

## Fluctuations d'échantillonnage et tests d'hypothèse

Yohann.Foucher@univ-nantes.fr

Equipe d'Accueil 4275 "Biostatistique, recherche clinique et mesures subjectives en  
santé", Université de Nantes

Master 2 - Bioinformatique, 26 Octobre 2011



UNIVERSITÉ DE NANTES



CENTRE HOSPITALIER  
UNIVERSITAIRE DE NANTES

itun

institut  
transplantation  
urologie  
néphrologie

INSERM - UNIV NANTES

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

1. Fluctuations d'échantillonnage
2. Quelques rappels
3. Intervalle de confiance
4. Les tests d'hypothèse

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

1. Fluctuations d'échantillonnage

2. Quelques rappels

3. Intervalle de confiance

4. Les tests d'hypothèse

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Le réel est monstrueux. Il est énorme, il est hors norme.\*
- Etudier totalement une population reviendrait à vouloir percevoir de manière simultanée et continue, toutes les caractéristiques de tous les individus de cette population.†
- Pour approcher et distinguer les choses, pour s'en faire une idée tangible, il faut s'en tenir à l'appréhension d'un nombre limité de caractéristiques.

---

\*. Edgar Morin

†. Y. Macé. Journal de la Société Française de Statistique, tome 147, 2006.

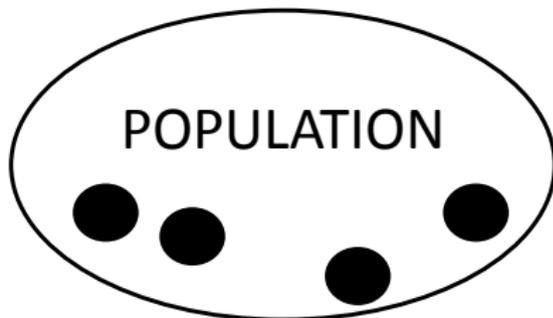
www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse



