

CHAPITRE 2 - STATISTIQUE DESCRIPTIVE

Marouane IL IDRISSI
il_idrissi.marouane@uqam.ca

STT 1000 - Automne 2025
Département de Mathématiques, Université du Québec à Montréal



Statistique descriptive

Définition (*Statistique descriptive*).

La **statistique descriptive** est un **ensemble de méthodes** (représentations graphiques et calculs de caractéristiques numériques) permettant de **faire une synthèse statistique** à partir de **données**.

On s'intéresse à un **phénomène**

Le cyclisme à Montréal, la fission nucléaire, l'érosion des côtes...

On réalise des **expériences** et on collecte des données

Sondages, capteurs, mesures...

On **analyse ces données** pour **mieux comprendre le phénomène**

☞ C'est là qu'entre en jeu la **statistique descriptive**

Illustration en utilisant R et l'environnement de développement RStudio

Notes

Plan du chapitre

1. Analyse de données
2. Mesures de position
3. Mesures de dispersion
4. Représentations graphiques univariées
5. Représenter deux variables

Sommaire

Sommaire

1. Analyse de données
 - 1.1 Glossaire
 - 1.2 Types de variables
 - 1.3 Illustration sur R
2. Mesures de position
3. Mesures de dispersion
4. Représentations graphiques univariées
5. Représenter deux variables

Notes

Analyse de données - Variable, population, individu

Un **individu** est l'objet décrit par une donnée
Une personne, un assuré, un arbre...

Une **variable** est une **caractéristique** d'un individu
L'âge, le revenu, la côte de crédit, la circonférence...

La **population** est l'**ensemble des individus** que l'on souhaite étudier
Les habitants du Canada, les clients d'une compagnie, les arbres du parc national de Forillon...

Un **échantillon observé** est une **sous-partie de la population** dont les **variables sont mesurées**

Dans ce cours:

- On notera n la **taille de l'échantillon**
- On supposera que l'**échantillon est aléatoire**
⇒ Échantillon tiré au hasard dans la population, chaque individu a la même chance d'être tiré

Un **jeu de données**, c'est l'**échantillon observé** sous forme de **tableau de données**
Individus en ligne et variables en colonne

4/36

Notes

Analyse de données - Types de variables

Les **variables catégorielles** (ou **qualitatives**) **partitionnent** les individus en plusieurs groupes

- **Variables qualitatives nominales** ★
Pas d'ordre dans les modalités
Par ex: Couleur de la voiture 🚗
- **Variables qualitatives ordinales** ★
Il y a un ordre dans les modalités
Par ex: Niveau de satisfaction d'un client 🗳️

Les **variables quantitatives** prennent des **valeurs numériques**
On peut les additionner, les multiplier...

- **Variables quantitatives discrètes** ★
Prend des valeurs dans un ensemble dénombrable
Par ex: Nombre d'enfant dans une famille 🗳️
- **Variables quantitatives continues** ★
Prend des valeurs dans un intervalle
Par ex: Diamètre (cm) d'un arbre 🗳️

5/36

Analyse de données - Illustration sur R

Pour illustrer ce chapitre, nous allons étudier le **jeu de données**:

Données Ouverte de la Ville de Montréal - Inventaire des arbres publics de la ville

☞ Qu'est-ce qu'un individu? 🚩

Quelle est la population?

Les variables sont:

- **arrondissement:** Arrondissement de la ville de Montréal
- **lieu:** Type de lieu (parc, parterre, trottoir)
- **type:** Type d'arbre
- **dhp:** Diamètre à hauteur de poitrine de l'arbre (cm)
- **remarquable:** Arbre remarquable ou non
- **age:** Âge (années) depuis la plantation
- **maturité:** Nombre de décennies depuis la plantation

☞ Quel est le type de chacune de ces variables? 📌

Voyons ce que ça donne sur le logiciel R! 6/36

6/36

Sommaire

- ## 1. Analyse de données

- ### 2.1 Moyenne

- ## 2.2 Fréquence et proportion

- ### 2.3 Médiane

- 2.5 Moyenne tronquée

- ## 2.6 Mode

Notes

7/36

Mesures de position - Médiane

Pour calculer la **médiane** $q_{0.5}$ d'un **échantillon observé** x_1, \dots, x_n :

