rangs signés de Wilcoxon</b> (K=2)</li> <li>• Friedman</li> <li>• Test de McNemar (K=2, variables binaires)</li> <li>• Test Q de Cochran (variables binaires)</li> </ul>
Tests multivariés pour K échantillons indépendants	<ul style="list-style-type: none"> <li>• <math>T^2</math> de Hotelling, comparaison de K=2 barycentres (vecteur des moyennes)</li> <li>• <b>MANOVA</b> (analyse de variance multivariée), comparaison de K barycentres : Lambda de Wilks, Trace de Pillai, Trace de Hotelling-Lawley, La plus grande valeur propre de Roy</li> </ul>	

# Données indépendantes ou appariées ?

- Données **indépendantes** : les observations sont indépendantes à l'intérieur de chaque échantillon et d'un échantillon à l'autre

*Ex: résultats scolaires filles et garçons, dosage d'un produit chez 2 groupes de patients ayant reçu une molécule ou un placebo...*

- Données **appariées** : les mêmes individus sont soumis à 2 mesures successives d'une même variable

*Ex: notes de copies soumises à une double correction, dosage d'un produit avant et après un traitement chez les mêmes individus...*

# Test paramétrique ou non paramétrique ?

- **Test paramétrique** : les hypothèses nulle et alternative du test portent sur un paramètre statistique (moyenne ou variance par exemple). Ces tests nécessitent généralement des conditions de validité (distribution normale des données par exemple).
- **Test non paramétrique** : un test non paramétrique porte globalement sur la répartition des données sans hypothèse sur leur distribution.

# Cas de 2 échantillons

## « Comparaison de moyennes »

Type de test	Test paramétrique	Test non paramétrique
Type de données		
<b>Données indépendantes</b>	Test de Student pour 2 échantillons	Test de Wilcoxon-Mann-Whitney <i>Rank-sum test</i>
<b>Données appariées</b>	Test de Student pour 1 échantillon (sur la différence)	Test de Wilcoxon <i>Signed-rank test</i>

# Le test de Wilcoxon-Mann-Whitney (1)

Exemple : la concentration d'un produit est mesurée sur 2 échantillons indépendants de taille respective  $n_1=5$  et  $n_2=6$ . Voici les mesures :

Ech 1 : 1.31 1.46 1.85 1.58 1.64

Ech 2 : 1.49 1.32 2.01 1.59 1.76 1.86

On souhaite savoir si les données sont significativement différentes dans les 2 groupes.

## Procédure du test de W-M-W

- 1) Classer toutes les observations par ordre croissant
- 2) Affecter son rang à chaque observation
- 3) Calculer la somme des rangs d'un échantillon (en général celui de plus petite taille)

## Mise en œuvre :

- 1) 1.31 1.32 1.46 1.49 1.58 1.59 1.64 1.76 1.85 1.86 2.01
- 2) 1 2 3 4 5 6 7 8 9 10 11
- 3) Somme des rangs en bleu :  $W_1 = 25$

L'hypothèse d'égalité des 2 distributions est rejetée si cette valeur s'éloigne « trop » d'une valeur « moyenne ».



# Le test de Wilcoxon-Mann-Whitney (2)

Une procédure alternative (et équivalente) consiste à utiliser la statistique de test U liée à la précédente par la relation :

$$U = W_1 - n_1(n_1 + 1)/2$$

Elle correspond au nombre total de fois où un élément de l'échantillon 1 dépasse un élément de l'échantillon 2.

## Illustration

1.31 1.32 1.46 1.49 1.58 1.59 1.64 1.76 1.85 1.86 2.01

- 1.85 est plus grand que 4 éléments de l'échantillon 2 (1.76 1.59 1.49 1.32)
- 1.64 est plus grand que 3 éléments de l'échantillon 2 (1.59 1.49 1.32)
- 1.58 est plus grand que 2 éléments de l'échantillon 2 (1.49 1.32)
- 1.46 est plus grand que 1 élément de l'échantillon 2 (1.32)

On obtient ainsi **U = 10** (ce qui est bien égal à  $25 - (5 \cdot 6)/2$ ).

C'est cette valeur que l'on retrouve dans la sortie de la fonction `wilcox.test()` du logiciel R.

La p-value obtenue ici (**0.4286**) indique qu'il n'y a pas de décalage (*shift*) entre les positions des 2 séries d'observations.

```
> x<-c(1.31,1.46,1.85,1.58,1.64)
> y<-c(1.49,1.32,2.01,1.59,1.76,1.86)
> wilcox.test(x,y)
           Wilcoxon rank sum test
data:  x and y
W = 10, p-value = 0.4286
alternative hypothesis: true location
shift is not equal to 0
```

# Le test de Student

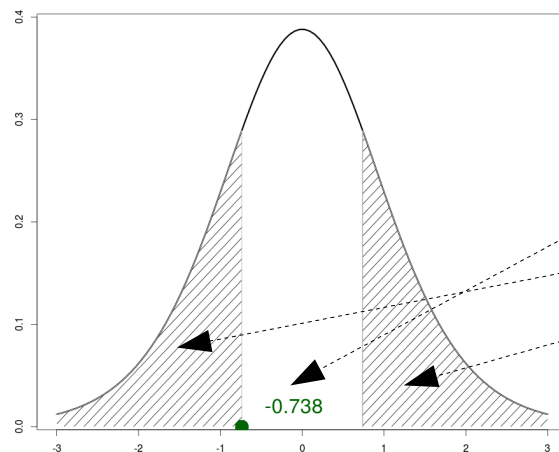
On considère le même problème que précédemment et on applique un test de Student pour comparer la moyenne des 2 échantillons même si les conditions d'application sont plus que discutables.

## Calculs

	1.31	1.49
	1.46	1.32
	1.85	2.01
	1.58	1.59
	1.64	1.76
		1.86
Moyenne	1.658	1.672
Variance	0.041	0.064
Var. Commune	0.054	

$$t = -0.738$$

Densité de la loi de Student à 9 ddl



**Formules** Sous  $H_0$ , hypothèse d'égalité des moyennes, on a :

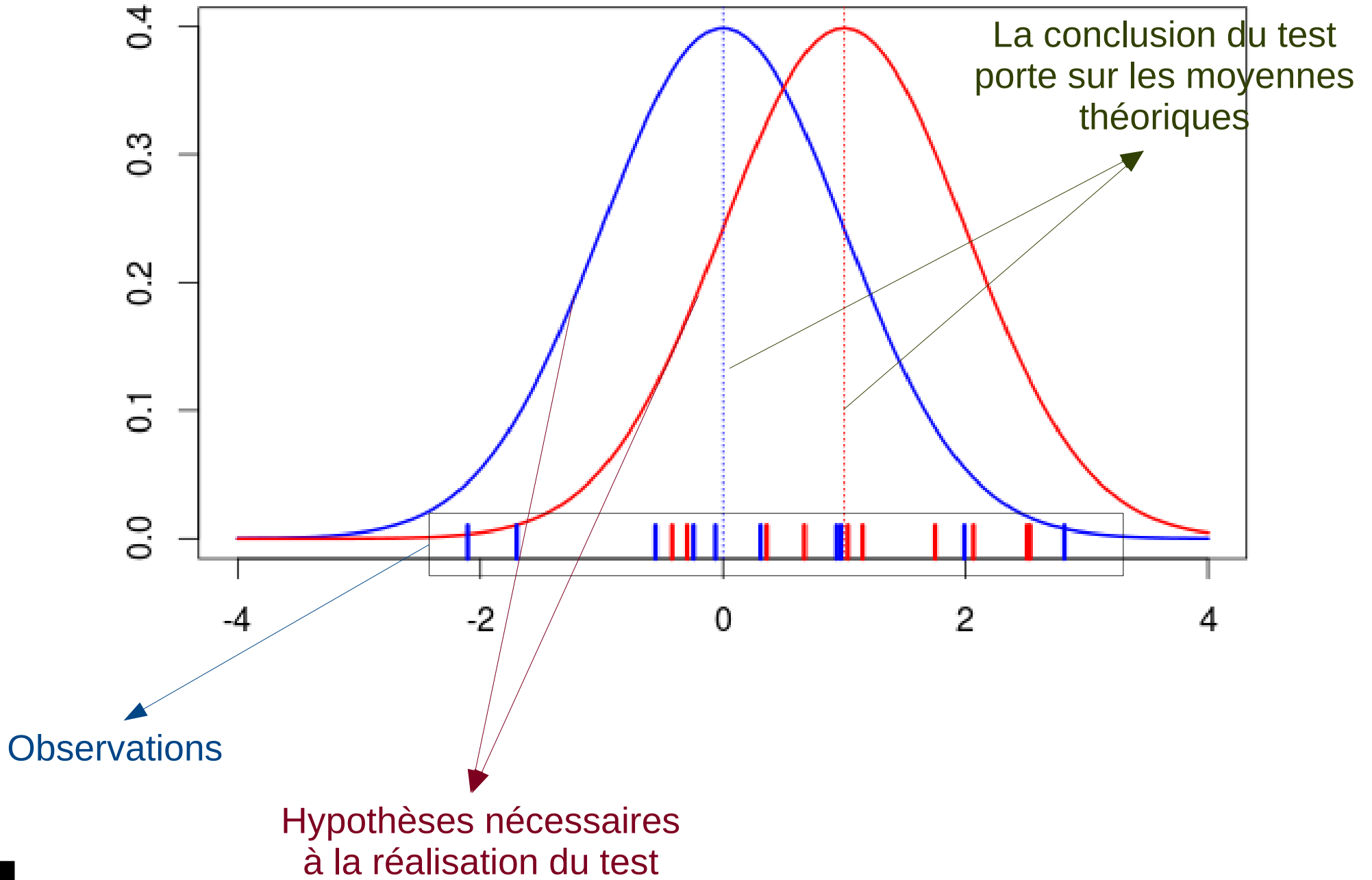
$$\frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim Student(n_1 + n_2 - 2)$$