Population non-observable



Echantillons observables

www.divat.fr

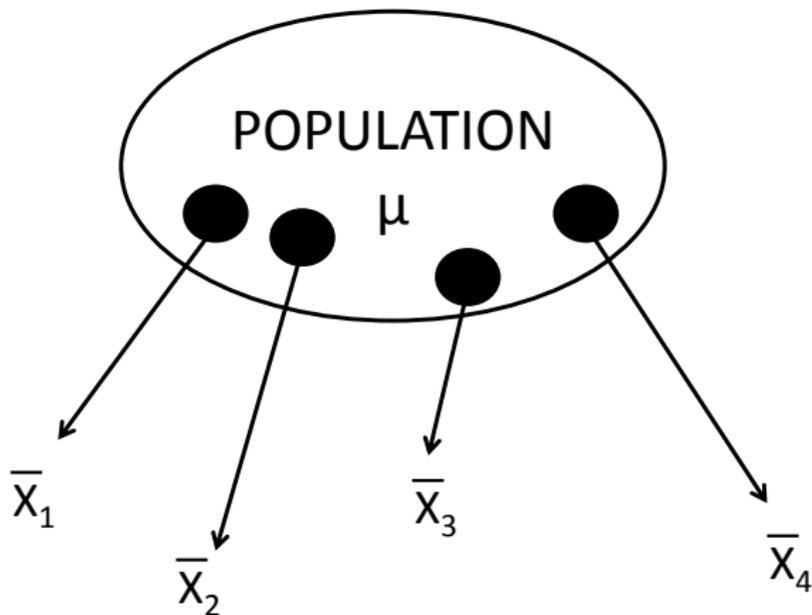
Fluctuations  
d'échantillon-  
nage

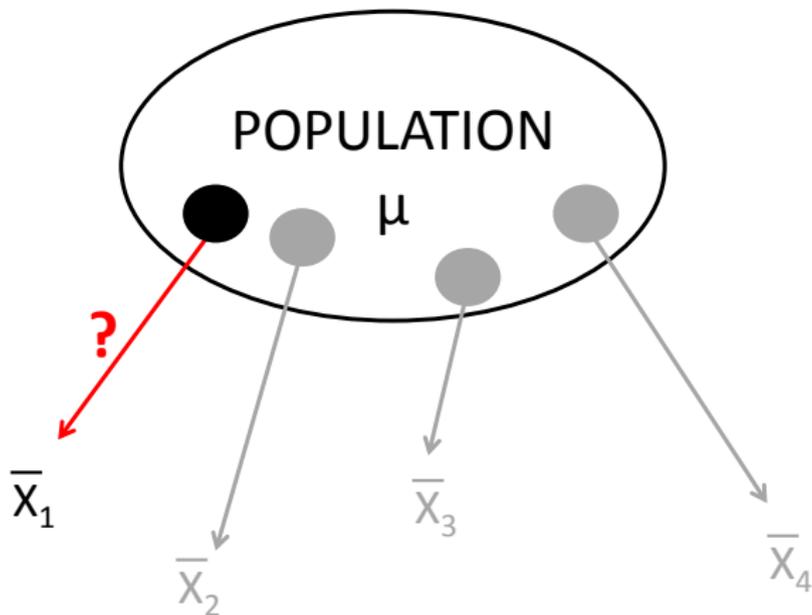
Quelques  
rappels

Intervalle de  
confiance

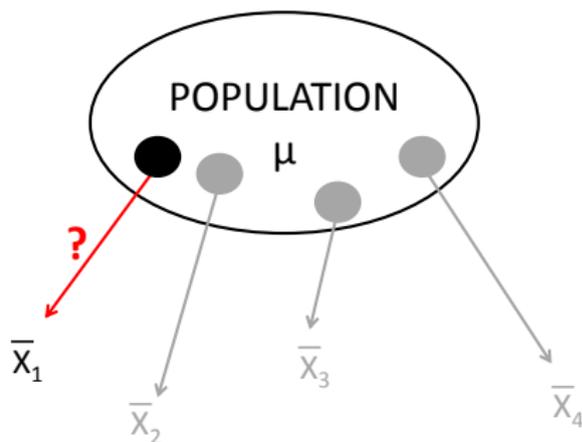
Les tests  
d'hypothèse

- Soit  $\mathcal{P}$  la population exhaustive d'intérêt.
  - Ex :  $\mathcal{P}$  = toutes les femmes atteintes d'un cancer du sein.
- Soit  $\mu$  l'âge moyen de cette population au diagnostic.
  - Ex :  $\mu = 55$  ans (supposition).
- Il n'est pas possible de mesurer  $\mu$  à partir de tous les patients.
- Soit  $E$ , un échantillon de  $N$  patientes à partir desquelles on observe la moyenne,  $\bar{x}$ . On espère que  $\bar{x}$  est proche de  $\mu$ .
- **Problème** : si plusieurs échantillons sont réalisés ( $E_1, E_2, E_3$ ), on observera autant de moyennes.
  - Ex :  $N = 100$  femmes,  $\bar{x}_1 = 53.1$ ,  $\bar{x}_2 = 58.0$ ,  $\bar{x}_3 = 56.3$ .

Fluctuations  
d'échantillon-  
nageQuelques  
rappelsIntervalle de  
confianceLes tests  
d'hypothèse



- Exemple : On souhaite mesurer l'association entre l'expression d'un gène et le succès thérapeutique d'une chimiothérapie chez les femmes avec un cancer du sein. On mesure chez 100 femmes l'expression par PCR quantitative et on suit les femmes. On sait que l'âge moyen des femmes atteintes d'un cancer du sein est proche de 55 ans. Soit  $\bar{x}$  l'âge moyen observé des 100 femmes. Peut-on considérer cet échantillon représentatif?
  - a. Si  $\bar{x} = 40$  ans : il est très probable que l'échantillon soit composé de femmes plus jeunes qu'en moyenne  
→ Résultats peu généralisables en pratique
  - b. Si  $\bar{x} = 70$  ans : il est très probable que l'échantillon soit composé de femmes plus âgées qu'en moyenne  
→ Résultats peu généralisables en pratique
  - c. Et si  $\bar{x} = 58$  ans ?



- L'inférence statistique ne conduit pas à une conclusion "stricte".
- Elle attache une probabilité d'erreur à chaque conclusion.

