

# **MAT 2777**

## **Probabilités et statistique pour ingénieur.e.s**

### **Chapitre 5**

#### **L'estimation et les intervalles de confiance**

P. Boily (uOttawa)

Hiver 2023

P.Boily (uOttawa)

## Aperçu

### 5.1 – L'inférence statistique (p.2)

- Les statistiques (p.4)
- La variance d'estimation et l'erreur-type (p.5)

### 5.2 – L'intervalle de confiance de $\mu$ lorsque $\sigma$ est connu (p.8)

- La règle du 68 – 96 – 99.7 et les intervalles de confiance (p.9)
- L'intervalle de confiance de  $\mu$  lorsque  $\sigma$  est connu (reprise) (p.14)

### 5.3 – Le choix de la taille de l'échantillon (p.26)

### 5.4 – L'intervalle de confiance de $\mu$ lorsque $\sigma$ est inconnu (p.30)

### 5.5 – L'intervalle de confiance d'une proportion (p.35)

### Annexe – Résumé (p.39)

## 5.1 – L'inférence statistique

L'un des objectifs de l'**inférence statistique** est de pouvoir tirer des conclusions sur une **population** à partir d'un échantillon aléatoire de cette population.

### Exemples :

- Peut-on évaluer la fiabilité du processus de fabrication d'un produit en sélectionnant au hasard un échantillon du produit final et en déterminant combien d'entre eux sont conformes à un certain schéma d'évaluation de la qualité ?
- Peut-on déterminer qui va remporter une élection en interrogeant un petit échantillon de répondants ?

Plus précisément, nous cherchons à estimer un **paramètre**  $\theta$  inconnu, disons, à l'aide d'une seule quantité, l'**estimé ponctuel**  $\hat{\theta}$ .

Cette estimation ponctuelle est obtenue à l'aide d'une **statistique**, une fonction d'un échantillon aléatoire. La distribution de probabilité de la statistique est sa **distribution d'échantillonnage**. Leur description est une des principales voies de recherche.

**Exemple :** Considérons un processus qui fabrique des roues dentées (dans un certain calibre). Soit  $X$  la v.a. qui enregistre le poids d'une roue dentée choisie au hasard. Quelle est la moyenne de la population  $\mu_X = E[X]$  ?

**Solution :** en l'absence d'une f.d.p.  $f(x)$ , nous pouvons estimer  $\mu_X$  à l'aide d'un échantillon aléatoire  $X_1, \dots, X_n$  de mesures du poids, *via* la statistique de la moyenne de l'échantillon :

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} \approx \mathcal{N}(\mu, \sigma^2/n) \text{ selon le TLC.}$$

## Les statistiques

Voici quelques exemples de statistiques :

- la moyenne et la médiane de l'échantillon
- la variance et l'écart-type d'échantillon
- les quantiles de l'échantillon (médiane, quartiles, percentiles)
- les statistiques de test ( $t$ ,  $\chi^2$ ,  $f$ , etc.)
- les statistiques d'ordre (le max./min. et l'étendue de l'échantillon, etc.)
- les moments de l'échantillon et leurs fonctions (l'asymétrie, l'aplatissement, etc.)

## La variance d'estimation et l'erreur-type

L'**erreur-type** d'une statistique est l'**écart-type de sa distribution d'échantillonnage**.

Par exemple, si les observations  $X_1, \dots, X_n$  proviennent d'une population de **moyenne inconnue**  $\mu$  et de **variance connue**  $\sigma^2$ , alors  $\text{Var}(\bar{X}) = \sigma^2/n$  et l'**erreur-type** de  $\bar{X}$  est

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}.$$

Si la variance de la population d'origine est **inconnue**, alors elle s'approche de la variance de l'échantillon  $S^2$  et l'erreur-type approximative de  $\bar{X}$  est

$$\hat{\sigma}_{\bar{X}} = \frac{S}{\sqrt{n}}, \quad \text{où} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

## Exemples :

1. Voici un échantillon de 20 tailles de joueurs de baseball (en pouces) :

74, 74, 72, 72, 73, 69, 69, 71, 76, 71, 73, 73, 74, 74, 69, 70, 72, 73, 75, 78.

Soit  $\bar{X}$  la moyenne d'échantillon. Alors ,

$$\bar{X} = \frac{X_1 + \cdots + X_{20}}{20} = 72.6$$

et la variance d'échantillon  $S^2$  est

$$S^2 = \frac{1}{20 - 1} \sum_{i=1}^{20} (X_i - 72.6)^2 \approx 5.6211.$$

L'erreur-type de  $\bar{X}$  est donc

$$\hat{\sigma}_{\bar{X}} = \frac{S}{\sqrt{20}} \approx \sqrt{\frac{5.6211}{20}} \approx 0.5301.$$

2. Considérons un échantillon  $\{X_1, \dots, X_{100}\}$  d'observations indépendantes prélevées d'une population normale  $\mathcal{N}(\mu, \sigma^2)$ , où  $\sigma = 50$  est connu mais  $\mu$  ne l'est pas. Quelle est la meilleure estimation de  $\mu$  ? Quelle est la distribution d'échantillonnage de cette estimation ?

**Solution :** la moyenne d'échantillon  $\bar{X} = \frac{X_1 + \dots + X_{100}}{100}$  fournit la meilleure estimation de  $\mu_X = \mu_{\bar{X}}$ .

L'erreur-type de  $\bar{X}$  est  $\sigma_{\bar{X}} = \frac{50}{\sqrt{100}} = 5$ . Étant donné que les observations sont prélevées indépendamment d'une population normale avec une moyenne de  $\mu$  et un écart-type de 50,  $\bar{X} \sim \mathcal{N}(\mu, 5^2) = \mathcal{N}(\mu, 25)$ , selon le TLC.



