

Introduction aux statistiques

Intervalles de confiance

L1 STE

Extraction de n échantillons d'une population P

Si l'on extrait plusieurs échantillons représentatifs de taille n fixée, les différences observées entre les résultats obtenus sont dues à des *fluctuations d'échantillonnage*. A partir d'un échantillon, on n'a donc pas de certitudes mais des *estimations de paramètres*.

L'estimation d'un paramètre peut être faite

- par un seul nombre: **estimation ponctuelle**
- par 2 nombres entre lesquels le paramètre peut se trouver: **estimation par intervalle**

Estimation ponctuelle d'une moyenne

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

x barre

$$s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

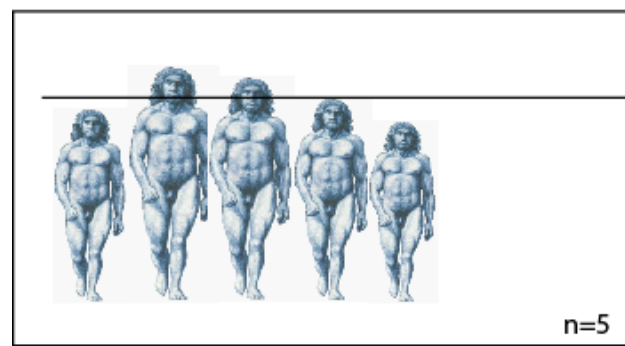
Estimateur sans biais

$$s_{\bar{x}} = \frac{s_x}{\sqrt{n}}$$

Ecart type de la moyenne

Echantillonnage – Estimation d'un paramètre

Echantillon

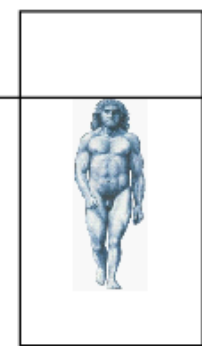


Néandertaliens males adultes

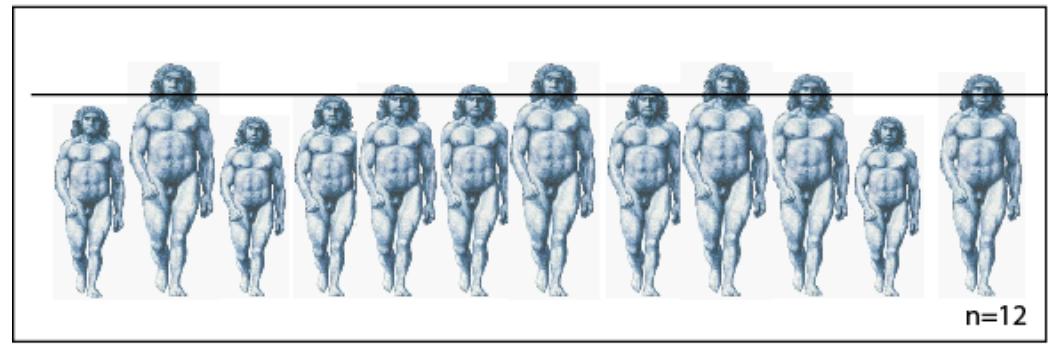
$$\bar{x} = 165 \text{ cm}$$

$$Sx = 10 \text{ cm}$$

Moyenne



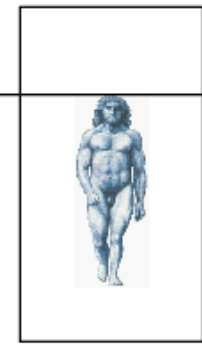
$$S\bar{x} = 10/\text{racine}(5)$$



Néandertaliens males adultes

$$\bar{x} = 164 \text{ cm}$$

$$Sx = 10,5 \text{ cm}$$



$$S\bar{x} = 10,5/\text{racine}(12)$$

Pour améliorer la connaissance de la moyenne, il faut augmenter la taille de l'échantillon

Intervalle de confiance de la moyenne

Cas des grands échantillons (>30) & variance connue:

Soit une population obéissant à une **loi normale** de moyenne μ et d'écart type σ .

$$\Pr\left(\bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

Exemple:

45 hommes de Neandertal males adultes

$$\bar{x} = 164 \text{ cm}$$

$$\sigma = 10 \text{ cm}$$

$$\mu \in \left[164 - 1.96 \cdot \frac{10}{\sqrt{45}} ; 164 + 1.96 \cdot \frac{10}{\sqrt{45}} \right]$$

$$\mu \in [161; 166.9] \quad \text{à 95\% de confiance}$$

$$\mu = 164 \pm 2.9$$



Echantillonnage – Estimation d'un paramètre

TABLE III – AIRES LIMITEES PAR LA COURBE NORMALE CENTRÉE RÉDUITE

La table fournit les valeurs de $\phi(z)$ pour z positif. Lorsque z est négatif il faut calculer le complément à l'unité de la valeur lue dans la table. La première colonne indique la première décimale de z et la première rangée fournit la deuxième décimale.