Avec  $s^2$  la variance commune aux 2 échantillons

$$s^2 = \frac{(n_1 - 1)V_1 + (n_2 - 1)V_2}{n_1 + n_2 - 2}$$

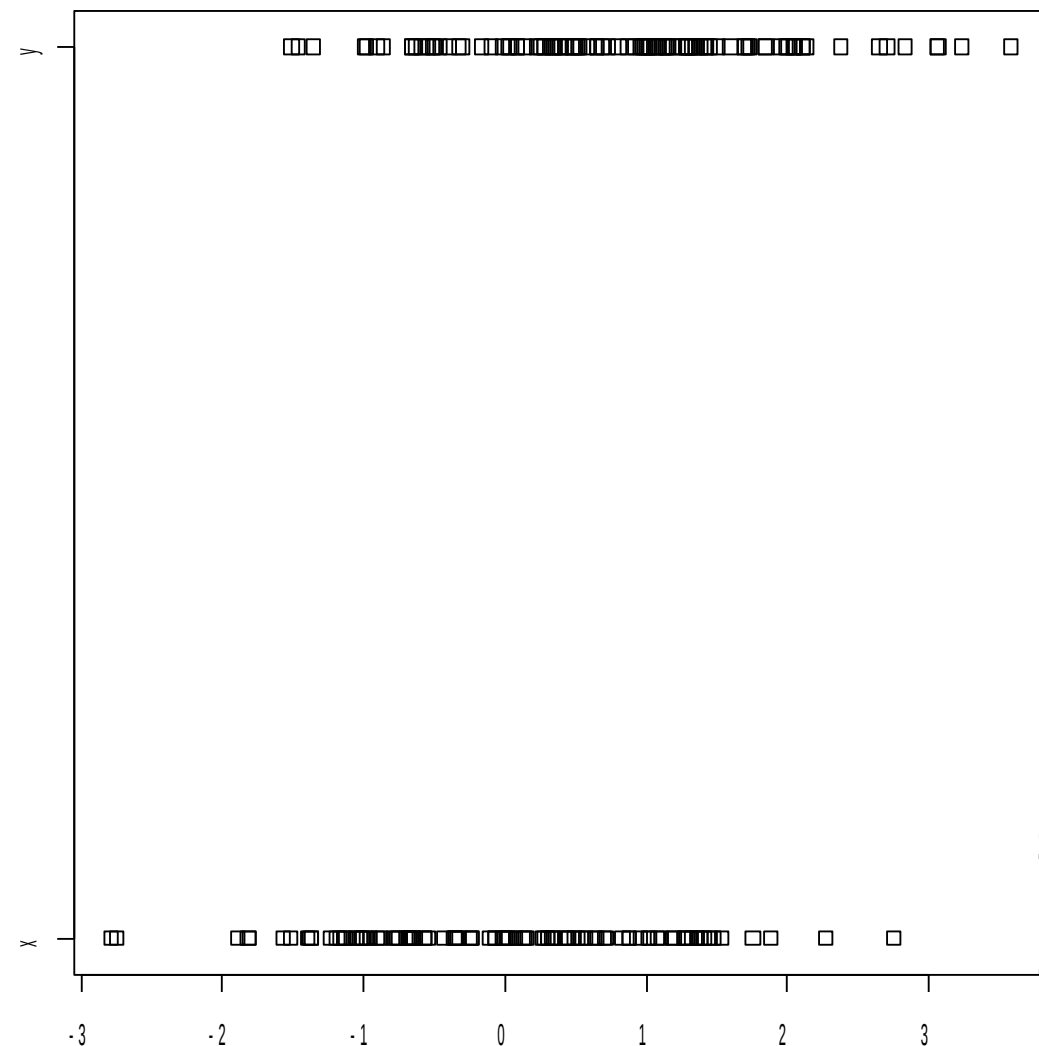
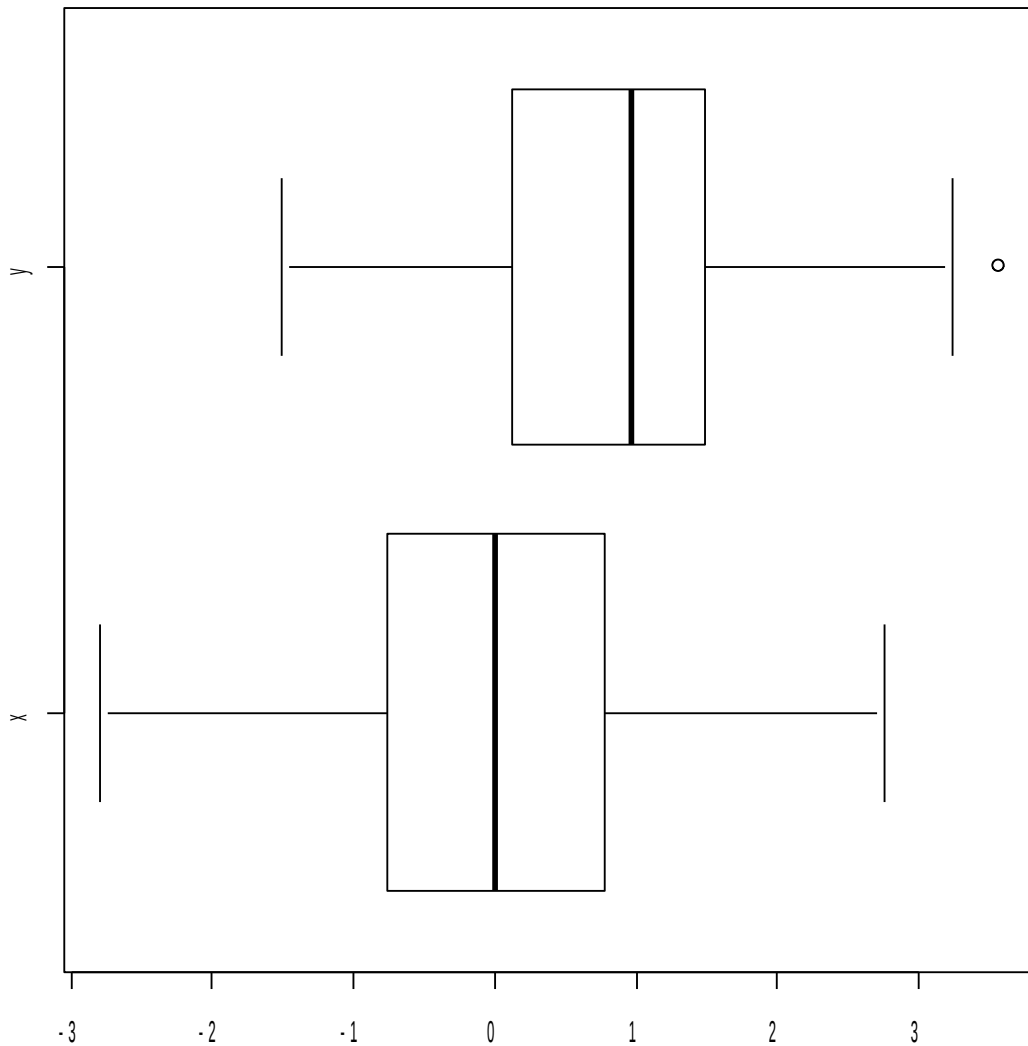
```
> t.test(x,y,var.equal=T)
      Two Sample t-test
data:  x and y
t = -0.7381, df = 9, p-value = 0.4792
alternative hypothesis: true difference
in means is not equal to 0
95 percent confidence interval:
 -0.4213783  0.2140450
sample estimates:
mean of x mean of y
 1.568000  1.671667
```

