

Seance 10: Comparaison de moyenne

Visseho Adjiwanou, PhD.

30 March 2022

Rappel

- 1 Comparaison de la moyenne (Chapitre 8 de Fox)
- 2 Test de comparaison de la moyenne: test t de student

Type de la variable indépendante			Type de la variable dépendante		
			Catégorielle		Continue
			Nominale	Ordinale	Intervalle/ratio
	Catégorielle	Nominale	Tableau croisé	Tableau croisé	Analyse de la variance
			Test du chi-carré		Test t
			V de Cramer		État carré
			Diagramme de barre empilée		Diagramme boxplot
	Continue	Ordinale	Tableau croisé	Tableau croisé	idem
		Intervalle/ratio	Transformer la VI en catégorielle		Corrélation / régression
					Test t
					R carré
					Diagramme de nuage de points

Figure 1

Comparaison de moyenne

- Exemple
- Différence des heures d'écoute de la télé entre Montréalais et Torontois
- Différence de revenus entre les hommes et les femmes
- Autres exemples?

Comparaison de moyenne

- Comme on le voit, on est en présence de deux variables:
- Une variable dépendante (intervalle/ratio)
 - Heure d'écoute de la télé
 - Revenus
- Une variable indépendante dichotomique
 - Résidence (Montréal, Toronto)
 - Sexe (Homme, Femme)

Comparaison de moyenne

1 Représentation

- Diagramme en boîte (à moustache) – boxplot
- Mais, attention, il ne nous donne pas la moyenne
- Il faut ajouter la moyenne de chaque groupe sur le graphique

2 Calcul

- Très simple: calculé la moyenne de la variable dépendante dans chaque catégorie de la variable indépendante

Comparaison de moyenne

- Comment savons-nous que la différence observée dans l'échantillon est valide au sein de la population?
- Il faut faire un **test statistique** pour conclure.
- Test est basé sur la distribution d'échantillonnage de la différence de la moyenne

Comparaison de moyenne

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories
- Prendre la différence de ces moyennes: vous vous retrouvez avec des millions de différences

Test t

- ① S'il y a **une grande différence** au sein de la **population**
 - On va s'attendre à ce que dans la très grande majorité des échantillons possibles, la différence est aussi grande
 - Quelque chose comme dans 19 fois sur 20, la différence va être significative, ou que 1 fois sur 20, elle ne va pas l'être
- ② En revanche, s'il **n'y a pas une grande différence** au sein de la **population**
 - On va s'attendre à ce que dans la très grande majorité des échantillons possibles, la différence est aussi faible, autour de 0

Test t

- ➊ Plus la différence entre les moyennes est faible dans la population, plus grande est la proportion des échantillons ayant une faible différence entre les moyennes
- ➋ Plus la différence entre les moyennes est forte dans la population, plus grande sera la proportion des échantillons ayant une forte différence entre les moyennes.

Distribution d'échantillonnage

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories

Distribution d'échantillonnage

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories
- Prendre la différence de ces moyennes: vous vous retrouvez avec des millions de différences

Distribution d'échantillonnage

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories
- Prendre la différence de ces moyennes: vous vous retrouvez avec des millions de différences
- La représentation graphique de cette différence est la **distribution d'échantillonnage** de la différence

Distribution d'échantillonnage

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories
- Prendre la différence de ces moyennes: vous vous retrouvez avec des millions de différences
- La représentation graphique de cette différence est la **distribution d'échantillonnage** de la différence
- Les statisticiens ont montré que cette distribution tend vers une **loi normale** au fur et à mesure que N devient grand

Distribution d'échantillonnage

- Prendre tous les échantillons possible de taille N de la population
- Calculer la moyenne pour les deux catégories
- Prendre la différence de ces moyennes: vous vous retrouvez avec des millions de différences
- La représentation graphique de cette différence est la **distribution d'échantillonnage** de la différence
- Les statisticiens ont montré que cette distribution tend vers une **loi normale** au fur et à mesure que N devient grand
- Plus concrètement, cette distribution suit une **loi de Student**

Distribution de Student et lecture

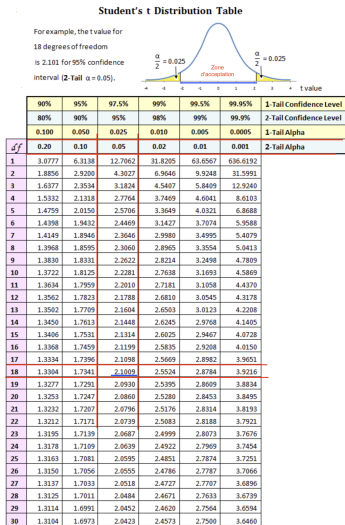


Figure 2

Test t

Cette statistique s'écrit :

