

Chapitre 1

Notations et notions utilisées

1.1 Bases de la Statistique

1.1.1 Définitions

- **Population et individus** : Une population est un ensemble fini d'objets sur lesquels une étude se porte. Ces objets sont appelés individus.
- **Caractère** : Un caractère est une qualité que l'on étudie chez des individus.
- **Échantillon** : Un échantillon est un ensemble d'individus issus d'une population.
- **Données** : Les données en notre possession sont les observations d'un caractère sur les individus d'un échantillon.
- **Estimation paramétrique** : L'enjeu de l'estimation paramétrique est d'évaluer/estimer avec précision un paramètre inconnu émanant d'un caractère à partir des données.
- **Moyenne et écart-type corrigé** : La moyenne et l'écart-type corrigé des données sont les principales mesures statistiques intervenant en estimation paramétrique.

En notant X un caractère numérique, n le nombre d'individus d'un échantillon et x_1, \dots, x_n les données associées, on définit :

- La moyenne de x_1, \dots, x_n :

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

C'est une estimation ponctuelle de la valeur moyenne de X .

- L'écart-type corrigé de x_1, \dots, x_n :

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

C'est une estimation ponctuelle de la variabilité de X autour de sa moyenne. La valeur obtenue a la même unité que X .

1.1.2 Exemple

- **Contexte :**

Population	Ensemble des kiwis d'une ferme												
Individu	Kiwi												
Caractère	Poids d'un kiwi (en grammes)												
Paramètre inconnu	Poids moyen d'un kiwi												
Échantillon	6 kiwis choisis au hasard ($n = 6$)												
Données	<table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td>x_1</td> <td>x_2</td> <td>x_3</td> <td>x_4</td> <td>x_5</td> <td>x_6</td> </tr> <tr> <td>72</td> <td>75</td> <td>78</td> <td>71</td> <td>71</td> <td>75</td> </tr> </table> <p>(par exemple, x_1 est le poids du premier kiwi de l'échantillon, soit 72 grammes)</p>	x_1	x_2	x_3	x_4	x_5	x_6	72	75	78	71	71	75
x_1	x_2	x_3	x_4	x_5	x_6								
72	75	78	71	71	75								
Objectif	Évaluer le poids moyen inconnu d'un kiwi à l'aide des données x_1, \dots, x_6												

- **Mesures statistiques :**

Moyenne	$\bar{x} = \frac{1}{6} \sum_{i=1}^6 x_i = 73.6666$
Écart-type corrigé	$s = \sqrt{\frac{1}{6-1} \sum_{i=1}^6 (x_i - \bar{x})^2} = 2.8047$

1.2 Éléments de probabilités

1.2.1 Notations

- $\mathbb{P}(A)$ = probabilité que A se réalise.
- var = variable aléatoire réelle.
- $X \sim \dots$ = la var X suit la loi \dots
- $\mathbb{E}(X)$ = espérance de X .
- $\mathbb{V}(X)$ = variance de X .
- $\sigma(X) = \sqrt{\mathbb{V}(X)}$ = écart-type de X .
- $\mathcal{N}(\mu, \sigma^2)$ = loi normale de moyenne μ et de variance σ^2 (écart-type σ).
- $\mathcal{B}(p)$ = loi de Bernoulli de paramètre p .
- $\mathcal{T}(\nu)$ = loi de Student à ν degrés de liberté.
Si $\nu \geq 31$, on a l'approximation $T(\nu) \approx \mathcal{N}(0, 1)$.
- $\chi^2(\nu)$ = loi du Chi-deux à ν degrés de liberté.
- $\mathcal{F}(\nu_1, \nu_2)$ = loi de Fisher à (ν_1, ν_2) degrés de liberté.

1.2.2 Modélisations

- **Loi normale** : On adopte la modélisation

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

quand X représente une grandeur sujette à une somme d'erreurs mineures indépendantes.

Par exemple, X peut être : poids, taille, temps, distance, masse, vitesse, température, indice, score, salaire, note, quantité ou teneur.

Dans ce cas, μ est la moyenne de X et σ^2 (ou σ) mesure la variabilité/dispersion de X autour de μ .