⇒ **AUCUNE CERTITUDE SUR LES CONCLUSIONS**

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

**Quelques  
rappels**

Intervalle de  
confiance

Les tests  
d'hypothèse

1. Fluctuations d'échantillonnage

2. Quelques rappels

3. Intervalle de confiance

4. Les tests d'hypothèse

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Les v.a. ont été définies à l'origine pour représenter un gain.
- Ex : Lancer d'une pièce de monnaie.  $X$  peut prendre les valeurs :
  - 1 euro si pile
  - -1 euro si face
- De façon plus générale une variable aléatoire est le résultat inconnu d'une expérience aléatoire.
- Une fois l'expérience réalisé, on observe la réalisation de la v.a.  $X$ . Cette valeur n'est plus aléatoire mais est certaine. On la note souvent  $x$ .  $x_1$  pour la 1ère expérience,  $x_2$  pour la 2nde expérience, etc.

- La moyenne de  $x_1, x_2, \dots, x_N$  :

$$\bar{x} = N^{-1} \sum_{i=1}^N x_i$$

- La variance de  $x_1, x_2, \dots, x_N$  :

$$s^2 = (N - 1)^{-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

- La proportion de  $k$  événements parmi  $N$  :

$$p = k/N$$

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Loi normale ou loi de Laplace-Gauss
- Caractérisée par sa moyenne ( $\mu$ ) et son écart-type ( $\sigma$ ) :

$$X \sim \mathcal{N}(\mu, \sigma)$$

- La distribution est symétrique et centrée autour de la moyenne.
- Moyenne = Mode = Médiane.
- L'écart-type représente la dispersion autour de la moyenne.

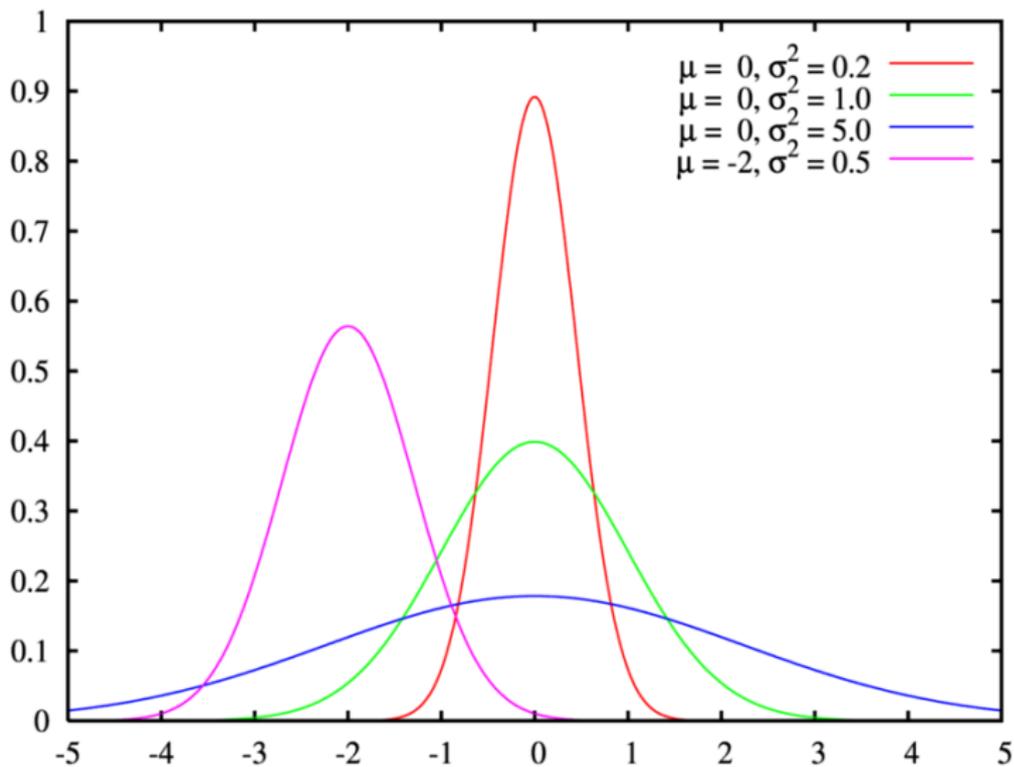
www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse



- Si  $X \sim \mathcal{N}(\mu, \sigma)$  alors :

$$U = (X - \mu)/\sigma \sim \mathcal{N}(0, 1)$$

- $\mathcal{N}(0, 1)$  est la loi normale centrée et réduite.
- En anglais : *standard normal distribution*.
- Intérêt de la transformation : la probabilité  $P(T < t)$  est connue.



