

## Vocabulaire des statistiques

<p>Une étude porte sur une « <b>population</b> »</p> <p><i>Par exemple :</i> Les hommes français de plus de 50 ans. Les sujets atteints du VIH.</p>	<p><b>Population</b> = ensemble des sujets étudiés, décrit avec précision.</p>
<p>Il est souvent impossible d'étudier la grandeur d'intérêt sur toute la « population », on l'observe alors sur une partie seulement de la population, appelée « échantillon »</p> <p><i>Par exemple :</i> les sujets atteints du VIH dépistés en un an dans un centre de santé donné.</p>	<p><b>Echantillon</b>=groupe issu de la population (sous-ensemble de la population).</p> <p><b>Echantillon représentatif</b> : échantillon constitué par tirage au sort dans la population d'étude.</p>

## Expériences aléatoires

<p>Résultats non mesurables sur une échelle numérique (modalités) : l'expérience aléatoire est <b>qualitative</b></p>	<p>Résultats mesurables sur une échelle numérique : l'expérience aléatoire est <b>quantitative</b>, appelée <b>variable aléatoire</b></p>		
<p><i>Pour savoir si une expérience est qualitative ou quantitative, se poser une question simple : « chaque résultat possible est-il un nombre ? »</i></p>			
<p>Deux types d'expériences <b>qualitatives</b></p>	<p>Deux types d'expériences <b>quantitatives</b></p>		
<p><b>qualitatives ordonnées</b></p>	<p><b>qualitatives catégoriques</b></p>	<p><b>quantitatives discrètes</b></p>	<p><b>quantitatives continues</b></p>
<p>Relation d'ordre évidente entre les résultats possibles</p>	<p>Pas de relation d'ordre évidente entre les résultats possibles</p>	<p>L'ensemble des résultats possibles est fini ou dénombrable</p>	<p>L'ensemble des résultats possibles forme un intervalle</p>
<p><i>Pour savoir si une expérience qualitative est ordinaire ou catégorique, se poser une question simple : « puis-je ordonner les résultats ? »</i></p>		<p><i>Pour savoir si une variable quantitative est discrète ou continue, se poser une question simple : « puis-je numéroter les résultats ? »</i></p>	
<p><b>Remarques :</b></p> <p>1. La distinction entre les variables discrètes ou continues dépend de l'échelle de mesure. En pratique, dès que le nombre de valeurs possibles est élevé, on peut traiter la variable comme continue.</p> <p>2. Une expérience aléatoire qualitative peut être transformée en variable quantitative discrète en associant des nombres aux résultats.</p>			

## Probabilités

Soit  $E$  l'ensemble des résultats possibles (ou issues) d'une expérience aléatoire,  $E$  est appelé **ensemble fondamental** ou **univers des possibles**.

### Événements

Un événement  $A$  est une partie de  $E$ .

Les événements peuvent se combiner entre eux pour former de nouveaux événements :

intersection de  $A$  et de  $B$  :

$$\ll A \text{ et } B \gg = A \cap B$$

réunion de  $A$  et de  $B$  :

$$\ll A \text{ ou } B \gg = A \cup B$$

contraire de  $A$  :

$$\ll \text{non } A \gg = \bar{A}$$

Remarquons que  $(A \cap B) \cup (A \cap \bar{B}) = A$

Le contraire d'une réunion est l'intersection des contraires :

Non ( $A$  ou  $B$ ) = ni  $A$ , ni  $B$

$$\overline{A \cup B} = \bar{A} \cap \bar{B}$$

Le contraire d'une intersection est la réunion des contraires :

Non ( $A$  et  $B$ ) = Non  $A$  ou Non  $B$

$$\overline{A \cap B} = \bar{A} \cup \bar{B}$$

### Probabilité d'un événement

$$P(E) = 1 \quad P(\emptyset) = 0 \quad 0 \leq P(A) \leq 1$$

$$P(\bar{A}) = 1 - P(A)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

### Événements incompatibles

$$\begin{aligned} A \text{ et } B \text{ sont incompatibles} &\Leftrightarrow A \cap B = \emptyset \\ &\Leftrightarrow P(A \cap B) = 0 \\ &\Leftrightarrow P(A \cup B) = P(A) + P(B) \end{aligned}$$

### Formule des probabilités totales

$$P(A) = P(A \cap B) + P(A \cap \bar{B})$$

### En situation d'équiprobabilité

$$P(A) = \frac{\text{card}A}{\text{card}E}$$

### Probabilités conditionnelles

$$P(A/B) = P_B(A) = \frac{P(A \cap B)}{P(B)}$$

le contraire de  $A/B$  est  $\bar{A}/B$

### Théorème de la multiplication

$$P(A \cap B) = P(A/B) \times P(B) = P(B/A) \times P(A)$$

### Formule de Bayes

$$P(A/B) = \frac{P(B/A) \times P(A)}{P(B)}$$

### Formule des probabilités totales (2<sup>ème</sup> forme)

$$P(A) = P(A/B) \times P(B) + P(A/\bar{B}) \times P(\bar{B})$$

### Événements indépendants

$$\begin{aligned} A \text{ et } B \text{ sont indépendants} &\Leftrightarrow P(A/B) = P(A) \\ &\Leftrightarrow P(B/A) = P(B) \\ &\Leftrightarrow P(A/\bar{B}) = P(A) \\ &\Leftrightarrow P(B/\bar{A}) = P(B) \\ &\Leftrightarrow P(A/B) = P(A/\bar{B}) \\ &\Leftrightarrow P(B/A) = P(B/\bar{A}) \end{aligned}$$

$$A \text{ et } B \text{ sont indépendants} \Leftrightarrow P(A \cap B) = P(A) \times P(B)$$

**Attention ! ne pas confondre**  
« indépendants » et « incompatibles »

## Signe et diagnostic

### Comment faire un diagnostic ?

A partir d'un signe, par exemple une douleur ou le résultat d'un examen.

$S = \text{signe présent}$     $\bar{S} = \text{signe absent}$

$T^+ = \text{test positif}$     $T^- = \text{test négatif}$

Si le test consiste à mesurer une grandeur continue, une valeur « seuil » est fixée :

Pour une maladie qui augmente la valeur de X  
Si  $X \geq \text{seuil}$  alors  $T^+$    Si  $X < \text{seuil}$  alors  $T^-$

Pour une maladie qui diminue la valeur de X  
Si  $X \leq \text{seuil}$  alors  $T^+$    Si  $X > \text{seuil}$  alors  $T^-$

### Les paramètres d'un signe, d'un test diagnostique

$$\text{Sensibilité} = Se = P(S / M) = P(T^+ / M)$$

$$\text{Spécificité} = Sp = P(\bar{S} / \bar{M}) = P(T^- / \bar{M})$$

- Quand on fait varier le seuil, Se et Sp varient en sens inverse.
- Se et Sp ne dépendent pas de la prévalence  $p$  de la maladie.

