

SYS865 Inférence statistique avec programmation R

Ornwipa Thamsuwan

24 janvier 2024

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Plan de la séance

- ▶ Récap: variables aléatoires
 - ▶ Espérance
 - ▶ Variance et covariance
 - ▶ Indépendance
- ▶ Échantillonnage
 - ▶ Méthodes d'échantillonnage
 - ▶ Taille d'échantillon
- ▶ Début de l'inférence statistique
 - ▶ Intervalle de confiance

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Récap et matière à réflexion

Lire des données

```
data <- read.csv("diabetes.csv")
```

Espérance

```
expectations <- sapply(data, mean)
```

Variance et covariance

```
covariances <- var(data)
```

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Comment savoir si deux variables sont indépendantes l'une de l'autre ?

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Comment savoir si deux variables sont indépendantes l'une de l'autre ?

- Inspection visuelle par graphiques de dispersion (“Scatter plot” en anglais)

Comment savoir si deux variables sont indépendantes l'une de l'autre ?

- ▶ Inspection visuelle par graphiques de dispersion (“Scatter plot” en anglais)
- ▶ Test de hypothèse
 - ▶ Test χ^2
 - ▶ Test de corrélation
 - ▶ Regression linéaire
 - ▶ Regression logistique

Graphique de dispersion

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

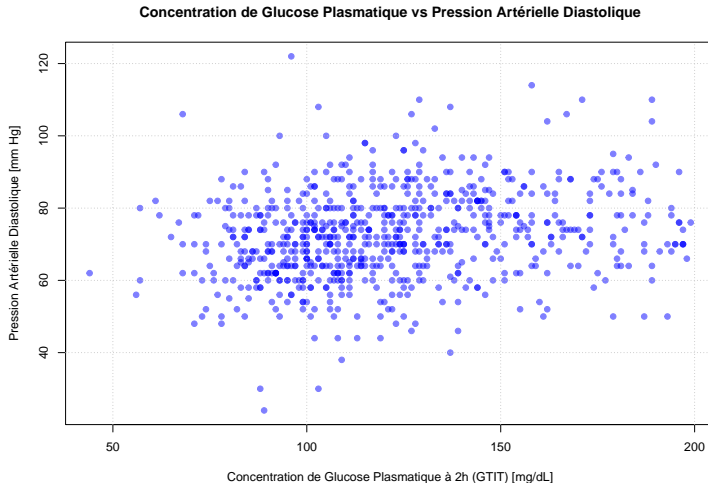
Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé



Graphique de dispersion

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

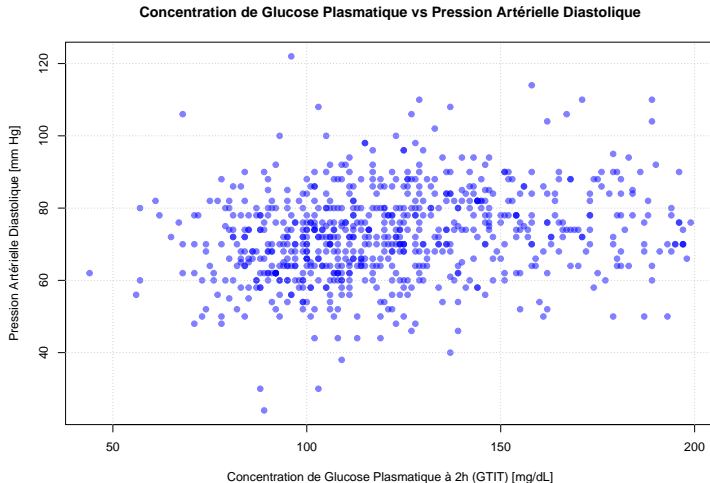
Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

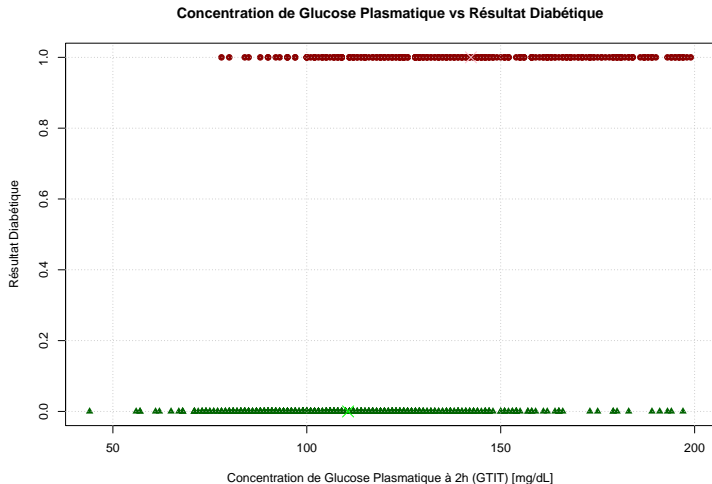


À votre avis, la covariance entre le glucose et la pression artérielle est positive, négative ou proche de zéro ?

