

Fondements statistiques de la méthode expérimentale en technologie éducative

Mattia A. Fritz

31/03/2023

Résumé

La méthode expérimentale est étroitement liée à la modélisation des données notamment avec des finalités d'inférence statistique. Le rapport entre les données, la modélisation de celles-ci, et les conclusions qu'on peut tirer des tests statistiques effectués est cependant complexe et régit par différentes approches philosophiques et méthodologiques. À complément du document sur les fondements empiriques de la méthode expérimentale, ce document s'intéresse donc aux fondements statistiques de cette méthode. Il propose d'abord une brève introduction à la modélisation des données et aux probabilités. Ensuite, les statistiques *fréquentistes* (ou *classiques*) et les statistiques Bayésiennes sont illustrées dans les grandes lignes. En guise de conclusion, le document propose une brève analyse sur les avantages des statistiques Bayésiennes dans un contexte pédagogique.

Introduction

Très peu des chercheurs en sciences sociales s'intéressent aux statistiques. Elles sont souvent vécues comme un *mal nécessaire* pour pouvoir bénéficier des avantages épistémologiques fournis par la méthode expérimentale. Cette perspective se reflète souvent dans les manuels de méthode expérimentale qui associent les différents design expérimentaux avec les tests statistiques correspondantes, afin d'illustrer comment les résultats de ces tests doivent être interprétés en relation aux hypothèses opérationnelles d'abord, et théoriques ensuite. Il en suit que les statistiques ressortent de cette démarche comme une liste de recettes à appliquer machinalement selon les *ingrédients* du plan expérimental: si j'ai n variables indépendantes de type t et m variables dépendantes de type u , alors il faut utiliser le test x . Cette démarche mécanique est d'ailleurs corroborée par des diagrammes qu'on peut facilement trouver dans des manuels ou en ligne, dont l'objectif est précisément d'accompagner les chercheurs à travers un nombre de bifurcations avant de trouver, enfin, le test nécessaire à leurs besoins. Malheureusement, cette approche ne semble pas donner ses fruits, car il existe désormais

plusieurs témoignages dans la littérature scientifique qui montrent comme les statistiques sont mal comprises et utilisées dans les contributions expérimentales (Bakker & Wicherts, 2011; Greenland et al., 2016; Nickerson, 2000; Nuijten et al., 2016; Singmann et al., 2023). Même les ressources pédagogiques qui sont censées former les nouveaux chercheurs présentent souvent des erreurs dans l'exposition de concepts clé (Cobb, 2007; Maurer et al., 2019; McElreath, 2020).

Cette contribution utilise une approche différente caractérisée en premier lieu par la séparation entre les fondements empiriques et les fondements statistiques de la méthode expérimentale. Ce document s'intéresse à ces derniers et le fait de manière à dévoiler des concepts importants qui, dans une approche machinale, sont cachés ou traités comme s'il n'existaient pas d'alternatives (Lakens, 2021; McElreath, 2020; Rodgers, 2010). L'objectif de cette contribution est de sensibiliser des étudiant-es en technologie éducative sans un background en méthode expérimentale à l'importance que les statistiques jouent dans la démarche expérimentale. Pour ce faire, le texte présente d'abord une brève introduction à la modélisation des données, c'est-à-dire à l'utilisation de techniques mathématiques et probabilistes pour manier et représenter des données empiriques. Ensuite, le document illustre deux approches et philosophies différentes aux statistiques: l'approche *fréquentiste* ou *classique*, dominante en sciences sociales depuis plusieurs décennies, et l'approche Bayésienne, qui commencent à s'introduire dans certaines contributions. Les deux approches sont illustrées dans les grandes lignes avec l'objectif de formuler qu'est-ce qu'on peut *vraiment* demander aux tests statistiques, et comment leur réponses peuvent nous aider dans la démarche explicative de la méthode expérimentale. En guise de conclusion, ce document prend position par rapport aux deux approches et suggère que les statistiques Bayésiennes devraient être enseignées en priorités dans les cours de méthodologie. Cette position est brièvement argumentée.

1 Modélisation des données

Les chercheurs récoltent des données empiriques dans la tentative de pouvoir en tirer des informations utiles qui, souvent, dépassent le cadre spécifique des données récoltées et visent plutôt à des affirmations ou principes avec une portée plus étendue. Ce mécanisme qui permet de passer du particulier à un cadre plus généralisé est souvent identifié avec le terme d'**inférence**. On peut identifier différents objectifs relatifs à l'inférence (Gelman et al., 2021; McElreath, 2020; Rodgers, 2010).