1. On **range par ordre croissant** les n observations $x_{(1)}, x_{(2)}, \dots, x_{(n)}$
 $\Leftrightarrow x_{(1)}$ est l'observation la plus petite et $x_{(n)}$ la plus grande
2. La **médiane** $q_{0.5}$ est égale à :
 - $x_{(\frac{n+1}{2})}$ si n est impair
 - $\frac{1}{2} (x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)})$ si n est pair

Interprétation: 50% des individus de l'échantillon ont une **valeur inférieure ou égale à la médiane** $q_{0.5}$ et 50% ont une valeur supérieure ou égale à la médiane.

Échantillon rangé par ordre croissant: 3, 8, 9, 9, 9, 9, 10, 11, 12, 14, 14, 15, 18

Quelle est la consommation médiane des voitures dans l'échantillon? 🚗

☞ Qu'arrive-t-il à la médiane dans l'exemple précédent si on ajoute une nouvelle observation égale à 6? 🐼

📌 La médiane est une mesure de position plus robuste que la moyenne ★📈10/36

Mesures de position - Quantile

La médiane est un cas particulier d'un quantile

Pour $\alpha \in (0, 1)$, le **quantile d'ordre α** , noté q_α , **de l'échantillon observé** est l'observation telle que $(\alpha \times 100)\%$ des observations **sont en dessous**

A partir d'un échantillon observé **rangé par ordre croissant** $x_{(1)}, x_{(2)}, \dots, x_{(n)}$, on a que ★

$$q_{\alpha} = (1 - \gamma) x_{(k)} + \gamma x_{(k+1)}$$

où $k = |(n+1)\alpha|$ et $\gamma = (n+1)\alpha - k$

Il y a d'autres méthodes pour calculer un quantile, **dans ce cours on utilisera celle-ci**

✎ Pour $\alpha = 0.5$, on retrouve **exactement la définition de la médiane** 📌

• Pour $\alpha = 1/4$, on parle de **premier Quartile** (souvent noté Q_1)

✎ Pour $\alpha = 3/4$, on parle de **troisième Quartile** (souvent noté Q_3)

☞ Échantillon observé: 3, 8, 9, 9, 9, 9, 10, 11, 12, 14, 14, 15, 18

Quelle est la valeur de Q_1 et Q_3 ? 📎

Notes

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Mesures de position - Moyenne tronquée

La **moyenne tronquée** à un niveau α (en général 0.1 ou 0.2) est la moyenne en omettant les valeurs les $(\alpha \times 100)\%$ les plus petites et les plus grandes:

$$\bar{X}_\alpha = \frac{X_{(\lfloor n\alpha+1 \rfloor)} + \cdots + X_{(n-\lfloor n\alpha \rfloor)}}{n-2\lfloor n\alpha \rfloor}$$

☛ Plus robuste que la moyenne de l'échantillon

On ne prend pas en compte les très petites et très grandes valeurs

Échantillon observé: 3, 8, 9, 9, 9, 9, 10, 11, 12, 14, 14, 15, 18

Quelle est la valeur de la moyenne tronquée à $\alpha = 0.1$? 🏠

12/36

Mesures de position - Mode

Le **mode** est l'observation qui **revient le plus souvent** dans l'échantillon

Échantillon observé: 3, 8, 9, 9, 9, 9, 10, 11, 12, 14, 14, 15, 18

Quelle est la valeur du mode?

Pour des des **variables aléatoires**:

- **Discrètes:** Le mode est la **valeur où la fonction de masse est la plus grande**
- **Continues:** Le mode est la **valeur où la densité est la plus grande**

Notes

This image shows a full page of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page, typical of notebook paper. There are no margins, text, or other markings on the page.