Exemples : pour $z = 1,21$, $\phi(z) = 0,8869$ et pour $z = -1,21$, $\phi(z) = 0,1131$

z	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	0,7291	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	0,9778	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
z	0,0	0,1	0,2	0,3	0,4	0,5	0,6	0,7	0,8	0,9
3	0,9987	0,9990	0,9993	0,9995	0,9997	0,9998	0,9998	0,9999	0,9999	1,0000
4	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000	1,0000

Intervalle de confiance de la moyenne

Cas des petits échantillons:

Quand $n < 30$ ou quand la variance est inconnue, on prend la loi de Student.

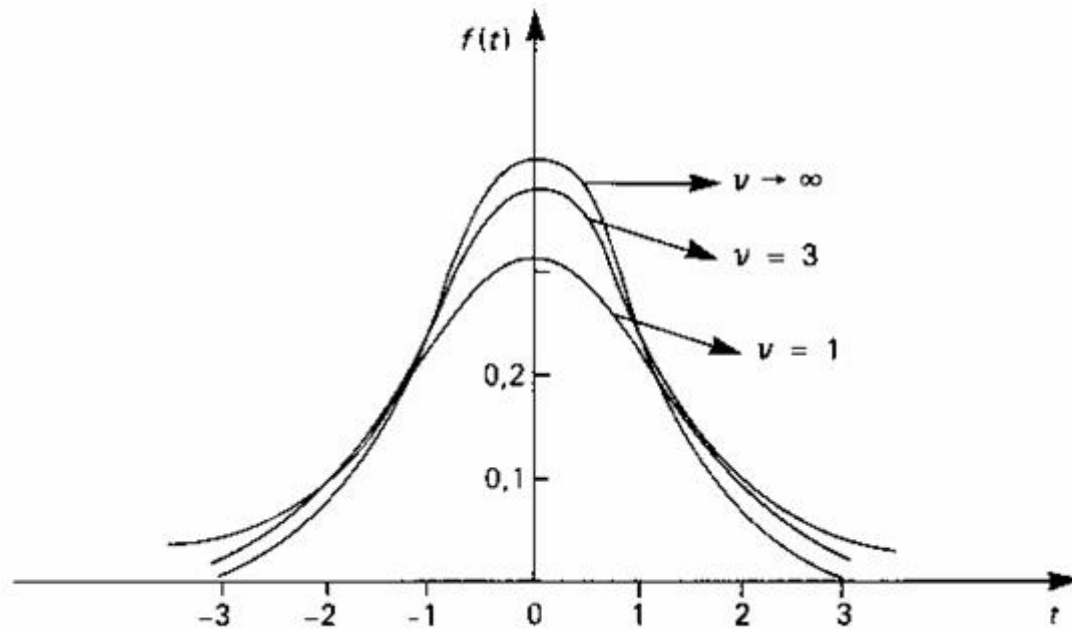
$$\Pr\left(\bar{x} - t_{\alpha/2} \cdot \frac{s_x}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2} \cdot \frac{s_x}{\sqrt{n}}\right) = 1 - \alpha$$

Pour $\nu = n - 1$ degrés de liberté

Enfin, on peut toujours utiliser la loi de Student puisque t tend vers la loi normale quand n est grand...

La loi de Student: $t(\nu)$

Famille de courbes de densité de probabilité obéissant à la loi de Student pour différents nombres de degrés de liberté



ν degrés de liberté

Converge vers **la loi Normale** quand ν augmente.