Réduction des données à des indicateurs représentatifs. Grâce à la représentation en forme d'indicateurs on peut *inférer* des caractéristiques d'un ensemble de données. Un exemple d'inférence de ce type sont la moyenne et l'écart type tirés d'un nombre

d'observations. Quand on indique que les notes obtenues dans un cours ont une moyenne de $M = 3.5$ ($SD = 0.8$) on peut par exemple inférer que la plupart des étudiant-es n'ont pas validé le cours, car la moyenne est inférieure au barème de 4. En même temps, l'écart type qui est représentatif de la dispersion des données indique que certain-es étudiant-es ont sûrement dépassé ce barème, car si on ajoute un écart type à la moyenne, on obtient 4.3, et en ajoutant deux écarts types on obtient 5.1. Avec ce calcul on peut également inférer que le système de notation de l'enseignant-e est assez sévère, car il y a peu de chances qu'une étudiant-e ait obtenu une très bonne note.

Évaluer le degré d'association entre des variables. On peut utiliser un jeu de données pour établir une mesure d'intérêt et utiliser d'autres variables récoltées pour faire des comparaisons depuis lesquelles inférer, par exemple, si les observations avec un certain type de valeurs sur une variable sont associées à des valeurs différentes sur la mesure d'intérêt. Par exemple, dans un jeu de données qui s'intéresse à la mesure du sentiment de bien être des étudiant-es dans un cursus de bachelor, on peut comparer si la perception de bien être diffère entre la première, la deuxième et la troisième année. Si c'est le cas, on peut inférer la présence d'une association entre les deux variables.

Prédire des événements futurs. Des données récoltées peuvent être utilisées pour produire un *algorithme* qui prend des variables comme Input et produisent la mesure d'intérêt comme Output. Une fois cet algorithme de conversion obtenu avec les données observées, on peut entrer des nouvelles données pour estimer la mesure d'intérêt. Dans ce mécanisme, on infère que le comportement des nouvelles données sera similaire au *comportement* des données observées. Cette inférence est traduite par l'algorithme lui-même. Ce principe est notamment à la base des techniques de *machine learning*.

Déterminer l'effet d'une intervention. À travers des données qui varient systématiquement sur une ou plusieurs variables dont les valeurs ont été fixées par les chercheurs, on peut inférer des mécanismes contre-factuels du type: que se serait-il passé si la personne avait été attribuée à une autre valeur/modalité de l'intervention? Dans ce contexte, l'objectif de l'inférence est en général double: (1) déterminer si l'effet est présent dans le *micro-monde* des données observées, et (2) estimer à quel point on peut être confiant que l'effet puisse se reproduire au *macro-monde*.

Comme il a été indiqué dans les fondements empiriques de la méthode expérimentale, les expériences appartiennent donc à ce dernier cas de figure. Cependant, l'inférence depuis des données empiriques s'appuie en général sur les mêmes instruments: des **modèles mathématiques**. Ceci est souvent source de confusion, surtout dans un contexte introductif aux méthodes dits *quantitatifs*. En effet, très concrètement, on utilise souvent les mêmes logiciels et les mêmes *fonctions/tests* à l'intérieur de ces logiciels indépendamment de comment les

données ont été créées (e.g. observation, simulation ou expérience). Il en résulte une compréhensible difficulté à cerner les différences qui ne résident en effet pas dans les modèles eux-mêmes, mais plutôt dans les connaissances scientifiques relatives à:

- Les relations causales entre les variables impliquées dans le modèle, qui peuvent notamment être explicitées avec un modèle structural de causalité sous forme de *Directed Acyclic Graph* (DAG) (Bareinboim & Pearl, 2016; Cinelli et al., 2020; Pearl, 2000; Pearl et al., 2016; Pearl & Mackenzie, 2018);
- Le processus génératif des données, c'est-à-dire sous quelles conditions les données ont été produites (Maxwell et al., 2017; McElreath, 2020).

Sur la base de ces informations, les modèles mathématiques et les indicateurs que ces modèles produisent doivent être interprétés de manière conforme à ce qu'ils permettent ou ne permettent pas d'inférer. À ce propos, cette section illustre d'abord en quoi consiste la modélisation des données et quels sont ces avantages. L'étape suivante introduit la *famille* de modèles la plus fréquente en science sociale: la modélisation linéaire. Ensuite, cette famille de modèle sera centrée plus spécifiquement dans le cadre des expériences, notamment en ce qui concerne la perspective contre-factuelle. Enfin, elle introduit le concept d'inférence statistique qui sera ensuite décliné dans les deux approches fréquentiste et Bayésien.

1.1 Pourquoi modéliser des données?

L'une des questions souvent inexplorée dans les manuels de méthodologie *quantitative* est la suivante: pourquoi avons-nous besoins des statistiques en premier lieu? La vie des chercheurs – et des étudiant-es! – seraient tellement plus simple sans elles. Il est donc légitime de s'attendre à ce que leur association pratiquement indissoluble avec la recherche *quantitative* soit justifiée par un apport exceptionnel en termes épistémologiques.

