

1. L'information statistique

Nature et présentation

Sam Gyetvay

ECO 2273 – Économétrie I

September 12, 2025

- ▶ Introduction
- ▶ Bref historique des données et de la statistique
- ▶ Qu'est-ce que les données ?
- ▶ Types de variables
- ▶ D'où viennent les données ?
- ▶ Exemples de bases de données
- ▶ Visualisation des données

Introduction

Professeur

Sam Gyetvay

Économiste du travail

Recherche principalement sur l'immigration

Nouveau prof dans le département (arrivé août 2025)

PhD UBC 2024, Postdoc Ohio State 2024-25

Première fois que j'enseigne

Je suis un anglophone montréalais, j'étais autrefois bilingue fluide mais mon français est maintenant rouillé

Auxiliaire d'enseignement

Oumar Djamaldiev

Étudiant MA en économie UQAM-ESG

Responsable de la démonstration lundi 14h00-15h00, où il enseignera STATA

Ce cours

Le nom du cours est Économétrie I, mais nous n'apprendrons pas beaucoup d'économétrie *à proprement parler* ce semestre

Nous allons plutôt construire une base de connaissances sur laquelle l'économétrie pourra se développer

$$\begin{aligned}\text{Économétrie} &= \mathbf{\text{Probabilités}} \\ &+ \mathbf{\text{Statistiques}} \\ &+ \text{Algèbre linéaire} \\ &+ \text{Théorie économique}\end{aligned}$$

Notre base se concentrera sur les deux premiers

Les prochaines conférences se concentreront sur les probabilités. Dans cette introduction, nous présentons quelques concepts de base des données et de la statistique

Statistique

- ▶ Latin *statisticum collegium* (« conseil de l'État »)
- ▶ Italien *statista* (« homme d'État »)
- ▶ Allemand *statistik* (« science de l'État »)

Données

- ▶ Latin *data / datum* (« choses supposées être des faits »)

Histoire de la statistique

Comme le montre son étymologie, l'origine de la statistique est étroitement liée à l'émergence des États centralisés.

Les États collectaient des informations systématiques pour connaître leur population et leur territoire

- ▶ Recensement de la population
- ▶ Comptes nationaux
- ▶ Relevés de température

Le produit de cette collecte systématique d'informations s'appelle **données**

Au cours de plusieurs siècles, des méthodes mathématiques ont été développées pour appliquer la **théorie des probabilités** aux données afin d'obtenir des enseignements. La **statistique** désigne cet ensemble de connaissances.

Dynastie Ming "Livre jaune" (1422)



Livre jaune des impôts et services de travail
Source : <https://iguoxue.ifeng.com/51700719/>

Premier recensement canadien, 1871

[illegible]

Page de retour du recensement d'Ascot, Québec, 1871

CENSUS OF 1871.											
PRINCE EDWARD ISLAND.											
TABLE IV.—Birth Places of the People.											
TABLEAU IV.—Population par Lieux de Naissances.											
Counties. Comtés.	Popula- tion.	E. Anglo- Irish.	Ireland. Irlande.	Brit. Isles.	Natives. Natijs.	British North- American Colonies.	Other Foreign Born.	Not Classed.			
					Various Other Foreign Born.	Various Other Foreign Born.	Various Other Foreign Born.				
<i>Prince.</i>											
In District.....	8,323	119	515	78	7,261	33	285	14	20		
2nd do.....	4,739	125	317	77	3,660	33	193	8	511		
3rd do.....	5,520	27	64	67	4,958	35	339	7	30		
4th do.....	4,495	166	302	142	3,655	39	341	14	29		
5th do.....	5,528	94	62	36	5,277	1	241	39			
Total.....	29,392	594	895	413	24,512	138	1,947	73	236		
<i>Queen's.</i>											
In District.....	9,947	551	524	794	5,742		339	33			
2nd do.....	4,495	399	381	438	35		440	11	25		
3rd do.....	4,052	194	232	315	4,439	28	96	7	39		
4th do.....	4,291	54	341	717	3,232	14	317	18	34		
Charlottetown.....	9,947	413	426	545	772		545	100			
Total.....	45,851	1,603	2,063	2,229	19,463	77	1,389	204	77		
<i>King's.</i>											
In District.....	5,514	30	124	229	5,141	4	264	31	4		
2nd do.....	5,539	45	126	273	5,149		266	38			
3rd do.....	5,506	196	254	411	4,911		246	33			
4th do.....	5,563	121	140	265	5,095	4	246	31	4		
Georgetown.....	1,056	13	41	44	947		179	5			
Total.....	23,688	338	744	1,313	20,584	8	879	102	8		
Grand Total.....	94,821	1,567	2,772	4,519	79,968	85	3,246	384	223		