## 5.2 – L'intervalle de confiance de $\mu$ lorsque $\sigma$ est connu

Soit un échantillon  $\{x_1, \dots, x_n\}$  prélevé d'une population normale dont la variance  $\sigma^2$  est **connue** et dont la moyenne  $\mu$  est **inconnue**. La moyenne de l'échantillon

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

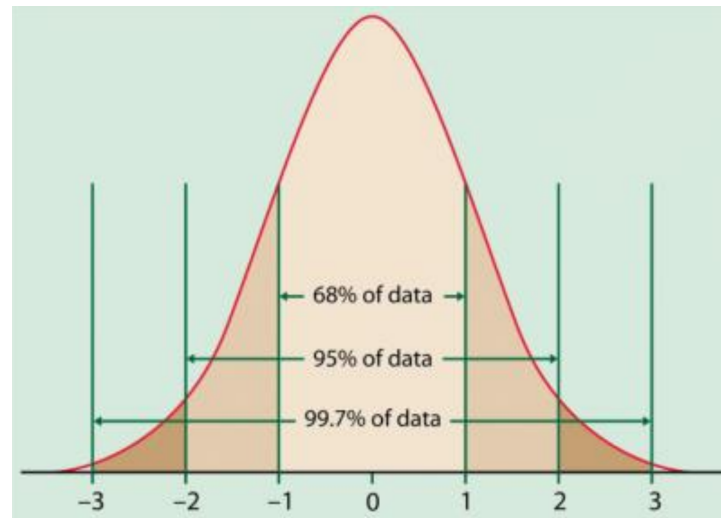
est une **estimation ponctuelle** de  $\mu$ .

Bien sûr, cette estimation n'est sans doute pas exacte, car  $\bar{x}$  est une valeur **observée** de  $\bar{X}$  ; il est peu probable que la valeur observée  $\bar{x}$  coïncide avec  $\mu$ .

Mais nous savons que  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ , et donc que

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

## La règle du 68 – 96 – 99.7 et les intervalles de confiance



$$P(-1 < Z < 1) \approx 0.683$$

$$P(-2 < Z < 2) \approx 0.955$$

$$P(-3 < Z < 3) \approx 0.997.$$

Chaque fois que nous observons une moyenne d'échantillon  $\bar{X}$  en provenance d'une population normale de moyenne  $\mu$ , nous nous attendons à ce que l'inégalité

$$-k < Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < k$$

se vérifie approximativement

$$g(k) = \begin{cases} 68.3\% \text{ du temps} & \text{si } k = 1 \\ 95.5\% \text{ du temps} & \text{si } k = 2 \\ 99.7\% \text{ du temps} & \text{si } k = 3 \end{cases}$$

De manière équivalente, l'**intervalle de confiance (symétrique)** de  $\mu$  à  $g(k)$  est

$$\bar{X} - k \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + k \frac{\sigma}{\sqrt{n}} \implies \text{IC}(\mu; g(k)) \equiv \bar{X} \pm k \frac{\sigma}{\sqrt{n}}.$$

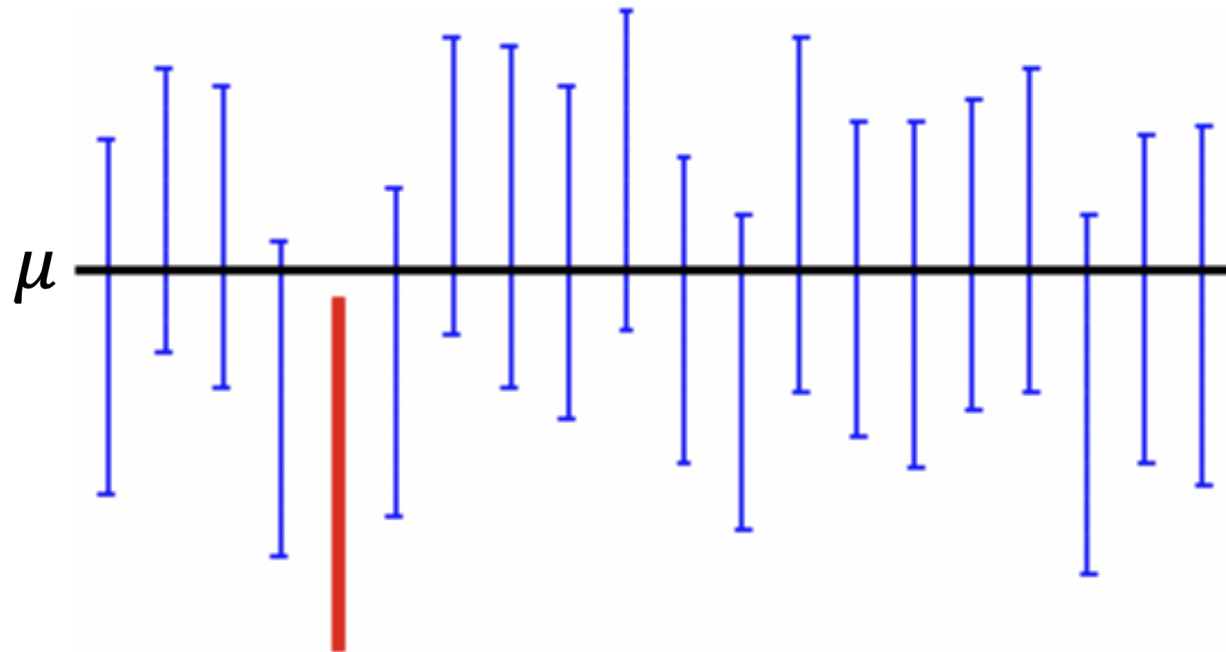
## Exemples :

1. Soit  $\{X_1, \dots, X_{64}\}$  un échantillon aléatoire prélevé d'une population normale avec un écart-type de  $\sigma = 72$  et une moyenne  $\mu$  inconnue. La moyenne de l'échantillon est  $\bar{X} = 375.2$ . Construisez un intervalle de confiance de  $\mu$  à 68.3%.