Graphique de dispersion

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan



Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

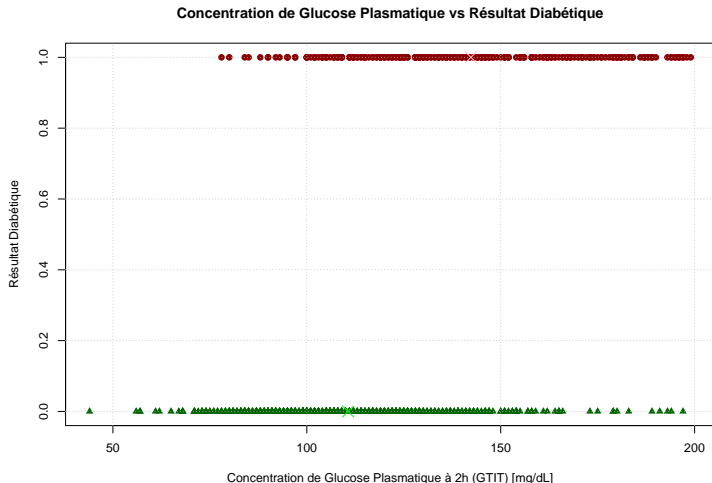
Intervalle de
confiance

Résumé

Graphique de dispersion

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan



Remarquez la différence dans la moyenne et dans la plage en comparant le cas des diabétiques et des non-diabétiques ?

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Récap et matière à réflexion (suite)

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Pouvons-nous utiliser les données fournies pour répondre à ces questions ?

Les données sont-elles représentatives de la population ?

Récap et matière à réflexion (suite)

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Pouvons-nous utiliser les données fournies pour répondre à ces questions ?

Les données sont-elles représentatives de la population ?

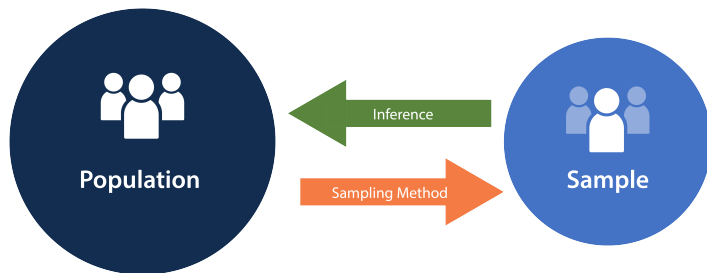


Figure 1: Relation entre population et échantillon

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Méthodes d'échantillonnage

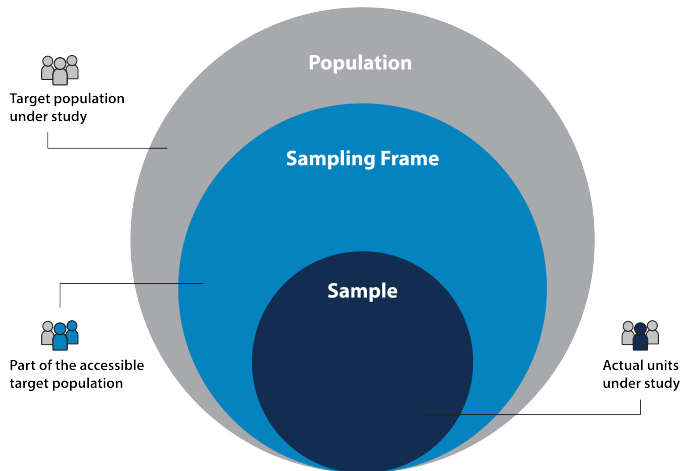


Figure 2: Échantillonnage

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Recherche sur la complexité des problèmes de statistiques dispensés par les professeurs aux différentes spécialisations à l'ÉTS :

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Recherche sur la complexité des problèmes de statistiques dispensés par les professeurs aux différentes spécialisations à l'ÉTS :

Population : Un ensemble complet d'éléments (personnes, objets ou sujets) ayant des caractéristiques spécifiques que vous souhaitez étudier et sur lesquelles vous souhaitez faire des inférences.