Sommaire

1. Analyse de données
2. Mesures de position
3. Mesures de dispersion
 - 3.1 Variance et Écart-type
 - 3.2 Étendue et écart interquartile
4. Représentations graphiques univariées
5. Représenter deux variables

Mesures de dispersion - Variance

Une **mesure de dispersion** est un indicateur de l'**étendue des valeurs** prises par une variable

La **variance** s^2 d'un **échantillon observé** x_1, \dots, x_n :

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

et l'**écart-type** d'un **échantillon observé** est $s = \sqrt{s^2}$.

D'où vient ce $n-1$? Réponse au Chapitre 3!

Définition (Centrer-réduire une variable). ★

Soit x_1, \dots, x_n un échantillon observé. **Centrer-réduire** cet échantillon c'est **retirer la moyenne, puis diviser chaque observation par l'écart-type**. C'est l'échantillon z_1, \dots, z_n , où, pour $i = 1, \dots, n$:

$$z_i = \frac{x_i - \bar{x}}{s} = \frac{x_i - \frac{1}{n} \sum_{j=1}^n x_j}{\sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}}$$

Quelle est la moyenne et la variance de l'échantillon z_1, \dots, z_n ? 15/36

D'où vient ce $n-1$? Réponse au Chapitre 3!

Définition (*Centrer-réduire une variable*). ★

Soit x_1, \dots, x_n un échantillon observé. **Centrer-réduire** cet échantillon c'est **retirer la moyenne, puis diviser chaque observation par l'écart-type**. C'est l'échantillon z_1, \dots, z_n , où, pour $i = 1, \dots, n$:

$$z_i = \frac{x_i - \bar{x}}{s} = \frac{x_i - \frac{1}{n} \sum_{j=1}^n x_j}{\sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}}$$

Définition (Centrer-réduire une variable). ★

Soit x_1, \dots, x_n un échantillon observé. **Centrer-réduire** cet échantillon c'est **retirer la moyenne, puis diviser chaque observation par l'écart-type**. C'est l'échantillon z_1, \dots, z_n , où, pour $i = 1, \dots, n$:

$$z_i = \frac{x_i - \bar{x}}{s} = \frac{x_i - \frac{1}{n} \sum_{j=1}^n x_j}{\sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}}$$

Soit x_1, \dots, x_n un échantillon observé. **Centrer-réduire** cet échantillon c'est **retirer la moyenne, puis diviser chaque observation par l'écart-type**. C'est l'échantillon z_1, \dots, z_n , où, pour $i = 1, \dots, n$:

$$z_i = \frac{x_i - \bar{x}}{s} = \frac{x_i - \frac{1}{n} \sum_{j=1}^n x_j}{\sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}}$$
$$z_i = \frac{x_i - \bar{x}}{s} = \frac{x_i - \frac{1}{n} \sum_{j=1}^n x_j}{\sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}}$$

Quelle est la moyenne et la variance de l'échantillon z_1, \dots, z_n ? 

Notes

Quelle est la moyenne et la variance de l'échantillon z_1, \dots, z_n ? 15/36

→ Quelle est la moyenne et la variance de l'échantillon z_1, \dots, z_n ?

Mesures de dispersion - Étendue et écart interquartile

L'**étendue** d'un échantillon observé x_1, \dots, x_n est la différence entre la plus grande et la plus petite valeur:

$$E = x_{(n)} - x_{(1)}$$

La **variance** et l'**étendue** ne sont pas robustes

Une très petite ou très grande valeur peut grandement les influencer

Pour résoudre ce problème, on peut utiliser l'écart interquartile

C'est la **différence** entre le troisième et le premier quartile:

$$I_{HQ} = Q_3 - Q_1$$

Échantillon observé: 3, 8, 9, 9, 9, 9, 10, 11, 12, 14, 14, 15, 18

Quelle est la valeur de la variance, de l'étendue et de l'écart interquartile? 🏷️

16/36

Sommaire

1. Analyse de données
2. Mesures de position
3. Mesures de dispersion
4. Représentations graphiques univariées
 - 4.1 Variables qualitatives
 - 4.2 Variables quantitatives
5. Représenter deux variables

Notes

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Jusqu'à présent, nous avons **résumé les informations contenues dans le jeu de données** à l'aide d'**indicateurs numériques**

Même si ces résumés sont **très précis**, il reste difficile de se faire une **idée globale de la manière dont sont distribuées les données**

Pour se faire une idée **plus générale**, on a recourt à des **représentations graphiques** pour représenter l'échantillon observé.