### Un sujet peut être :

VP =  $T^+$  et malade   FP =  $T^+$  et sain   VN =  $T^-$  et sain   FN =  $T^-$  et malade

### Taux de dépistage

« Etre dépisté » signifie « présenter le signe » ou « avoir un test positif »

$$P(T^+) = P(VP) + P(FP)$$

$$P(T^+) = P(T^+ \cap M) + P(T^+ \cap \bar{M})$$

$$P(T^+) = P(T^+ / M) \times P(M) + P(T^+ / \bar{M}) \times P(\bar{M})$$

$$P(T^+) = p \times Se + (1 - p) \times (1 - Sp)$$

### Probabilité d'erreur

$$P(\text{erreur}) = P(FP) + P(FN)$$

$$P(\text{erreur}) = P(T^+ \cap \bar{M}) + P(T^- \cap M)$$

$$P(\text{erreur}) = P(T^+ / \bar{M}) \times P(\bar{M}) + P(T^- / M) \times P(M)$$

$$P(\text{erreur}) = (1 - p) \times (1 - Sp) + p(1 - Se)$$

### Valeurs prédictives

$$VPP = P(M / T^+) = \frac{P(T^+ / M) \times P(M)}{P(T^+)} = \frac{pSe}{pSe + (1 - p)(1 - Sp)}$$

$$VPN = P(\bar{M} / T^-) = \frac{P(T^- / \bar{M}) \times P(\bar{M})}{P(T^-)} = \frac{(1 - p)Sp}{(1 - p)Sp + p(1 - Se)}$$

- VPP et VPN dépendent de la prévalence  $p$ .
- Quand  $p \nearrow$ , la VPP  $\nearrow$  et la VPN  $\searrow$  et inversement.

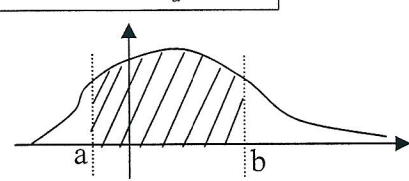


## Variables aléatoires

Une variable aléatoire est une fonction à valeurs dans l'ensemble des réels qui transforme tout événement élémentaire en un nombre (ou observation).

Par extension, un phénomène aléatoire quantitatif est appelé « variable aléatoire »

Soit  $X$  une variable aléatoire et  $E$  l'ensemble de ses valeurs.

<u>Variable aléatoire discrète</u>	<u>Variable aléatoire continue</u>										
<p><math>X</math> est <b>discrète</b> si <math>E</math> est fini ou dénombrable  <math>E = \{x_1, x_2, x_3, \dots, x_i, \dots\}</math></p>	<p><math>X</math> est <b>continue</b> si <math>E</math> n'est pas dénombrable (en général, c'est un intervalle)</p>										
<p><b>Loi de probabilité</b> de <math>X</math> (ou distribution de <math>X</math>) :            ensemble des probabilités des valeurs prises par <math>X</math></p> <table border="1" style="display: inline-table; margin-right: 20px;"> <tr> <td><math>x_1</math></td> <td><math>x_2</math></td> <td>...</td> <td><math>x_i</math></td> <td>...</td> </tr> <tr> <td><math>p_1</math></td> <td><math>p_2</math></td> <td></td> <td><math>p_i</math></td> <td></td> </tr> </table> <p><math>P(X = x_i) =</math></p>	$x_1$	$x_2$	...	$x_i$	...	$p_1$	$p_2$		$p_i$		<p>La <b>distribution</b> de <math>X</math> est déterminée à partir d'une fonction qui est la <b>densité de probabilité de <math>X</math></b>.</p> <p><i>Définition</i> : Une fonction <math>f</math> est une densité de probabilité si :</p> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 5px auto;"> <math display="block">\forall x \in \mathbb{R}, f(x) \geq 0 \text{ et } \int_{-\infty}^{+\infty} f(x) dx = 1</math> </div> <p><i>Autrement dit</i> : si la courbe représentative de <math>f</math> est située au-dessus de l'axe des abscisses et l'aire sous cette courbe est égale à 1. On a alors</p> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 5px auto;"> <math display="block">P(a &lt; X &lt; b) = \int_a^b f(t) dt \text{ et } P(X = a) = 0</math> </div> <div style="text-align: center; margin: 10px 0;">  <p>The diagram shows a bell-shaped curve representing a probability density function f(x) plotted against x. The x-axis has two points, a and b, marked with vertical dashed lines. The area under the curve between x=a and x=b is shaded with diagonal lines, representing the probability P(a &lt; X &lt; b). The curve is above the x-axis, and the total area under the curve is 1.</p> </div> <p><i>Remarque</i> : changer <math>&lt;</math> en <math>\leq</math> ou <math>&gt;</math> en <math>\geq</math> ne change pas la probabilité.</p>
$x_1$	$x_2$	...	$x_i$	...							
$p_1$	$p_2$		$p_i$								
<b>Fonction de répartition</b>											
<p><math>F(x) = P(X \leq x) = \sum_{x_i \leq x} P(X = x_i)</math>  <math>F</math> est une fonction en escalier, croissante</p>	<p><math>F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt</math>  <math>F</math> est continue, croissante  <math>F</math> est une primitive de la densité.</p>										
<b>Espérance (ou moyenne théorique ou moyenne vraie)</b>											
<p><math>\mu_X = E(X) = \sum_i x_i P(X = x_i) = \sum_i p_i x_i</math></p>	<p><math>\mu_X = E(X) = \int_{-\infty}^{+\infty} x f(x) dx</math></p>										
<b>Variance (vraie ou théorique) et écart-type</b>											
$Var(X) = E(X - E(X))^2 = E(X^2) - (E(X))^2$ $\sigma_X = \sqrt{Var(X)}$											
<p><math>Var(X) = \sum_i p_i (x_i - E(X))^2</math>  <math>Var(X) = \sum_i p_i x_i^2 - (E(X))^2</math></p>	<p><math>Var(X) = \int_{-\infty}^{+\infty} (x - E(X))^2 f(x) dx</math>  <math>Var(X) = \int_{-\infty}^{+\infty} x^2 f(x) dx - (E(X))^2</math></p>										



<u>Règles de calcul avec les espérances et les variances</u>	
X et Y sont des variables aléatoires discrètes ou continues, a est une constante réelle	
L'espérance est linéaire	$\text{Var}(aX) = a^2 \text{Var}(X)$ $\text{Var}(a) = 0$ $\text{Var}(X+a) = \text{Var}(X)$
$E(X+Y) = E(X) + E(Y)$ $E(aX) = aE(X)$ $E(a) = a$ $E(X+a) = E(X) + a$	
<p><b>Si X et Y sont INDEPENDANTES</b> ce qui signifie :</p> $P((X = x_i) \cap (Y = y_j)) = P(X = x_i) \times P(Y = y_j) \text{ pour tous les couples } (i, j)$ <p><b>Alors :</b></p> $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ <p><i>La réciproque est fausse.</i></p>	

**Lois discrètes usuelles à connaître.**

<p><b>Loi de Bernoulli de paramètre <math>\pi</math></b></p>	<p>Expérience aléatoire à deux résultats possibles (épreuve de Bernoulli)  <math>X=1</math> si succès      <math>X=0</math> si échec</p> <p><math>P(X=1) = \pi</math>      <math>P(X=0) = 1-\pi</math></p>	<p><math>E(X) = \pi</math></p> <p><math>Var(X) = \pi(1-\pi)</math></p>
<p><b>Loi binomiale de paramètres n et <math>\pi</math></b></p>	<p>La répétition n fois de manière indépendante de la même épreuve de Bernoulli définit n variables de Bernoulli <math>X_1, X_2, \dots, X_n</math> de même paramètre <math>\pi</math></p> <p>La variable aléatoire <math>S_n</math> égale au nombre de succès sur les n résultats prend ses valeurs dans <math>\{0, 1, 2, \dots, n\}</math> et est égale à :</p> <p><math>S_n = X_1 + X_2 + \dots + X_n</math></p> <p><math>P(S_n = k) = C_n^k \pi^k (1-\pi)^{n-k}</math>      <math>C_n^k = \binom{n}{k} = \frac{n!}{k!(n-k)!}</math></p>	<p><math>E(S_n) = n\pi</math></p> <p><math>Var(S_n) = n\pi(1-\pi)</math></p>
<p><b>Loi de Poisson de paramètre <math>\lambda</math></b> (<math>\lambda</math> est un réel strictement positif)</p>	<p><math>X =</math> nombre d'événements observés pendant une période de temps donnée dans le cas où ces événements sont indépendants et faiblement probables (nombre de colonies bactériennes dans une boîte de Petri, nombre d'accidents ...)</p> <p>X prend ses valeurs dans <math>\mathbb{N}</math></p> <p><math>P(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}</math></p> <p><i>Remarque : Si <math>X_1</math> et <math>X_2</math> sont indépendantes et distribuées selon des lois de Poisson de paramètres <math>\lambda_1</math> et <math>\lambda_2</math> alors <math>X_1 + X_2</math> suit une loi de Poisson de paramètre <math>\lambda_1 + \lambda_2</math></i></p>	<p><math>E(X) = Var(X) = \lambda</math></p>