- **Loi de Bernoulli** : On adopte la modélisation

$$X \sim \mathcal{B}(p)$$

quand X prend deux valeurs : 0 ou 1, correspondant souvent à un codage binaire.

Par exemple, $X = 1$ peut caractériser :

- le succès à une épreuve,
- la présence d'un élément caractéristique.

Le paramètre p est la probabilité que $X = 1$ se réalise, laquelle peut aussi s'interpréter en terme de proportion d'individus dans la population vérifiant $X = 1$.

- **Exemple** :

Population	Ensemble des plaquettes de beurre d'une production
Individu	Plaquette de beurre
Caractère 1	$X =$ Poids d'une plaquette (en grammes)
Modélisation	$X \sim \mathcal{N}(\mu, \sigma^2)$
Paramètres : μ et σ^2	$\mu =$ Poids moyen d'une plaquette σ^2 mesure la dispersion du poids d'une plaquette autour de μ
Caractère 2	$Y = 1$ si la plaquette présente un défaut de conditionnement et $Y = 0$ sinon
Modélisation	$Y \sim \mathcal{B}(p)$
Paramètre p	$p =$ Proportion de plaquettes présentant un défaut de conditionnement

Chapitre 2

Intervalles de confiance

2.1 Méthodes

2.1.1 Bases des intervalles de confiance

- **Objectif** : On veut évaluer avec précision un paramètre inconnu émanant d'un caractère à partir des données. On note θ ce paramètre.
- **Intervalle de confiance ; première approche** : Un intervalle de confiance pour θ est une "fourchette de valeurs", déterminée à partir des données, qui a de fortes chances de contenir θ .
- **Niveau** : Le niveau est le pourcentage de chances que l'intervalle de confiance contienne θ .
On veut que ce niveau soit aussi élevé que possible.
Il s'écrit sous la forme : $100(1 - \alpha)\%$, avec $\alpha \in]0, 1[$ (par exemple, 95%, soit $\alpha = 0.05$, ou 99%, soit $\alpha = 0.01$).
- **Intervalle de confiance ; bilan** : Un intervalle de confiance pour θ au niveau $100(1 - \alpha)\%$, $\alpha \in]0, 1[$, est un intervalle :

$$i_\theta = [a, b],$$

où a et b sont des réels calculés à partir des données de sorte qu'il y ait $100(1 - \alpha)\%$ de chances que i_θ contienne θ .

2.1.2 Intervalles de confiance pour une moyenne : Z-IntConf, T-IntConf et Z-IntConf-Lim

Contexte

On étudie un caractère représenté par une $\text{var } X \sim \mathcal{N}(\mu, \sigma^2)$.

Les données sont constituées de la valeur de X pour chacun des n individus d'un échantillon. Ces valeurs sont notées x_1, \dots, x_n .

Le paramètre μ , représentant la moyenne de X , est inconnu.

Objectif

On veut déterminer un intervalle de confiance pour μ au niveau $100(1 - \alpha)\%$, $\alpha \in]0, 1[$.

On distingue alors 2 cas :

- le cas où σ est connu,
- le cas où σ est inconnu.

La connaissance ou non de σ va impacter sur la construction théorique de l'intervalle de confiance.

Cas où σ est connu : Z-IntConf

Dans le cas où σ est connu, on utilise le Z-IntConf.

Sa construction repose sur la loi normale centrée réduite $\mathcal{N}(0, 1)$ (symbolisée par le "Z" dans "Z-IntConf").

Pour obtenir le Z-IntConf, il faut calculer :

- \bar{x} : la moyenne de x_1, \dots, x_n ,
- le réel z_α vérifiant :

$$\mathbb{P}(|Z| \geq z_\alpha) = \alpha,$$

où $Z \sim \mathcal{N}(0, 1)$.

Ce réel est évaluable dans la table 1 page 223.

Le Z-IntConf est alors donné par :

$$i_\mu = \left[\bar{x} - z_\alpha \frac{\sigma}{\sqrt{n}}, \bar{x} + z_\alpha \frac{\sigma}{\sqrt{n}} \right].$$

Il y a donc $100(1 - \alpha)\%$ de chances que i_μ contienne μ .