# Le test de Student



# Mise en œuvre de quelques tests

Données simulées : génération aléatoire selon une loi normale de 2 échantillons de longueur 100 :  $x \sim N(0,1)$  et  $y \sim N(1,1)$

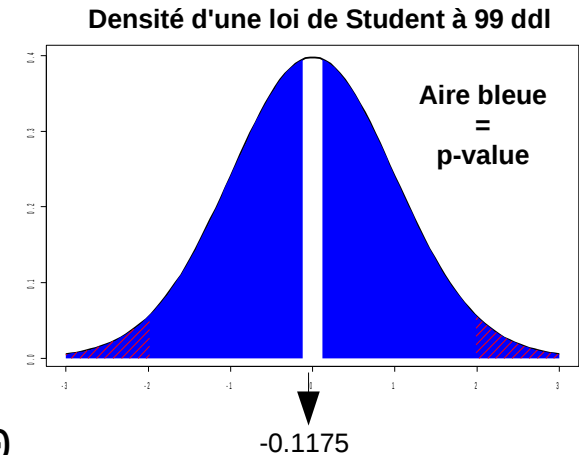


# Mise en œuvre de quelques tests

## Test de Student pour un échantillon

### One Sample t-test

On ne peut pas rejeter  $H_0$ , la moyenne est probablement nulle.

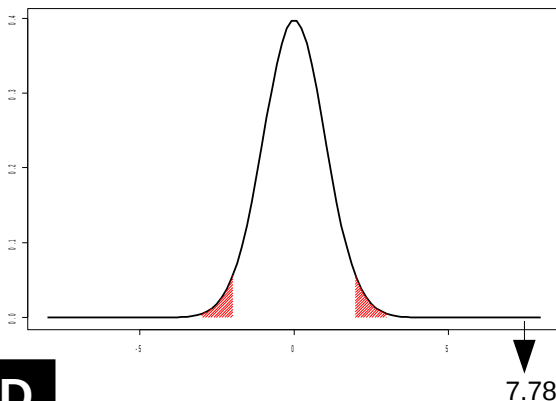


```
data: x
t = -0.1175, df = 99, p-value = 0.9067
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.2239679  0.1989233
sample estimates:
 mean of x
-0.01252230
```

Rejet de  $H_0$  avec une très faible probabilité de se tromper.

### One Sample t-test

Densité d'une loi de Student à 99 ddl



```
data: y
t = 7.78, df = 99, p-value = 7.082e-12
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 0.6291375  1.0599157
sample estimates:
 mean of x
0.8445266
```

# Mise en œuvre de quelques tests

## Test de Fisher d'égalité des variances

### F test to compare two variances

```
data:  x and y
F = 0.9637, num df = 99, denom df = 99, p-value = 0.8545
alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 0.6484291 1.4323091
sample estimates:
ratio of variances
 0.9637173
```

On ne peut pas rejeter  $H_0$ , les 2 variances sont très probablement égales.

# Mise en œuvre de quelques tests

## Test de Student pour 2 échantillons

**Two Sample t-test** (*variances supposées égales*)

data: x and y

t = -5.6342, df = 198, p-value = **5.982e-08**

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

-1.1570238 -0.5570741

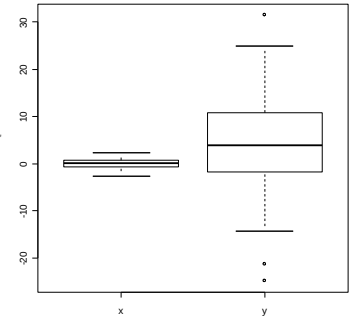
sample estimates:

mean of x mean of y

-0.01252230 0.84452662

On rejette  $H_0$ , les 2 moyennes sont très probablement différentes.

*Pour effectuer ce test, on suppose les 2 variances égales. Cela peut être contrôlé par un test de Fisher d'égalité des variances. Dans le cas ci-contre, la comparaison des moyennes n'a pas vraiment de sens.*



## Welch Two Sample t-test

adaptation du test de Student sans l'hypothèse de variances égales

data: x and y

t = -5.6342, df = 197.932, p-value = **5.985e-08**

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

-1.1570244 -0.5570734

sample estimates:

mean of x mean of y

-0.01252230 0.84452662

# Mise en œuvre de quelques tests

## Test sur le coefficient de corrélation

### Pearson's product-moment correlation

data: x and y

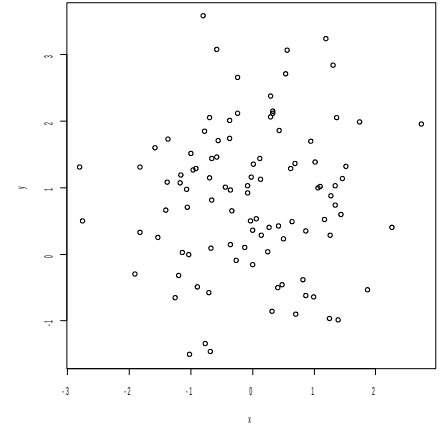
t = 0.5464, df = 98, p-value = 0.586

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

-0.1428544 0.2488346

sample estimates: cor = 0.05511005



### Pearson's product-moment correlation

data: x and z1

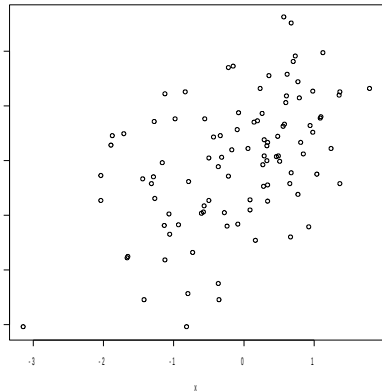
t = 5.5115, df = 98, p-value = 2.88e-07

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

0.3206572 0.6233025

sample estimates: cor = 0.486438



### Pearson's product-moment correlation

data: x and z2

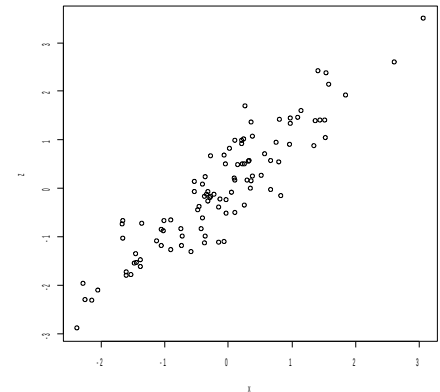
t = 22.3231, df = 98, p-value < 2.2e-16

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

0.8748002 0.9415099

sample estimates: cor = 0.914144





# Mise en œuvre de quelques tests

## Test de normalité Kolmogorov-Smirnov

> `ks.test(x,y)` # *x et Y sont-ils des échantillons d'une même distribution ?*  
 Two-sample Kolmogorov-Smirnov test

data: x and y  
 D = 0.33, p-value = 3.729e-05  
 alternative hypothesis: two-sided