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_{\bar{X}_1 - \bar{X}_2}}$$

Test t

\bar{X}_1 = la moyenne de la variable dépendante pour la catégorie 1 de la variable indépendante de l'**échantillon**

\bar{X}_2 = la moyenne de la variable dépendante pour la catégorie 2 de la variable indépendante de l'**échantillon**

μ_1 = la moyenne de la variable dépendante pour la catégorie 1 de la variable indépendante de la **population**

μ_2 = la moyenne de la variable dépendante pour la catégorie 2 de la variable indépendante de la **population**

$s_{\bar{X}_1 - \bar{X}_2}$ = erreur-type de la différence entre les moyennes

Test t

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_{\bar{X}_1 - \bar{X}_2}}$$

$$s_{\bar{X}_1 - \bar{X}_2} = \sqrt{\left(\frac{N_1 * s_1^2 + N_2 * s_2^2}{N_1 + N_2 - 2}\right) \left(\frac{N_1 + N_2}{N_1 * N_2}\right)}$$

s_1^2 = variance de la variable dépendante pour la catégorie 1 de la variable indépendante

s_2^2 = variance de la variable dépendante pour la catégorie 2 de la variable indépendante

N_1 = Le nombre de cas dans la catégorie 1 de la variable indépendante

N_2 = Le nombre de cas dans la catégorie 2 de la variable indépendante

Exemple de calcul

Hommes (Catégorie 1)	Femmes (Catégorie 2)
$\bar{X}_1 = 2.75$	$\bar{X}_2 = 3.01$
$s_1 = 2.030$	$s_2 = 2.225$
$N_1 = 855$	$N_2 = 1085$

Test statistique

- ❶ Hypothèse nulle : $\mu_1 = \mu_2$ (la moyenne dans le groupe 1 de la population est la même que dans le groupe 2)
 - Ou Hypothèse nulle : $\mu_1 - \mu_2 = 0$
- ❷ Postule qu'il n'y a pas de différence entre la moyenne pour les femmes et la moyenne pour les hommes dans la population
- ❸ Dans ces conditions, le t devient $t = \frac{(\bar{X}_1 - \bar{X}_2)}{s_{\bar{X}_1 - \bar{X}_2}}$
- ❹ Le degré de liberté vaut $N_1 + N_2 - 2$
- ❺ Nous décidons de rejeter ou de ne pas rejeter H_0 , en déterminant où se trouve cette valeur dans la distribution t

En résumé (1e manière)

Pour faire un test de comparaison de deux moyennes:

- ➊ Posez votre hypothèse nulle
- ➋ Choisissez votre niveau de significativité
- ➌ Trouvez votre degré de liberté
- ➍ Trouvez votre zone d'acceptabilité et de rejet (où trouver la valeur MINIMALE de t pour rejeter l'hypothèse nulle)
- ➎ Calculez votre t selon les formules ci-dessus
- ➏ Prendre une décision:
 - Si votre $|t|$ calculé est supérieur ou égal au $t(lu) \implies$ Rejeter l'hypothèse nulle
 - Si votre $|t|$ calculé est inférieur au $t(lu) \implies$ Vous ne pouvez pas rejeter l'hypothèse nulle

En résumé (2e manière)

OU

- ⑥ Prendre une décision:
 - Si votre t calculé se trouve entre la zone d'acceptation : Ne rejeter pas l'hypothèse nulle
 - Si votre t calculé se trouve à l'extérieur de la zone d'acceptation : rejeter l'hypothèse nulle

Exemple

- Est ce que les hommes regardent la télé plus que les femmes?
- Vous obtenez les résultats suivants à partir d'un échantillon de ... (à vous de trouver la valeur) hommes et femmes

Hommes (Catégorie 1)	Femmes (Catégorie 2)
$\bar{X}_1 = 2.75$	$\bar{X}_2 = 3.01$
$s_1 = 2.030$	$s_2 = 2.225$
$N_1 = 855$	$N_2 = 1085$

Exemple (1e manière)

- ① Posez votre hypothèse nulle
 - H_0 : La moyenne d'heures d'écoute de la télé est la même pour les hommes et les femmes

Exemple (1e manière)

- 1 Posez votre hypothèse nulle
 - H_0 : La moyenne d'heures d'écoute de la télé est la même pour les hommes et les femmes
- 2 Choisissez votre niveau de significativité
 - $\alpha = 0.05$

Exemple (1e manière)