La loi de Student: $t(\nu)$

La probabilité d'obtenir une valeur de t à l'extérieur de l'intervalle $(-t_{\alpha/2} \text{ et } t_{\alpha/2}) \rightarrow$ TABLES.

$$P(|t| > t_{\alpha/2}) = \alpha$$

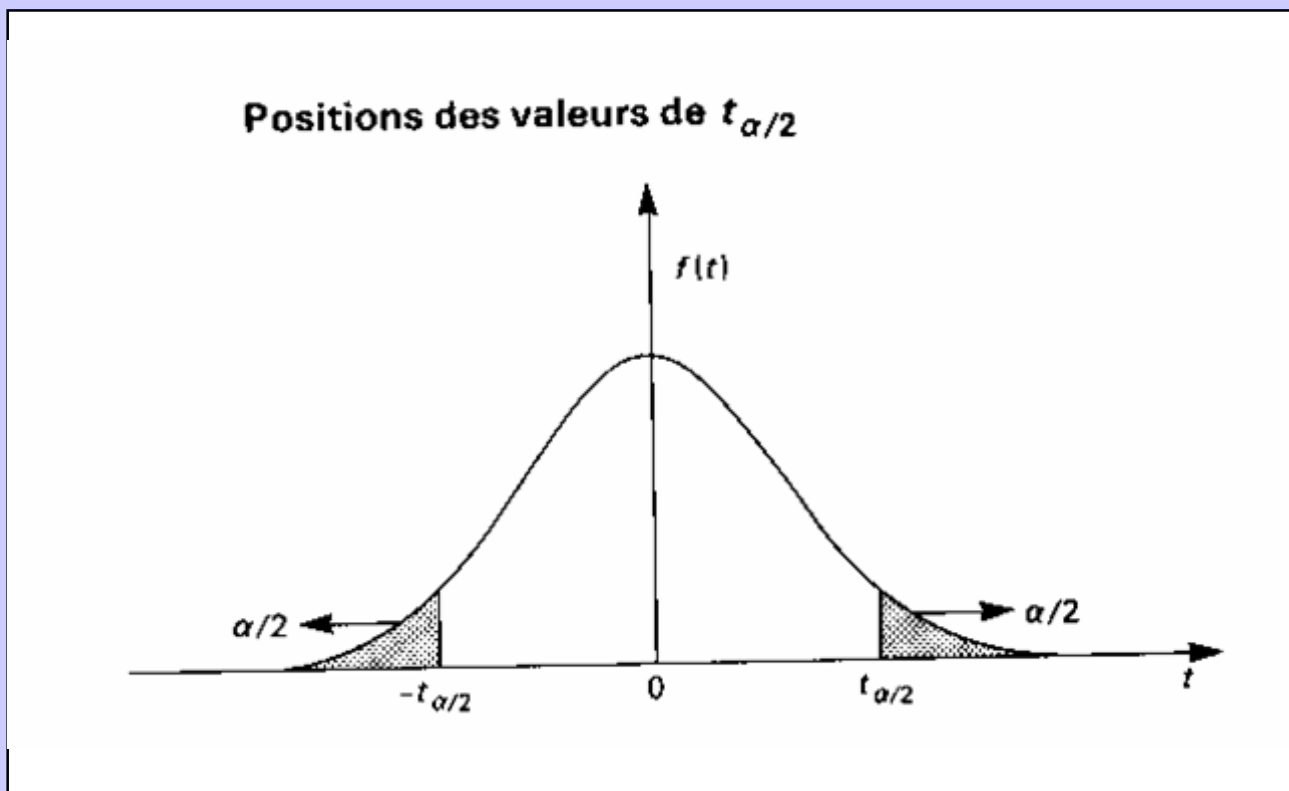


Table de Student t

ν	α					
	0.100	0.050	0.025	0.010	0.005	0.001
1	3.078	6.314	12.706	31.821	63.657	318.309
2	1.886	2.920	4.303	6.965	9.925	22.327
3	1.638	2.353	3.182	4.541	5.841	10.215
4	1.533	2.132	2.776	3.747	4.604	7.173
5	1.476	2.015	2.571	3.365	4.032	5.893
6	1.440	1.943	2.447	3.143	3.707	5.208
7	1.415	1.895	2.365	2.998	3.499	4.785
8	1.397	1.860	2.306	2.896	3.355	4.501
9	1.383	1.833	2.262	2.821	3.250	4.297
10	1.372	1.812	2.228	2.764	3.169	4.144
11	1.363	1.796	2.201	2.718	3.106	4.025
12	1.356	1.782	2.179	2.681	3.055	3.930
13	1.350	1.771	2.160	2.650	3.012	3.852
14	1.345	1.761	2.145	2.624	2.977	3.787
15	1.341	1.753	2.131	2.602	2.947	3.733
16	1.337	1.746	2.120	2.583	2.921	3.686
17	1.333	1.740	2.110	2.567	2.898	3.646
18	1.330	1.734	2.101	2.552	2.878	3.611
19	1.328	1.729	2.093	2.539	2.861	3.579
20	1.325	1.725	2.086	2.528	2.845	3.552
21	1.323	1.721	2.080	2.518	2.831	3.527
22	1.321	1.717	2.074	2.508	2.819	3.505
23	1.319	1.714	2.069	2.500	2.807	3.485
24	1.318	1.711	2.064	2.492	2.797	3.467
25	1.316	1.708	2.060	2.485	2.787	3.450
26	1.315	1.706	2.056	2.479	2.779	3.435
27	1.314	1.703	2.052	2.473	2.771	3.421
28	1.313	1.701	2.048	2.467	2.763	3.408
29	1.311	1.699	2.045	2.462	2.756	3.396
30	1.310	1.697	2.042	2.457	2.750	3.385
100	1.290	1.660	1.984	2.365	2.626	3.174
∞	1.282	1.645	1.960	2.326	2.576	3.090