### Approximation de la loi binomiale par la loi de Poisson

Si  $X$  suit une loi **binomiale** de paramètres  $n$  et  $\pi$   
avec  $n$  grand ( $n > 50$ ) et  $\pi$  petit ( $\pi < 0,1$ )

alors  $X$  suit **approximativement** une loi de **Poisson** de paramètre  $n\pi$

### Lois continues usuelles à connaître

<p><b>Loi uniforme sur <math>[a ; b]</math></b> <i>modélise les tirages de nombres au hasard dans un intervalle.</i></p>	<p><math>E = [a ; b]</math> Densité de probabilité  <math display="block">\begin{cases} f(x) = \frac{1}{b-a} \text{ si } x \in [a; b] \\ f(x) = 0 \text{ si } x \notin [a; b] \end{cases}</math></p>	<p><math>E(X) = \frac{a+b}{2}</math> <math>Var(X) = \frac{(b-a)^2}{12}</math></p>
<p><b>Loi exponentielle de paramètre <math>\lambda &gt; 0</math></b> <i>modélise les durées de vie lorsque le vieillissement n'intervient pas.</i></p>	<p><math>E = \mathbb{R}^+</math> Densité de probabilité  <math display="block">\begin{cases} f(x) = \lambda e^{-\lambda x} \text{ si } x \geq 0 \\ f(x) = 0 \text{ si } x &lt; 0 \end{cases}</math> Fonction de répartition : Pour <math>x \geq 0</math>, <math>F(x) = P(X \leq x) = 1 - e^{-\lambda x}</math> <b>Propriété importante :</b> Si <math>X</math> suit une loi exponentielle Alors <math>P(X &gt; a + b \mid X &gt; a) = P(X &gt; b)</math> La loi exponentielle est dite « sans mémoire ».</p>	<p><math>E(X) = \frac{1}{\lambda}</math> <math>Var(X) = \frac{1}{\lambda^2}</math> <math>med(X) = \frac{\ln 2}{\lambda}</math></p>
<p><b>Loi normale (ou de Gauss) de paramètres <math>\mu</math> et <math>\sigma &gt; 0</math></b> <b>notée <math>N(\mu ; \sigma^2)</math></b> <i>modélise de nombreuses mesures biologiques.</i></p>	<p><math>E = \mathbb{R}</math> Densité de probabilité  <math display="block">f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}</math> la courbe de la densité est une courbe en cloche symétrique par rapport à la droite d'équation <math>x = \mu</math> <i>Dans la pratique :</i> Si <math>X</math> suit <math>N(\mu ; \sigma^2)</math> Alors <math>Z = \frac{X - \mu}{\sigma}</math> suit <math>N(0 ; 1)</math> <i>Autrement dit :</i> <math>Z = \frac{X - \mu}{\sigma}</math> est normale <b>centrée réduite</b>  <i>puis on utilise la table de la loi normale centrée réduite pour le calcul des probabilités</i></p>	<p><math>E(X) = \mu</math> <math>Var(X) = \sigma^2</math></p>



### Intervalle de pari d'une variable aléatoire Normale X de moyenne $E(X)=\mu$ et d'écart-type $\sigma$

L'intervalle de pari (appelé aussi intervalle de fluctuation) de X (Normale) de niveau  $1 - \alpha$  (ou de risque  $\alpha$ ) est :

$$[\mu - \varepsilon_\alpha \sigma ; \mu + \varepsilon_\alpha \sigma] \quad \text{C'est l'intervalle } [a, b] \text{ centré sur } \mu \text{ tel que } P(a \leq X \leq b) = 1 - \alpha$$

Cas particulier classique : L'IP de X de niveau 95% (de risque 5 %) est  $[\mu - 2\sigma ; \mu + 2\sigma]$  (avec l'arrondi habituel de 1,96 à 2)

95% des réalisations de X sont dans cet IP

Il y a 95 chances sur 100 que X prenne sa valeur dans cet IP

Remarque : Il existe d'autres intervalles que l'IP de niveau  $1 - \alpha$  (non centrés sur  $\mu$ ) ayant pour probabilité  $1 - \alpha$

### Variable aléatoire Somme, variable aléatoire Moyenne, variable aléatoire Proportion

Soit  $n$  réalisations d'une variable aléatoire X, par exemple la mesure d'une grandeur biologique chez  $n$  sujets (variable aléatoire continue) ou l'étude d'une propriété en codant 1 la présence de cette propriété et 0 son absence (variable aléatoire de Bernoulli).

Ces  $n$  répétitions indépendantes de X sur  $n$  sujets sont  $n$  variables aléatoires  $X_1, X_2, \dots, X_n$  de même loi que X (iid : indépendantes identiquement distribuées), donc de moyenne (espérance)  $\mu$  et d'écart-type  $\sigma$ .

A partir de  $n$  réalisations de X ( $n$  valeurs prises par X) dans un échantillon de  $n$  sujets, on peut en calculer la somme qui

est la valeur prise par la variable aléatoire Somme  $S_n = \sum_{i=1}^n X_i$

ou la moyenne arithmétique (moyenne observée) qui est la valeur prise par la variable aléatoire Moyenne

$$M_n = \frac{\sum_{i=1}^n X_i}{n} \quad \text{S/C}$$

$$E(S_n) = nE(X) = n\mu$$

$$Var(S_n) = nVar(X)$$

écart type

$$E(M_n) = E(X) = \mu$$

$$Var(M_n) = \frac{Var(X)}{n}$$

$$\sigma(M_n) = \frac{\sigma(X)}{\sqrt{n}}$$

#### Remarque importante :

Dans le cas où X est une variable aléatoire de Bernoulli,

$$M_n = \frac{\sum_{i=1}^n X_i}{n} = \frac{S_n}{n} = \frac{\text{v.a. Nombre de succès dans un échantillon de taille } n}{n}$$

= v.a. Proportion de succès dans un échantillon de taille  $n$

$M_n$  est alors notée alors  $P_n$

La variable aléatoire Proportion  $P_n$  est la variable aléatoire Moyenne  $M_n$  de  $n$  répétitions d'une variable de Bernoulli.

Ne pas confondre  $P_n$  et  $S_n$  :  $P_n = \frac{S_n}{n}$   $S_n$  suit une loi binomiale, mais pas  $P_n$  !

$$E(P_n) = E(X) = \pi$$

$$Var(P_n) = \frac{Var(X)}{n} = \frac{\pi(1-\pi)}{n}$$

$$\sigma(P_n) = \sqrt{\frac{\pi(1-\pi)}{n}}$$



Loi de  $S_n$ , Loi de  $M_n$ , Loi de  $P_n$ 

**Théorème central limite** : Quelle que soit la loi de  $X$ , variable aléatoire de moyenne  $\mu$  et d'écart-type  $\sigma$ , les lois de  $S_n$  et de  $M_n$  (donc aussi de  $P_n$ ) convergent, pour  $n$  suffisamment grand, vers une loi normale.