Malheureusement, ce n'est pas vraiment le cas, au moins selon cette contribution qui adopte une attitude désenchantée. Le rôle prépondérant des statistiques dans la recherche s'explique principalement par deux raisons:

- La complexité des phénomènes étudiées qui sont souvent dépendantes d'un large éventail de facteurs qui s'influencent mutuellement. Dans cette complexité, même des *patterns stables* peuvent ne pas se produire à chaque fois, mais seulement *la plupart des fois*. Par exemple, il existe un lien de causalité stable entre les heures passées à étudier et la réussite à un examen. Mais il se peut que, parfois, un-e étudiant-e qui a beaucoup étudié rencontre un échec, et un-e étudiant-e qui n'a pas beaucoup étudié réussisse néanmoins l'examen.

- La tendance intrinsèque aux êtres humains à voir des liens de cause à effet lorsqu'en réalité il ne s'agit que d'épiphénomènes circonstanciels, dont l'occurrence est indépendante des causes supposées par la personne.

Ces deux raisons sont brièvement développées par la suite. En guise de conclusion de cette partie introductive sur la modélisation, une définition formelle d'un modèle applicable dans le contexte statistique sera fournie.

1.1.1 Vivre – et faire de la recherche – dans un monde pseudo-déterministe

1.1.2 Se protéger de ses propres biais

1.1.3 Qu'est-ce qu'un modèle statistique

1.2 La modélisation linéaire

Il existe différentes manières pour modéliser des données. Dans les sciences sociales, les modèles les plus utilisés appartiennent à la *famille* des modélisations linéaires. Ces modèles se caractérisent par le fait que la variable d'intérêt (le outcome ou la mesure) peut être représenté par ce qu'on appelle une équation de régression. L'équation de régression dans les modèles linéaires correspond à l'équation représentant une pente dans un plan cartésien:

$$Y = \text{Intercepte} + \text{Pente} \times X$$

Dans les manuels statistiques cette équation est plus souvent représentée de la manière suivante:

$$y_i = \beta_0 + \beta_1(x_i) + \epsilon_i$$

L'explication des composantes de cette version de l'équation est la suivante:

- y_i est la mesure de la variable Y pour l'observation i dans le jeu de données, c'est-à-dire y_1, y_2, \dots, y_n . Malheureusement, en statistique on fait souvent la distinction entre notations qui pour les non-statisticiens ne sont pas très marquée, comme la distinction entre minuscule et majuscule ou entre lettre latine et grecque.
- β_0 correspond à l'intercepte, c'est-à-dire la valeur de Y lorsque toutes les éventuelles variables sur la droite de l'équation (dans ce cas seulement X) équivalent à 0.
- β_1 correspond au paramètre (ou coefficient de régression) pour la variable X . Ce coefficient est commun à toutes les valeurs observées pour X , c'est-à-dire x_i .

- ϵ_i correspond à ce qu'on appelle le résidu de l'équation. Il s'agit d'une valeur de *compensation* dû au fait que les paramètres β_0 et β_1 sont communs à toutes les observations, mais pratiquement jamais les paramètres du modèles permettent d'arriver *exactement* à la valeur de y_i correspondante.

Contrairement à ce qui est souvent indiqué, une régression linéaire ne veut pas forcément dire *rectilinéaire*. En effet, on peut par exemple créer des courbes en ajoutant des valeurs exponentielles aux variables prédictives:

$$\checkmark y_i = \beta_0 + \beta_1(x_i^2) + \epsilon_i$$

Par contre, il n'est pas possible d'utiliser des exponentiels pour les coefficients de régression. Par exemple, cette équation ne serait pas considérée comme un modèle linéaire:

$$\times y_i = \beta_0 + \beta_1^{x_i} + \epsilon_i$$

Ce préambule très technique, qui sera développé davantage dans les exemples plus bas, sert à ce point pour indiquer que les modèles linéaires sont des modèles qu'on peut facilement accommoder pour prendre en ligne de compte plusieurs types de relations entre les variables prédictives et la mesure de outcome. Cette flexibilité se traduit malheureusement dans la littérature scientifique avec des noms d'analyses différentes qui sont en réalité de cas spéciaux de la modélisation linéaire. Ce tableau propose une liste de ces *cas spéciaux*.

Table 1: Cas spéciaux de la modélisation linéaire présents dans la littérature scientifique

| Nom de l'analyse | Outome | Variable prédictive |
|-------------------------------|------------|--|
| Régression simple | 1 continue | 1 continue |
| Régression multiple | 1 continue | 1 continue ou plus |
| t-test à groupes indépendants | 1 continue | 1 catégorielle avec 2 modalités |
| ANOVA simple | 1 continue | 1 catégorielle avec plus de 2 modalités |
| ANOVA factorielle | 1 continue | 2 catégorielles ou plus |
| ANCOVA | 1 continue | 1 catégorielle ou plus et 1 continue ou plus |

Table 1: Cas spéciaux de la modélisation linéaire présents dans la littérature scientifique (*continued*)

| Nom de l'analyse | Outome | Variable prédictive |
|---------------------------------|------------------------------|--|
| Corrélation de Pearson | 1 continue standardisée | 1 continue standardisée |
| t-test avec un seul groupe | 1 continue | Intercepte seulement |
| t-test apparié | Différence entre 2 continues | Intercepte seulement |
| t-test de Hotelling | 2 continues ou plus | 1 catégorielle avec 2 modalités |
| MANOVA | 2 continues ou plus | 1 catégorielle avec plus de 2 catégorielles |
| Régression multiple multivariée | 2 continues ou plus | 1 catégorielle ou plus et 1 continue ou plus |