La table de Student donne les valeurs $t_{(\alpha,\nu)}$ telles que

$$P\{T > t_{(\alpha,\nu)}\} = \alpha$$

Exemple:

6 hommes de Neandertal males adultes

$$\bar{x} = 165 \text{ cm}$$

$$s_x = 11 \text{ cm}$$

$$\mu \in \left[165 - 2.57 \cdot \frac{11}{\sqrt{6}}; 165 + 2.57 \cdot \frac{11}{\sqrt{6}} \right]$$

$$\mu \in [153; 177] \quad \text{à 95\% de confiance}$$

$$\mu = 165 \pm 12$$



Enfin, finalement on peut toujours utiliser la loi de Student puisque t tend vers la loi normale quand n est grand...

Un problème très fréquent!

Un quotidien publie tous les mois la cote du chef du gouvernement à partir d'un sondage réalisé sur un échantillon représentatif de 1000 personnes. En janvier, la cote publiée était de 38% d'opinions favorables, en février de 36%. Un journaliste commente alors ces valeurs par "Le chef du gouvernement perd 2 points !!"

En fait: On construit un intervalle de confiance autour des proportions. Avec un seuil de 95%, on obtient respectivement [35;41] et [33;39] pour les valeurs 38% et 36%. Les deux intervalles ayant une intersection non vide, on ne peut pas conclure qu'il y ait eu baisse ou augmentation de la cote du chef de gouvernement.

Estimation ponctuelle d'un pourcentage

La population est formée d'individus ayant ou non un caractère A. Soit p la probabilité pour qu'un individu pris au hasard dans la population présente le caractère A.

$$p = a / n$$

$$s_p^2 = \frac{p(1-p)}{n-1}$$

Quand on dispose d'un seul échantillon de taille n , la meilleure estimation ponctuelle de P est donc la fréquence p observée sur l'échantillon.

Intervalle de confiance d'un pourcentage

Grands échantillons ($n > 30$), p ni voisin de 0, ni voisin de 1, ($np > 5$, $n(1-p) > 5$)

La variable fréquence obéit alors à une loi normale centrée réduite.
Théorème Central Limit

$$\Pr\left(p - Z_{\alpha/2} \cdot \sqrt{\frac{p(1-p)}{n}} < P < p + Z_{\alpha/2} \cdot \sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$

Introduction aux statistiques

Premiers tests statistiques

L1 STE

Quel est le problème...?

On sait qu'un homme de Neandertal mesure en moyenne 165 cm.

Sur un site on trouve 40 hommes avec une moyenne de 167 et un écart type de 8 cm (e.t. échantillon).

Comparaison de la moyenne avec la valeur théorique de 165 cm



Possibilités:

Moyenne très élevée: Nous pourrions être amenés à croire que ces hommes ont des tailles différentes de 165 cm

Moyenne faiblement plus élevée: on ne pourra pas conclure si c'est significativement supérieur à la norme ou si c'est l'effet du hasard.

Question: à partir de quelle limite pouvons nous raisonnablement conclure à une différence?

H_0 : $\mu=165$ (il n'y pas de différence)

H_1 : $\mu \neq 165$

Calcul de

$$s_{\bar{x}} = \frac{s_x}{\sqrt{n}} = \frac{8}{\sqrt{40}} = 1.265$$

On mesure en fait 167 +/- 2.48 à 95% de confiance, ce qui n'est pas différent de 165 cm!

Les deux risques d'erreur dans un test.

Erreur de 2^{nde} espèce (compliquée)

Décision	H_0 est vraie	H_1 est vraie
H_0 acceptée	Bonne décision	Erreur β
H_0 rejetée	Erreur α	Bonne décision

Erreur de 1^{ere} espèce

A priori on ne sait pas à quel type d'erreur on sera confronté:

Le résultat de l'échantillon a révélé 167 cm probablement par pur hasard. On conclue que la moyenne pourrait être 165 cm alors qu'en fait elle est mesurée à 167 cm.