CENSUS OF 1871.											
ISLE DU PRINCE-EDOUARD.											
TABLE V.—Land and Cattle.											
TABLEAU V.—Terres et Bétail.											
Counties. Comtés.	Land improved with Crops.	Land improved with Pasture.	Land improved with Timber.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.	Land improved with Other Crops.
<i>Prince.</i>											
In District.....	1,396	95,062	55,778	2,708
2nd do.....	678	61,126	32,761	1,369
3rd do.....	738	64,274	39,383	5,067
4th do.....	857	93,365	52,667	4,834	1,565
5th do.....	146	13,582	7,233	1,339	144

Tableau imprimé contenant des statistiques sur le lieu de naissance des individus à l'Î.-P.-É., la superficie de terres et le bétail possédés

John Graunt *Natural and Political Observations Made upon the Bills of Mortality* (1662)

Collecte et analyse des Bills of Mortality de Londres et des registres de baptêmes paroissiaux (naissances)

Utilisation des registres de baptême et d'enterrement pour estimer la population

Construction d'une table de vie :

- Probabilité de survie par âge
- Forte mortalité infantile et déclin régulier ensuite

[illegible]

Les données à l'ère moderne

Aujourd'hui, les données proviennent de nombreuses sources

Sondages: recueillent des opinions ou des informations sur la population

- ▶ *L'Enquête sur la population active* utilisée pour calculer le taux de chômage

Expérimentations: études contrôlées pour identifier des effets causaux

- ▶ Les essais de Pfizer pour un vaccin contre la COVID-19 recueillent des données sur les participants à l'étude

Données administratives: collectées dans le cadre des opérations gouvernementales

- ▶ Les dossiers fiscaux basés sur le T1, T4

Données en ligne: générées par les plateformes numériques

- ▶ Les données d'Amazon sur chaque transaction et avis produit

La disponibilité de ces sources de données a transformé la recherche en économie

Qu'est-ce qu'une base de données ?

Les données sont des faits

- ▶ Sam est un homme de 34 ans né à Montréal, Canada
- ▶ La population du Québec en 2015 était de 8,25 M

Avant d'appliquer la statistique à des données, celles-ci doivent être organisées en **base de données**

Une base de données est une manière de collecter un ensemble de faits dans une **matrice** ou un **tableau**

Chaque ligne représente une **observation**

Chaque colonne représente une **variable** différente

Exemple d'un base de données 1

Nom	Sexe	Âge	Lieu de naissance	Éducation
Sam	M	34	Montréal, Canada	PhD
Alex	M	29	Toronto, Canada	BA
Marie	F	31	Québec, Canada	MA

Chaque ligne = une observation (personne)

Chaque colonne = une variable (Nom, Sexe, Âge, Lieu de naissance, Éducation)

La première ligne n'est pas une observation, elle donne seulement les noms des variables

Exemple d'un base de données 2

Année	Province	Population
1985	Québec	6,665,800
1990	Québec	6,997,000
1995	Québec	7,219,200
2000	Québec	7,357,000
2005	Québec	7,581,200
2010	Québec	7,929,400
2015	Québec	8,254,900

Chaque ligne = une observation (année-province)

Chaque colonne = une variable (Année, Province, Population)

Types de variables

Comme nous l'avons déjà vu, il existe différents types de variables.

Certaines variables sont **numériques**, comme la population, l'âge ou l'année.

Certaines variables sont **non numériques**, comme le nom, le sexe, la province, le lieu de naissance ou le niveau d'éducation.

Parfois, les variables non numériques ont un **ordre**. Par exemple, on peut classer les niveaux d'éducation : PhD > MA > BA.

Certaines variables non numériques n'ont pas d'ordre ; elles représentent juste différentes catégories. Nom, sexe, lieu de naissance, province sont ainsi.

Visualisation des données

La plupart des bases de données sont assez grands : ils contiennent des centaines, des milliers, parfois des millions de lignes.

Il n'est donc pas pratique de lire un base de données ligne par ligne.

La visualisation des données transforme les données en image. Ces images sont faciles à interpréter et peuvent révéler des motifs intéressants dans les données.

Aujourd'hui, nous discuterons de trois visualisations classiques, bien que beaucoup d'autres existent

1. Séries temporelles
2. Nuage de points
3. Histogram

Séries temporelles

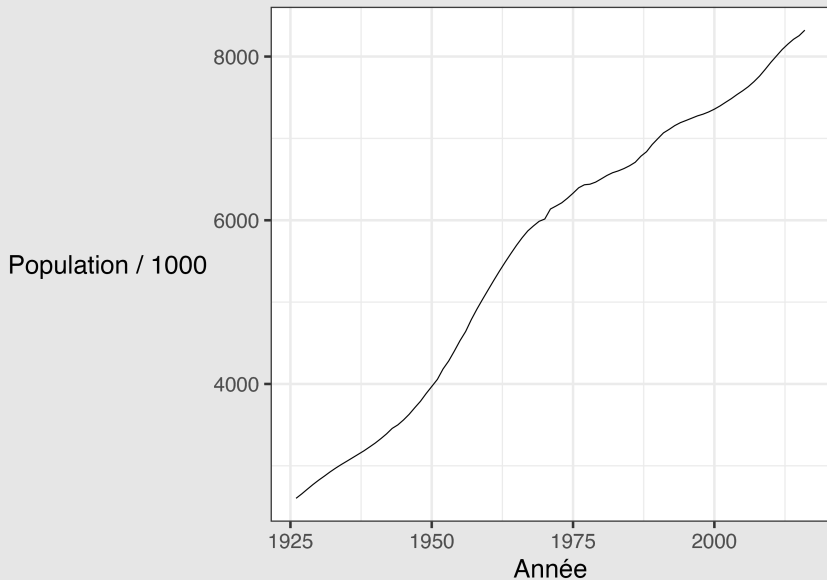
Un graphique en séries temporelles montre comment une variable change dans le temps.