**Solution :** c'est le cas  $k = 1$ . Selon la formule,  $IC(\mu; 68.3\%)$  est

$$375.2 \pm 1 \cdot \frac{72}{\sqrt{64}} \implies IC(\mu; 68.3\%) \equiv (366.2, 384.2).$$

**TRÈS IMPORTANT:** ceci ne dit pas que nous sommes certains à 68.3% que  $\mu$  se situe entre 366.2 et 384.2 – lorsqu'un échantillon de taille 64 est prélevé dans une population normale  $\mathcal{N}(\mu, 72^2)$  et qu'on construit l'intervalle  $IC(\mu; 68.3\%)$ ,  $\mu$  se retrouve entre les extrémités de l'intervalle environ 68.3% du temps.



Dans un IC à 95%, nous nous attendons à ce que 19 échantillons sur 20, prélevés d'une population unique, produisent des intervalles de confiance qui contiennent le paramètre de population d'intérêt, en moyenne.

2. Construisez  $IC(\mu; 95.5\%)$  avec les données du problème précédent.

**Solution :** c'est le cas  $k = 2$ , alors

$$375.2 \pm 2 \cdot \frac{72}{\sqrt{64}} \implies IC(\mu; 95.5\%) \equiv (357.2, 393.2).$$

3. Construisez  $IC(\mu; 99.7\%)$  avec les données du problème précédent.

**Solution :** c'est le cas  $k = 3$ , alors

$$375.2 \pm 3 \cdot \frac{72}{\sqrt{64}} \implies IC(\mu; 99.7\%) \equiv (348.2, 402.2).$$

## L'IC de $\mu$ lorsque $\sigma$ est connu (reprise)

Une autre approche de construction de l'IC consiste à spécifier la proportion d'intérêt de l'aire sous la f.d.p.  $\phi(z)$ , puis à déterminer les valeurs critiques correspondantes (les extrémités de l'intervalle).

Soit  $\{X_1, \dots, X_n\}$  un échantillon prélevé de  $\mathcal{N}(\mu, \sigma^2)$  ; on rappelle que  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1)$ .

Pour un IC à 95%, disons, on doit trouver la **valeur critique**  $z^* > 0$  telle que  $P(-z^* < Z < z^*) = 0.95$ .

Mais le côté gauche peut être ré-écrit comme suit :

$$\begin{aligned} P(-z^* < Z < z^*) &= \Phi(z^*) - \Phi(-z^*) \\ &= \Phi(z^*) - (1 - \Phi(z^*)) = 2\Phi(z^*) - 1. \end{aligned}$$

Nous cherchons donc un  $z^* > 0$  tel que

$$0.95 = 2\Phi(z^*) - 1 \implies \Phi(z^*) = \frac{0.95 + 1}{2} = 0.975.$$

D'après la table de la f.d.p.  $\phi(z)$ , nous voyons que  $\Phi(1.96) \approx 0.9750$ , d'où

$$P(-1.96 < Z < 1.96) = P\left(-1.96 < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < 1.96\right) \approx 0.95.$$

Autrement dit, l'inégalité

$$-1.96 < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < 1.96$$

est valide avec une probabilité de 0.95 (sous l'interprétation préalable).



De façon équivalente,

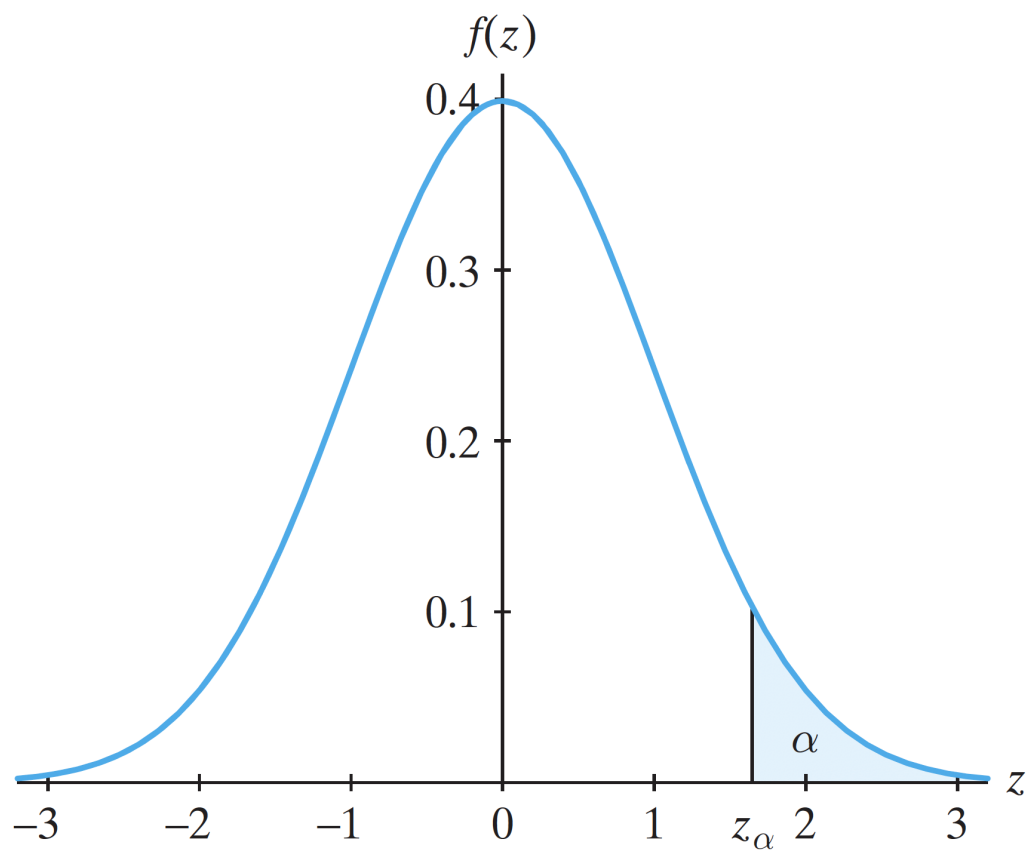
$$\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}} \implies \text{IC}(\mu; 95\%) \equiv \bar{X} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