H_0 : hypothèse nulle ou principale

Ex: Les haches de type A présentent les mêmes teneurs en Sn que les haches de type B.

H_1 : hypothèse alternative ou contraire ...

Soumission à une épreuve de vérité!

Conclusion : différence attribuable aux fluctuations d'échantillonnage???

Niveau de signification : un peu arbitraire...

significatif : 0.05

hautement significatif : 0.01

très hautement significatif : 0.001.

Test bilatéral / unilatéral :

bilatéral : différence sans se préoccuper du sens.

Unilatéral : $>$ ou $<$. Zone de rejet d'un seul côté de la distribution de probabilité de référence.

Echantillons indépendants

Indépendants : aucune influence du 1^{er} ech sur le 2nd.

Comparaison des moyennes de 2 grands échantillons indépendants (n_1 et $n_2 > 30$):

Deux échantillons qui suivent des **lois normales**: $\mu_1, \sigma^2_1; \mu_2, \sigma^2_2$

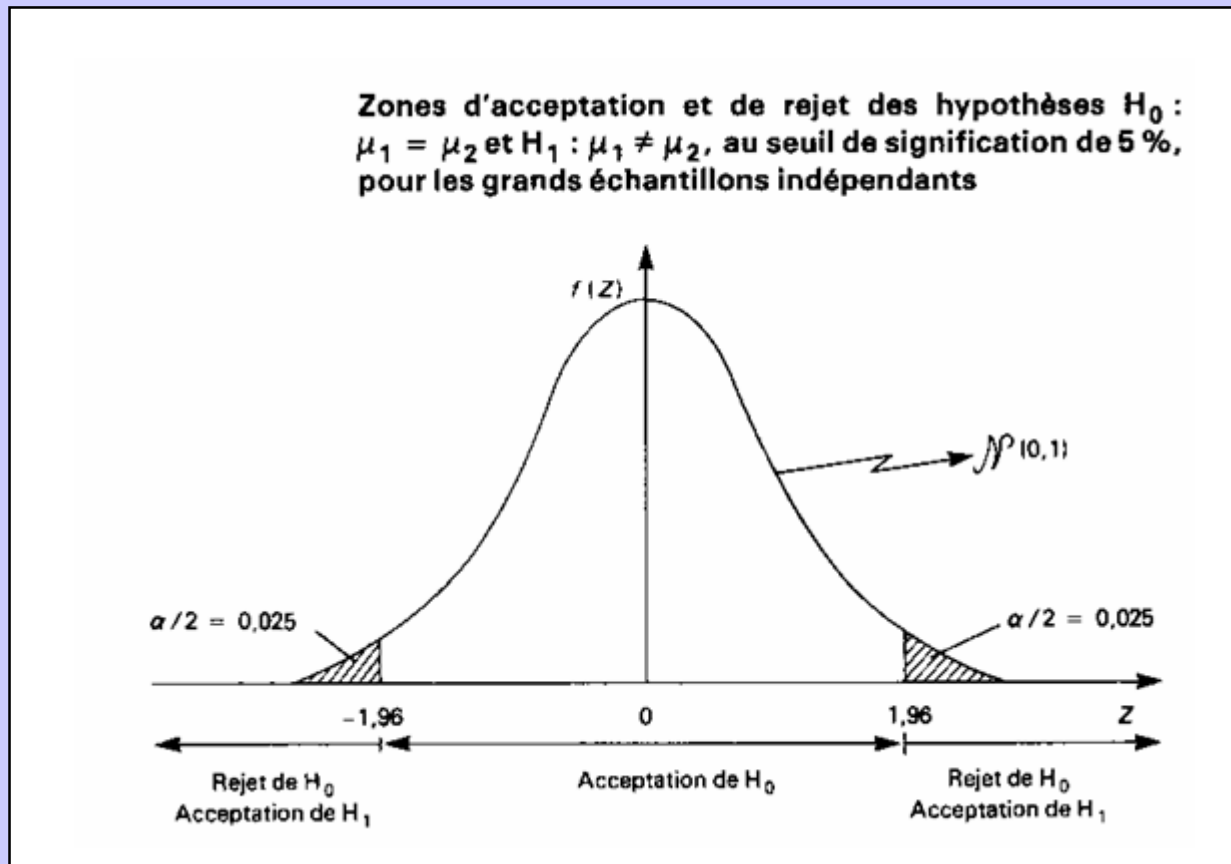
$$H_0 : \mu_1 = \mu_2$$

$$Z_c = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_{x_1}^2}{n_1} + \frac{s_{x_2}^2}{n_2}}}$$