- ❶ Posez votre hypothèse nulle
 - H_0 : La moyenne d'heures d'écoute de la télé est la même pour les hommes et les femmes
- ❷ Choisissez votre niveau de significativité
 - $\alpha = 0.05$
- ❸ Trouvez votre degré de liberté
 - $df = 855 + 1085 - 2 = 1938$

Exemple (1e manière)

- ➊ Posez votre hypothèse nulle
 - H_0 : La moyenne d'heures d'écoute de la télé est la même pour les hommes et les femmes
- ➋ Choisissez votre niveau de significativité
 - $\alpha = 0.05$
- ➌ Trouvez votre degré de liberté
 - $df = 855 + 1085 - 2 = 1938$
- ➍ Trouvez la valeur MINIMALE de t pour rejeter l'hypothèse nulle
 - $t(lu) = 1.960$

Exemple (1e manière)

- 5 Calculez votre t selon les formules ci-dessus

Rappelons-nous la formule de t:

$$t = \frac{(\bar{X}_1 - \bar{X}_2)}{s_{\bar{X}_1 - \bar{X}_2}}$$

$$s_{\bar{X}_1 - \bar{X}_2} = \sqrt{\left(\frac{N_1 * s_1^2 + N_2 * s_2^2}{N_1 + N_2 - 2}\right) \left(\frac{N_1 + N_2}{N_1 * N_2}\right)}$$

```
numérateur <- 2.75 - 3.01
denominateur <- sqrt((855*2.030^2 + 1085*2.225^2)/(855 + 1085))
t_calcule <- numérateur/denominateur

t_calcule

## [1] -2.653868
```

Exemple (1e manière)

6 Prendre une décision:

- Si votre $|t|$ calculé est supérieur ou égal au $t(lu) \implies$ Rejeter l'hypothèse nulle
- Si votre $|t|$ calculé est inférieur au $t(lu) \implies$ Vous ne pouvez pas rejeter l'hypothèse nulle

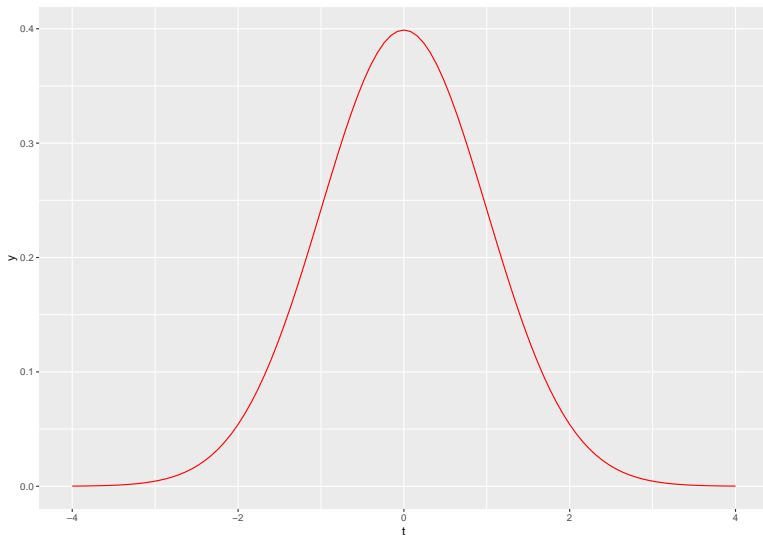
• Notre décision

- $|-2.653868| = 2.653868 > 1.960$,
- décision: On rejette l'hypothèse nulle.

Exemple (2e manière)

- H_0 : La moyenne d'heures d'écoute de la télé est la même pour les hommes et les femmes
- Sous cette hypothèse, la statistique t prend la valeur : $t = \frac{(\bar{X}_1 - \bar{X}_2)}{s_{\bar{X}_1 - \bar{X}_2}}$
- Le degré de liberté vaut $N_1 + N_2 - 2 = 855 + 1085 - 2 = 1938$
- Avec ce niveau de degré de liberté, la distribution de Student aura l'allure ci-dessous:

Exemple



Exemple

- Si dans la population, il n'y a pas relation entre le sexe et la durée d'écoute de la télé (Hypothèse nulle), on s'entend que si on tire par exemple, 20 échantillons aléatoires de taille N , la grande majorité va avoir la valeur de t autour de 0.