*Probablement pas, avec une faible chance de se tromper*

> `ks.test(x,"pnorm")` # *x est-il un échantillon d'une loi normale  $N(0,1)$  ?*  
 One-sample Kolmogorov-Smirnov test

data: x  
 D = 0.0718, p-value = 0.6803  
 alternative hypothesis: two-sided

*Les données ne permettent pas de dire le contraire.*

> `ks.test(y,"pnorm")` # *y est-il un échantillon d'une loi normale  $N(0,1)$  ?*  
 One-sample Kolmogorov-Smirnov test

data: y  
 D = 0.3408, p-value = 1.641e-10  
 alternative hypothesis: two-sided

*Probablement pas, avec une très faible chance de se tromper.*

> `ks.test(y,"pnorm",1)` # *y est-il un échantillon d'une loi normale  $N(1,1)$  ?*  
 One-sample Kolmogorov-Smirnov test

data: y  
 D = 0.0923, p-value = 0.3614  
 alternative hypothesis: two-sided

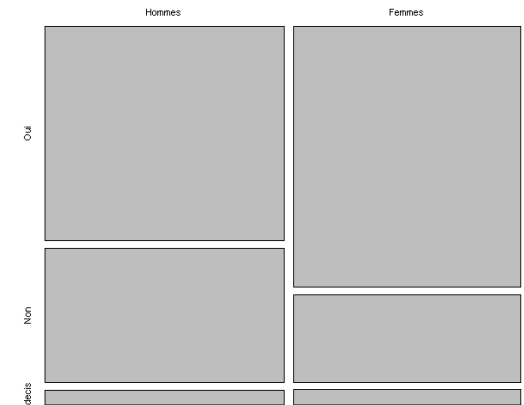
*Les données ne permettent pas de dire le contraire.*

# Quelques tests classiques

## Test du $\chi^2$ d'indépendance

- **Données** : effectifs recueillis dans une table de contingence (tableau croisé pour 2 variables qualitatives)
- **Question** : les 2 variables qualitatives sont-elles indépendantes ?
- **Exemple** : 1250 personnes ont répondu à la question « Êtes-vous satisfaits des programmes TV ? ». On souhaite savoir si la satisfaction dépend du sexe.

	OUI	NON	Indécis	Somme
Hommes	378	237	26	<b>641</b>
Femmes	438	146	25	<b>609</b>
Somme	<b>816</b>	<b>383</b>	<b>51</b>	<b>1250</b>



**Hypothèse H0** : Satisfaction et sexe sont indépendants

### Effectifs théoriques sous l'hypothèse d'indépendance

(effectif d'une case = effectif de la ligne \* effectif de la colonne / effectif total)

	OUI	NON	Indécis
Hommes	418	196	26
Femmes	398	187	25

Statistique de test :  $\chi^2_{\text{obs}} = \sum (\text{Obs} - \text{Théo})^2 / \text{Théo}$

Reflète l'écart entre les données observées et les effectifs théoriques en cas d'indépendance

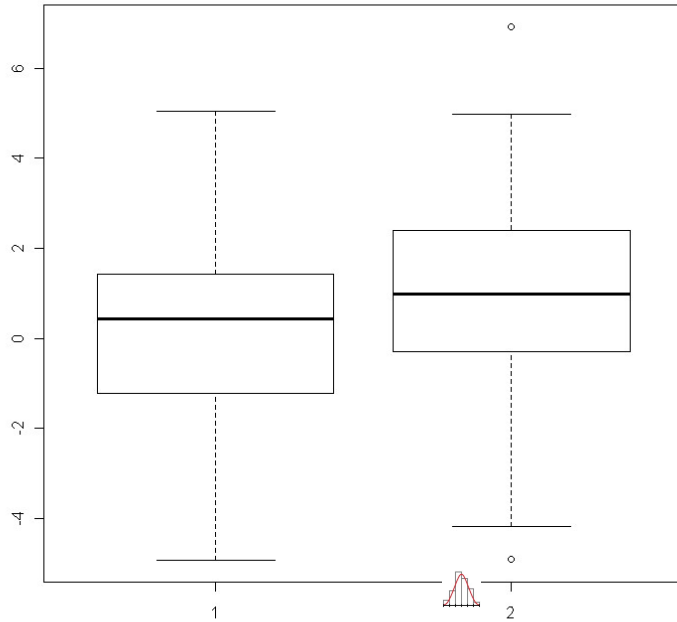
### Pearson's Chi-squared test

data: data.chisq

X-squared = 25.2501, df = 2, p-value = 3.289e-06

Les 2 caractères ne semblent pas indépendants.

# P-value et \*

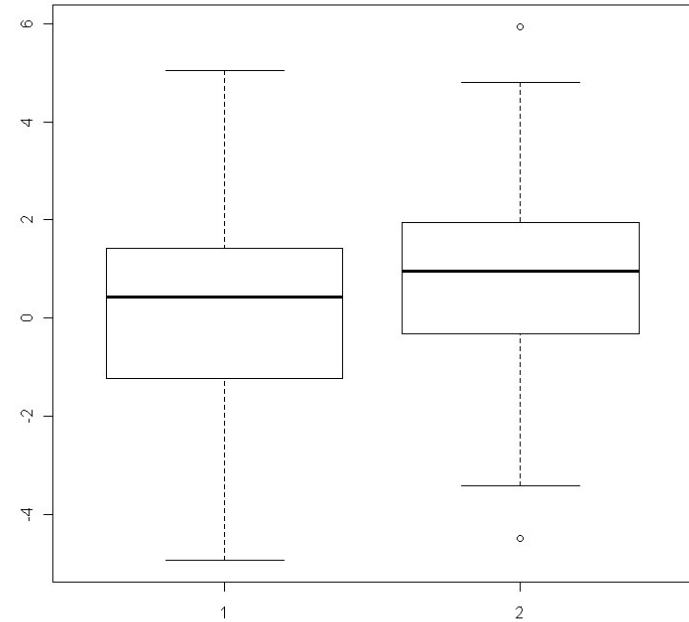
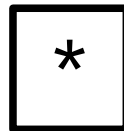


Two Sample t-test : p-value = **0.08284**

P-value > 5%

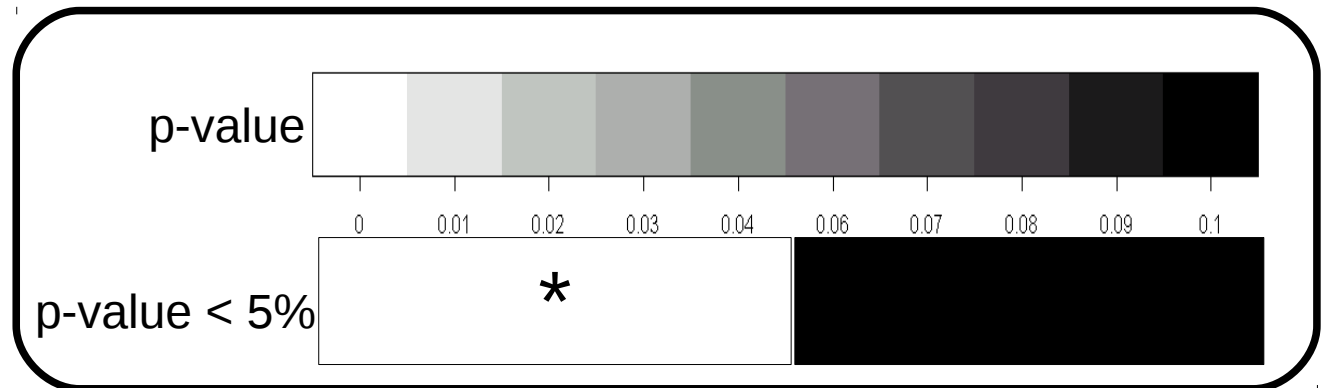
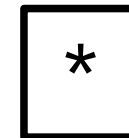


P-value < 10%



Two Sample t-test : p-value = **0.03556**

P-value < 5%



# Une petite liste de tests « courants »

Extrait de l'aide en ligne du logiciel R (sans package additionnel)

stats::ansari.test	Ansari-Bradley Test
stats::bartlett.test	Bartlett Test of Homogeneity of Variances
stats::binom.test	Exact Binomial Test
stats::Box.test	Box-Pierce and Ljung-Box Tests
stats::chisq.test	Pearson's Chi-squared Test for Count Data
stats::cor.test	Test for Association/Correlation Between Paired Samples
stats::fisher.test	Fisher's Exact Test for Count Data
stats::fligner.test	Fligner-Killeen Test of Homogeneity of Variances
stats::friedman.test	Friedman Rank Sum Test
stats::kruskal.test	Kruskal-Wallis Rank Sum Test
stats::ks.test	Kolmogorov-Smirnov Tests
stats::mantelhaen.test	Cochran-Mantel-Haenszel Chi-Squared Test for Count Data
stats::mauchly.test	Mauchly's Test of Sphericity
stats::mcnemar.test	McNemar's Chi-squared Test for Count Data
stats::mood.test	Mood Two-Sample Test of Scale
stats::oneway.test	Test for Equal Means in a One-Way Layout
stats::pairwise.prop.test	Pairwise comparisons for proportions
stats::pairwise.t.test	Pairwise t tests
stats::pairwise.wilcox.test	Pairwise Wilcoxon rank sum tests
stats::poisson.test	Exact Poisson tests
stats::power.anova.test	Power calculations for balanced one-way analysis of variance tests
stats::power.prop.test	Power calculations two sample test for proportions
stats::power.t.test	Power calculations for one and two sample t tests
stats::PP.test	Phillips-Perron Test for Unit Roots
stats::prop.test	Test of Equal or Given Proportions
stats::prop.trend.test	Test for trend in proportions
stats::quade.test	Quade Test
stats::shapiro.test	Shapiro-Wilk Normality Test
stats::t.test	Student's t-Test
stats::var.test	F Test to Compare Two Variances
stats::wilcox.test	Wilcoxon Rank Sum and Signed Rank Tests