www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Symétrie  $\rightarrow P(T > t) = 1 - P(T \leq t) = P(T < -t)$ .
- Valeurs importantes à retenir :
  - $P(-1,96 < T < 1,96) = 0,95$
  - $P(T < -1,64) = P(T > 1,64) = 0,95$
- Si  $X \sim \mathcal{N}(\mu, \sigma)$  alors  $\bar{X} \sim \mathcal{N}(\mu, \sigma/\sqrt{n})$

- Soit un échantillon de taille  $N$  et  $X_1, X_2, \dots, X_N$  une suite de variables aléatoires indépendantes et identiquement distribuées.
- Supposons que cette loi possède une espérance  $\mu$  et un écart-type  $\sigma$  (avec  $\sigma \neq 0$ ).
- Considérons la somme  $S_N = X_1 + X_2 + \dots + X_N = \sum_{i=1}^N X_i$ . Alors l'espérance de  $S_N$  vaut  $N\mu$  et son écart-type vaut  $\sigma\sqrt{N}$ .
- Quand  $N$  est assez grand ( $N > 30$ ), la loi normale  $\mathcal{N}(N\mu, N\sigma^2)$  est une bonne approximation de la loi de  $S_N$ .
- Cette approximation s'écrit :

$$U = \frac{S_N - N\mu}{\sigma\sqrt{N}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1)$$

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

**Quelques  
rappels**

Intervalle de  
confiance

Les tests  
d'hypothèse

- Conséquence : Une donnée influencée par une multitude de phénomènes aléatoires indépendants qui s'additionnent est approximativement décrite par une loi normale même si les phénomènes qui la composent ne suivent pas des lois normales.

- Exemple

Reprenons l'exemple du cancer du sein. On sait que  $\mu$  est proche de 55 ans. Supposons l'écart-type est égal à 30 ans.<sup>‡</sup> On observe une moyenne  $\bar{x}$  de 58 ans, soit 3 ans d'écart avec la valeur théorique.

Quelle est la probabilité (A) d'observer un échantillon de 100 femmes avec une moyenne d'âge de plus de 3 ans d'écart avec la moyenne attendue ?

$$\begin{aligned} A &= Pr(\bar{X} > 58) + Pr(\bar{X} < 52) \\ &= Pr\left(\frac{\bar{X} - 55}{30/\sqrt{100}} > \frac{58 - 55}{30/\sqrt{100}}\right) + Pr\left(\frac{\bar{X} - 55}{30/\sqrt{100}} < \frac{52 - 55}{30/\sqrt{100}}\right) \\ &= Pr(U > 1) + Pr(U < -1) \\ &= 2 \times Pr(U > 1) \end{aligned}$$

---

<sup>‡</sup> Valeur le plus souvent inconnue dans la population et estimée à partir de l'échantillon



[www.divat.fr](http://www.divat.fr)Fluctuations  
d'échantillon-  
nageQuelques  
rappelsIntervalle de  
confianceLes tests  
d'hypothèse

$$\begin{aligned} A &= 2 \times Pr(U > 1) \\ &= 2 \times (1 - 0.8413) \\ &= 0.3174 \end{aligned}$$

La probabilité d'observer un échantillon de 100 femmes avec une moyenne d'âge observée écartée de plus de 3 ans avec la moyenne théorique est égale à 0.32. Autrement dit, il y a environ une chance sur trois d'observer un échantillon encore plus éloigné de la population.

- Remarque : On peut penser que l'échantillon observé soit assez vraisemblable et que sa représentativité ne soit pas remise en question.

[www.divat.fr](http://www.divat.fr)Fluctuations  
d'échantillon-  
nageQuelques  
rappelsIntervalle de  
confianceLes tests  
d'hypothèse

Soit un échantillon de taille  $N$  et  $X_1, X_2, \dots, X_N$  une suite de variables aléatoires binaires indépendantes et identiquement distribuées selon une loi de Bernoulli de paramètre  $\pi$ .

La somme  $S_N$  suit une loi binomiale de paramètre  $N$  et  $\pi$ .