forme un **intervalle de confiance de  $\mu$  à 95% lorsque  $\sigma$  est connu.**

On peut aussi montrer que

$$\bar{X} - 2.575 \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + 2.575 \frac{\sigma}{\sqrt{n}} \implies \text{IC}(\mu; 99\%) \equiv \bar{X} \pm 2.575 \frac{\sigma}{\sqrt{n}}$$

forme un **intervalle de confiance de  $\mu$  à 99% lorsque  $\sigma$  est connu.**



$$P(Z > z_\alpha) = \alpha$$

$$P(Z > z) = 1 - \Phi(z) = \Phi(-z)$$

Le **niveau de confiance**  $1 - \alpha$  est généralement exprimé en termes de **petits**  $\alpha$ , par exemple,  $\alpha = 0.05 \implies 1 - \alpha = 0.95$ .

Pour  $\alpha = 0.01, 0.02, \dots, 0.98, 0.99$ , les  $z_\alpha$  correspondants sont appelés les **pourcentiles** (ou tout simplement les **centiles**) de la loi normale centrée réduite. En général,

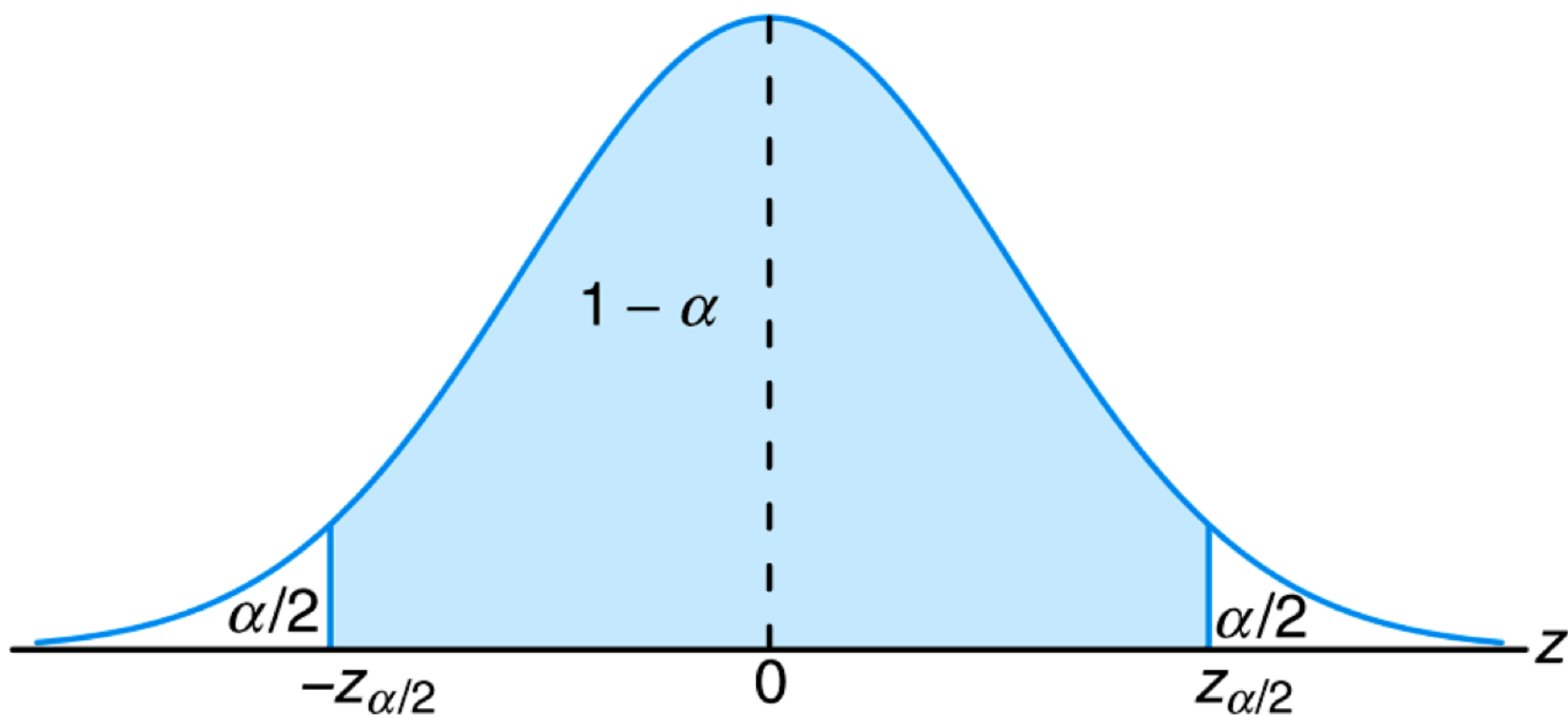
$$P(Z > z_\alpha) = \alpha \implies z_\alpha \text{ est le } 100(1 - \alpha) \text{ centile.}$$