- *Tous les enseignants de l'ÉTS qui dispensent des cours en statistiques aux différentes spécialisations.*

Recherche sur la complexité des problèmes de statistiques dispensés par les professeurs aux différentes spécialisations à l'ÉTS :

Population : Un ensemble complet d'éléments (personnes, objets ou sujets) ayant des caractéristiques spécifiques que vous souhaitez étudier et sur lesquelles vous souhaitez faire des inférences.

- *Tous les enseignants de l'ÉTS qui dispensent des cours en statistiques aux différentes spécialisations.*

Cadre d'échantillonnage : le matériel source ou la liste complète à partir de laquelle un échantillon est tiré. C'est une compilation exhaustive de tous les éléments de votre population.

- *Le registre de l'ÉTS qui liste tous les enseignants des cours en statistiques.*

Échantillon : Un sous-ensemble d'une population. Il s'agit de l'ensemble spécifique d'éléments à partir desquels vous collecterez des données.

- *Sous-ensemble des enseignants de l'ÉTS que vous sélectionnez pour votre étude.*

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillon : Un sous-ensemble d'une population. Il s'agit de l'ensemble spécifique d'éléments à partir desquels vous collecterez des données.

- *Sous-ensemble des enseignants de l'ÉTS que vous sélectionnez pour votre étude.*

Taille de l'échantillon : le nombre de membres de la population enquêtés, mesurés ou observés.

La taille de l'échantillon détermine la quantité de données, ce qui influence davantage la précision de votre étude et la fiabilité de vos résultats.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage aléatoire simple : Chaque membre de la population a une chance égale d'être sélectionné.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage aléatoire simple : Chaque membre de la population a une chance égale d'être sélectionné.

Échantillonnage systématique : On sélectionne des membres de la population à intervalles réguliers, par exemple, choisir chaque 10ème personne de la liste.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage aléatoire simple : Chaque membre de la population a une chance égale d'être sélectionné.

Échantillonnage systématique : On sélectionne des membres de la population à intervalles réguliers, par exemple, choisir chaque 10ème personne de la liste.

Échantillonnage stratifié : La population est divisée en sous-groupes (strates) qui partagent des caractéristiques similaires. Un échantillon aléatoire est ensuite prélevé dans chacune de ces strates. Cette méthode garantit une représentation de chaque sous-groupe.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage aléatoire simple : Chaque membre de la population a une chance égale d'être sélectionné.

Échantillonnage systématique : On sélectionne des membres de la population à intervalles réguliers, par exemple, choisir chaque 10^{ème} personne de la liste.

Échantillonnage stratifié : La population est divisée en sous-groupes (strates) qui partagent des caractéristiques similaires. Un échantillon aléatoire est ensuite prélevé dans chacune de ces strates. Cette méthode garantit une représentation de chaque sous-groupe.

Échantillonnage par grappes : La population est divisée en grappes, généralement basées sur des zones géographiques, et un échantillon aléatoire de ces grappes est choisi. Tous les individus des grappes sélectionnées sont dans l'échantillon.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage probabiliste (suite)

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

Récap et matière
à réflexion

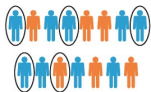
Méthodes
d'échantillonnage

Taille de
l'échantillon

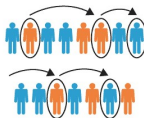
Intervalle de
confiance

Résumé

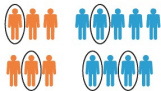
Simple random sample



Systematic sample



Stratified sample



Cluster sample

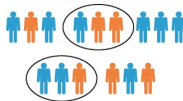


Figure 3: Échantillonnage probabiliste

Échantillonnage de convenance : Les participants choisis sont les plus faciles à atteindre. Ce n'est pas un échantillon aléatoire et est souvent utilisé pour les tests pilotes.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage de convenance : Les participants choisis sont les plus faciles à atteindre. Ce n'est pas un échantillon aléatoire et est souvent utilisé pour les tests pilotes.

Échantillonnage intentionnel : Les participants sont sélectionnés en fonction du but de l'étude et du jugement du chercheur.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage de convenance : Les participants choisis sont les plus faciles à atteindre. Ce n'est pas un échantillon aléatoire et est souvent utilisé pour les tests pilotes.

Échantillonnage intentionnel : Les participants sont sélectionnés en fonction du but de l'étude et du jugement du chercheur.

Échantillonnage à réponse volontaire : C'est les sujets qui choisissent de participer, souvent en réponse à une invitation générale.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Échantillonnage de convenance : Les participants choisis sont les plus faciles à atteindre. Ce n'est pas un échantillon aléatoire et est souvent utilisé pour les tests pilotes.