Comme le tableau l'indique, ces cas spéciaux présupposent que la mesure soient toujours représentée sur une échelle continue potentiellement infinie sur les côtés. Ceci n'est cependant pas toujours le cas, par exemple lorsque la mesure est exprimée sur des échelles de Lickert ou sur une échelle avec des limites inférieurs et/ou supérieurs. Dans le reste de cette partie, le texte propose des exemples concrets de différentes typologies de modèles linéaires dans lesquels le modèle suivante est une version plus flexible du précédent et peut donc être appliqué à des données plus complexes. Les modèles sont dans l'ordre:

- Modèle linéaire simple
- Modèle linéaire multiple
- Modèle linéaire généralisée
- Modèle linéaire généralisée mixte

1.2.1 Modèle linéaire simple

Comme indiqué plus haut, le modèle linéaire simple présuppose que la variable outcome est le résultat de l'addition entre l'intercepte, la pente d'une variable prédictive, et le résidu:

$$y_i = \beta_0 + \beta_1(x_i) + \epsilon_i$$

Voici un jeu de données avec 10 observations. Chaque observation se compose simplement de la variable prédictive x et de la variable outcome y .

Table 2: Jeu de données pour une régression linéaire simple.

| i | x | y |
|----|-------|--------|
| 1 | 16.64 | 101.99 |
| 2 | 13.69 | 107.25 |
| 3 | 16.61 | 125.49 |
| 4 | 20.00 | 130.52 |
| 5 | 15.70 | 96.45 |
| 6 | 12.81 | 80.77 |
| 7 | 10.00 | 50.00 |
| 8 | 11.78 | 63.24 |
| 9 | 14.84 | 102.25 |
| 10 | 19.31 | 200.00 |

Lorsqu'on demande à un logiciel d'analyse statistique comme par exemple R de mener une régression linéaire simple sur ces données, le résultat qu'on obtient sera le suivant.

Table 3: Tableau des paramètres d'une régression linéaire simple dans une perspective fréquentiste.

| Predictor | b | 95% CI | t | df | p |
|-----------|--------|------------------|-------|------|------|
| Intercept | -63.44 | [-148.49, 21.61] | -1.72 | 8 | .124 |
| X | 11.18 | [5.67, 16.69] | 4.68 | 8 | .002 |

Le résultat se réfère à une analyse de type fréquentiste, mais pour l'instant cet aspect n'a pas d'importance. Ce qui est à noter depuis les résultats sont les deux coefficients attribués aux paramètres $\beta_0 = -63.44$ et $\beta_1 = 11.18$. Il en résulte qu'on peut écrire les résultats de notre modèle également sous forme d'équation, mais cette fois-ci avec les coefficients intégrés.

$$\hat{y} = -63.44 + 11.18(x) \quad (1)$$

Cette équation diffère de celle générique présentée plus haut de deux manières. En premier lieu, y a été remplacé par \hat{y} . Lorsque un élément à un *chapeau* en statistique, cela signifie qu'il s'agit d'une estimation liée à l'utilisation d'un modèle. Donc \hat{y} n'est pas une valeur observée, mais une valeur calculée sur la base des paramètres inférés par le modèle. La deuxième différence consiste dans la disparition du résidu ϵ . En effet, en s'agissant d'une prédiction, on ne peut pas savoir à quel point ce valeur s'éloigne de la *vraie* valeur. Par contre, on peut

calculer cette distance si on compare les valeurs observées du jeu de données avec les valeurs qui seraient prédites par le modèle lorsque x a exactement la même valeur de l'observation dans le jeu de données. La figure suivante montre ce principe graphiquement.