Plus précisément :

X $\neq$ Bernoulli	Si $X \sim N(\mu ; \sigma^2)$ <i>pas besoin du TCL</i> $\longrightarrow$	$S_n \sim N(n\mu ; n\sigma^2)$ (quel que soit $n$ ) <b>Loi exacte</b>	$M_n \sim N(\mu ; \frac{\sigma^2}{n})$ (quel que soit $n$ ) <b>Loi exacte</b>
	Si $X \sim \text{Poisson}(\lambda)$ <i>pas besoin du TCL</i> $\longrightarrow$	$S_n \sim \text{Poisson}(n\lambda)$ (quel que soit $n$ ) <b>Loi exacte</b>	$M_n$ <b>Loi exacte inconnue</b>
	Si de plus $n \geq 30$ <b>TCL</b> $\longrightarrow$	$S_n \sim N(n\lambda ; n\lambda)$ <b>Loi approchée</b>	$M_n \sim N(\lambda ; \frac{\lambda}{n})$ <b>Loi approchée</b>
	Si $X$ ne suit pas une loi Normale (autre loi ou loi inconnue, moy $\mu$ et var $\sigma^2$ ) et $n \geq 30$ <b>TCL</b> $\longrightarrow$	$S_n \sim N(n\mu ; n\sigma^2)$ <b>Loi approchée</b>	$M_n \sim N(\mu ; \frac{\sigma^2}{n})$ <b>Loi approchée</b>
X de Bernoulli de paramètre $\pi$	<i>pas besoin du TCL</i> $\longrightarrow$	$S_n \sim$ loi binomiale de paramètres $(n, \pi)$ <b>Loi exacte</b>	$P_n$ <b>Loi exacte inconnue</b>
	Si $n\pi \geq 5$ et $n(1-\pi) \geq 5$ <b>TCL</b> $\longrightarrow$	$S_n \sim N(n\pi ; n\pi(1-\pi))$ <b>Loi approchée</b>	$P_n \sim N(\pi ; \frac{\pi(1-\pi)}{n})$ <b>Loi approchée</b>

Remarque : On dispose finalement de deux approximations possibles pour une loi binomiale :

## Approximations de la loi binomiale

1) par la loi de Poisson (voir séance 2)

Si  $X$  suit une loi binomiale de paramètres  $n$  et  $\pi$

avec  $n$  grand ( $n > 50$ ) et  $\pi$  petit ( $\pi < 0,1$ )

alors  $X$  suit approximativement une loi de Poisson de paramètre  $n\pi$

2) par la loi normale (voir ci-dessus)

Si  $X$  suit une loi binomiale de paramètres  $n$  et  $\pi$

avec  $n\pi \geq 5$  et  $n(1-\pi) \geq 5$

alors  $X$  suit approximativement une loi normale  $N(n\pi ; n\pi(1-\pi))$

## Approximation de la loi de Poisson par la loi Normale:

Si  $X$  suit une loi de Poisson avec paramètre grand  $\lambda > 30$

alors  $X$  suit approximativement une loi Normale d'espérance et variance égales à  $\lambda$ .

*Remarque* : Quand  $S_n$ ,  $M_n$  ou  $P_n$  suit une loi normale (exacte ou approchée par le TCL, voir conditions dans tableau), on peut écrire des intervalles de pari de ces variables aléatoires, avec leur propre moyenne et leur propre écart-type bien sûr !

$$IP_{niv1-\alpha}(S_n) = [\mu(S_n) - \varepsilon_\alpha \sigma(S_n); \mu(S_n) + \varepsilon_\alpha \sigma(S_n)] = [n\mu_X - \varepsilon_\alpha \sqrt{n}\sigma_X; n\mu_X + \varepsilon_\alpha \sqrt{n}\sigma_X]$$

$$IP_{niv1-\alpha}(M_n) = [\mu(M_n) - \varepsilon_\alpha \sigma(M_n); \mu(M_n) + \varepsilon_\alpha \sigma(M_n)] = [\mu_X - \varepsilon_\alpha \frac{\sigma_X}{\sqrt{n}}; \mu_X + \varepsilon_\alpha \frac{\sigma_X}{\sqrt{n}}]$$

$$IP_{niv1-\alpha}(P_n) = [\mu(P_n) - \varepsilon_\alpha \sigma(P_n); \mu(P_n) + \varepsilon_\alpha \sigma(P_n)] = [\pi - \varepsilon_\alpha \sqrt{\frac{\pi(1-\pi)}{n}}; \pi + \varepsilon_\alpha \sqrt{\frac{\pi(1-\pi)}{n}}]$$

## Estimation

Problème de l'estimation : on cherche à connaître un paramètre  $\theta$  d'une variable aléatoire (en général la moyenne vraie  $\mu$  ou la variance vraie  $\sigma^2$ ) grâce à des observations réalisées sur un échantillon de taille  $n$  extrait de la population.

### Estimation ponctuelle d'un paramètre $\theta$

Pour approcher un paramètre **inconnu**  $\theta$ , on utilise un **estimateur**, c'est-à-dire une variable aléatoire  $T_n$  qui prend, dans des échantillons de taille  $n$ , des valeurs voisines du paramètre  $\theta$  que l'on cherche à estimer.

Une valeur prise par  $T_n$  est une **estimation ponctuelle** de  $\theta$ .

Biais de l'estimateur  $T_n$  :  $B(T_n) = E(T_n) - \theta$

L'estimateur  $T_n$  est **non biaisé** si son biais est nul, autrement dit si  $E(T_n) = \theta$

Dans le cas où  $T_n$  est biaisé, si  $\lim_{n \rightarrow +\infty} B(T_n) = 0$  on dit que l'estimateur  $T_n$  est **asymptotiquement non biaisé**

Soit X une variable aléatoire de moyenne vraie $\mu$ et de variance vraie $\sigma^2$ inconnues	
<p style="text-align: center;">Estimateur usuel</p> <p style="color: red; font-style: italic;">Tous les estima. usuel sont non biaisés</p>	<p>Sur un échantillon de <math>n</math> sujets on peut calculer la valeur prise par l'estimateur</p>
<p>De <math>\mu</math> :</p> $M_n = \frac{\sum_{i=1}^n X_i}{n} \quad \text{non biaisé}$	$m = \frac{x_1 + x_2 + \dots + x_n}{n} \quad \text{moyenne observée, estimation ponctuelle de la moyenne vraie } \mu.$
<p>De <math>\pi</math>, proportion vraie (paramètre d'une loi de Bernoulli, autrement dit moyenne vraie d'une loi de Bernoulli) :</p> $P_n = \frac{\sum_{i=1}^n X_i}{n} \quad \text{non biaisé}$	$p = \frac{\text{nombre de succès dans l'échantillon}}{n}$ <p>proportion observée de succès, estimation ponctuelle de la proportion vraie <math>\pi</math>.</p>
<p>De <math>\sigma^2</math> :</p> $V_n \text{ ou } S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_n)^2$ $= \frac{1}{n-1} \left[ \sum_{i=1}^n X_i^2 - nM_n^2 \right] \quad \text{non biaisé}$	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - n \times m^2 \right]$ <p>variance observée, estimation ponctuelle de la variance vraie <math>\sigma^2</math>.</p>

Remarque : Des échantillons différents fournissent des estimations ponctuelles différentes.



### Estimation par intervalle d'une moyenne vraie ou d'une proportion vraie

Il s'agit de donner une fourchette de valeurs pour  $\mu$  ou  $\pi$  (paramètre inconnu) : un intervalle dans lequel se trouve le paramètre inconnu avec une probabilité fixée  $1 - \alpha$ , appelé **intervalle de confiance de niveau  $1 - \alpha$  ou de risque  $\alpha$** .