Pour les intervalles de confiance **symétriques (à 2 côtés)**, on trouve les valeurs appropriées en résolvant  $P(|Z| > z^*) = \alpha$  pour  $z^*$ . Par les propriétés de  $\mathcal{N}(0, 1)$ ,

$$\alpha = P(|Z| > z^*) = 1 - P(-z^* < Z < z^*) = 1 - (2\Phi(z^*) - 1) = 2(1 - \Phi(z^*)),$$

de sorte que

$$\Phi(z^*) = 1 - \alpha/2 \implies z^* = z_{\alpha/2}.$$



Par exemple,

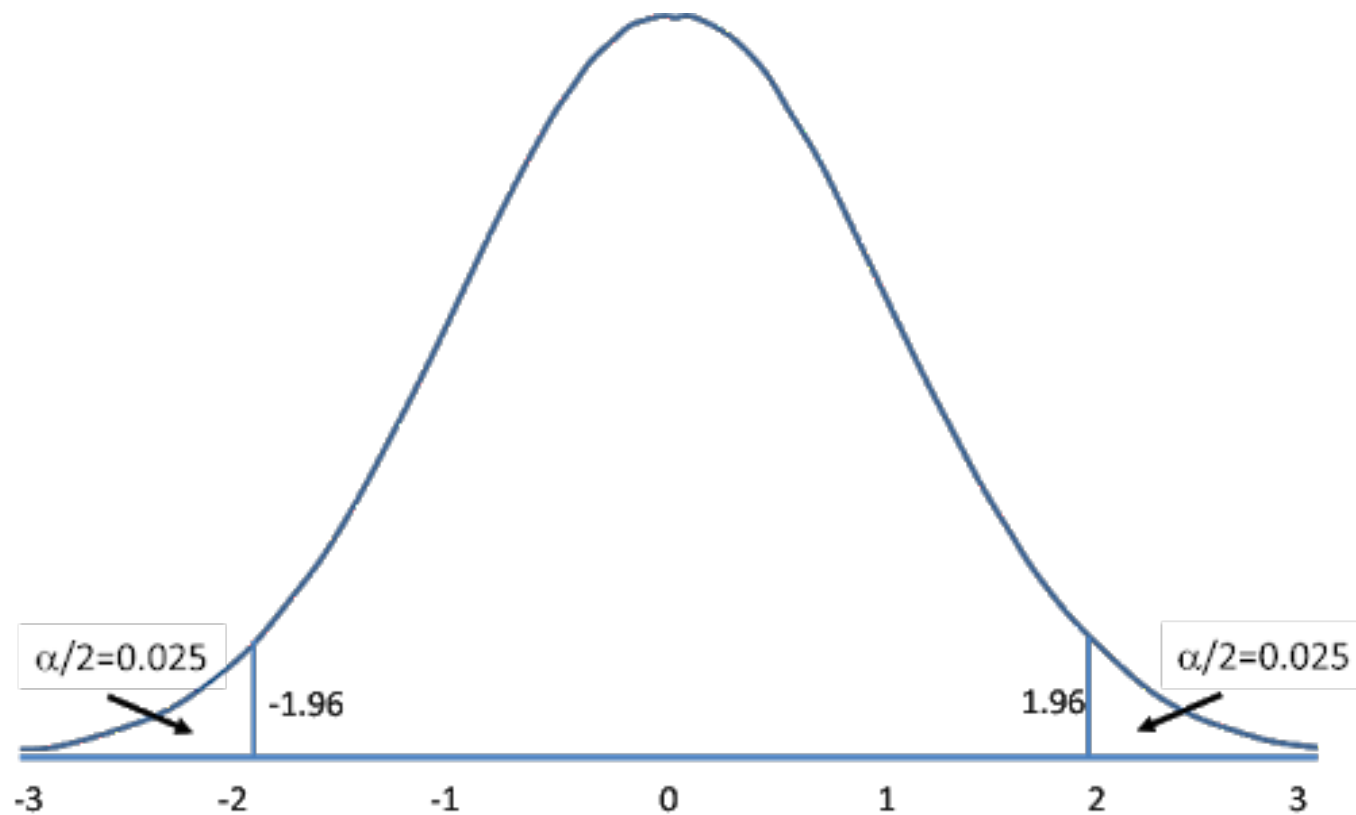
$$P(|Z| > z_{0.025}) = 0.05 \implies z_{0.025} = 1.96$$

$$P(|Z| > z_{0.005}) = 0.01 \implies z_{0.005} = 2.575.$$

Dans ce même contexte ( $X_i \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\sigma$  connu), l'intervalle de confiance de  $\mu$  pour un  $\alpha$  donné prend généralement la forme

$$\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \implies \text{IC}(\mu; 100(1 - \alpha)\%) \equiv \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Il importe alors de savoir calculer les centiles  $z_{\alpha/2}$  (soit à l'aide de tables, soit à l'aide d'un logiciel).



Pour un niveau de confiance  $\alpha$  donné, les intervalles de confiance **plus étroits** sont primés par rapport à l'estimation de la moyenne :

- les estimations deviennent “**meilleures**” lorsque la taille  $n$  de l'échantillon **augmente** ;
- les estimations deviennent “**meilleures**” lorsque  $\sigma$  diminue.