# Problème de la multiplicité

La **roulette russe** est un jeu consistant à mettre une ou plusieurs cartouches (suivant la probabilité souhaitée) dans le barillet d'un revolver, à tourner ce dernier de manière aléatoire puis à pointer le revolver sur sa tempe avant d'actionner la détente. Si une cartouche se trouve alors dans la chambre placée dans l'axe du canon elle sera percutée et le « joueur » mourra ou sera blessé. (Wikipedia)

Supposons que le barillet contienne **100** emplacements de cartouche et que le « joueur » mette **5** cartouches aléatoirement. Sa probabilité d'être blessé (soyons optimiste) est de 5%.

Supposons maintenant que l'individu « joue » plusieurs fois de suite. Intuitivement, on sent bien qu'il va finir par « se faire mal » et que la probabilité de « perdre » au moins une fois augmente à chaque tentative.

Nous allons traduire mathématiquement cette petite histoire.

# Problème de la multiplicité

Calculons la probabilité qu'à le « joueur » de rester bien portant après n tentatives

- Pour  $n=1$ ,  $P[\text{bp}(1)]=1-0.05=\mathbf{0.95}$
- Pour  $n=2$ , il faut que le « joueur » soit bien portant après le 1er essai **ET** après le second

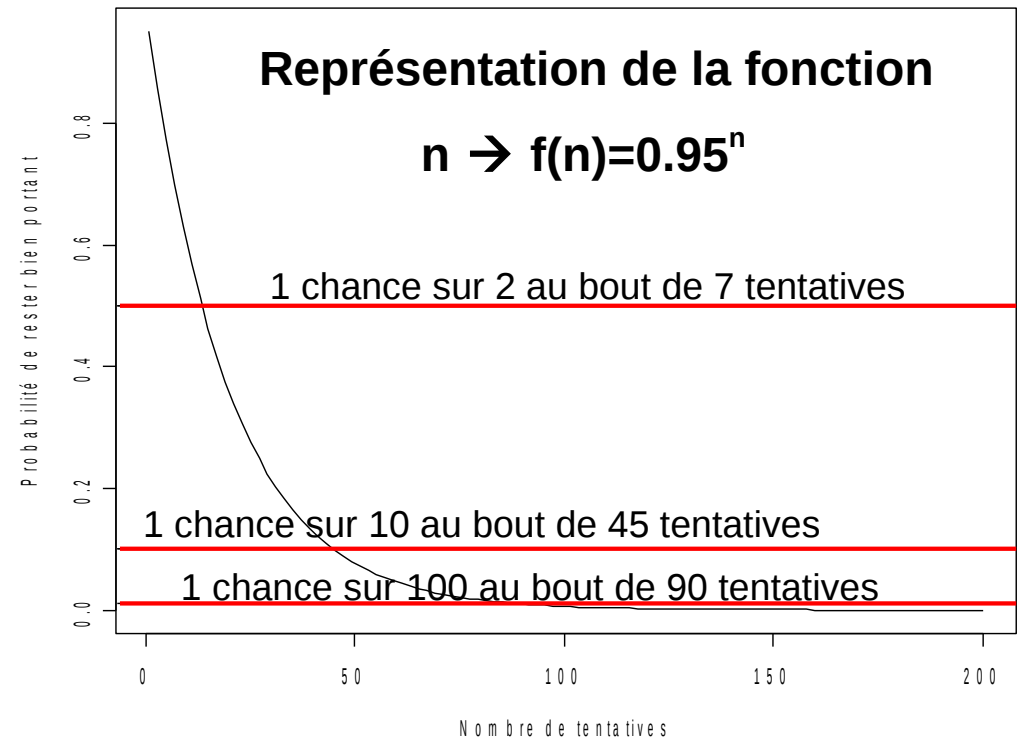
$$P[\text{bp}(2)]=(1-0.05)*(1-0.05)=0.95^2=\mathbf{0.9025}$$

- Et ainsi de suite...

$$\text{Règle générale : } P[\text{bp}(n)]=(1-0.05)^n$$

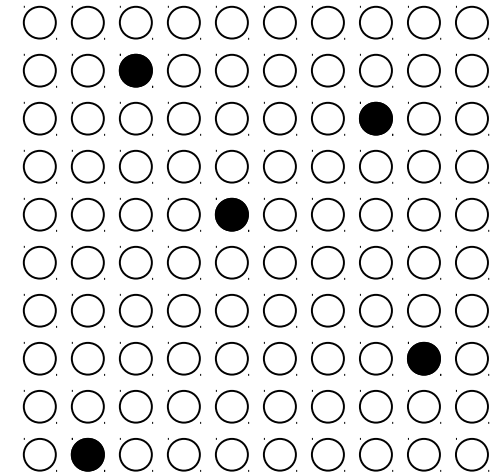
Quelques valeurs

n	1	2	5	10	100
<b>P[bp(n)]</b>	0.95	0.902	0.774	0.599	0.006



# Problème de la multiplicité

Une autre façon de considérer le problème consiste à envisager plusieurs « joueurs » jouant simultanément. Si le risque est toujours fixé à 5% (avec 5 cartouches sur 100 emplacements), et si 100 joueurs actionne la détente, on peut s'attendre à ce que 5 joueurs perdent. Ils seront 50 pour 1000 joueurs...



Cette façon de voir le problème de multiplicité est plus proche de la problématique de la détection de gènes différentiellement exprimés. En effet, un test statistique effectué pour l'ensemble des tests d'une biopuce (plusieurs milliers) au risque de 5% implique nécessairement la détection de gènes non pertinents.

Illustration : dans le cas de l'analyse des biopuces, la conséquence de l'erreur est moins radicale que pour la roulette russe.

Si les tests sont effectués pour détecter les gènes différentiellement exprimés (DE) entre 2 conditions alors le risque consiste à déclarer un gène DE alors qu'il ne l'est pas : on parle de gène faux-positif (FP). Au contraire, un gène « réellement » DE non détecté par le test est dit faux-négatif (FN).

Les conséquences de ces 2 types d'erreur sont de natures diverses. Dans le premier cas, on va engager des moyens (humains et financiers) pour valider l'hypothèse que le gène FP est effectivement intéressant alors qu'il ne l'est pas. Dans le second cas (FN), on passe peut-être à côté de la découverte du siècle en négligeant un gène important !!!

# Exemple de tests multiples

- 1 comparaison :  $\alpha = P[\text{Rejeter } H_0 // H_0 \text{ vraie}] = 5\%$

→ Probabilité de ne pas commettre d'erreur :  $1 - 0.05 = 0.95$

- 3 comparaisons :  $\left\{ \begin{array}{l} \alpha_1 = P[\text{Rejeter } H_{0_1} // H_{0_1} \text{ vraie}] = 5\% \\ \alpha_2 = P[\text{Rejeter } H_{0_2} // H_{0_2} \text{ vraie}] = 5\% \\ \alpha_3 = P[\text{Rejeter } H_{0_3} // H_{0_3} \text{ vraie}] = 5\% \end{array} \right.$

→ Probabilité de ne pas commettre d'erreur = produit des probabilités de ne pas commettre d'erreur à chacune des 3 comparaisons =  $(1 - 0.05) * (1 - 0.05) * (1 - 0.05) = 0.86$

Le risque (global) de commettre au moins une erreur est :

$$1 - 0.86 = 0.14$$

**Il faut donc diminuer le risque associé à chaque comparaison pour contrôler le risque global.**



# Exemple de tests multiples

## Méthode de Bonferroni

La méthode de Bonferroni consiste à corriger le risque associé à chaque comparaison ( $\alpha_i$ ) en le divisant par le nombre de comparaisons à effectuer pour contrôler le risque global ( $\alpha_g$ )

En prenant  $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_g/3 = 0.05/3$ , on a la probabilité de ne pas commettre d'erreur :

$$(1-0.05/3) * (1-0.05/3) * (1-0.05/3) = 0.9508$$

Ce qui correspond à un risque global inférieur à 5%.