**IC d'une moyenne vraie  $\mu$  (d'une variable aléatoire  $X \neq$  Bernoulli) établi à partir d'un échantillon de  $n$  sujets**

On calcule  $m$  la moyenne observée et  $s^2$  la variance observée dans l'échantillon

**Condition de validité :**

**Soit  $X$  est Normale (alors pas de condition sur  $n$ )**

**Soit  $n \geq 30$**

$$IC_{niv 1-\alpha}(\mu) = \left[ m - \varepsilon_{\alpha} \frac{s}{\sqrt{n}} ; m + \varepsilon_{\alpha} \frac{s}{\sqrt{n}} \right]$$

précision de l'IC :  $\varepsilon_{\alpha} \frac{s}{\sqrt{n}}$     largeur de l'IC :  $2 \times \varepsilon_{\alpha} \frac{s}{\sqrt{n}}$

**IC d'une proportion vraie  $\pi$  établi à partir d'un échantillon de  $n$  sujets**

On calcule  $p$  la proportion observée

$$IC_{niv 1-\alpha}(\pi) = \left[ p - \varepsilon_{\alpha} \sqrt{\frac{p(1-p)}{n}} ; p + \varepsilon_{\alpha} \sqrt{\frac{p(1-p)}{n}} \right] = [\pi_1 ; \pi_2]$$

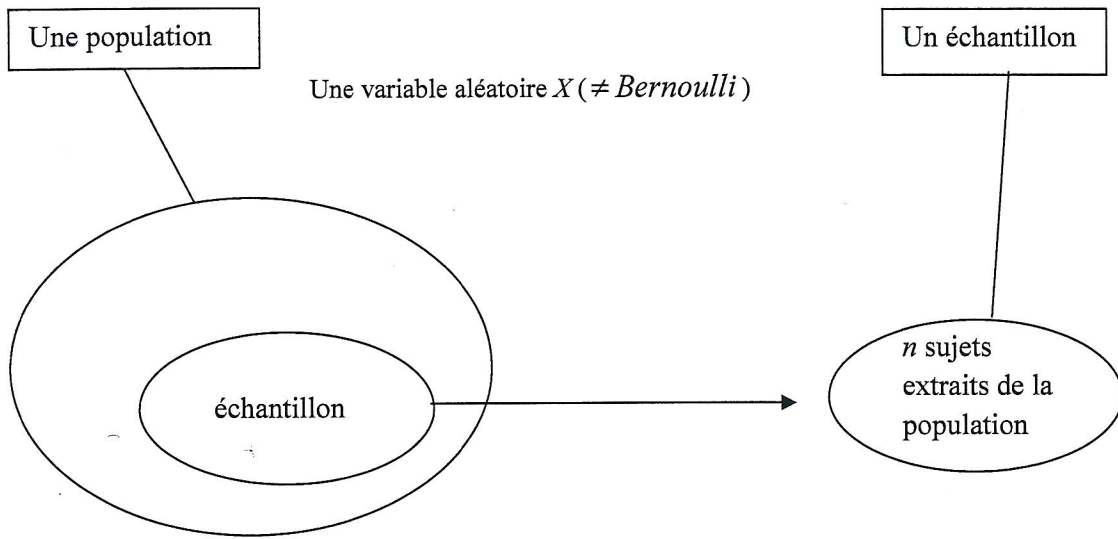
**Attention ! Conditions de validité à vérifier aux bornes de l'intervalle, après son calcul :**

$$n\pi_1 \geq 5 \quad n(1-\pi_1) \geq 5 \quad n\pi_2 \geq 5 \quad n(1-\pi_2) \geq 5$$

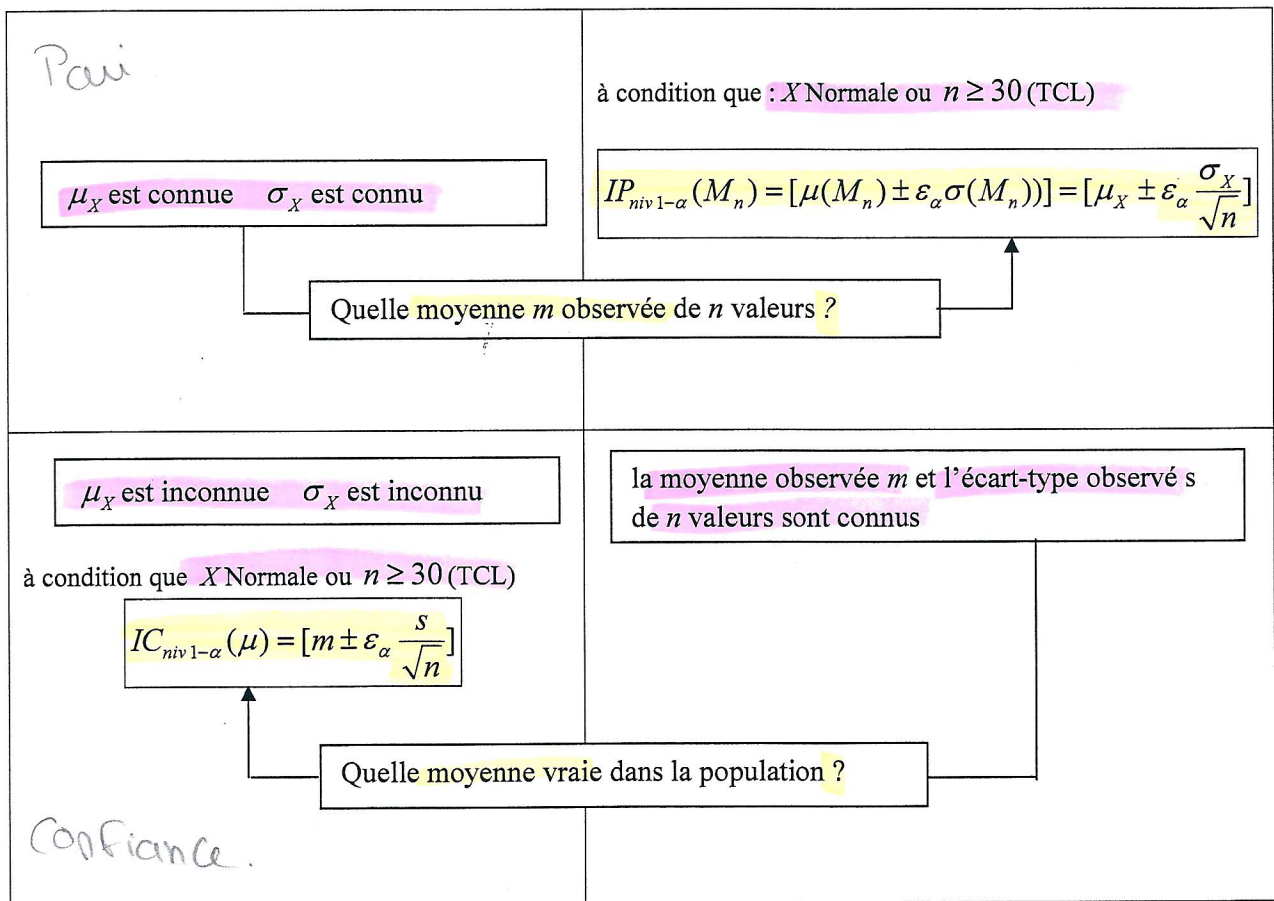
Précision de l'IC :  $\varepsilon_{\alpha} \sqrt{\frac{p(1-p)}{n}}$     largeur de l'IC :  $2 \times \varepsilon_{\alpha} \sqrt{\frac{p(1-p)}{n}}$