Cas où σ est inconnu : T-IntConf

Dans la pratique, le σ est souvent inconnu. On utilise alors le T-IntConf.

Sa construction repose sur la loi de Student $\mathcal{T}(\nu)$ (symbolisée par le "T" dans "T-IntConf").

Pour obtenir le T-IntConf, il faut calculer :

- \bar{x} : la moyenne de x_1, \dots, x_n ,
- s : l'écart-type corrigé de x_1, \dots, x_n ,
- le réel $t_\alpha(\nu)$ vérifiant :

$$\mathbb{P}(|T| \geq t_\alpha(\nu)) = \alpha,$$

où $T \sim \mathcal{T}(\nu)$, $\nu = n - 1$.

Ce réel est évaluable dans la table 3 page 225.

Le T-IntConf est alors donné par :

$$i_\mu = \left[\bar{x} - t_\alpha(\nu) \frac{s}{\sqrt{n}}, \bar{x} + t_\alpha(\nu) \frac{s}{\sqrt{n}} \right].$$

Il y a donc $100(1 - \alpha)\%$ de chances que i_μ contienne μ .

Approximation : Lorsque $\nu \geq 31$, on peut utiliser l'approximation $\mathcal{T}(\nu) \approx \mathcal{N}(0, 1)$ qui entraîne

$$t_\alpha(\nu) \simeq z_\alpha,$$

où z_α est le réel vérifiant :

$$\mathbb{P}(|Z| \geq z_\alpha) = \alpha,$$

$Z \sim \mathcal{N}(0, 1)$ (celui-ci est évaluable dans la table 1 page 223).

Complément : Z-IntConf-Lim

Dans le cas où $n \geq 1000$, l'hypothèse : " $X \sim \mathcal{N}(\mu, \sigma^2)$ " peut être omise.

On représente alors toujours le caractère par une *var* X et sa moyenne est donnée par $\mathbb{E}(X)$. Elle est inconnue.

On veut déterminer un intervalle de confiance pour $\mathbb{E}(X)$ au niveau $100(1 - \alpha)\%$, $\alpha \in]0, 1[$.

Dans ce cas, on utilise le Z-IntConf-Lim.

Sa construction repose sur la loi normale centrée réduite $\mathcal{N}(0, 1)$, laquelle est obtenue en tant que loi limite (d'où le "Z" et le "Lim").

Pour obtenir le Z-IntConf-Lim, il faut calculer :

- \bar{x} : la moyenne de x_1, \dots, x_n ,
- s : l'écart-type corrigé de x_1, \dots, x_n ,
- le réel z_α vérifiant :

$$\mathbb{P}(|Z| \geq z_\alpha) = \alpha,$$

où $Z \sim \mathcal{N}(0, 1)$.

Ce réel est évaluable dans la table 1 page 223.

Le Z-IntConf-Lim est alors donné par :

$$i_{\mathbb{E}(X)} = \left[\bar{x} - z_\alpha \frac{s}{\sqrt{n}}, \bar{x} + z_\alpha \frac{s}{\sqrt{n}} \right].$$

Il y a donc $100(1 - \alpha)\%$ de chances que $i_{\mathbb{E}(X)}$ contienne $\mathbb{E}(X)$.

2.1.3 Intervalles de confiance pour la variance : Chi2-IntConf

Contexte

On étudie un caractère représenté par une $\text{var } X \sim \mathcal{N}(\mu, \sigma^2)$.

Les données sont constituées de la valeur de X pour chacun des n individus d'un échantillon. Ces valeurs sont notées x_1, \dots, x_n .

Les paramètres μ et σ sont inconnus. On rappelle que σ^2 (ou σ) représente la dispersion de X autour de μ .

Objectif

On veut déterminer un intervalle de confiance pour σ^2 au niveau $100(1 - \alpha)\%$, $\alpha \in]0, 1[$.

Chi2-IntConf

On utilise un Chi2-IntConf.

Sa construction repose sur la loi du Chi-deux $\chi^2(\nu)$.

Pour obtenir le Chi2-IntConf, il faut calculer :

- s : l'écart-type corrigé de x_1, \dots, x_n ,