On ne peut utiliser les séries temporelles que si la base de données contient une variable temporelle comme l'année, le trimestre, le mois ou la date.

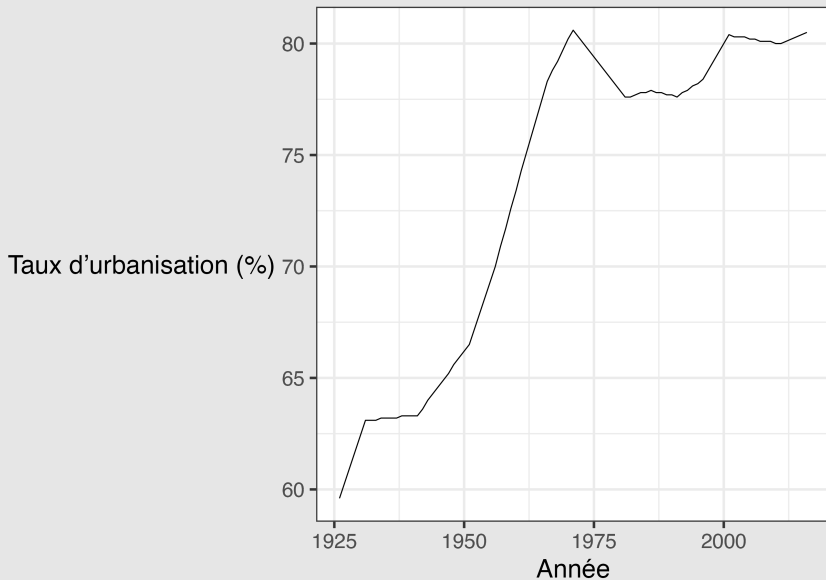
Dans un graphique en séries temporelles, la variable temps (ex. année) est sur l'axe X, et l'autre variable sur l'axe Y.

Les valeurs à gauche sont plus anciennes, et celles à droite plus récentes.

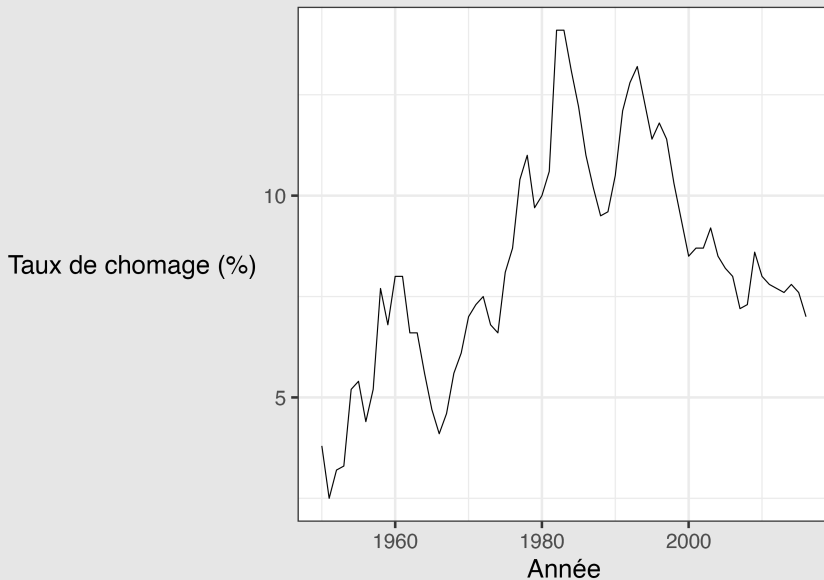
Population du Québec, 1926-2016



Taux d'urbanisation du Québec, 1925-2016



Taux de chômage du Québec, 1950-2016



Nuage de points

Un nuage de points montre comment deux variables sont liées, généralement à un instant donné.

Un nuage de points a une variable sur l'axe Y et une sur l'axe X.

On place ensuite un point à la coordonnée (y, x) correspondant à chaque observation.