**Remarque :** Des échantillons différents fournissent des intervalles de confiance différents.

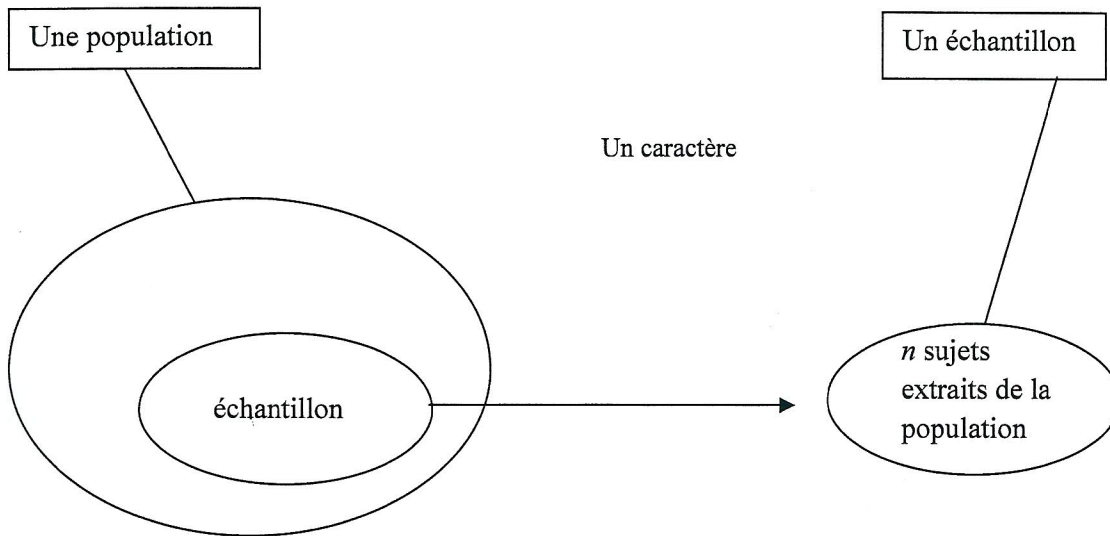
**MOYENNE : Ne pas confondre Intervalle de Pari et Intervalle de Confiance**



La moyenne $\mu_X$ de $X$ dans la population (dite vraie ou théorique) est une constante	La moyenne observée $m$ de $n$ valeurs de $X$ est variable d'un échantillon à l'autre ( $m$ est une réalisation de $M_n$ )
--	--



**PROPORTION : Ne pas confondre Intervalle de Pari et Intervalle de Confiance**



La proportion  $\pi$  de sujets de la population (dite vraie ou théorique) possédant un certain caractère est une constante

La proportion observée  $p$  de sujets possédant ce caractère parmi  $n$  sujets est variable d'un échantillon à l'autre ( $p$  est une réalisation de  $P_n$ )

<p><math>\pi</math> est connue</p>	<p>à condition que : <math>n\pi \geq 5</math> et <math>n(1-\pi) \geq 5</math></p>
<p>Quelle proportion observée <math>p</math> parmi <math>n</math> sujets ?</p>	$IP_{niv\ 1-\alpha}(P_n) = [\pi \pm \varepsilon_\alpha \sqrt{\frac{\pi(1-\pi)}{n}}]$
<p><math>\pi</math> est inconnue</p>	<p>La proportion <u>observée</u> <math>p</math> de sujets possédant ce caractère parmi <math>n</math> sujets est connue</p>
$IC_{niv\ 1-\alpha}(\pi) = [p \pm \varepsilon_\alpha \sqrt{\frac{p(1-p)}{n}}] = [\pi_1; \pi_2]$	<p>Quelle proportion vraie dans la population ?</p>
<p>valide à condition qu'après calcul les bornes vérifient :</p> $n\pi_1 \geq 5 \quad n(1-\pi_1) \geq 5 \quad n\pi_2 \geq 5 \quad n(1-\pi_2) \geq 5$	



## Test de comparaison de deux moyennes sur échantillons APPARIÉS

### Appariés = Non Indépendants

- Un seul groupe de sujets d'effectif  $n$
- Condition de validité :  $n \geq 30$
- Deux variables  $X_A$  et  $X_B$  prenant deux valeurs pour chaque sujet du groupe

On effectue pour la variable aléatoire différence  $Y = X_A - X_B$

le test de comparaison d'une moyenne observée à la valeur zéro (test de comparaison d'une moyenne à une norme, à une valeur donnée)

$\mu = \mu_A - \mu_B$

$$H_0 : \mu = 0 \quad H_1 : \mu \neq 0$$

On calcule  $z = \frac{m}{\sqrt{\frac{s^2}{n}}}$

$Y \sim N(\mu; \frac{\sigma^2}{n})$

Avec un risque de 1<sup>ère</sup> espèce fixé à  $\alpha$ :

Si  $|z| > \varepsilon_\alpha$  alors **RRH<sub>0</sub>** au risque  $\alpha$ , et on dit que la moyenne observée des différences est *significativement* différente de zéro et on conclut : les moyennes vraies  $\mu_A$  et  $\mu_B$  sont différentes.

### Degré de signification =

$$P(|Z| > |z_{calculé}|)$$

sinon **NRH<sub>0</sub>** et on dit que la moyenne observée des différences n'est pas *significativement* différente de zéro,

mais attention, on ne conclut pas  $H_0$ , on ne conclut pas que les moyennes vraies  $\mu_A$  et  $\mu_B$  sont égales !

*Pourquoi une procédure différente pour comparer deux moyennes quand les échantillons ne sont pas indépendants ?*

$X_A$  et  $X_B$  ne sont pas indépendantes,

donc  $M_A$  (variable Moyenne de  $n$  répétitions de  $X_A$ ) et

$M_B$  (variable Moyenne de  $n$  répétitions de  $X_B$ ) ne sont pas indépendantes et dans ce cas :

$$Var(M_A - M_B) = Var(M_A) + Var(M_B) - 2cov(M_A, M_B)$$

$$E(M_A - M_B) = \mu_A - \mu_B$$

La différence entre les deux valeurs relevées est calculée pour chaque sujet : on dispose de  $n$  réalisations de la variable différence  $Y = X_A - X_B$  dont la moyenne vraie

est  $\mu = \mu_A - \mu_B$

$m$  et  $s^2$  sont la moyenne observée et la variance observée de ces différences

**Test de comparaison d'une distribution observée à une distribution théorique :  $\chi^2$  d'adéquation**

Soit une variable qualitative à k modalités.

On observe sur un échantillon de taille n l'effectif de chaque modalité : les effectifs observés sont  $O_1, O_2, \dots, O_k$

$H_0$  : la distribution vraie de la variable est égale à la distribution théorique de référence.

$H_1$  : la distribution de la variable diffère de la distribution théorique de référence.

On calcule les effectifs attendus sous  $H_0$  :  $C_1, C_2, \dots, C_k$

Condition de validité : tous les  $C_i \geq 5$

Sous  $H_0$ ,  $K = \sum_{i=1}^k \frac{(O_i - C_i)^2}{C_i} \sim \chi^2 \text{ à } (k-1) \text{ ddl}$

Avec  $\alpha = 5\%$

Si  $K > \chi_{0,05}^2(k-1)$  (valeur lue dans la table du  $\chi^2$ ) alors  $RH_0$

Degré de signification :  $P(\chi^2(k-1) > K \text{ calculé})$

Sinon  $NRH_0$

Dans le cas d'une variable à 2 modalités ( $k=2$ , nombre de ddl=1) test équivalent au test de comparaison d'une proportion observée à une proportion théorique, même conclusion.