Si  $\alpha_1 > \alpha_2$ , l'intervalle de confiance de  $\mu$  à  $100(1 - \alpha_1)\%$  est plus étroit que l'intervalle de confiance de  $\mu$  à  $100(1 - \alpha_2)\%$ .

En particulier,

$$IC(\mu; 95\%) \subseteq IC(\mu; 99\%).$$

Si l'échantillon provient d'une population normale, alors l'intervalle de confiance est **exact**. Autrement, nous pouvons utiliser le TLC afin d'obtenir un intervalle de confiance **approximatif**, lorsque  $n$  est suffisamment élevé.

**Exemples :**

1. Un échantillon de  $n = 9$  observations provenant d'une population normale ayant un écart-type  $\sigma = 5$  connu a une moyenne d'échantillon  $\bar{X} = 19.93$ . Donnez  $IC(\mu; 0.95)$  et  $IC(\mu; 0.99)$  sur la base de cet échantillon.

**Solution :** l'estimation ponctuelle de  $\mu$  est  $\bar{X} = 19.93$ , et

$$IC(\mu; 100(1 - \alpha)\%) \equiv \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Ainsi,

$$95\% : \bar{X} \pm z_{0.025} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 1.96 \cdot \frac{5}{\sqrt{9}} \Rightarrow (16.66, 23.20)$$

$$99\% : \bar{X} \pm z_{0.005} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 2.575 \cdot \frac{5}{\sqrt{9}} \Rightarrow (15.64, 24.22)$$



2. Un échantillon de 25 observations provenant d'une population normale ayant un écart-type  $\sigma = 5$  connu a une moyenne d'échantillon  $\bar{X} = 19.93$ . Donnez  $IC(\mu; 0.95)$  et  $IC(\mu; 0.99)$  sur la base de cet échantillon.

**Solution :** l'estimation ponctuelle de  $\mu$  est toujours  $\bar{X} = 19.93$ ; les  $IC(\mu; 100(1 - \alpha)\%)$  recherchés sont

$$95\% : \bar{X} \pm z_{0.025} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 1.96 \cdot \frac{5}{\sqrt{25}} \Rightarrow (17.97, 21.89)$$

$$99\% : \bar{X} \pm z_{0.005} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 2.575 \cdot \frac{5}{\sqrt{25}} \Rightarrow (17.35, 22.51)$$

3. Un échantillon de 25 observations provenant d'une population normale ayant un écart-type  $\sigma = 10$  connu a une moyenne d'échantillon  $\bar{X} = 19.93$ . Donnez  $IC(\mu; 0.95)$  et  $IC(\mu; 0.99)$  sur la base de cet échantillon.

**Solution :** l'estimation ponctuelle de  $\mu$  est toujours  $\bar{X} = 19.93$ ; les  $IC(\mu; 100(1 - \alpha)\%)$  recherchés sont

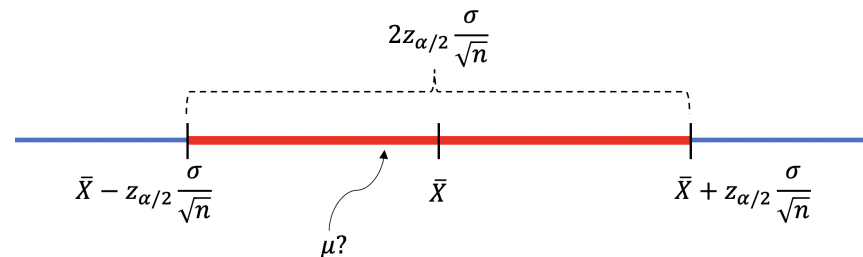
$$95\% : \bar{X} \pm z_{0.025} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 1.96 \cdot \frac{10}{\sqrt{25}} \Rightarrow (16.01, 23.85)$$

$$99\% : \bar{X} \pm z_{0.005} \frac{\sigma}{\sqrt{n}} \Rightarrow 19.93 \pm 2.575 \cdot \frac{10}{\sqrt{25}} \Rightarrow (14.78, 25.08)$$

Notez comment les intervalles de confiance sont affectés par  $\alpha$ ,  $n$ , et  $\sigma$ .

## 5.3 – Le choix de la taille de l'échantillon

Lorsque l'on prélève un échantillon de taille  $n$  d'une population normale dont l'écart-type  $\sigma$  est connu, l'erreur que l'on commet en estimant la moyenne  $\mu$  à l'aide de la moyenne empirique  $\bar{X}$  est en général bornée par  $E = z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ , à un niveau de confiance de  $100(1 - \alpha)\%$ .



Afin de contrôler l'erreur, il faut contrôler la **taille de l'échantillon** :

$$E > \frac{z_{\alpha/2}\sigma}{\sqrt{n}} \implies n > \left( \frac{z_{\alpha/2}\sigma}{E} \right)^2 = \frac{z_{\alpha/2}^2 \sigma^2}{E^2}.$$

## Exemples :

1. Un échantillon  $\{X_1, \dots, X_n\}$  est prélevé d'une population normale dont l'écart-type  $\sigma = 100$  est connu. Quelle taille d'échantillon  $n$  est requise afin de s'assurer que l'erreur sur l'estimation est au plus  $E = 10$ , à un niveau de confiance  $\alpha = 0.05$ ?

**Solution :** tant que

$$n > \left( \frac{z_{\alpha/2}\sigma}{E} \right)^2 = \left( \frac{z_{0.025} \cdot 100}{10} \right)^2 = (19.6)^2 = 384.16,$$

alors l'erreur d'estimation commise en utilisant  $\bar{X}$  pour approcher  $\mu$  sera d'au plus  $E = 10$ , avec une probabilité de 95%.