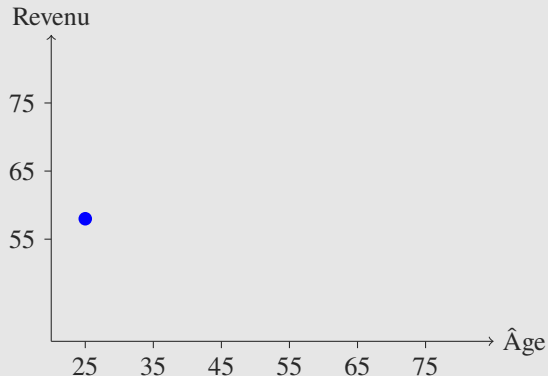
Nuage de points

Un nuage de points montre comment deux variables sont liées, généralement à un instant donné.

Un nuage de points a une variable sur l'axe Y et une sur l'axe X.

On place ensuite un point à la coordonnée (y, x) correspondant à chaque observation.

Nom	Âge	Revenu
Alice	25	58



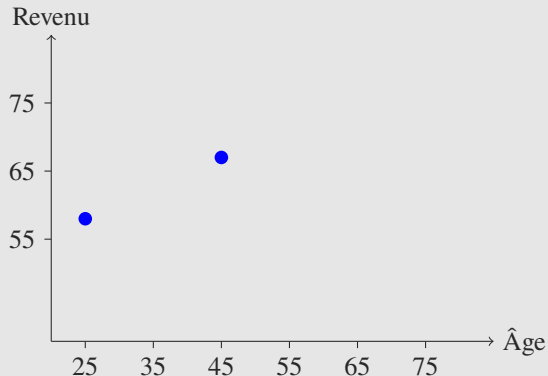
Nuage de points

Un nuage de points montre comment deux variables sont liées, généralement à un instant donné.

Un nuage de points a une variable sur l'axe Y et une sur l'axe X.

On place ensuite un point à la coordonnée (y, x) correspondant à chaque observation.

Nom	Âge	Revenu
Alice	25	58
Bob	45	67



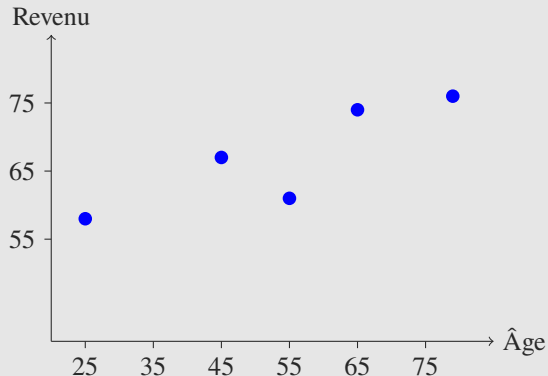
Nuage de points

Un nuage de points montre comment deux variables sont liées, généralement à un instant donné.

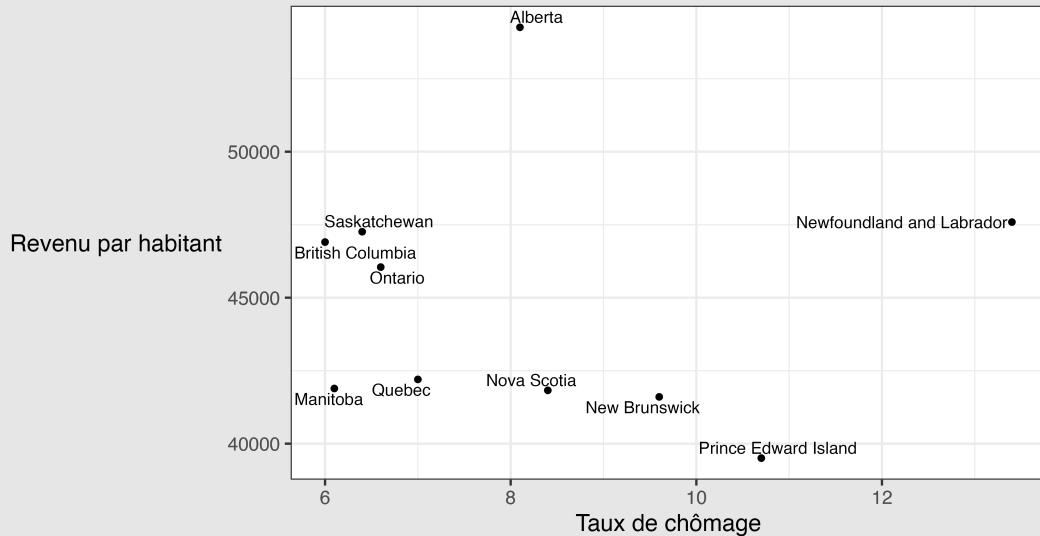
Un nuage de points a une variable sur l'axe Y et une sur l'axe X.

On place ensuite un point à la coordonnée (y, x) correspondant à chaque observation.

Nom	Âge	Revenu
Alice	25	58
Bob	45	67
Carol	55	61
Dave	65	74
Eve	70	76



Nuage de points : taux de chômage et revenu par habitant, 2016



Nuage de points : notes des étudiants dans deux cours



Histogramme

Un histogramme montre la fréquence à laquelle chaque valeur d'une variable apparaît dans les données.