Exemple

- Si dans la population, il n'y a pas relation entre le sexe et la durée d'écoute de la télé (Hypothèse nulle), on s'entend que si on tire par exemple, 20 échantillons aléatoires de taille N , la grande majorité va avoir la valeur de t autour de 0.
- Si nous souhaitons nous tromper 1 fois sur 20 (niveau de significativité de 5%), alors, il faut trouver :
 - la borne inférieure et
 - la borne supérieure

Exemple

- Si dans la population, il n'y a pas relation entre le sexe et la durée d'écoute de la télé (Hypothèse nulle), on s'entend que si on tire par exemple, 20 échantillons aléatoires de taille N , la grande majorité va avoir la valeur de t autour de 0.
- Si nous souhaitons nous tromper 1 fois sur 20 (niveau de significativité de 5%), alors, il faut trouver :
 - la borne inférieure et
 - la borne supérieure
- A partir de cette courbe, on peut ainsi trouver la valeur de t pour laquelle
 - la distribution est inférieure à 0,025 (valeur inférieure)
 - la distribution est inférieure à 0,975 (valeur supérieure)

Exemple

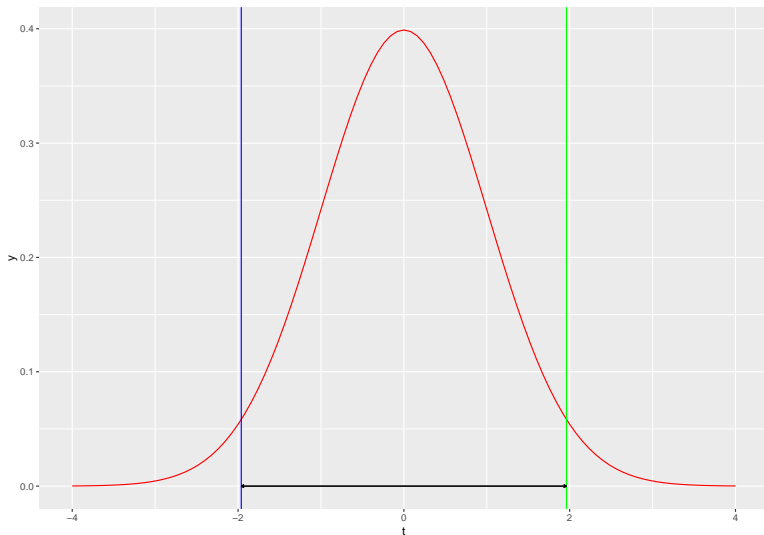
```
bi <- qt(0.025, df=1938)
bs <- qt(0.975, df=1938)

c(bi, bs)
```

```
## [1] -1.961189  1.961189
```

Plaçons ces informations sur le graphique

Plaçons ces informations sur le graphique

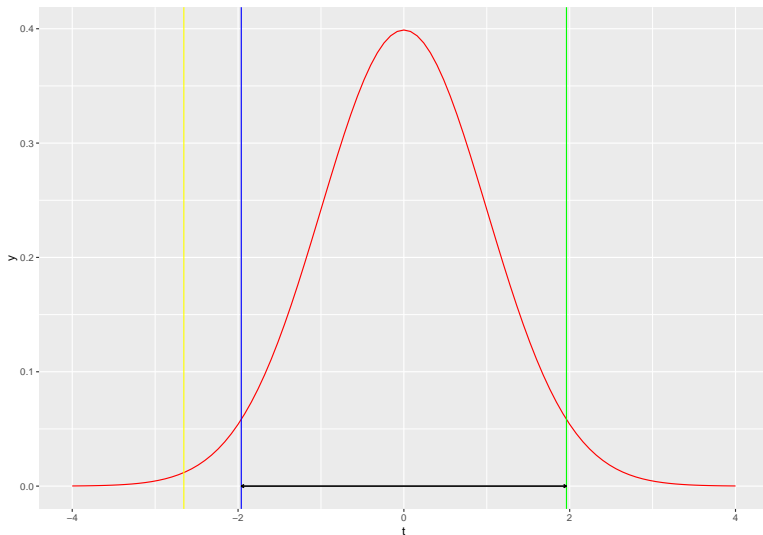


- Maintenant, si vous calculez votre t issu de votre échantillon et qu'il se trouve
 - entre $-1,96$ et $1,96$ (entre les lignes bleue et verte), alors vous ne pouvez pas rejeter votre hypothèse nulle
 - plus petit que $-1,96$ ou plus grand que $1,96$, alors on rejette l'hypothèse nulle.

Calculons alors le t

- On avait calculé le t qui vaut: -2.653868
- On voit ainsi que la valeur du t calculé se trouve en dehors de la zone d'acceptation:
- Conclusion: On rejette l'hypothèse nulle.
- En classe, je vais vous montrer comment lire la table de distribution du t de Student

Décision



Conclusion

- Pour la semaine prochaine
 - Faire Labo 10 sur la comparaison de moyenne
 - Lire chapitre 9
- Le Quiz va porter sur tout ce qu'on a vu jusqu'au chapitre 8