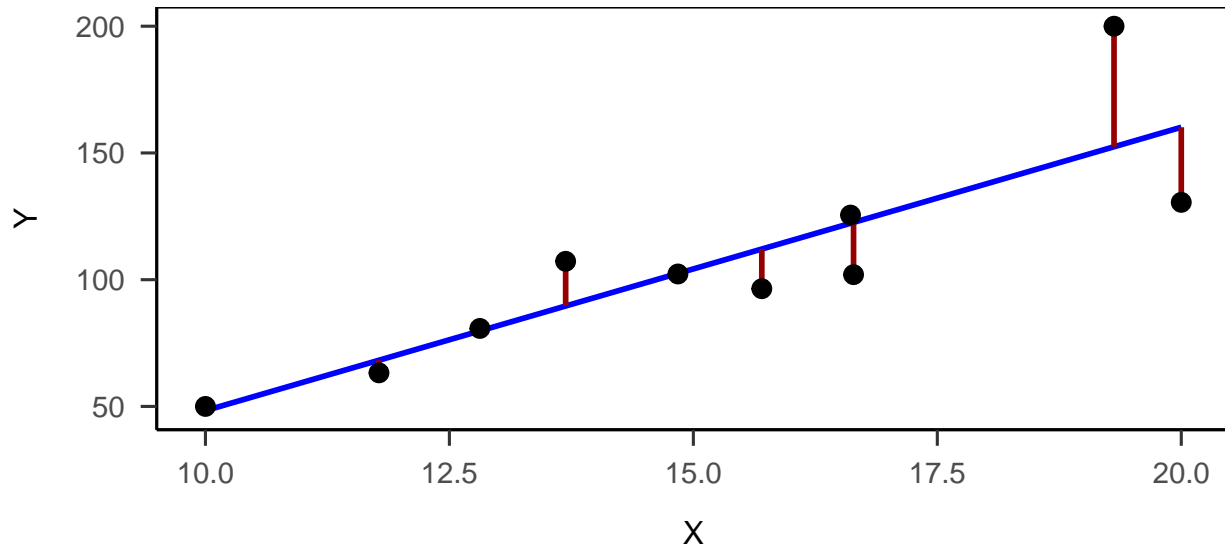


Figure 1: Représentation graphique d'une régression linéaire simple. La ligne bleu représente les valeurs prédites par le modèle. Les segments rouges entre les points et la ligne bleu sont les résidus des données observées.

En fait, c'est exactement avec ce type de processus que les paramètres du modèle linéaire sont calculés en premier lieu. En effet, la régression linéaire simple consiste à trouver la *ligne* qui passe à travers les observations et minimise la distance avec toutes les observations du jeu de données. Le tableau suivant reprend à ce propos les observations, mais ajoute trois autres colonnes:

- \hat{y} : la valeur prédite par le modèle pour une observation x avec la même valeur du jeu de donnée
- $y - \hat{y}$: la différence entre la valeur observée et la valeur prédite
- $(y - \hat{y})^2$: la différence élevée au carré. L'utilisation de l'exponentiel sert deux objectifs: (1) éviter que dans la somme de toutes les différences les valeurs positives et négatives s'annulent, et (2) donner plus de poids aux distances plus extrêmes, comme par exemple la dernière, dont la différence de 47.56 unités devient de 2262.41.

Table 4: Jeu de données avec valeur prédite et différence (normale et au carré) entre la valeur observée et prédite

| i | x | y | \hat{y} | $y - \hat{y}$ | $(y - \hat{y})^2$ |
|----|-------|--------|-----------|---------------|-------------------|
| 1 | 16.64 | 101.99 | 122.60 | -20.60 | 424.44 |
| 2 | 13.69 | 107.25 | 89.61 | 17.64 | 311.16 |
| 3 | 16.61 | 125.49 | 122.25 | 3.25 | 10.54 |
| 4 | 20.00 | 130.52 | 160.14 | -29.62 | 877.52 |
| 5 | 15.70 | 96.45 | 112.08 | -15.64 | 244.56 |
| 6 | 12.81 | 80.77 | 79.79 | 0.98 | 0.96 |
| 7 | 10.00 | 50.00 | 48.35 | 1.65 | 2.73 |
| 8 | 11.78 | 63.24 | 68.23 | -4.99 | 24.91 |
| 9 | 14.84 | 102.25 | 102.48 | -0.23 | 0.05 |
| 10 | 19.31 | 200.00 | 152.44 | 47.56 | 2262.41 |

La somme de $(y - \hat{y})^2$ est notamment utilisée comme mesure pour établir à quel point le modèle s'encastre avec les données observées. Le plus cette somme est élevée, le moins le modèle est en adéquation avec les données. Est-ce que la somme de 4,159.27 est élevée? C'est difficile à dire sans une mesure de comparaison. À cet effet, on peut par exemple utiliser la moyenne $\bar{Y} = 105.80$ comme modèle linéaire alternative. La moyenne peut être tout à fait considérée un modèle linéaire qui a seulement l'intercepte. L'intercepte correspond justement à la moyenne de la variable considérée. Le tableau suivant reprend la même structure du précédent, mais en utilisant la valeur fixe de la moyenne pour calculer la distance des observations.

Table 5: Jeu de données avec valeur prédite et différence (normale et au carré) entre la valeur observée et la moyenne de Y

| i | x | y | \hat{y} | $y - \bar{Y}$ | $(y - \bar{Y})^2$ |
|---|-------|--------|-----------|---------------|-------------------|
| 1 | 16.64 | 101.99 | 105.8 | -3.80 | 14.46 |
| 2 | 13.69 | 107.25 | 105.8 | 1.46 | 2.12 |
| 3 | 16.61 | 125.49 | 105.8 | 19.70 | 388.05 |
| 4 | 20.00 | 130.52 | 105.8 | 24.72 | 611.06 |
| 5 | 15.70 | 96.45 | 105.8 | -9.35 | 87.41 |
| 6 | 12.81 | 80.77 | 105.8 | -25.03 | 626.41 |
| 7 | 10.00 | 50.00 | 105.8 | -55.80 | 3113.18 |
| 8 | 11.78 | 63.24 | 105.8 | -42.56 | 1811.11 |