Si H_0 est vraie, Z_c suit une loi normale $N(0,1)$

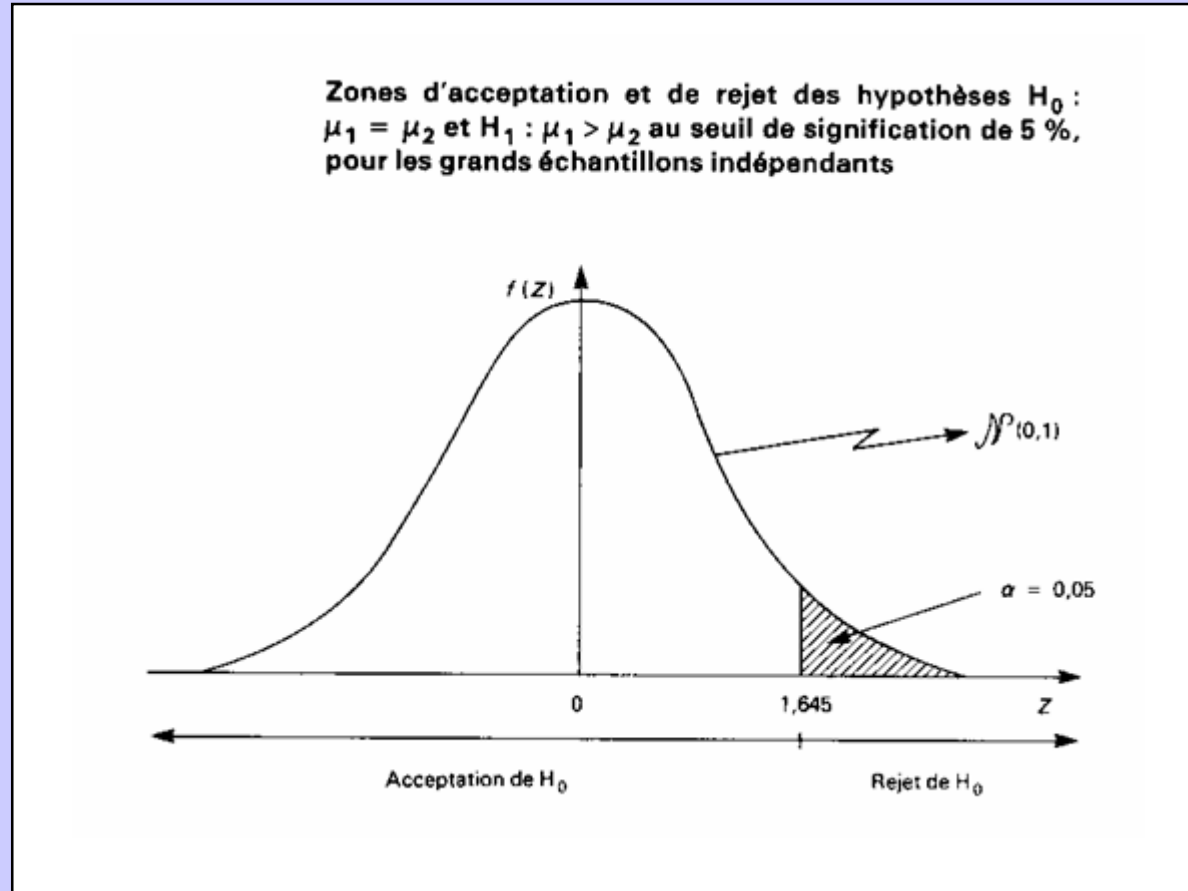
Comparaison de deux moyennes – grands échantillons -