- ▶ Axe X : valeurs de la variable
- ▶ Axe Y : nombre de fois que la valeur apparaît dans les données

Par exemple, si la variable est une note d'examen, la hauteur de chaque barre montre combien d'étudiants ont obtenu cette note.

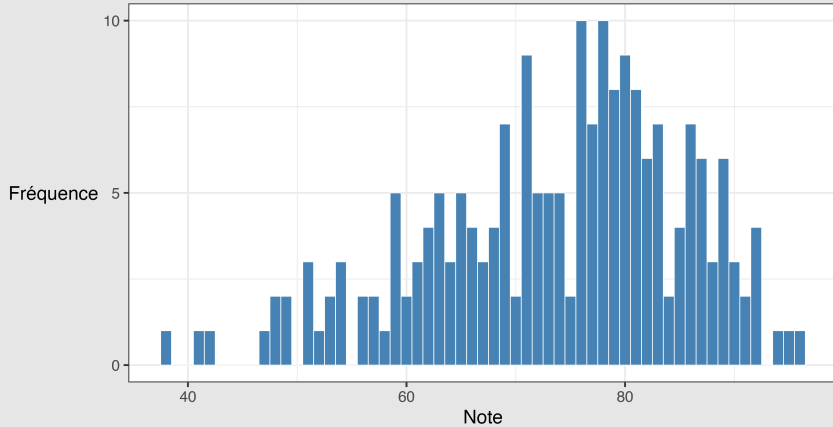
- ▶ Si 9 étudiants ont obtenu la note de 80, l'histogramme place une barre de hauteur 9 à 80

Parfois, on regroupe des valeurs proches en **classes**.

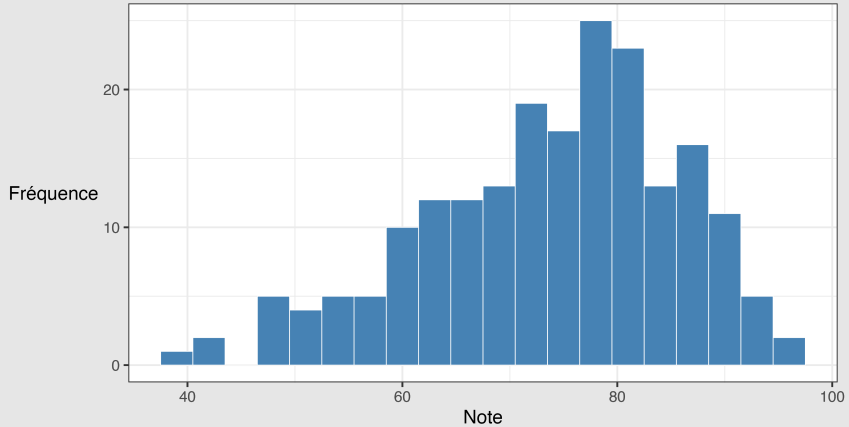
La **largeur de classe** détermine la largeur de chaque regroupement

- ▶ Si la largeur de classe = 3, on peut regrouper les notes 80, 81, 82 ensemble

Histogramme : distribution des notes dans une classe



Histogramme : distribution des notes dans une classe (largeur de classe = 3)



Caractéristiques d'une distribution

Comment décrieriez-vous la distribution des notes dans cette classe ?

En mots, on pourrait décrire

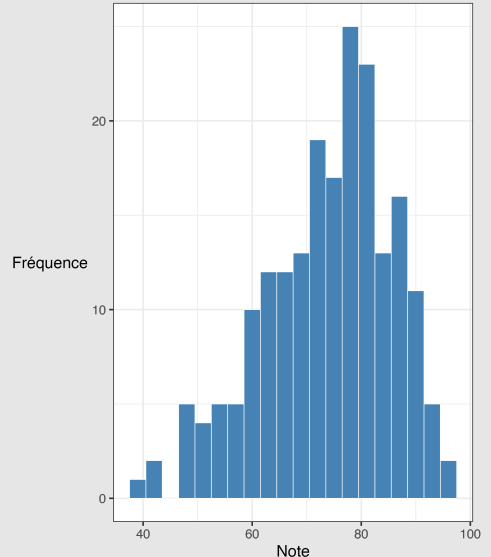
1. la tendance centrale

- ▶ La plupart obtiennent une note proche de 75

2. la dispersion

- ▶ Certains obtiennent une note > 90
- ▶ Certains obtiennent une note < 60

Comment être plus précis ?



Résumé des données avec statistiques descriptives

Une autre façon d'analyser rapidement les données est de calculer les **statistiques descriptives**.

Certaines statistiques, comme la **moyenne**, la **médiane** et le **mode**, renseignent sur la **tendance centrale** d'une variable.

- ▶ Quelle est la valeur typique ou centrale de la variable ?

D'autres, comme la **variance**, l'**écart-type** ou l'**intervalle interquartile**, renseignent sur la dispersion de la variable.