Si  $N > 30$ ,  $N\pi > 5\%$  et  $N(1 - \pi) > 5\%$ , alors

$$S_N \sim \mathcal{N}(\mu = N\pi, \sigma = \sqrt{N\pi(1 - \pi)})$$

Si  $P$  est la proportion estimée à partir de cet échantillon, alors :

$$P = S_N/N \sim \mathcal{N}\left(\mu = \pi, \sigma = \sqrt{\frac{\pi(1 - \pi)}{N}}\right)$$

- Exemple

On souhaite réaliser un échantillon de taille  $N = 100$  individus. Ces individus sont issus d'une population où la prévalence d'un certain phénotype est de 20%.  $x_i = 1$  si le sujet  $i$  possède ce phénotype et  $x_i = 0$  sinon.  $S_N$  représente le nombre d'individus avec ce phénotype. **Quelle est la probabilité d'observer moins de 10% des individus avec ce phénotype dans l'échantillon ?**

$$S_N \sim \mathcal{B}(100, 0.2)$$

Comme  $N = 100 > 30$  et  $Np = 20 > 5$ ,

$$S_n \sim \mathcal{N}(20, 4)$$



$$P = \frac{S_N}{N} \sim \mathcal{N}(0.20, 0.04)$$

$$P \sim \mathcal{N}(0.20, 0.04)$$



$$U = \frac{P - 0.20}{0.04} \sim \mathcal{N}(0, 1)$$



$$\begin{aligned} Pr(P < 0.10) &= Pr\left(\frac{P - 0.20}{0.04} < \frac{0.10 - 0.20}{0.04}\right) \\ &= Pr\left(U < \frac{0.10 - 0.20}{0.04}\right) \\ &= Pr(U < -2.50) \end{aligned}$$



www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

$$\begin{aligned}Pr(P < 0.10) &= Pr(U < -2.50) \\ &= 1 - 0.9938 \\ &= 0.0062\end{aligned}$$

La probabilité d'observer un échantillon de 100 patients parmi lesquels moins de 10% possède le phénotype est égale à 6 pour 1000.

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

**Intervalle de  
confiance**

Les tests  
d'hypothèse

1. Fluctuations d'échantillonnage

2. Quelques rappels

3. Intervalle de confiance

4. Les tests d'hypothèse

- Soit un échantillon de taille  $N$  et  $x_1, x_2, \dots, x_N$  les observations d'une variable aléatoire continue.
- La moyenne de la population ( $\mu$ ) est estimée par  $\bar{x}$ .
- Si  $N > 30$  (TCL), alors on sait aussi que

$$\bar{X} \sim \mathcal{N}(\mu, s/\sqrt{N})$$

- En centrant et réduisant, on obtient :

$$\frac{\bar{X} - \mu}{s/\sqrt{N}} \sim \mathcal{N}(0, 1)$$

- Pour obtenir l'intervalle de confiance à 95% de la moyenne estimée, on cherche l'intervalle à l'intérieur duquel on a 95% de chance de retrouver la moyenne théorique de la population :

$$Pr(-1.96 < \frac{\bar{X} - \mu}{s/\sqrt{N}} < 1.96) = 0.95$$

[www.divat.fr](http://www.divat.fr)Fluctuations  
d'échantillon-  
nageQuelques  
rappelsIntervalle de  
confianceLes tests  
d'hypothèse

$$Pr(-1.96 < \frac{\bar{X} - \mu}{s/\sqrt{N}} < 1.96) = 0.95$$

 $\Leftrightarrow$ 

$$Pr(\bar{X} - 1.96 \frac{s}{\sqrt{N}} < \mu < \bar{X} + 1.96 \frac{s}{\sqrt{N}}) = 0.95$$

 $\Leftrightarrow$ 

$$IC_{95\%}(\mu) = [\bar{x} \pm 1.96 \frac{s}{\sqrt{N}}]$$

- Exemple

Reprenons l'exemple du cancer du sein. L'écart-type estimé est de 30 ans. On observe une moyenne  $\bar{x}$  de 58 ans. **Quelle est l'intervalle de confiance de la moyenne ?**

$$\begin{aligned} IC_{95\%}(\mu) &= \left[ 58 \pm 1.96 \frac{30}{\sqrt{100}} \right] \\ &= [52.12; 63.88] \end{aligned}$$

Si on réalise 100 échantillons comme celui-ci ( $N = 100$ ), on attend 95 moyennes estimées comprises entre 52 et 63 ans.