Échantillonnage intentionnel : Les participants sont sélectionnés en fonction du but de l'étude et du jugement du chercheur.

Échantillonnage à réponse volontaire : C'est les sujets qui choisissent de participer, souvent en réponse à une invitation générale.

Échantillonnage boule de neige : Les sujets actuels recrutent de futurs sujets parmi leurs connaissances. Cela est particulièrement utile pour atteindre des populations difficiles d'accès.

Échantillonnage non-probabiliste (suite)

SYS865 Inférence
statistique avec
programmation R

Ornwipa
Thamsuwan

Plan de la séance

Récap et matière
à réflexion

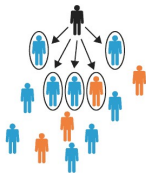
Méthodes
d'échantillonnage

Taille de
l'échantillon

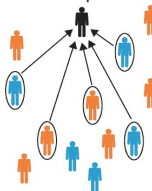
Intervalle de
confiance

Résumé

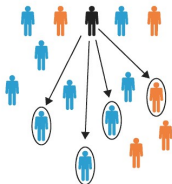
Convenience sample



Voluntary response sample



Purposive sample



Snowball sample

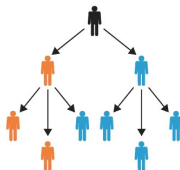


Figure 4: Échantillonnage non-probabiliste

En groupe de 3-4, discutez de quelle serait la situation dans laquelle chacune des méthodes d'échantillonnage serait utilisée ?

Pour une étude sur une maladie rare, un chercheur choisi délibérément les patients connus pour souffrir de cette maladie à partir des dossiers médicaux ou des établissements de santé spécialisés.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Pour une étude sur une maladie rare, un chercheur choisi délibérément les patients connus pour souffrir de cette maladie à partir des dossiers médicaux ou des établissements de santé spécialisés.

- Échantillonnage intentionnel

Pour une étude sur une maladie rare, un chercheur choisi délibérément les patients connus pour souffrir de cette maladie à partir des dossiers médicaux ou des établissements de santé spécialisés.

► Échantillonnage intentionnel

Pour évaluer l'efficacité d'une nouvelle campagne de santé, une organisation sélectionne au hasard cinq quartiers d'une ville. Ils enquêtent ensuite sur chaque foyer de ces quartiers.

Pour une étude sur une maladie rare, un chercheur choisi délibérément les patients connus pour souffrir de cette maladie à partir des dossiers médicaux ou des établissements de santé spécialisés.

- Échantillonnage intentionnel

Pour évaluer l'efficacité d'une nouvelle campagne de santé, une organisation sélectionne au hasard cinq quartiers d'une ville. Ils enquêtent ensuite sur chaque foyer de ces quartiers.

- Échantillonnage par grappes

Dans une étude sur la santé mentale des personnes sans abri, les premiers participants sans abri recommandent d'autres personnes sans abri qu'ils connaissent.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Dans une étude sur la santé mentale des personnes sans abri, les premiers participants sans abri recommandent d'autres personnes sans abri qu'ils connaissent.

- Échantillonnage boule de neige

Dans une étude sur la santé mentale des personnes sans abri, les premiers participants sans abri recommandent d'autres personnes sans abri qu'ils connaissent.

► Échantillonnage boule de neige

Une équipe de recherche sur la prévalence de l'hypertension divise la population en catégories ethniques, puis sélectionne au hasard un nombre proportionné d'individus dans chaque groupe pour garantir que tous soient représentés.

Dans une étude sur la santé mentale des personnes sans abri, les premiers participants sans abri recommandent d'autres personnes sans abri qu'ils connaissent.

- Échantillonnage boule de neige

Une équipe de recherche sur la prévalence de l'hypertension divise la population en catégories ethniques, puis sélectionne au hasard un nombre proportionné d'individus dans chaque groupe pour garantir que tous soient représentés.

- Échantillonnage stratifié

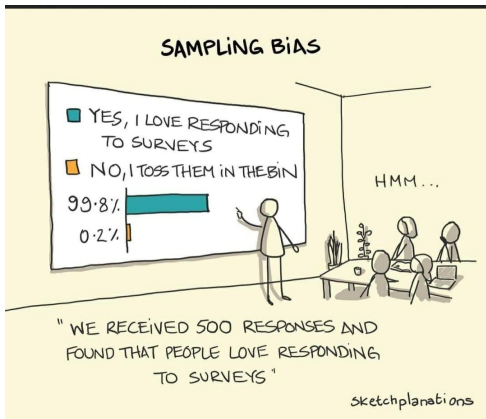


Figure 5: Biais de l'échantillonnage

Taille de l'échantillon

Le Théorème Central Limite (TCL) stipule que, pour une taille d'échantillon suffisamment grande, la distribution des moyennes d'échantillons se rapprochera d'une distribution normale, indépendamment de la distribution originale de la population.