- ▶ Quelle est l'étendue des valeurs ?
- ▶ La variable prend-elle souvent des valeurs extrêmes ou reste-t-elle dans un intervalle étroit ?

Notation vectorielle

Nous ferons souvent référence à une colonne d'un ensemble de données par un **vecteur** x

$$x = (x_1, x_2, \dots, x_n)$$

x_i : valeur de x dans la i -ème ligne de l'ensemble de données

n : nombre de lignes dans l'ensemble de données

La **somme** des éléments de x s'écrit $\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$.

Dans l'ensemble de données ci-dessous $n = 3$ et nous pourrions écrire l'âge comme $x = (x_1, x_2, x_3) = (34, 29, 31)$

Nom	Sexe	Âge	Lieu de naissance	Éducation
Sam	M	34	Montréal, Canada	PhD
Alex	M	29	Toronto, Canada	BA
Marie	F	31	Québec, Canada	MA

$$\sum_{i=1}^n x_i = \sum_{i=1}^3 x_i = x_1 + x_2 + x_3 = 34 + 29 + 31 = 94$$

Mesures de tendance centrale

La tendance centrale décrit la **valeur typique** d'une variable.

Trois mesures courantes :

- ▶ **Moyenne** — moyenne arithmétique
- ▶ **Médiane** — valeur du milieu une fois triée
- ▶ **Mode** — valeur la plus fréquente

Calculer moyenne, médiane et mode

Moyenne (arithmétique)

La moyenne d'une variable x_1, x_2, \dots, x_n est

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Médiane (valeur centrale)

- ▶ Trier les données par ordre croissant.
- ▶ Si n est impair : valeur centrale.
- ▶ Si n est pair : moyenne des deux valeurs centrales.

Mode (valeur la plus fréquente)

- ▶ Valeur apparaissant le plus souvent.
- ▶ Si aucune répétition \rightarrow pas de mode.

Exemple : calculer moyenne, médiane, mode

Nom	Âge
Alice	25
Bob	45
Carol	55
Dave	65
Eve	70
Fran	70

Moyenne : $(25 + 45 + 55 + 65 + 70 + 70)/6 = 55$

Médiane : valeur centrale = $(55 + 65)/2 = 60$

Mode : 70 apparaît deux fois, toutes les autres valeurs une fois

Mesures de dispersion

La dispersion décrit à quel point les données sont **étalées**.

Mesures courantes :

- ▶ **Variance** — écart moyen au carré par rapport à la moyenne
- ▶ **Écart-type** — racine carrée de la variance
- ▶ **Intervalle interquartile (IQR)** — distance entre le 25^e percentile (Q_1) et le 75^e percentile (Q_3)

Variance ou écart-type élevé → données très dispersées ; faible → données proches.

Calculer variance, écart-type et IQR

Variance :

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

Écart-type :

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

Intervalle interquartile (IQR) :

$$\text{IQR} = Q_3 - Q_1$$

- ▶ Q_1 (premier quartile) = médiane de la moitié inférieure
- ▶ Q_3 (troisième quartile) = médiane de la moitié supérieure

Exemple : calculer variance et écart-type

Nom	Âge
Alice	25
Bob	45
Carol	55
Dave	65
Eve	70
Fran	70

Comme précédemment, $\bar{x} = 55$.

Variance

$$\sigma^2 = \frac{(25 - 55)^2 + (45 - 55)^2 + (55 - 55)^2 + (65 - 55)^2 + (70 - 55)^2 + (70 - 55)^2}{6} = 258.33$$

Écart-type

$$\sigma = \sqrt{258.33} \approx 16.07$$

Exemple : calculer IQR

Nom	Âge
Alice	25
Bob	45
Carol	55
Dave	65
Eve	70
Fran	70

- ▶ Q_1 = médiane de la moitié inférieure (25, 45, 55) = 45
- ▶ Q_3 = médiane de la moitié supérieure (65, 70, 70) = 70

$$\text{IQR} = Q_3 - Q_1 = 70 - 45 = 25$$

Quelle mesure rapporter ?

Moyenne :

- ▶ Très standard – toujours rapporter la moyenne de vos variables
- ▶ Mais la moyenne peut être influencée par des valeurs extrêmes

Médiane :

- ▶ Non influencée par les valeurs extrêmes

Mode :

- ▶ À utiliser quand la variable prend un nombre fini de valeurs

Variance et Écart-type :

- ▶ Écart-type dans les mêmes unités que les données, variance en unités au carré

Intervalle interquartile (IQR) :

- ▶ Mesure la dispersion du milieu 50% des données
- ▶ Moins affecté par les valeurs extrêmes ou distributions asymétriques

Prochain cours...