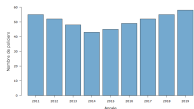
Ressource additionnelle pour vos révisions: [Chapitre 5 de la Formation Statistique Canada](#)

18/36

Représenter une variable qualitative - Diagramme à barres

Un **diagramme à barre** permet de représenter les **fréquences ou les proportions** des modalités d'une **variable qualitative**

Graphique 5.6.1
Nombre de policiers à Montréal, 2011 à 2019



Source: [Formation Statistique Canada](#)

❏ Que représente un individu dans ce jeu de données?

❏ Quelle est la variable qualitative que l'on étudie?

❏ Quelles sont ses modalités?

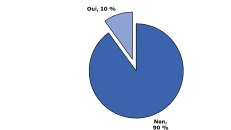
Les **diagramme à barre** permettent de **facilement comparer les occurrences des modalités dans un jeu de données**

19/36

Représenter une variable qualitative - Diagramme circulaire

Un **diagramme circulaire** remplit exactement la même fonction qu'un diagramme à barre
Comparer les **proportions des modalités** d'une **variable qualitative**

Graphique 9.4.1
Réponse des élèves et de la faculté à la question « Est-ce que les élèves de l'école Avenue devraient adopter l'uniforme? »



Source: Formation Statistique Canada

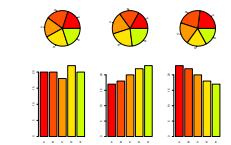
- ☞ Que représente un individu dans ce jeu de données?
- ☞ Quelle est la variable qualitative que l'on étudie?
- ☞ Quelles sont ses modalités?

Les **diagramme circulaire** permet une **comparaison rapide des modalités d'une variable qualitative dans un jeu de données**

20/36

Diagramme circulaire - Règles d'utilisation

Il est difficile pour le cerveau humain de jauger les proportions représentées dans un cercle



- Pas plus de 5 modalités (idéalement 2 ou 3)
- Représenter toutes les modalités (somme égale à 100)
- Les écarts en % entre les modalités ne doivent pas être négligeables
- Ranger les modalités de la moins fréquente à la plus fréquente dans le sens horaire
- Ne pas comparer deux diagrammes circulaires

Notes

[illegible]

21/36

Un **histogramme** permet de représenter la distribution des données d'une variable quantitative

Pour **créer un histogramme**:

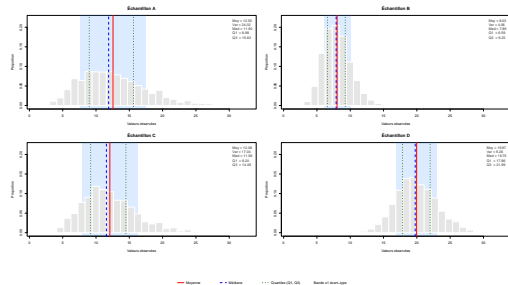
1. On **partitionne l'intervalle des valeurs de la variable**
Généralement en **sous-intervalle de même taille**, mais ce n'est pas obligatoire
2. On compte **combien d'observations tombent dans chaque partition**
3. On trace des **rectangles pour chaque partition**, dont la largeur est la **largeur de la partition**, et la hauteur est la **proportion d'observations dans la partition**

On peut aussi choisir comme hauteur des rectangles la **fréquence** ou la **densité**

☞ **Plus on a de données observées, plus le partitionnement peut être fin (partitions petites)**

☞ **L'histogramme approche la densité de la variable aléatoire que l'on a observé**

Plus de détails dans les cours de statistique plus avancés



Histogramme - Interprétation

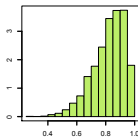
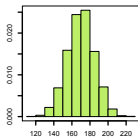
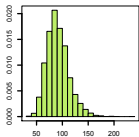
On dit que:

- Une distribution est **unimodale** si elle ne possède qu'un seul "pic" (mode)
- Une distribution est **symétrique** si les queues de distribution ont la même longueur
- Une distribution est **asymétrique à droite** si la queue droite de la distribution est plus longue que celle de gauche

Une distribution est **asymétrique à gauche** si la queue gauche de la distribution est plus longue que celle de droite

24/36

Histogramme - Symétrie



Commentez

Notes

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

25/36

Représenter une variable quantitative - Boîte à moustache

La **boîte à moustache** (boxplot) est une **représentation graphique** de la dispersion d'un échantillon observé

Figure 4.5.2.1 Construction d'une boîte à moustaches

The diagram illustrates the construction of a box plot. A central box is labeled 'BOÎTE' at the top. The box is divided into two equal horizontal sections by a vertical line. Below this line is the label 'MÉDIANE, Q₂'. The left edge of the box is labeled 'QUARTILE INFÉRIEUR, Q₁' and the right edge is labeled 'QUARTILE SUPÉRIEUR, Q₃'. Extending from the left edge of the box is a horizontal line labeled 'MOUSTACHE' at its end. Below this line is the label 'MINIMUM'. Similarly, extending from the right edge of the box is a horizontal line labeled 'MOUSTACHE' at its end. Below this line is the label 'MAXIMUM'. Arrows point from the labels 'QUARTILE INFÉRIEUR, Q₁', 'MÉDIANE, Q₂', and 'QUARTILE SUPÉRIEUR, Q₃' to their respective parts of the box plot.

Source: Formation Statistique Canada

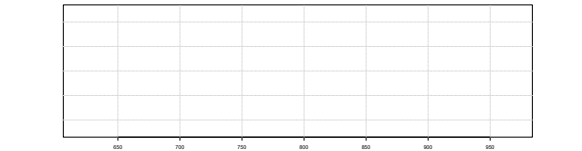
On peut **représenter les boîtes à moustache** verticalement ou horizontalement

Les points qui sont plus loin qu'une fois et demi l'écart interquartile sont considérés comme potentiellement aberrante ou extrêmes (outliers), et sont représentés en dehors de la boîte

Boîte à moustache - Exemple

Mesure de la vitesse de la lumière par Albert Michelson (1879) en m/s (moins 299 000):

880, 880, 880, 860, 720, 720, 620, 860, 970, 950, 880, 910, 850, 870, 840, 840, 850, 840, 840, 840

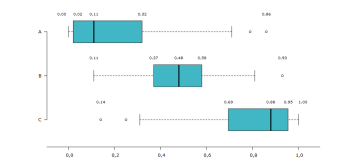


Notes

[illegible]

Boîte à moustache - Interprétation

Graphique 4.5.2.1
Boîtes à moustaches et résumés en cinq nombres des distributions A, B et C



Source: Formation Statistique Canada

Boîte à moustache:

Moins d'information qu'un histogramme, mais suffisamment pour pouvoir analyser la distribution de l'échantillon et propice à la comparaison entre variables ou échantillons

28/36

Sommaire

1. Analyse de données
2. Mesures de position
3. Mesures de dispersion
4. Représentations graphiques univariées
5. Représenter deux variables
 - 5.1 Nuage de point (quantit x quant)
 - 5.2 Tableaux croisés (qualit x qualit)
 - 5.3 Conditionnement (quantit x qualit)

Notes

This image shows a full page of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page, typical of notebook paper. There are no margins, text, or other markings on the page.

Jusqu'à présent, on a analysé les variables **une par une**

Mais on peut être intéressés à exhiber **le lien entre deux variables**

Lorsque l'on met en relation **deux variables** on parle de **statistiques descriptives bivariées**

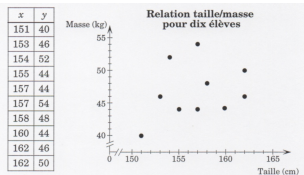
30/36

Représentation bivariées - Nuage de points

Un **nuage de point** permet de représenter la relation entre **deux variables quantitatives**

Cette représentation permet de **caractériser la nature du lien** entre deux variables

Linéaire, quadratique...