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Le Théorème Central Limite (TCL) stipule que, pour une taille d'échantillon suffisamment grande, la distribution des moyennes d'échantillons se rapprochera d'une distribution normale, indépendamment de la distribution originale de la population.

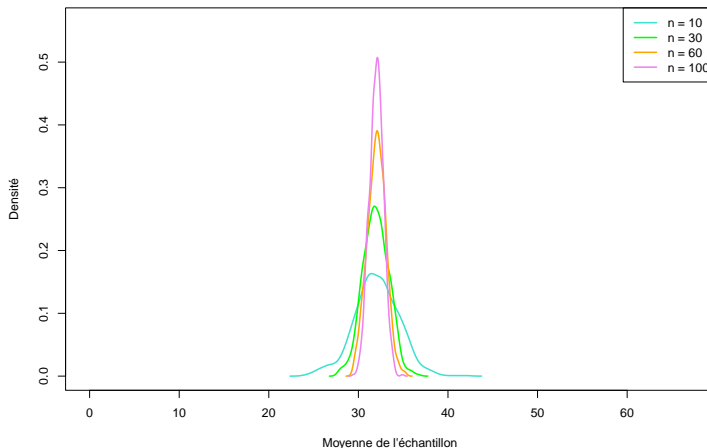
Définition

Si $X_1, X_2, X_3, \dots, X_n$ sont des échantillons aléatoires pris d'une population avec une moyenne générale μ et une variance finie σ^2 , la moyenne de l'échantillon $\bar{X} = \frac{1}{n}(X_1 + X_2 + X_3 + \dots + X_n)$ sera approximativement distribuée normalement avec une moyenne μ et une variance $\frac{\sigma^2}{n}$, à mesure que n devient grand.

La distribution normale est notée $N(\mu, \frac{\sigma^2}{n})$.

BMI dans la base de données de “Pima Indian Diabetes Dataset”

Distribution des moyennes d'échantillons pour différentes tailles d'échantillons



Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Le TCL justifie l'utilisation de la distribution normale dans l'inférence statistique et les tests d'hypothèses, même lorsque la population sous-jacente n'est pas normalement distribuée.

À mesure que n (la taille de l'échantillon) augmente, la forme de la distribution de la moyenne de l'échantillon \bar{X} devient de plus en plus en cloche ("bell-shaped" en anglais) ou normale.

L'erreur standard de la moyenne, qui mesure la variabilité des moyennes d'échantillons, est donnée par $SE = \frac{\sigma}{\sqrt{n}}$.

À mesure que la taille de l'échantillon n augmente, l'erreur standard SE diminue. Cela indique que des échantillons plus grands fournissent des estimations plus précises de la moyenne de la population, réduisant ainsi le risque d'erreur d'échantillonnage.

À noter que . . .

La taille d'échantillon “suffisamment grande” pour le TCL est généralement considérée comme étant 30 ou plus, mais cela peut varier en fonction de la population.

Plan de la séance

**Récap et matière
à réflexion**

**Méthodes
d'échantillonnage**

**Taille de
l'échantillon**

**Intervalle de
confiance**

Résumé

Intervalle de confiance

Plan de la séance

Récap et matière
à réflexion

Méthodes
d'échantillonnage

Taille de
l'échantillon

Intervalle de
confiance

Résumé

Résumé

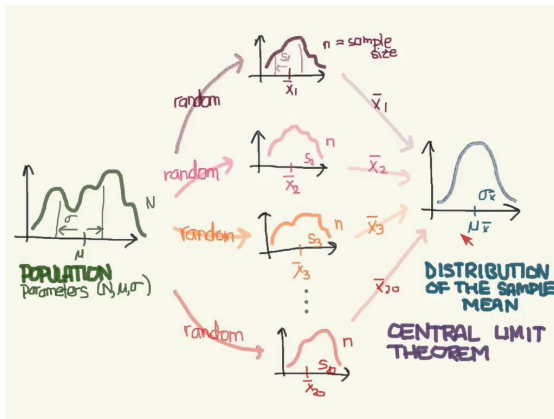


Figure 6: Théorème Central Limite