Dans le cas des biopuces comportant plusieurs centaines voire plusieurs milliers de gènes, cette correction peut devenir trop conservatrice et conduire à la détection d'aucun gène. D'autres méthodes plus sophistiquées existent : Sidak, Holm, Westfall et Young, Hochberg...

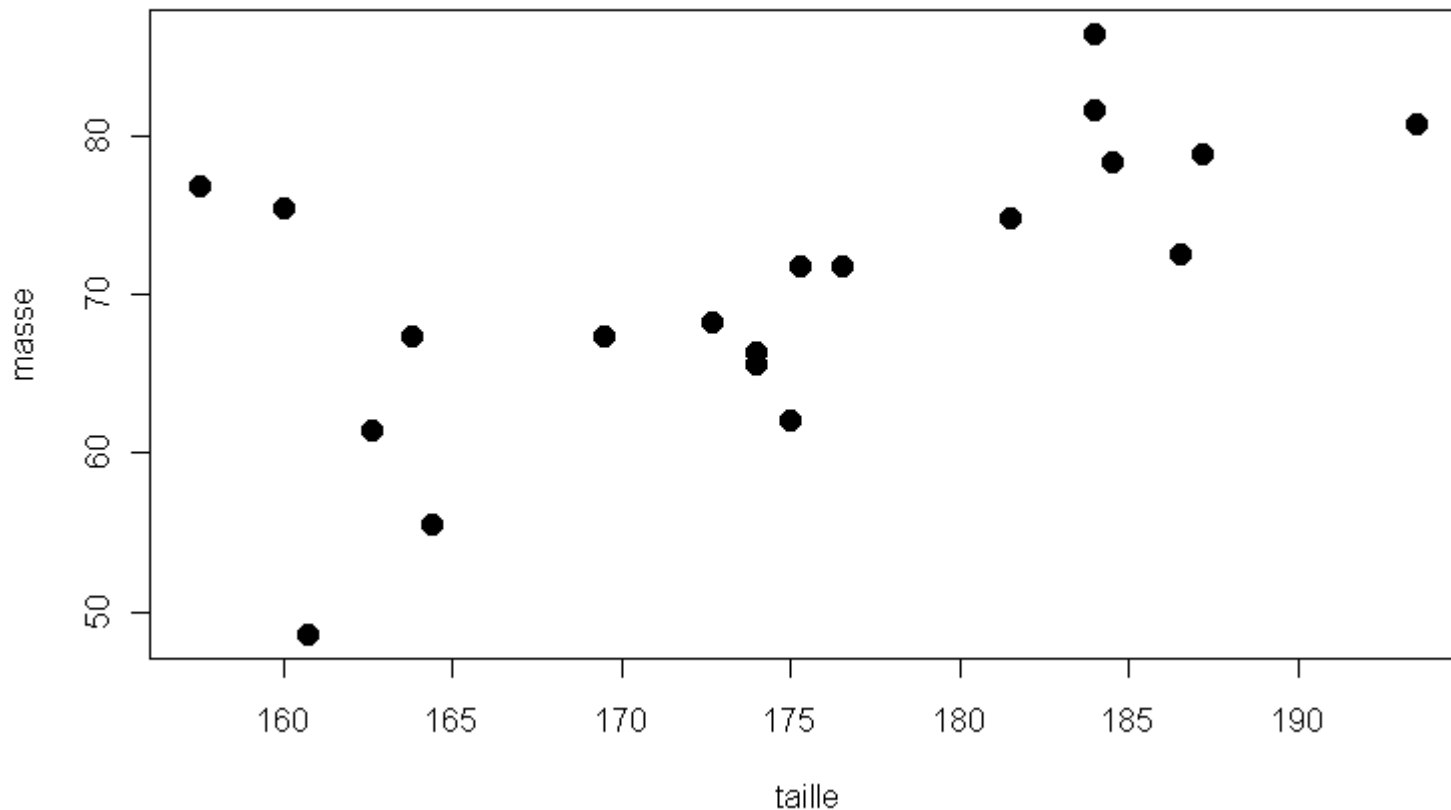
# Modélisation

1 ou plusieurs variables à **expliquer**  
par  
1 ou plusieurs variables  
**explicatives**

- Régression linéaire simple :
  - **1** variable à expliquer **quantitative** et **1** variable explicative **quantitative**
- ANOVA 1 facteur
  - **1** variable à expliquer **quantitative**, **1** variable explicative **qualitative** (facteur)

# Régression linéaire simple

**Taille (cm) :** 174.0 175.3 193.5 186.5 187.2 181.5 184.0 184.5 175.0 184.0 169.5 160.0 172.7 162.6 157.5 176.5 164.4 160.7 174.0 163.8  
**Masse (kg) :** 65.6 71.8 80.7 72.6 78.8 74.8 86.4 78.4 62.0 81.6 67.3 75.5 68.2 61.4 76.8 71.8 55.5 48.6 66.4 67.3



Peut-on « modéliser » « correctement » par une droite la masse des individus en fonction de la taille ?

# Régression linéaire simple

Équation d'une droite :  $Y = aX + b + \varepsilon$

Comment déterminer a et b ?

Par exemple, critère des moindres carrés : trouver a et b qui minimisent

$$\sum_i (y_i - ax_i - b)^2 = \sum_i \varepsilon_i^2$$

Sur l'exemple :  $a = 0.5445$  ;  $b = -24.37$

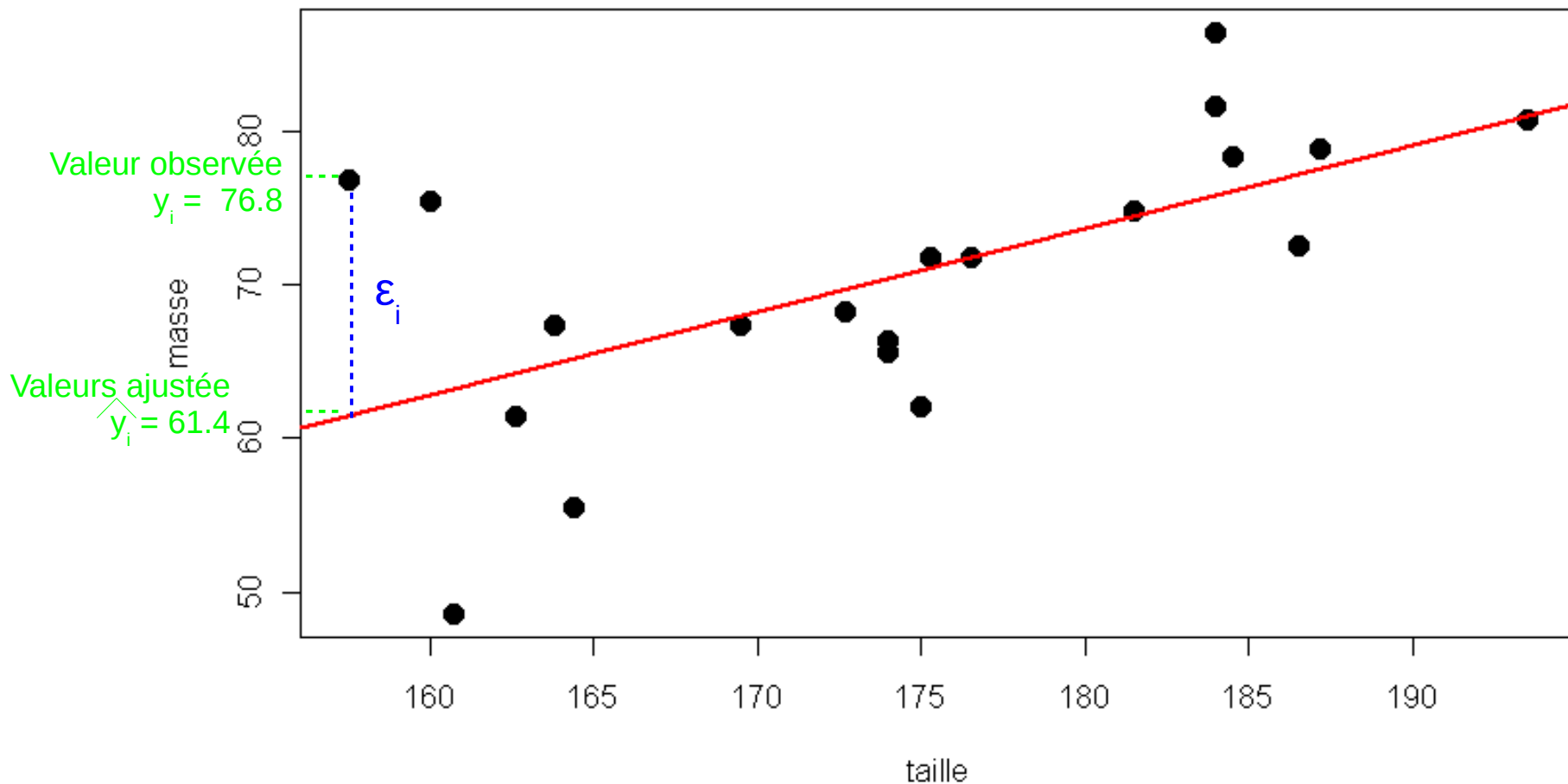
Un individu qui mesure 180cm pèse, selon le modèle,  
 $180 * 0.5445 - 24.37 = 73.6\text{kg}$