Table 5: Jeu de données avec valeur prédite et différence (normale et au carré) entre la valeur observée et la moyenne de Y (*continued*)

| i | x | y | \hat{y} | $y - \bar{Y}$ | $(y - \bar{Y})^2$ |
|----|-------|--------|-----------|---------------|-------------------|
| 9 | 14.84 | 102.25 | 105.8 | -3.55 | 12.58 |
| 10 | 19.31 | 200.00 | 105.8 | 94.20 | 8874.42 |

La somme des distances au carré en utilisant \bar{Y} comme modèle linéaire est de 15,540.79, c'est à dire 11,381.52 plus du modèle dérivé de la régression linéaire simple. On peut notamment inférer depuis cette différence qu'en connaissant la valeur de x , on peut avoir une meilleure idée de la valeur de y . En termes formelles: $\mathbb{P}(Y|X) \neq \mathbb{P}(Y)$.

1.2.2 Modèle linéaire multiple

Le modèle linéaire multiple, connu aussi comme modèle linéaire général, est simplement une extension de la régression linéaire simple à plusieurs variables prédictives. Dans l'exemple précédent, le modèle consistait dans une seule variable numérique, mais le modèle linéaire multiple accepte également des variables prédictives binaires ou catégorielles. Par exemple, dans le jeu de données suivante, la variable outcome Y est calculée sur la base de l'effet additive entre les variables numériques X et W , plus la variable catégorielle Z qui possède trois modalités: *Faible*, *Moyenne* et *Forte*. L'équation de la régression linéaire multiple est donc la suivante:

$$y_i = \beta_0 + \beta_1(x_i) + \beta_2(w_i) + \beta_3(z_{i\text{Moyenne}}) + \beta_4(z_{i\text{Forte}}) + \epsilon_i$$

Il est utile de remarquer comme dans cette équation n'apparaît pas la modalité faible de la variable Z . Ceci s'explique par un mécanisme adopté souvent en régression linéaire qui consiste à attribuer à la première modalité d'une variable catégorielle une sorte de valeur de base. Ensuite, on attribue aux autres modalités une *dummy* variable, c'est-à-dire une variable qui assume les valeurs 0 ou 1. Dans ce cas, si l'observation appartient à la modalité *Moyenne*, alors $z_{i\text{Moyenne}}$ sera 1 et $z_{i\text{Forte}}$ sera 0. Si la variable appartient au contraire à *Forte*, les valeurs seront inversés. Ceci à la conséquence d'additionner dans le calcul seulement le coefficient de la variable avec valeur 1 et d'annuler le coefficient de la variable avec valeur 0. En effet, multiplier n'importe quel coefficient par 1 signifie ajouter une fois ce coefficient au calcul, tandis que le multiplier par 0 signifie ajouter 0 au calcul. Le tableau suivant montre les premières 10 observations d'un jeu de données qui compte $N = 200$.

Table 6: Extrait des premières 10 observations d'un jeu de données avec plusieurs variables prédictives.

| i | x | w | z | y |
|----|--------|--------|---------|--------|
| 1 | 11.917 | 9.929 | Moyenne | 51.840 |
| 2 | 17.232 | 10.799 | Moyenne | 97.255 |
| 3 | 13.968 | 14.466 | Moyenne | 72.491 |
| 4 | 6.645 | 14.047 | Faible | 24.949 |
| 5 | 15.063 | 3.767 | Moyenne | 79.507 |
| 6 | 17.450 | 11.484 | Forte | 93.480 |
| 7 | 10.855 | 5.691 | Forte | 71.177 |
| 8 | 2.428 | 8.581 | Forte | 17.317 |
| 9 | 14.297 | 1.634 | Faible | 47.803 |
| 10 | 14.671 | 7.361 | Forte | 80.809 |

Le tableau de la régression multiple doit à ce moment proposer les coefficients pour X , Y , ainsi que pour les deux modalités *Moyenne* et *Forte* de Z .

Table 7: Tableau des paramètres d'une régression linéaire multiple dans une perspective fréquentiste.

| Predictor | b | 95% CI | t | df | p |
|-----------|-------|---------------|-------|------|--------|
| Intercept | -3.25 | [-8.84, 2.35] | -1.14 | 195 | .254 |
| X | 4.73 | [4.36, 5.10] | 25.20 | 195 | < .001 |
| W | 0.11 | [-0.23, 0.44] | 0.62 | 195 | .537 |
| ZMoyenne | 5.13 | [1.79, 8.48] | 3.02 | 195 | .003 |
| ZForte | 7.69 | [4.40, 10.98] | 4.62 | 195 | < .001 |