- Remarque : comme la moyenne de la population générale est proche de 55 ans dans l'énoncé, la représentativité de l'échantillon n'est pas remise en question.

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

1. Fluctuations d'échantillonnage

2. Quelques rappels

3. Intervalle de confiance

4. Les tests d'hypothèse

- Principe de la statistique : Est ce que la différence que j'observe est due à la fluctuation d'échantillonnage ?
- J'observe une moyenne  $\bar{x}$  (et sa variance) à partir d'un échantillon. Je veux la comparer à une valeur théorique  $\mu$  (connue dans  $\mathcal{P}$ ). 2 hypothèses :
  - La différence observée est minime. Elle est due au fait que tous les sujets de la population n'ont pas été inclus. Si on avait inclus tout le monde, la moyenne observée  $\bar{x}$  serait égale à la moyenne théorique  $\mu$ .
    - Cette hypothèse est appelée **hypothèse nulle** :  $H_0$
  - La différence observée est importante. Elle ne peut pas être due au fait que tous les sujets n'ont pas été inclus. Il est évident que si on avait inclus tout le monde, la moyenne observée  $\bar{x}$  aurait été différente de la moyenne théorique  $\mu$ .
    - Cette hypothèse est appelée **hypothèse alternative** :  $H_1$
- Probabilité de me tromper si je rejette  $H_0$  (risque de 1ère espèce).
- Probabilité de me tromper si je rejette  $H_1$  (risque de 2nd espèce).

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Calculer la probabilité de se tromper en rejetant  $H_0$ .
  - probabilité critique, notée  $p_c$ .
  - *p-value* en anglais, notée  $p$ .
- C'est la probabilité d'observer l'échantillon en supposant  $H_0$  comme vraie.
  - Si cette probabilité est forte, tendance au non-rejet de  $H_0$ .
  - Si cette probabilité est faible, tendance au rejet de  $H_0$ .
- On choisi a priori une probabilité maximale d'erreur ( $\alpha$ ), au delà de laquelle on ne rejettera pas  $H_0$ 
  - Si  $p_c \geq \alpha$  : non-rejet de  $H_0$ .
  - Si  $p_c < \alpha$  : rejet de  $H_0$ .

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Pour illustrer cette logique, choisissons le cas où l'on souhaite savoir si l'échantillon observé est issu d'une population où la moyenne théorique est égale à  $\mu$ .
- Exemple : un chercheur souhaite savoir si son échantillon de 100 femmes atteintes d'un cancer du sein est représentatif des femmes atteintes de cette pathologie concernant l'âge au diagnostic ( $\mu = 55$  ans). La moyenne estimée à partir de l'échantillon est de 58 ans et la variance estimée est de 900 ans<sup>2</sup> (soit un écart-type de 30 ans).

- Définir la variable aléatoire  $X$  dans la population d'étude.
  - Age au diagnostic de cancer : v.a. continue.
- Définir la population de référence  $\mathcal{P}$  et la moyenne théorique correspondante  $\mu$ .
  - Femmes atteintes d'un cancer du sein et  $\mu = 55$  ans.
- Choisir les hypothèses à tester :
  - $H_0$  :  $\bar{X} = \mu$ , l'échantillon est issu d'une population où l'âge moyen est  $\mu$  (55 ans).
  - $H_1$  :  $\bar{X} \neq \mu$ , l'échantillon est issu d'une population où l'âge moyen n'est pas égal à  $\mu$  (55 ans).

- Définir la statistique de test sous  $H_0$ . Comme  $N > 30$  :

$$\bar{X} \sim \mathcal{N}(\mu, \sigma/\sqrt{N}) \text{ donc } U = \frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1)$$

- Définir du risque de 1ère espèce maximum
  - ex :  $\alpha = 5\%$
- Définir la région critique (RC). Il s'agit des valeurs de la statistique de test  $U$  qui ont moins de  $\alpha$  chance d'être observées sous  $H_0$ .
  - $\alpha = 5\% \Rightarrow$  Quelle est la valeur de  $u_{2.5\%}$  pour que  $P(-u_{2.5\%} < U < u_{2.5\%}) = 0.95$  ?
  - Comme  $U \sim \mathcal{N}(0, 1) \rightarrow u_{2.5\%} = 1.96 \rightarrow RC : U \notin [-1.96, 1.96]$

Si la statistique de test est inférieure à -1.96 ou si elle est supérieure à 1.96, l'échantillon observé a moins de 5% de chance de provenir d'une population où  $H_0$  est vraie.