$K = z^2$

**Test de comparaison de plusieurs distributions observées :  $\chi^2$  d'homogénéité**

**Test d'indépendance entre variables qualitatives :  $\chi^2$  d'indépendance**

Il s'agit de comparer les répartitions d'une variable qualitative à k modalités, observées dans m échantillons.

$H_0$  : les distributions vraies de la variable sont les mêmes dans les m populations.

$H_1$  : les distributions vraies sont différentes.

Ou bien il s'agit de tester l'indépendance de deux variables qualitatives, l'une à k modalités, l'autre à m modalités observées dans un échantillon

$H_0$  : les deux variables sont indépendantes.

$H_1$  : les deux variables sont liées.

	$x_1$	...	$x_j$	...	$x_k$	Total ligne
$y_1$						$l_1$
⋮						
$y_i$			$O_{i,j}$			$l_i$
⋮						
$y_m$						$l_m$
Total colonne	$c_1$		$c_j$			$n$

On calcule les effectifs attendus sous  $H_0$  :

Pour chaque case,  $C_{i,j} = \frac{l_i \times c_j}{n}$

Condition de validité :

Tous les  $C_{i,j} \geq 5$

Sous  $H_0$ ,  $K = \sum_{i,j} \frac{(O_{i,j} - C_{i,j})^2}{C_{i,j}} \sim \chi^2 \text{ à } (k-1)(m-1) \text{ ddl}$

Avec  $\alpha = 5\%$

Si  $K > \chi_{0,05}^2(k-1)(m-1)$  (valeur lue dans la table du  $\chi^2$ ) alors  $RH_0$

Degré de signification :  $P(\chi^2(k-1)(m-1) > K \text{ calculé})$

Sinon  $NRH_0$

Dans le cas de variables à 2 modalités ( $k=m=2$ , nombre de ddl=1)

Test équivalent au test de comparaison de deux proportions observées, même conclusion.

$K = z^2$



## Tests d'hypothèses : généralités

On teste une hypothèse  $H_0$  (hypothèse nulle) contre une hypothèse  $H_1$  (hypothèse alternative).

Une règle de décision étant fixée, soit on rejette  $H_0$  ( $RH_0$ ), soit on ne rejette pas  $H_0$  ( $NRH_0$ ).

### Erreurs et risques d'erreur

**Erreur de 1<sup>ère</sup> espèce** (de type I) = erreur de décision sous  $H_0$  (si  $H_0$  est vraie) =  $RH_0/H_0$   
ne peut se produire que lorsqu'on rejette  $H_0$

$\alpha$  = **Risque de 1<sup>ère</sup> espèce** = probabilité de commettre l'erreur de 1<sup>ère</sup> espèce =  $P(RH_0/H_0)$

**Erreur de 2<sup>nde</sup> espèce** (de type II) = erreur de décision sous  $H_1$  (si  $H_1$  est vraie) =  $NRH_0/H_1$   
ne peut se produire que lorsqu'on ne rejette pas  $H_0$

$\beta$  = **Risque de 2<sup>nde</sup> espèce** = probabilité de commettre l'erreur de 2<sup>nde</sup> espèce =  $P(NRH_0/H_1)$

**Erreur totale** = ( $RH_0$  et  $H_0$  vraie) ou ( $NRH_0$  et  $H_1$  vraie)

$P(\text{erreur totale}) = P(RH_0 \cap H_0) + P(NRH_0 \cap H_1) = P(RH_0/H_0)P(H_0) + P(NRH_0/H_1)P(H_1)$

### Puissance

**Puissance** =  $P(RH_0/H_1) = 1 - P(NRH_0/H_1) = 1 - \beta$

### Degré de signification

Ne se calcule que si on rejette  $H_0$

En cas de rejet de  $H_0$ , le degré de signification, noté  $p$ , est la plus petite valeur de risque  $\alpha$  avec lequel on aurait encore rejeté  $H_0$ .

Le degré de signification mesure la force avec laquelle on rejette  $H_0$  : plus il est petit, plus confortable est le rejet.



# 1 Cas usuels sur des proportions avec des grands échantillons

<p><b>Test de comparaison d'une proportion observée à une proportion théorique donnée</b> <math>\pi_0</math></p>	<p><b>Un seul échantillon d'effectif <math>n</math></b></p>	<p><b>Conditions de validité :</b> <math>n\pi_0 \geq 5</math> et <math>n(1-\pi_0) \geq 5</math></p>	<p><math>H_0 : \pi = \pi_0</math>    <math>H_1 : \pi \neq \pi_0</math></p> <p>Sous <math>H_0</math>, <math>P_n \sim N(\pi_0; \frac{\pi_0(1-\pi_0)}{n})</math> donc <math>Z = \frac{P_n - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}} \sim N(0;1)</math></p> <p>On calcule <math>z = \frac{p - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}}</math> avec <math>p</math> la proportion observée dans l'échantillon</p> <p>Avec un risque de 1<sup>ère</sup> espèce fixé à <math>\alpha</math>: Si <math> z  &gt; \varepsilon_\alpha</math> alors <b>RRH<sub>0</sub></b> au risque <math>\alpha</math>, et on dit que la proportion observée est <b>significativement</b> différente de la valeur donnée et on conclut : la proportion vraie <math>\pi</math> est différente de la valeur donnée <math>\pi_0</math>.</p> <p><b>Degré de signification</b> = <math>P( Z  &gt;  z_{calculé} )</math></p> <p>sinon <b>NRH<sub>0</sub></b> et on dit que la proportion observée n'est pas <b>significativement</b> différente de la valeur donnée, mais attention, on ne conclut pas <math>H_0</math>, on ne conclut pas que la proportion vraie est égale à la valeur donnée <math>\pi_0</math>.</p>
<p><b>Test de comparaison de deux proportions observées</b></p>	<p><b>Deux échantillons d'effectifs <math>n_A</math> et <math>n_B</math></b></p> <p><b>Conditions de validité :</b> on calcule d'abord la proportion observée <math>p</math> commune dans les deux échantillons <math>p = \frac{n_A p_A + n_B p_B}{n_A + n_B}</math> avec <math>p_A</math> et <math>p_B</math> les proportions observées dans chaque groupe</p> <p>puis on vérifie que : <math>n_A p \geq 5</math> et <math>n_A(1-p) \geq 5</math> <math>n_B p \geq 5</math> et <math>n_B(1-p) \geq 5</math></p>	<p><math>H_0 : \pi_A = \pi_B (= \pi)</math>    <math>H_1 : \pi_A \neq \pi_B</math></p> <p>Sous <math>H_0</math>, <math>P_{n_A} - P_{n_B} \sim N(0; \frac{\pi(1-\pi)}{n_A} + \frac{\pi(1-\pi)}{n_B}) = \pi(1-\pi)(\frac{1}{n_A} + \frac{1}{n_B})</math></p> <p>Donc <math>Z = \frac{P_{n_A} - P_{n_B}}{\sqrt{\pi(1-\pi)(\frac{1}{n_A} + \frac{1}{n_B})}} \sim N(0;1)</math> On calcule <math>z = \frac{p_A - p_B}{\sqrt{p(1-p)(\frac{1}{n_A} + \frac{1}{n_B})}}</math></p> <p>Avec un risque de 1<sup>ère</sup> espèce fixé à <math>\alpha</math>: Si <math> z  &gt; \varepsilon_\alpha</math> alors <b>RRH<sub>0</sub></b> au risque <math>\alpha</math>, et on dit que les proportions observées sont <b>significativement</b> différentes et on conclut : les proportions vraies <math>\pi_A</math> et <math>\pi_B</math> sont différentes.</p> <p><b>Degré de signification</b> = <math>P( Z  &gt;  z_{calculé} )</math></p> <p>sinon <b>NRH<sub>0</sub></b> et on dit que les proportions observées ne sont pas <b>significativement</b> différentes, mais attention, on ne conclut pas <math>H_0</math>, on ne conclut pas que les proportions vraies <math>\pi_A</math> et <math>\pi_B</math> sont égales.</p>	