Encore une fois, les résultats de ce tableau font référence à l'approche fréquentiste, mais ces sont seulement les coefficients qui nous intéressent pour le moment. L'équation de régression avec les coefficients estimés sera donc la suivante:

$$\hat{y} = -3.25 + 4.73(x) + 0.11(w) + 5.13(z_{\text{Moyenne}}) + 7.69(z_{\text{Forte}}) \quad (2)$$

Un élément d'intérêt particulier dans la régression linéaire multiple concerne l'interprétation des coefficients les uns par rapport aux autres. En effet, dans ce type de modèle, la contribution additive de chaque coefficient est à interpréter comme une contribution unique par

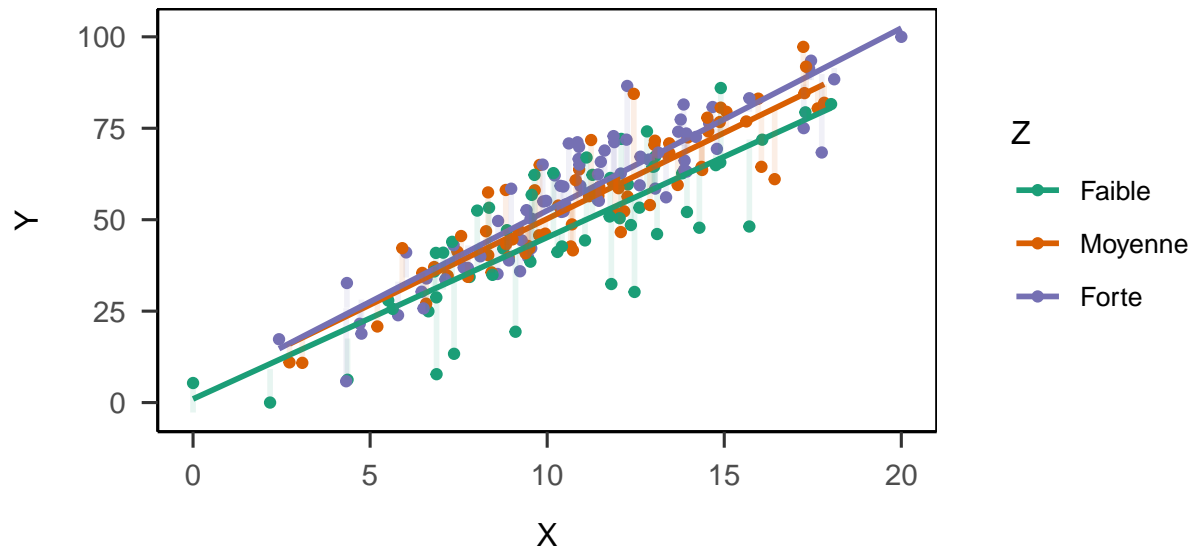
rapport aux autres coefficients. Voici comment interpréter chaque coefficient en termes de changements sur la variable outcome Y :

- β_0 Intercept(-3.25) : Lorsque toutes les variables prédictives (X , W et Z) sont égales à zéro, la valeur prédite de la variable Y est de -3,25. Cependant, cette interprétation peut ne pas être saillante si la valeur zéro pour l'une des variables prédictives n'est pas plausible dans le contexte des données, ou même si une variable négative de l'outcome Y ne fait pas de sens (e.g. temps négatif).
- β_1 x(4.73) : Deux observations qui diffèrent d'une unité sur la variable X , tout en maintenant les autres variables (w et z) constantes, diffèrent sur la variable outcome Y de 4,73 unités.
- β_2 w(0.11) : Deux observations qui diffèrent d'une unité sur la variable W , les autres variables (x et z) restant constantes, diffèrent de 0,11 unité sur la variable outcome Y .
- β_3 zMoyenne(5.13) : Il s'agit d'une variable catégorielle avec comme catégorie de référence la modalité *Faible*. Dans ce cas, une observation avec modalité *Moyenne* diffère d'une observation avec modalité *Faible* de 5,13 unités sur la variable outcome Y , tout en maintenant les autres variables (X et W) constantes.
- β_4 zForte(7.69) : De même, le coefficient de 7,69 représente la différence sur la variable outcome Y entre une observation avec modalité *Forte* et une avec modalité *Faible*, tout en gardant les autres variables (X et W) constantes.
- On peut récupérer le coefficient qui détermine la différence entre une observation avec modalité *Forte* et une avec modalité *Moyenne* à travers la subtraction 7.69 (zForte) - 5.13 (zMoyenne) = 2.56 . En d'autres termes, deux observations qui gardent X et W constantes diffèrent de 2.56 unités sur la variable outcome Y lorsque l'une est avec modalité *Forte* et l'autre modalité *Moyenne*.

La comparaison entre les modalités de Z est de quelque sorte plus compliquée à comprendre, car la valeur de *Faible* n'est pas vraiment explicitée. En effet, la valeur de Y pour les observations avec modalité *Faible* correspond tout simplement à l'effet additive de X et W , car zMoyenne et zForte sont annulés par la multiplication par 0.