# Régression linéaire simple

Équation de la droite :

$$Y = 0.5445 X - 24.37$$

$$R^2 = 0,3782 = \frac{\text{var}(y_i)_{\text{observés}}}{\text{var}(\hat{y}_i)_{\text{ajustés}}}$$



# ANOVA 1 facteur

## Notations

Modalité	Répétitions				Moyenne
	1	2	...	J	
1	$Y_{11}$	$Y_{12}$	...	$Y_{1J}$	$\bar{Y}_{1\cdot}$
2	$Y_{21}$	$Y_{22}$	...	$Y_{2J}$	$\bar{Y}_{2\cdot}$
...	...	...		...	
I	$Y_{I1}$	$Y_{I2}$	...	$Y_{IJ}$	$\bar{Y}_{I\cdot}$
					$\bar{Y}_{..}$

- $Y_{ij}$  :  $j^{\text{ème}}$  observation pour la  $i^{\text{ème}}$  modalité du facteur
- $\bar{Y}_{i\cdot}$  : moyenne des observations de la  $i^{\text{ème}}$  modalité
- $\bar{Y}_{..}$  : moyenne générale
- I modalités, J observations par modalité (plan équilibré)
- $N=IJ$  observations

# ANOVA 1 facteur

## Modèle et sommes de carrés

Modèle d'ANOVA : 
$$Y_{ij} = \mu_i + \varepsilon_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

$\mu$  : effet moyen général

$\alpha_i$  : effet différentiel de la  $i^{\text{ème}}$  modalité du facteur

Test de l'effet du facteur : 
$$H_0 : \forall i=1, \dots, I \quad \alpha_i = 0$$

$$\sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{..})^2 = J \sum_{i=1}^I (\bar{y}_{i.} - \bar{y}_{..})^2 + \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.})^2$$

**SC**<sub>totale</sub>
**SC**<sub>facteur</sub>
**SC**<sub>résidu</sub>

# ANOVA 1 facteur

## Table d'ANOVA

Source de variation	Degrés de liberté	Somme de carrés	Carrés moyens	Statistique de test F	P value
Variation « facteur »	$I-1$	$SC_{\text{facteur}}$	$CM_f = SC_f / (I-1)$	$F = CM_f / CM_r$	...
Variation « résidu »	$N-I$	$SC_{\text{résidu}}$	$CM_r = SC_r / (N-I)$		
Variation totale	$N-1$	$SC_{\text{totale}}$			

Règle de décision :

$$\text{rejet de } H_0 \iff F > F_{I-1, N-I, \alpha}$$

Plus la variabilité due au facteur est importante par rapport à celle due au résidu, plus le test est susceptible de conclure à un effet du facteur.



# Degré de liberté

Dans la table d'ANOVA présentée précédemment, on retrouve une colonne « degré de liberté ». Cette information est essentielle dans tous les aspects de la statistique liés à la modélisation et aux tests. Elle traduit la quantité d'information dont on dispose pour estimer une quantité inconnue. Dans le calcul des degrés de liberté, on se retrouve systématiquement avec une valeur **-1**.

→ Imaginons **3** personnes (Pierre, Paul, Jacques) face à **3** vêtements **Rouge**, **Vert** ou **Bleu**. Chacune des personnes doit choisir à son tour un vêtement.

- Pierre commence et a le choix entre **3** vêtements **R**, **V**, **B**. Il choisit **R**.
- Paul vient ensuite. Il ne reste que **2** choix possibles : **V** et **B**. Il choisit **B**.
- Jacques arrive et **n'a plus le choix**, il prend **V** !

Ainsi, dans une situation où **3** unités sont initialement présentes, on constate que seulement **2** choix sont possibles ; le troisième étant nécessairement imposé.

→ Autre exemple : au moment de remplir le formulaire d'évaluation de la formation, vous aurez probablement la possibilité de rester anonyme. Si tel est votre souhait, assurez-vous qu'une autre personne le souhaite également. Sinon avec l'identité des **n-1** autres personnes, votre anonymat ne sera pas trop difficile à lever.

# Mise en œuvre de l'ANOVA

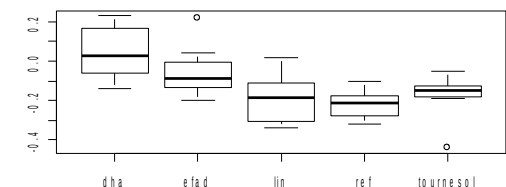
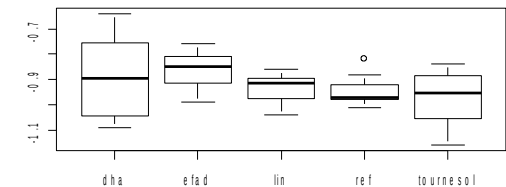
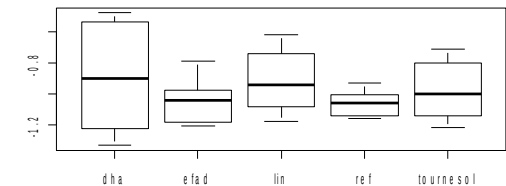
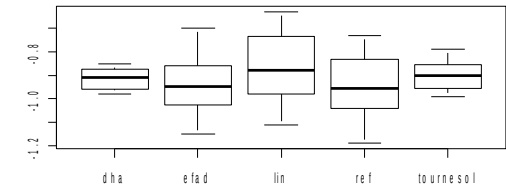
Exemple : mesure de 4 variables quantitatives selon les 5 modalités d'un facteur

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
facteur	4	0.03353	0.00838	0.5652	0.6895
Residuals	35	0.51918	0.01483		

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
facteur	4	0.18796	0.04699	0.9322	0.4566
Residuals	35	1.76434	0.05041		

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
facteur	4	0.06372	0.01593	1.4155	0.2492
Residuals	35	0.39391	0.01125		

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
facteur	4	0.39496	0.09874	7.0847	0.0002741 ***
Residuals	35	0.48780	0.01394		



*L'effet du facteur est significatif seulement dans le dernier cas.  
Au moins 1 des 5 moyennes est différente des autres.*

# ANOVA 2 facteurs

## Notations

Facteur A (I modalités)	Facteur B (J modalités)				Moyenne
	1	2	...	J	
1	$Y_{111}$ $\bar{Y}_{.11.}$ $Y_{11K}$	$Y_{121}$ $\bar{Y}_{.12.}$ $Y_{12K}$	...	$Y_{1J1}$ $\bar{Y}_{.1J.}$ $Y_{1JK}$	$\bar{Y}_{1..}$
2	$Y_{211}$ $\bar{Y}_{.21.}$ $Y_{21K}$	$Y_{221}$ $\bar{Y}_{.22.}$ $Y_{22K}$	...	$Y_{2J1}$ $\bar{Y}_{.2J.}$ $Y_{2JK}$	$\bar{Y}_{2..}$
...	...	...		...	
I	$Y_{I11}$ $\bar{Y}_{.I1.}$ $Y_{I1K}$	$Y_{I21}$ $\bar{Y}_{.I2.}$ $Y_{I2K}$	...	$Y_{IJ1}$ $\bar{Y}_{.IJ.}$ $Y_{IJK}$	$\bar{Y}_{I..}$
Moyenne	$\bar{Y}_{.1.}$	$\bar{Y}_{.2.}$	...	$\bar{Y}_{.J.}$	$\bar{Y}_{...}$