On verra (si le temps) au Chapitre 7 comment **mesurer la force d'un lien linéaire**

31/36

Nuage de point - Lien linéaire ou non-linéaire

On dit qu'il y a un **lien linéaire** entre deux variables si **on discerne une ligne droite** entre les points du nuage

Graphique 5.6.2
Relation linéaire ou relation non linéaire

The figure contains two side-by-side scatter plots, labeled A and B. Both plots have 'Variable X' on the x-axis and 'Variable Y1' on the y-axis, with scales from 0 to 100. Plot A, titled 'A. Linéaire', shows a clear positive linear trend where data points follow a straight line. Plot B, titled 'B. Non linéaire', shows a non-linear relationship where the data points form a curve that starts flat and then rises steeply.

A. Linéaire

Variable X	Variable Y1
10	5
20	20
25	22
35	35
40	25
45	45
50	42
55	58
60	58
70	72
80	78
85	88
90	92
95	100

B. Non linéaire

Variable X	Variable Y1
15	8
20	3
25	10
35	10
40	10
45	9
55	22
60	25
65	48
70	70
75	100

☞ Si on discerne **autre chose qu'une droite**, on dit que le **lien est non-linéaire**

32/36

Nuage de point - Lien linéaire positif ou négatif

Lorsque les deux variables croissent ensemble, on parle de **lien linéaire positif**

Si une variable diminue lorsque l'autre augmente, on parle de **lien linéaire négatif**

A. Positive

Variable X	Variable Y1
5	15
10	20
15	12
25	35
30	28
40	45
45	65
50	55
60	50
65	62
75	88

B. Négative

Variable X	Variable Y2
5	82
10	92
15	75
20	78
25	82
30	62
35	65
40	72
45	50
50	30
55	28
60	42
65	40
70	45
75	42
80	20
85	12

33/36

Notes

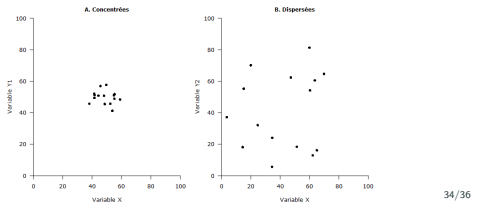
[illegible]

Nuage de point - Données concentrées et dispersées

Lorsque **les points sont rapprochés**, on dit que **les données sont concentrées**

Si les points sont étalés, on dit que les données sont dispersées

Graphique 5.6.4
Données concentrées ou données dispersées



34/36

Deux variables qualitatives - Tableau croisé

Lorsque l'on souhaite **comparer deux variables qualitatives**, on peut dresser un **tableau de fréquences croisées**

Pour une variable x à 3 modalités $\tilde{M}_1, \dots, \tilde{M}_3$, et une variable y à 4 modalités M_1, \dots, M_4 :

	M_1	M_2	M_3	M_4	Total
M_1					
M_2					
M_3					
Total					

Dans la **cellule relative** à la modalité \tilde{M}_j et M_k , on indique

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\tilde{M}_j}(x_i) \mathbb{1}_{M_k}(y_i)$$

C'est-à-dire le nombre d'individus ayant \tilde{M}_j et M_k comme modalités aux variables x et y

Que contiendront les cellules "Total"? 

Notes

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

35/36

Pour étudier le lien entre une **variable qualitative et quantitative**, on peut **conditionner ou stratifier** l'échantillon observé

Cela revient à:

1. Créer un **nouvel échantillon conditionnel/stratifié** par modalité observée
Sexe, profession...
2. Calculer les indicateurs ou représentation que l'on a vu jusqu'à présent **pour chaque nouvel échantillon conditionnel/stratifié**
3. Analyser les différences **entre les modalités**

This image shows a full page of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page, typical of notebook paper. There are no margins, text, or other markings on the page.