Tests usuels sur des moyennes avec des grands échantillons

<p>Test de comparaison d'une moyenne observée à une valeur donnée</p>	<p>Un seul échantillon d'effectif <math>n</math></p>	<p>Conditions de validité : <math>n \geq 30</math></p>	<p><math>H_0 : \mu = \mu_0</math>      <math>H_1 : \mu \neq \mu_0</math></p> <p>Sous <math>H_0</math>, <math>M_n \sim N(\mu_0, \frac{\sigma^2}{n})</math> donc <math>Z = \frac{M_n - \mu_0}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0;1)</math> On calcule <math>z = \frac{m - \mu_0}{\sqrt{\frac{s^2}{n}}}</math></p> <p>avec <math>m</math> la moyenne observée dans l'échantillon et <math>s^2</math> la variance observée dans l'échantillon.</p> <p>Avec un risque de 1<sup>ère</sup> espèce fixé à <math>\alpha</math>: Si <math> z  &gt; \varepsilon_\alpha</math> alors <b>RH<sub>0</sub></b> au risque <math>\alpha</math>, et on dit que la moyenne observée est <i>significativement</i> différente de la valeur donnée et on conclut : la moyenne vraie <math>\mu</math> est différente de la valeur donnée <math>\mu_0</math>.</p> <p><b>Degré de signification</b> = <math>P( Z  &gt;  z_{calculé} )</math></p> <p>sinon <b>NRH<sub>0</sub></b> et on dit que la moyenne observée n'est pas <i>significativement</i> différente de la valeur donnée, mais attention, on ne conclut pas <math>H_0</math>, on ne conclut pas que la moyenne vraie est égale à la valeur donnée <math>\mu_0</math>.</p>
<p>Test de comparaison de deux moyennes observées</p>	<p>Deux échantillons d'effectifs <math>n_A</math> et <math>n_B</math></p>	<p>Conditions de validité : <math>n_A \geq 30</math> et <math>n_B \geq 30</math></p>	<p><math>H_0 : \mu_A = \mu_B</math>      <math>H_1 : \mu_A \neq \mu_B</math></p> <p>Sous <math>H_0</math>, <math>M_{n_A} - M_{n_B} \sim N(0; \frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B})</math> donc <math>Z = \frac{M_{n_A} - M_{n_B}}{\sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}} \sim N(0,1)</math></p> <p>On calcule <math>z = \frac{m_A - m_B}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}</math> avec <math>m_A</math> et <math>m_B</math> les moyennes observées dans les échantillons et <math>s_A^2</math> et <math>s_B^2</math> les variances observées dans les échantillons.</p> <p>Avec un risque de 1<sup>ère</sup> espèce fixé à <math>\alpha</math>: Si <math> z  &gt; \varepsilon_\alpha</math> alors <b>RH<sub>0</sub></b> au risque <math>\alpha</math>, et on dit que les moyennes observées sont <i>significativement</i> différentes et on conclut : les moyennes vraies <math>\mu_A</math> et <math>\mu_B</math> sont différentes. <b>Degré de signification</b> = <math>P( Z  &gt;  z_{calculé} )</math></p> <p>sinon <b>NRH<sub>0</sub></b> et on dit que les moyennes observées ne sont pas <i>significativement</i> différentes, mais attention, on ne conclut pas <math>H_0</math>, on ne conclut pas que les moyennes vraies <math>\mu_A</math> et <math>\mu_B</math> sont égales.</p>



## Calcul de la puissance

### Test de comparaison de deux proportions observées

$$H_0 : \pi_A = \pi_B (= \pi) \quad H_1 : \pi_A - \pi_B = \Delta$$

$$\text{Sous } H_1, Z' = \frac{P_{n_A} - P_{n_B}}{\sqrt{\pi(1-\pi)\left(\frac{1}{n_A} + \frac{1}{n_B}\right)}} \sim N\left(\frac{\Delta}{\sqrt{\pi(1-\pi)\left(\frac{1}{n_A} + \frac{1}{n_B}\right)}}; 1\right)$$

$$\alpha = 5\%$$

$$\text{Puissance} = P(RH_0/H_1)$$

$$= P(|Z'| > 1,96)$$

$$= P(Z' < -1,96) + P(Z' > 1,96)$$

$$= P\left(Z < -1,96 - \frac{\Delta}{\sqrt{p(1-p)\left(\frac{1}{n_A} + \frac{1}{n_B}\right)}}\right) + P\left(Z > 1,96 - \frac{\Delta}{\sqrt{p(1-p)\left(\frac{1}{n_A} + \frac{1}{n_B}\right)}}\right)$$

L'un des deux termes est en général négligeable.

$\Delta > 0$   $p_1$  négl.  
 $\Delta < 0$   $p_2$  négl.

### Test de comparaison de deux moyennes observées

$$H_0 : \mu_A = \mu_B \quad H_1 : \mu_A - \mu_B = \Delta$$

$$\text{Sous } H_1, Z' = \frac{M_{n_A} - M_{n_B}}{\sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}} \sim N\left(\frac{\Delta}{\sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}}; 1\right)$$

$$\alpha = 5\%$$

$$\text{Puissance} = P(RH_0/H_1)$$

$$= P(|Z'| > 1,96)$$

$$= P(Z' < -1,96) + P(Z' > 1,96)$$

$$= P\left(Z < -1,96 - \frac{\Delta}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}\right) + P\left(Z > 1,96 - \frac{\Delta}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}\right)$$

L'un des deux termes est en général négligeable.

### Test de comparaison d'une proportion observée à une proportion théorique donnée $\pi_0$

$$H_0 : \pi = \pi_0 \quad H_1 : \pi - \pi_0 = \Delta$$

$$\text{Sous } H_1, Z' = \frac{P_n - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}} \sim N\left(\frac{\Delta}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}}; 1\right)$$

$$\alpha = 5\%$$

$$\text{Puissance} = P(RH_0/H_1)$$

$$= P(|Z'| > 1,96)$$

$$= P(Z' < -1,96) + P(Z' > 1,96)$$

$$= P\left(Z < -1,96 - \frac{\Delta}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}}\right) + P\left(Z > 1,96 - \frac{\Delta}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}}\right)$$

L'un des deux termes est en général négligeable.

### Test de comparaison d'une moyenne observée à une valeur donnée

$$H_0 : \mu = \mu_0 \quad H_1 : \mu - \mu_0 = \Delta$$

$$\text{Sous } H_1, Z' = \frac{M_n - \mu_0}{\sqrt{\frac{\sigma^2}{n}}} \sim N\left(\frac{\Delta}{\sqrt{\frac{\sigma^2}{n}}}; 1\right)$$

$$\alpha = 5\%$$

$$\text{Puissance} = P(RH_0/H_1)$$

$$= P(|Z'| > 1,96)$$

$$= P(Z' < -1,96) + P(Z' > 1,96)$$

$$= P\left(Z < -1,96 - \frac{\Delta}{\sqrt{\frac{s^2}{n}}}\right) + P\left(Z > 1,96 - \frac{\Delta}{\sqrt{\frac{s^2}{n}}}\right)$$

L'un des deux termes est en général négligeable.