2. On répète le premier exemple, mais avec  $\sigma = 10$ .

**Solution :** on doit avoir

$$n > \left( \frac{z_{\alpha/2}\sigma}{E} \right)^2 = \left( \frac{z_{0.025} \cdot 10}{10} \right)^2 = (1.96)^2 = 3.8416.$$

3. On répète le premier exemple, mais avec  $E = 1$ .

**Solution :** on doit avoir

$$n > \left( \frac{z_{\alpha/2}\sigma}{E} \right)^2 = \left( \frac{z_{0.025} \cdot 100}{1} \right)^2 = (196)^2 = 38416.$$

4. On répète le premier exemple, mais avec  $\alpha = 0.01$ .

**Solution :** on doit avoir

$$n > \left( \frac{z_{\alpha/2\sigma}}{E} \right)^2 = \left( \frac{z_{0.005} \cdot 100}{10} \right)^2 = (25.75)^2 = 663.0625.$$

5. On répète le premier exemple, mais avec  $\sigma = 10$ ,  $E = 1$ , et  $\alpha = 0.01$ .

**Solution :** on doit avoir

$$n > \left( \frac{z_{\alpha/2\sigma}}{E} \right)^2 = \left( \frac{z_{0.005} \cdot 10}{1} \right)^2 = (25.75)^2 = 663.0625.$$

La relation entre  $\alpha$ ,  $\sigma$ ,  $E$ , et  $n$  est simple, mais pas toujours intuitive !

## 5.4 – L'intervalle de confiance de $\mu$ lorsque $\sigma$ est inconnu

Jusqu'à présent, nous nous sommes trouvés dans la situation heureuse d'échantillonner à partir d'une population dont la variance  $\sigma^2$  est connue.

Que faire lorsque cette dernière est **inconnue** ?

Nous donnons un estimé de  $\sigma$  par l'entremise de la **variance d'échantillon**

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

(rappelez-vous que la moyenne de la population  $\mu$  est également inconnue... c'est ce que nous cherchons !) et l'**écart-type empirique**  $S = \sqrt{S^2}$ .

Si l'on connaît  $\sigma$ , nous savons grâce au TLC que  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  est approximativement  $\mathcal{N}(0, 1)$ .

Si  $\sigma$  est inconnu, on peut montrer que  $\frac{\bar{X} - \mu}{S/\sqrt{n}}$  suit approximativement une **loi  $t$  de Student** avec  $n - 1$  **degrés de liberté**,  $t(n - 1)$ .

Par conséquent, pour un niveau de confiance  $\alpha$ ,

$$P\left(-t_{\alpha/2}(n-1) < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{\alpha/2}(n-1)\right) \approx 1 - \alpha,$$

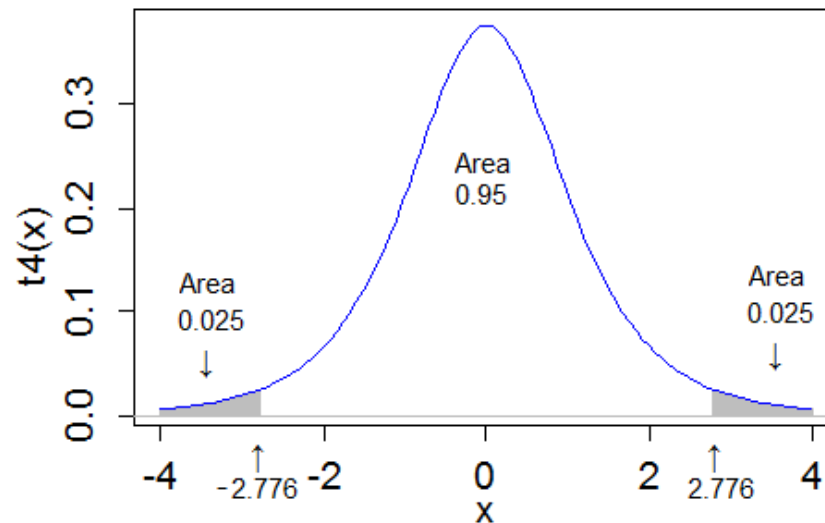
où  $t_{\alpha/2}(n-1)$  est le  $100(1 - \alpha/2)^{\text{e}}$  centile de  $t(n-1)$  (on peut les lire dans les tables). L'égalité est atteinte si la population sous-jacente est **normale**:

$$\text{IC}_S(\mu; 100(1 - \alpha)\%) \equiv \bar{X} \pm t_{\alpha/2}(n-1) \cdot \frac{S}{\sqrt{n}}.$$



Par exemple, si  $\alpha = 0.05$  et  $\{X_1, X_2, X_3, X_4, X_5\}$  est un échantillon prélevé d'une normale dont la variance  $\sigma^2$  est inconnue, alors

$$t_{0.025}(5 - 1) = 2.776 \quad \text{et} \quad P\left(-2.776 < \frac{\bar{X} - \mu}{S/\sqrt{5}} < 2.776\right) = 0.95.$$



## Exemples :

1. Pour une année donnée, on obtient  $n = 9$  mesures de la concentration d'ozone :

3.5 5.1 6.6 6.0 4.2 4.4 5.3 5.6 4.4

En supposant que les concentrations d'ozone mesurées suivent une loi normale avec  $\sigma^2 = 1.21$ , construisez  $IC(\mu; 95\%)$ . Notez que  $\bar{X} = 5.01$  et que  $S = 0.97$ .