Aujourd'hui, tout ce que nous avons vu était très **concret**

Nous avons vu des exemples de jeux de données, des méthodes pour les visualiser, et des statistiques pour les résumer

Pour approfondir nos connaissances, nous devons d'abord explorer le monde **abstrait** de la **théorie des probabilités**

La théorie des probabilités a été initialement développée pour étudier les jeux de hasard, tels que les dés, pile ou face, ou le poker

Des mathématiciens comme **Blaise Pascal** et **Pierre de Fermat** étaient fascinés par les jeux aléatoires et les ont étudiés formellement

Cela était considéré comme une trivialité jusqu'à ce que, des siècles plus tard, on découvre que cette théorie pouvait également s'appliquer à l'analyse de phénomènes réels

La semaine prochaine : nous ferons une pause avec les données et étudierons ces jeux !

Annexe : Au-delà des mesures de tendance centrale et de dispersion

Moyenne pondérée

- ▶ Comme la moyenne, mais en utilisant des poids pour tenir compte de l'importance ou de la fréquence différente des observations
- ▶ Exemple : pondérer les ménages par le nombre d'habitants

Asymétrie (skewness)

- ▶ Mesure si la variable est symétrique et, sinon, dans quelle direction
- ▶ Asymétrie à droite : la plupart ont une épargne modeste mais certains sont milliardaires
- ▶ Asymétrie à gauche : le revenu du travail de la plupart des gens augmente de 2-3% par an, mais certaines personnes (qui perdent leur emploi et deviennent chômeurs) perdent 100%

Covariance et corrélation

- ▶ Décrit comment deux variables « évoluent ensemble »
- ▶ Corrélation positive : les personnes ayant plus d'expérience gagnent plus
- ▶ Corrélation négative : lorsque le PIB diminue, le chômage augmente généralement

Quantiles

- ▶ Comme la médiane ou les quartiles, mais pour n'importe quelle partie de la distribution
- ▶ 99e percentile du revenu = plus riche que 99% des gens et plus pauvre que 1% des gens

Moyenne pondérée

Parfois, toutes les observations ne sont pas également importantes. On peut attribuer des **poids** à chaque observation.

Moyenne pondérée :

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

où w_i est le poids de l'observation i .

- ▶ Si tous les $w_i = 1$, la moyenne pondérée revient à la moyenne classique.

Exemple : Moyenne pondérée

Ménage	Revenu	Nombre d'habitants
1	80	1
2	90	2
3	70	1

$$\bar{x}_w = \frac{1 \cdot 80 + 2 \cdot 90 + 1 \cdot 70}{1 + 2 + 1} = \frac{330}{4} = 82.5$$

Quelle est l'interprétation de la moyenne pondérée \bar{x}_w ?

Exemple : Moyenne pondérée

Ménage	Revenu	Nombre d'habitants
1	80	1
2	90	2
3	70	1

$$\bar{x}_w = \frac{1 \cdot 80 + 2 \cdot 90 + 1 \cdot 70}{1 + 2 + 1} = \frac{330}{4} = 82.5$$

Quelle est l'interprétation de la moyenne pondérée \bar{x}_w ?

- Revenu moyen par personne

Quelle est l'interprétation de la moyenne non pondérée $\bar{x} = (80 + 90 + 70)/3 = 80$?

Exemple : Moyenne pondérée

Ménage	Revenu	Nombre d'habitants
1	80	1
2	90	2
3	70	1

$$\bar{x}_w = \frac{1 \cdot 80 + 2 \cdot 90 + 1 \cdot 70}{1 + 2 + 1} = \frac{330}{4} = 82.5$$

Quelle est l'interprétation de la moyenne pondérée \bar{x}_w ?

- Revenu moyen par personne

Quelle est l'interprétation de la moyenne non pondérée $\bar{x} = (80 + 90 + 70)/3 = 80$?

- Revenu moyen par ménage

Asymétrie (skewness)

L'asymétrie mesure la **non-symétrie** d'une distribution autour de sa moyenne.¹

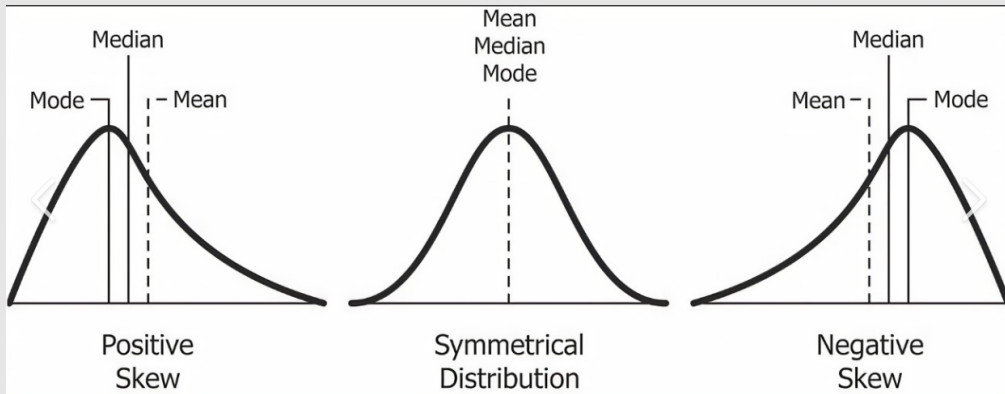
$$s = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3$$

- ▶ $s > 0$: asymétrie à droite (longue queue à droite)
- ▶ $s < 0$: asymétrie à gauche (longue queue à gauche)
- ▶ $s = 0$: distribution symétrique

Pourquoi cette mesure reflète-t-elle l'asymétrie ?