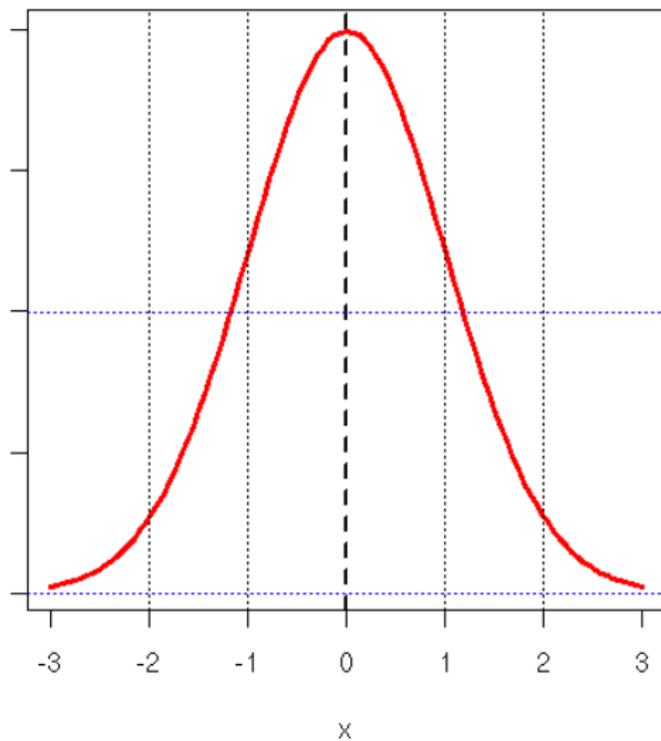
[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse



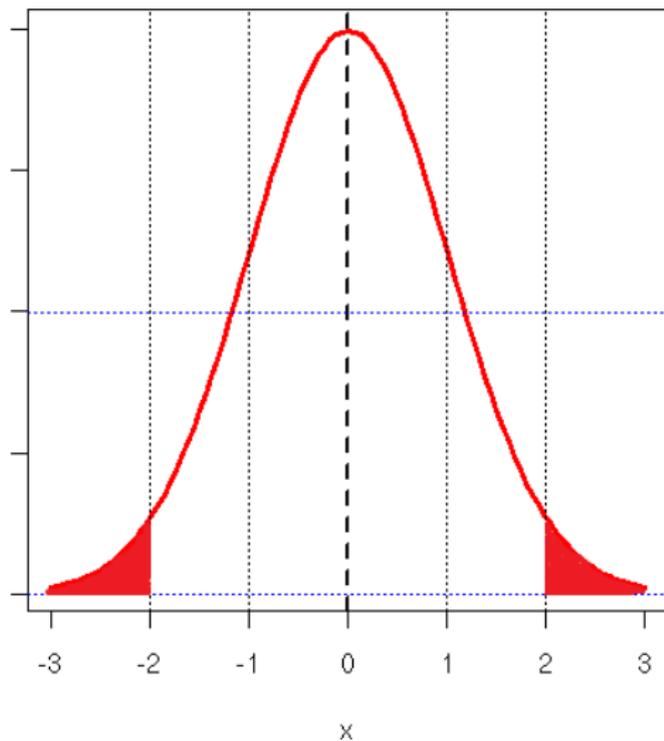
[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse



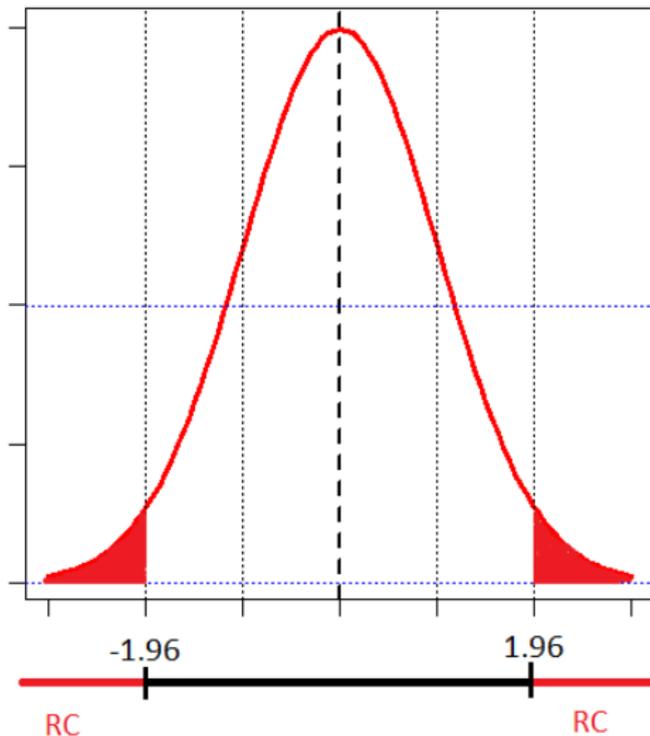
www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

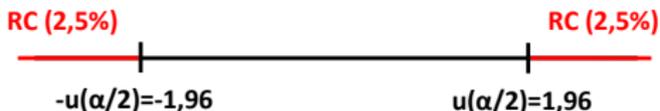
Intervalle de  
confiance

Les tests  
d'hypothèse



- Application numérique

- $u = (\bar{x} - \mu)/(s/\sqrt{N})$
- $N = 100$
- $\bar{x} = N^{-1} \sum_i x_i = 58$
- $s^2 = (N - 1)^{-1} \sum_i (x_i - \bar{x})^2 = 900$
- $u = (58 - 55)/(\sqrt{900}/\sqrt{40}) = 1.00$



- **Conclusion** : La statistique de test n'appartient pas à la région critique. L'hypothèse nulle ne peut donc pas être rejetée. Il n'est pas possible de conclure que l'échantillon ne provienne pas de la population de référence ( $p_c > 5\%$ ).

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Le risque d'erreur précédent est de rejeter à tort l'hypothèse alternative
  - Risque de 2<sup>de</sup> espèce.
  - Seule interprétation possible : On ne peut pas montrer que l'échantillon n'est pas représentatif.
  - On ne peut pas conclure sur la représentativité.
- Les conclusions ne sont jamais déterministes.
- Incertitude toujours présente et doit ressortir dans la conclusion.
- Si la conclusion avait été le rejet de  $H_0$  car moins de 5% de chance qu'elle soit vraie, la conclusion aurait pu être du type :
  - Il semble que l'échantillon observé ne soit pas représentatif de l'âge moyen de diagnostic du cancer du sein ( $p_c < 5\%$ ).

www.divat.fr

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Calculer directement la probabilité critique.
- Quelle est la probabilité d'observer cet échantillon ou un échantillon encore moins probable si  $H_0$  est vrai.
- Exemple précédent :  $u = 1.00$ .
- $p_c = P(U > 1.00) + P(U < -1.00) = 2 * P(U > 1.00) = 2 * (1 - P(U < 1.00))$ .
- Dans la table :  $P(U < 1.00) \approx 0.8413 \Rightarrow p_c = 0.3174$
- $p_c > \alpha$  : On ne rejette pas  $H_0$  car l'observation de cet échantillon ou d'un échantillon encore moins probable sous cette hypothèse est assez vraisemblable (une chance sur trois).

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

Nous venons de voir le cas d'un seul échantillon de grande taille.  
Mais il existe d'autres situations :

- Petits échantillons (TCL non-adapté)
- Comparaison de deux échantillons indépendants
- Comparaison de deux échantillons appariés
- La v.a. est catégorielle (comparaison de proportions)

[www.divat.fr](http://www.divat.fr)

Fluctuations  
d'échantillon-  
nage

Quelques  
rappels

Intervalle de  
confiance

Les tests  
d'hypothèse

- Définir la variable aléatoire  $X$ .
- Définir la population de référence.
- Choisir les hypothèses à tester.
- Définir la statistique de test sous  $H_0$ .
- Définir le risque de 1ère espèce maximum.
- Définir la région critique ( $RC$ ).
- Application numérique.
- Conclusion.