# ANOVA 2 facteurs

## Modèle et sommes de carrés

Modèle  
e

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk}$$

$\mu$  : effet moyen général

$\alpha_i$  : effet différentiel de la  $i^{\text{ème}}$  modalité du facteur A

$\beta_j$  : effet différentiel de la  $j^{\text{ème}}$  modalité du facteur B

$\gamma_{ij}$  : effet d'interaction de la  $i^{\text{ème}}$  modalité du facteur A avec la  $j^{\text{ème}}$  du facteur B

Décompositio  
n

SC<sub>totale</sub>

$$\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (y_{ijk} - \bar{y}_{...})^2 =$$

$$JK \sum_{i=1}^I (\bar{y}_{i..} - \bar{y}_{...})^2 + IK \sum_{i=1}^I (\bar{y}_{.j.} - \bar{y}_{...})^2 + K \sum_{i=1}^I \sum_{j=1}^J (\bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...})^2 + \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (y_{ijk} - \bar{y}_{ij.})^2$$

SC<sub>facteur A</sub>

SC<sub>facteur B</sub>

SC<sub>interaction</sub>

SC<sub>résidu</sub>

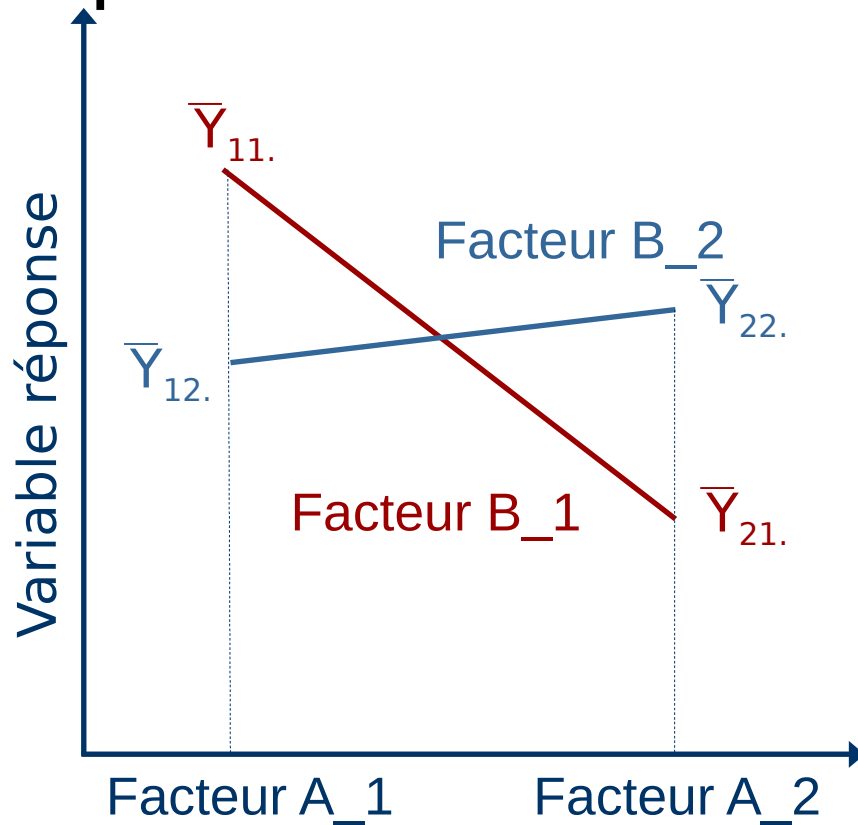
# ANOVA 2 facteurs

## Table d'ANOVA

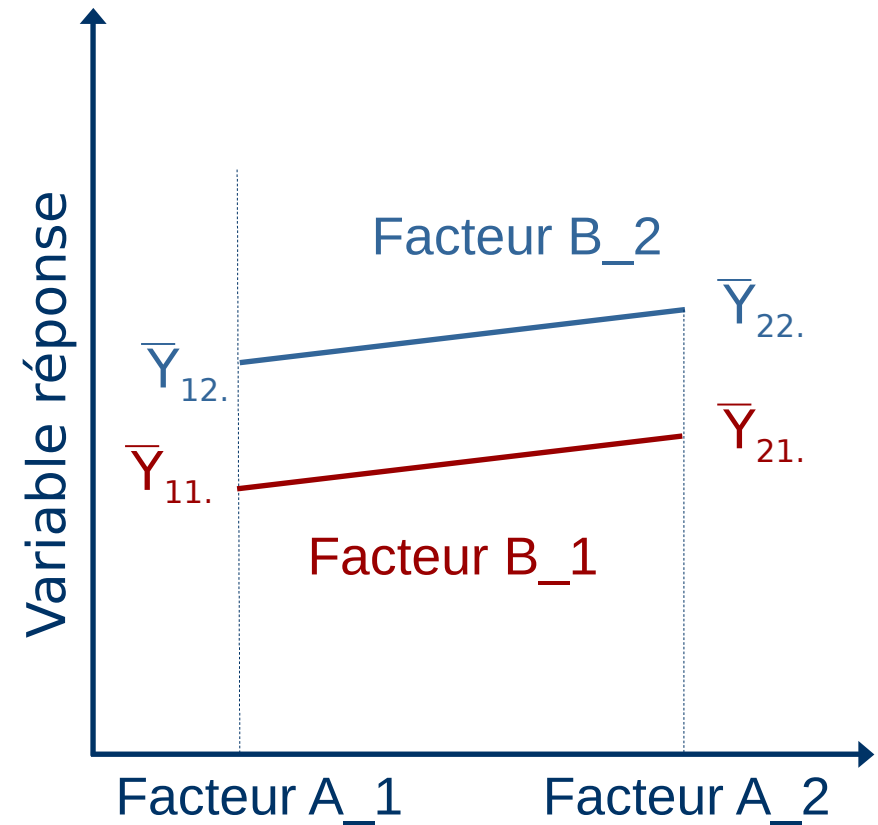
Sources de variation	Degrés de liberté	SCE (Somme de carrés d'écart)	CM (carré moyen)	F	p-value
Facteur A	<b>I-1</b>	<b>SCE<sub>A</sub></b>	<b>CM<sub>A</sub> = SCE<sub>A</sub> / (I-1)</b>	<b>CM<sub>A</sub> / CM<sub>résidu</sub></b>	...
Facteur B	<b>J-1</b>	<b>SCE<sub>B</sub></b>	<b>CM<sub>B</sub> = SCE<sub>B</sub> / (J-1)</b>	<b>CM<sub>B</sub> / CM<sub>résidu</sub></b>	...
Interaction	<b>(I-1)(J-1)</b>	<b>SCE<sub>A*B</sub></b>	<b>CM<sub>A*B</sub> = SCE<sub>A*B</sub> / (I-1)(J-1)</b>	<b>CM<sub>A*B</sub> / CM<sub>résidu</sub></b>	...
Résiduelle	<b>IJ(K-1)</b>	<b>SCE<sub>résidu</sub></b>	<b>CM<sub>resid</sub> = SCE<sub>résidu</sub> / IJ(K-1)</b>		

# ANOVA 2 facteurs

## Graphe d'interaction



L'effet d'un facteur dépend de la modalité de l'autre facteur → **l'interaction** sera probablement déclarée **significative** lors du test statistique



L'effet d'un facteur ne dépend pas de la modalité de l'autre facteur → **l'interaction** sera probablement déclarée **non significative** lors du test statistique

# Modèle linéaire

$$Y = X\beta + \varepsilon$$

Régression linéaire

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} X_1 & 1 \\ X_2 & 1 \\ \dots & \dots \\ X_n & 1 \end{pmatrix} \quad \beta = \begin{pmatrix} a \\ b \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_n \end{pmatrix}$$

ANOVA 1 facteur

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & 0 \\ \dots & \dots \\ 1 & 0 \\ 0 & 1 \\ \dots & \dots \\ 0 & 1 \end{pmatrix} \quad \beta = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_n \end{pmatrix}$$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

À suivre...