- ▶ $(x_i - \bar{x})^3$ est positif si $x_i > \bar{x}$ et négatif si $x_i < \bar{x}$
- ▶ Le cube amplifie les grandes différences entre x_i et \bar{x}
- ▶ Si $\bar{x} = 5$, $(7 - 5)^3 = 8$, $(10 - 5)^3 = 125$, $(125 - 5)^3 = 1728000$

¹Techniquement, il s'agit de la formule du « troisième moment centré ». L'asymétrie se réfère à une valeur standardisée où l'on divise s par $\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)^{3/2}$.



Exemple : Asymétrie (skewness)

Nom	Âge
Alice	25
Bob	45
Carol	55
Dave	65
Eve	70

Moyenne $\bar{x} = 52$, écart-type $\sigma \approx 16.24$

$$\begin{aligned}s &= \frac{1}{5}((25 - 52)^3 + (45 - 52)^3 + (55 - 52)^3 + (65 - 52)^3 + (70 - 52)^3) \\&= \frac{1}{5}((-27)^3 + (-7)^3 + 3^3 + 13^3 + 18^3) \\&= \frac{1}{5}(-19683 - 343 + 27 + 2197 + 5832) \\&= \frac{1}{5}(-11970) \\&= -2394\end{aligned}$$

Covariance et corrélation

La covariance et la corrélation mesurent la **relation entre deux variables** x et y .

Covariance :

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- ▶ Positive $\rightarrow x$ et y ont tendance à évoluer dans le même sens
- ▶ Négative $\rightarrow x$ et y ont tendance à évoluer dans des sens opposés
- ▶ Unités = unités de $x \times$ unités de y

Corrélation :

$$\rho_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

- ▶ Version standardisée de la covariance
- ▶ $-1 \leq \rho_{xy} \leq 1$
- ▶ 1 \rightarrow relation linéaire positive parfaite, -1 \rightarrow relation négative parfaite

Exemple : Covariance et corrélation

X (Heures d'étude)	Y (Score)
2	70
3	75
5	80
6	85

$$\bar{X} = 4, \quad \bar{Y} = 77.5$$

$$\sigma_X \approx 1.58, \quad \sigma_Y \approx 5.59$$

$$\text{Cov}(X, Y) = \frac{1}{4} \left((2 - 4)(70 - 77.5) + \dots + (6 - 4)(85 - 77.5) \right) = 8.75$$

$$\rho_{XY} = \frac{8.75}{1.82 \cdot 5.59} \approx 0.99$$

- Relation positive forte. Alors, étudiez bien !

Quantiles

Les **quantiles** sont des points qui divisent un jeu de données en intervalles de taille égale.

- ▶ Le p -ième quantile Q_p est la valeur en dessous de laquelle se trouvent $100 \cdot p\%$ des observations.
- ▶ Types courants de quantiles :
 - ▶ **Terciles** : divisent les données en 3 parties égales ($Q_{1/3}, Q_{2/3}$)
 - ▶ **Quartiles** : divisent les données en 4 parties égales ($Q_{1/4}, Q_{2/4}, Q_{3/4}$)
 - ▶ **Déciles** : divisent les données en 10 parties égales ($Q_{1/10}, Q_{2/10}, \dots, Q_{9/10}$)
 - ▶ **Percentiles** : divisent les données en 100 parties égales ($Q_{1/100}, Q_{2/100}, \dots, Q_{99/100}$)

Les quantiles sont une manière très flexible de décrire une distribution.

Plus de quantiles → description plus précise de la distribution

Mais la flexibilité a un coût : plus de nombres, moins facile à interpréter.

Plus simple de regarder la moyenne et l'écart-type (seulement deux nombres) que de regarder 10 déciles ou 100 percentiles, surtout pour comparer deux variables.

Exemple : Quantiles

ID	Score
1	42
2	45
3	48
4	50
5	52
6	53
7	55
8	56
9	57
10	58
11	60
12	61
13	62
14	64
15	66
16	68
17	70
18	72
19	75
20	78

► Quintiles :

- $Q_{1/5} = 50$
- $Q_{2/5} = 56$
- $Q_{3/5} = 61$
- $Q_{4/5} = 68$

► Déciles :

- $D_{(1/10)} = 45$
- $D_{(2/10)} = 50$
- $D_{(3/10)} = 53$
- $D_{(4/10)} = 56$
- $D_{(5/10)} = 58$
- $D_{(6/10)} = 61$
- $D_{(7/10)} = 64$
- $D_{(8/10)} = 68$
- $D_{(9/10)} = 72$