$H_1 : \mu_1 \neq \mu_2$ bilatéral



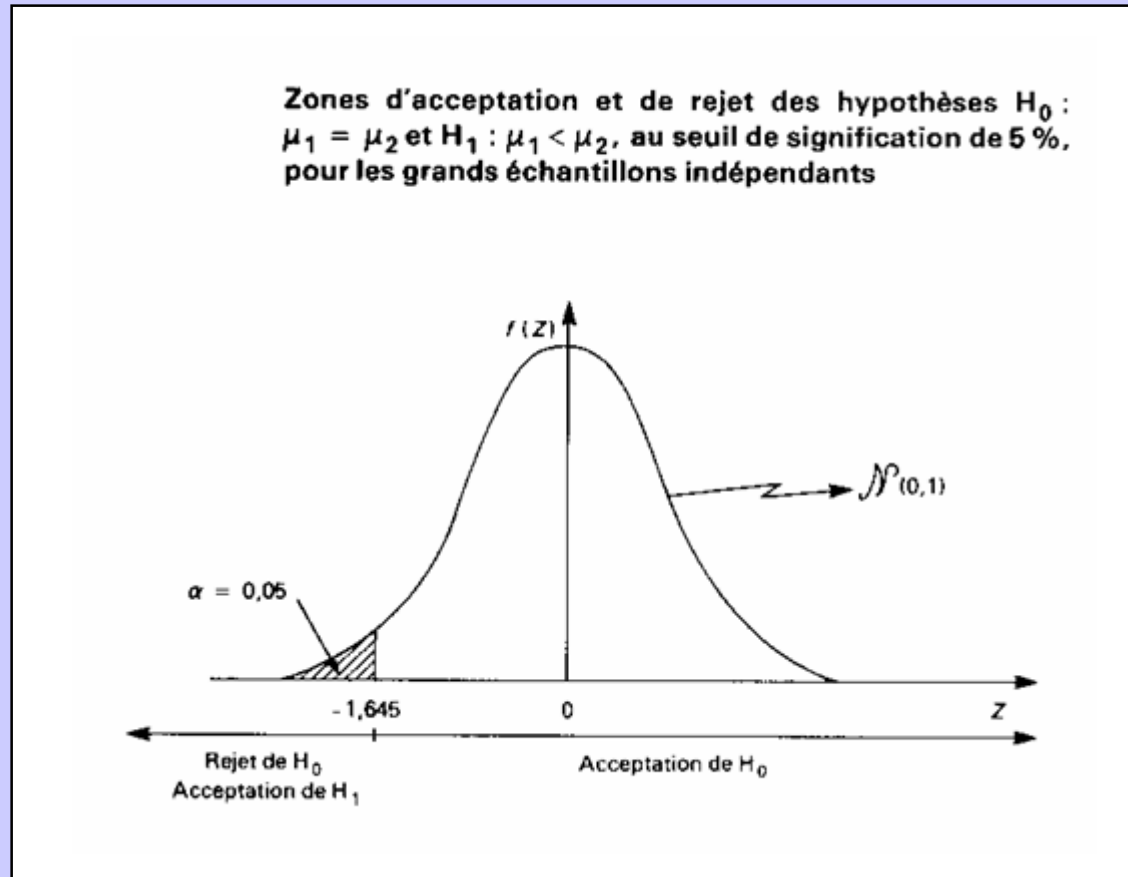
Comparaison de deux moyennes – grands échantillons -

$H_1 : \mu_1 > \mu_2$ unilatéral



Comparaison de deux moyennes – grands échantillons -

$H_1 : \mu_1 < \mu_2$ unilatéral



Pour résumer:

H_0	H_1	Rejet de H_0 si	$\alpha = 0.05$	$\alpha = 0.01$
$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$	$ Z_c \geq z_{\alpha/2} $	$ z_{\alpha/2} = 1.96$	$ z_{\alpha/2} = 2.57$
	$\mu_1 > \mu_2$	$Z_c \geq z_\alpha$	$z_\alpha = 1.64$	$z_\alpha = 2.33$
	$\mu_1 < \mu_2$	$Z_c \leq -z_\alpha$	$z_\alpha = 1.64$	$z_\alpha = 2.33$

Maintenant un exemple...

Taille des silex sur deux sites

$$n_1 = 50$$

$$\bar{x}_1 = 158,86mm$$

$$s_{x_1}^2 = 37,18mm^2$$

$$s_{x_1} = 6,09mm$$

$$n_2 = 67$$

$$\bar{x}_2 = 134,46mm$$

$$s_{x_2}^2 = 25,92mm^2$$

$$s_{x_2} = 5,09mm$$



Les moyennes de ces deux échantillons prélevés indépendamment l'un de l'autre diffèrent-elles d'une façon hautement significative?