L'interprétation des coefficients d'une régression linéaire multiple nécessite de beaucoup de pratique avant d'être maîtrisée. Mais ces coefficients jouent un rôle fondamental dans l'inférence et il était donc nécessaire d'en illustrer le mécanisme. D'ailleurs, une régression linéaire multiple est également plus difficile à visualiser graphiquement par rapport à la régression linéaire simple. Ici nous proposons un graphique qui ne prend pas en compte la variable W juste pour des propos illustratifs. Le graphique propose trois lignes de régression

en correspondance aux trois modalités de la variable Z .



Nous verrons par la suite que l'interprétation de ces coefficients sera plus saillante lorsque les données seront issues d'une expérience, tandis que dans cette partie nous nous référons de manière générale aux modèles linéaires indépendamment de leur application concrète. Avec des exemples dans lesquels les variables assument des connotations plus concrètes (e.g. heures d'études, attribution à des interfaces différentes d'un logiciel, etc.), aussi l'interprétation des coefficients en résultera facilitée.

1.2.3 Modèle linéaire généralisée

1.2.4 Modèle linéaire généralisée mixte

1.3 La modélisation d'outcomes potentielles

1.4 L'inférence statistique

2 Statistiques *fréquentistes*

3 Statistiques Bayésiennes

4 Conclusion

Références

- Bakker, M., & Wicherts, J. M. (2011). The (Mis)Reporting of Statistical Results in Psychology Journals. *Behavior Research Methods*, 43(3), 666-678. <https://doi.org/10.3758/s13428-011-0089-5>
- Bareinboim, E., & Pearl, J. (2016). Causal Inference and the Data-Fusion Problem. *Proceedings of the National Academy of Sciences*, 113(27), 7345-7352. <https://doi.org/10.1073/pnas.1510507113>
- Cinelli, C., Forney, A., & Pearl, J. (2020). A Crash Course in Good and Bad Controls. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.3689437>
- Cobb, G. W. (2007). The Introductory Statistics Course: A Ptolemaic Curriculum. *Technology Innovations in Statistics Education*, 1(1), 1-16.
- Gelman, A., Hill, J., & Vehtari, A. (2021). *Regression and Other Stories*.
- Greenland, S., Senn, S. J., Rothman, K. J., Carlin, J. B., Poole, C., Goodman, S. N., & Altman, D. G. (2016). Statistical Tests, P Values, Confidence Intervals, and Power: A Guide to Misinterpretations. *European Journal of Epidemiology*, 31(4), 337-350. <https://doi.org/10.1007/s10654-016-0149-3>
- Lakens, D. (2021). The Practical Alternative to the p Value Is the Correctly Used p Value. *Perspectives on Psychological Science*, 1745691620958012. <https://doi.org/10.1177/1745691620958012>
- Maurer, K., Hudiburgh, L., Werwinski, L., & Bailer, J. (2019). Content Audit for P-Value Principles in Introductory Statistics. *The American Statistician*, 73(sup1), 385-391. <https://doi.org/10.1080/00031305.2018.1537890>
- Maxwell, S. E., Delaney, H. D., & Kelley, K. (2017). *Designing experiments and analyzing data: a model comparison perspective* (Third edition). Routledge.
- McElreath, R. (2020). *Statistical rethinking: A Bayesian course with examples in R and Stan* (2 éd.). Taylor; Francis, CRC Press.
- Nickerson, R. S. (2000). Null Hypothesis Significance Testing: A Review of an Old and Continuing Controversy. *Psychological Methods*, 5(2), 241-301. <https://doi.org/10.1037/1082-989X.5.2.241>
- Nuijten, M. B., Hartgerink, C. H. J., van Assen, M. A. L. M., Epskamp, S., & Wicherts, J. M. (2016). The Prevalence of Statistical Reporting Errors in Psychology (1985–2013). *Behavior Research Methods*, 48(4), 1205-1226. <https://doi.org/10.3758/s13428-015-0664-2>
- Pearl, J. (2000). *Causality: models, reasoning, and inference*. Cambridge University Press.
- Pearl, J., Glymour, M., & Jewell, N. P. (2016). *Causal Inference in Statistics: A Primer*. Wiley.
- Pearl, J., & Mackenzie, D. (2018). *The book of why: the new science of cause and effect*. Basic

Books.

- Rodgers, J. L. (2010). The epistemology of mathematical and statistical modeling: A quiet methodological revolution. *American Psychologist*, 65(1), 112. <https://doi.org/10.1037/a0018326>
- Singmann, H., Kellen, D., Cox, G. E., Chandramouli, S. H., Davis-Stober, C. P., Dunn, J. C., Gronau, Q. F., Kalish, M. L., McMullin, S. D., Navarro, D. J., & Shiffrin, R. M. (2023). Statistics in the Service of Science: Don't Let the Tail Wag the Dog. *Computational Brain & Behavior*, 6(1), 64-83. <https://doi.org/10.1007/s42113-022-00129-2>