**Solution:** puisque nous connaissons la variance, nous devons utiliser la valeur critique  $z_{\alpha/2} = z_{0.025} = 1.96$ :

$$IC(\mu; 95\%) \equiv \bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = 5.01 \pm 1.96 \cdot \frac{\sqrt{1.21}}{\sqrt{9}} \equiv (4.29, 5.73).$$

2. On aborde le même problème, mais en supposant cette fois que la variance de la population sous-jacente est inconnue.

**Solution:** Comme nous ne connaissons pas la variance, nous devons utiliser valeur critique  $t_{\alpha/2}(n-1) = t_{0.025}(8) = 2.306$  (assurez-vous de comprendre comment obtenir cette valeur à partir du tableau) :

$$\text{IC}(\mu; 95\%) \equiv \bar{X} \pm t_{\alpha/2}(n-1) \cdot \frac{S}{\sqrt{n}} = 5.01 \pm 2.306 \cdot \frac{0.97}{\sqrt{9}} \equiv (4.26, 5.76).$$

Lorsque nous connaissons la variance, l'IC est **plus étroit (plus petit)**, ce qui est naturel puisque nous sommes plus confiants lorsque nous avons plus d'informations.

## 5.5 – L'intervalle de confiance d'une proportion

Si  $X \sim \mathcal{B}(n, p)$  (le nombre de réussites dans  $n$  épreuves de Bernouilli), alors l'estimateur ponctuel de  $p$  est  $\hat{P} = \frac{X}{n}$ .

Rappelons que  $E[X] = np$  et  $\text{Var}[X] = np(1 - p)$ .

Nous pouvons normaliser toute variable aléatoire, normale ou non :

$$Z = \frac{X - \mu}{\sigma} = \frac{n\hat{P} - np}{\sqrt{np(1 - p)}} = \frac{\hat{P} - p}{\sqrt{\frac{p(1-p)}{n}}}$$

est approximativement  $\mathcal{N}(0, 1)$  si  $n$  est suffisamment élevé.

Ainsi,

$$P \left( -z_{\alpha/2} < \frac{\hat{P} - p}{\sqrt{\frac{p(1-p)}{n}}} < z_{\alpha/2} \right) \approx 1 - \alpha.$$

En utilisant l'approche précédente, on construit un IC( $p$ ;  $100(1 - \alpha)\%$ ) **approximatif** :

$$\hat{P} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} < p < \hat{P} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}},$$

mais ce n'est pas bien utile  $p$  est inconnu ! Au lieu, nous utilisons :

$$\hat{P} - z_{\alpha/2} \sqrt{\frac{\hat{P}(1 - \hat{P})}{n}} < p < \hat{P} + z_{\alpha/2} \sqrt{\frac{\hat{P}(1 - \hat{P})}{n}}.$$

## Exemples :

1. Deux candidats ( $A$  et  $B$ ) se présentent aux élections. Dans un sondage, 1000 électeurs sont choisis au hasard : 52% soutiennent  $A$ , tandis que 48% soutiennent  $B$ . Donnez un  $IC(p; 0.95)$  pour le soutien de chaque candidat.

**Solution :** on utilise  $\alpha = 0.05$  et  $\hat{P} = 0.52$ . L'intervalle de confiance recherché du candidat  $A$  est ainsi

$$\hat{P} \pm z_{\alpha/2} \sqrt{\frac{\hat{P}(1 - \hat{P})}{n}} = 0.52 \pm 1.96 \sqrt{\frac{0.52 \cdot 0.48}{1000}} \approx 0.52 \pm 0.031.$$

L'intervalle de confiance pour  $B$  a la même largeur:  $0.48 \pm 0.031$ .

2. Sur la base de ce résultat de sondage, un quotidien publie la une suivante :  
“La candidate  $A$  devance le candidat  $B$  !”. Ce titre est-il justifié ?

**Solution** : bien qu'il y ait un écart de 4 point de pourcentage dans les sondages, le soutien réel pour la candidate  $A$  se retrouve dans l'intervalle de confiance

$(48.9\%, 55.1\%)$ ,

avec une probabilité de 95%. De même, le soutien réel pour le candidat  $B$  se situe dans l'intervalle

$(44.9\%, 51.1\%)$ ,

avec la même probabilité.

Puisque les intervalles de confiance se **chevauchent**, il est probable que la une **n'est pas justifiée**.

## Annexe – Résumé

**Échantillon** :  $\{X_1, \dots, X_n\}$ . **Objectif** : prédire  $\mu$  avec confiance  $\alpha$ .

- Si la population est **normale**, de variance  $\sigma^2$  **connue**,  $IC(\mu; 100(1-\alpha)\%)$  est **exact** et

$$IC(\mu; 100(1-\alpha)\%) \equiv \bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}.$$

- Si la population est **non-normale**, de variance  $\sigma^2$  **connue**, et  $n$  est suffisamment élevé, alors  $IC(\mu; 100(1-\alpha)\%)$  est **approximatif** et

$$IC(\mu; 100(1-\alpha)\%) \equiv \bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}.$$



- Si la population est **normale**, de variance **inconnue**,  $IC_S(\mu; 100(1-\alpha)\%)$  est **exact** et

$$IC_S(\mu; 100(1-\alpha)\%) \equiv \bar{X} \pm t_{\alpha/2}(n-1) \cdot \frac{S}{\sqrt{n}}.$$

- Si la population est **non-normale**, de variance **inconnue**, et  $n$  est suffisamment élevé, alors  $IC_S(\mu; 100(1-\alpha)\%)$  est **approximatif** et

$$IC_S(\mu; 100(1-\alpha)\%) \equiv \bar{X} \pm z_{\alpha/2} \cdot \frac{S}{\sqrt{n}}.$$

- Si la variance de la population est **inconnue** et si  $n$  est '**trop petit**', il n'y a rien à faire...