Comparaison de deux moyennes – grands échantillons -

n_1 et n_2 grands \rightarrow test sur la **loi normale**

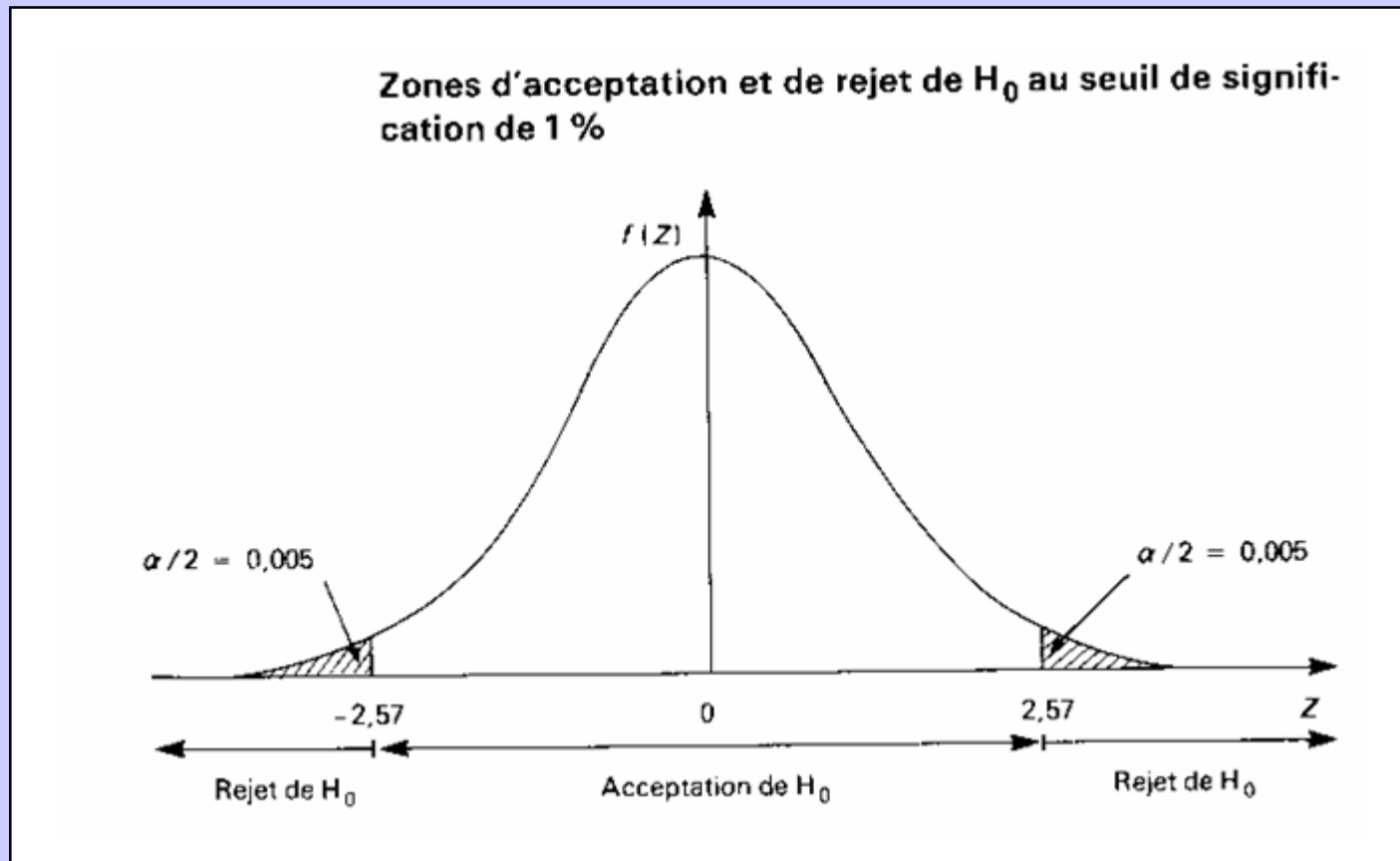
$$H_0 : \mu_a = \mu_b$$

$$H_1 : \mu_a \neq \mu_b \text{ (bilatéral)}$$

$$Z_c = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_{x1}^2}{n_1} + \frac{s_{x2}^2}{n_2}}} \quad Z_c = \frac{158.86 - 134.66}{\sqrt{\frac{37.18}{50} + \frac{25.92}{67}}} = 22.9$$

$$\alpha = 0.01, Z_{\alpha/2} = 2.57$$

Comparaison de deux moyennes – grands échantillons -



H_0 rejetée au seuil de signification de 1%

Même principe que précédemment (quand n est grand):

$$H_0: \mu = \mu_0$$

$$Z_c = \frac{\bar{x} - \mu_0}{\frac{s_x}{\sqrt{n}}}$$

que l'on teste sur la **loi normale $N(0,1)$**