

Chapitre 2 : Tableaux statistiques et représentation graphiques

Statistiques 1 : Statistiques descriptives

ISA NUM 1

Enseignante: Elene ANTON - BALERDI

elene.anton-balerdi@univ-pau.fr

Campus Côte Basque – ISA

2024/2025



Hasard :

- Phénomènes déterministes : sont reliées aux expériences telles que si elles sont renouvelées dans des conditions totalement identiques, elles produiront le même résultat, qui devient donc prévisible.
- Phénomènes aléatoires : sont reliées aux expériences telles qu'effectuées dans des conditions totalement identiques elles donneront des résultats différents. Le résultat est non prévisible et on dit qu'il est dû au hasard. Exemples : lancement d'une pièce de monnaie, ...
- La nature du hasard est inconnue. L'objectif de la statistique est de maîtriser au mieux cette incertitude pour extraire des informations utiles des données, par l'intermédiaire de l'analyse des variations dans les observations.

1. MOTS CLÉS

- **Population** : Le groupe ou l'ensemble d'objets équivalents dans lequel on a les résultats de l'expérience.
- **Individu** : Un objet de la population est appelé individu.
- **Variables** : Les caractéristique sur lesquelles on a enquêté les individus de la populations. Les variables, comme dans le cours de probabilités prennent des valeurs dans un ensemble fondamental.
- **Échantillon** : En général, la population est trop vaste pour pouvoir être observée exhaustivement. On étudie alors la variable sur une sous partie de la population que on l'appelle échantillon.

1. MOTS CLÉS

- Étant donnée une expérience, on s'intéresse à analyser une variable X prenant des valeurs sur l'ensemble fondamental Ω , sur une population \mathcal{P} .
 - Exemple de : population, échantillon, individu et variable.
- Normalement, on ne peut pas observer les caractères de toute la population, mais d'un sous-ensemble de taille n :
 - La sous-population ou échantillon de \mathcal{P} de taille n est donnée par : $\{i_1, \dots, i_j, \dots, i_n\}$ où chaque individu est choisi au hasard.
 - l'échantillon des données ou des observations de la variable X pour les n individus est donné par $\{x_1, \dots, x_j, \dots, x_n\}$.

- **Echantillonnage :**

- Si l'échantillon est constitué de tous les individus de la population, on dit que l'on fait un recensement.
- Quand l'échantillon n'est qu'une partie de la population, on parle de sondage. Le principe des sondages est d'étendre à l'ensemble de la population les enseignements tirés de l'étude de l'échantillon. Pour que cela ait un sens, il faut que l'échantillon soit représentatif de la population. Il existe des méthodes pour y parvenir. Ça n'est pas l'objectif de ce cours.

1. MOTS CLÉS

- **Les variables** : Chaque individu est décrit par un ensemble des variables X . Ces variables peuvent être clarifiés en deux catégories selon leur nature :
 - Variable qualitative : s'exprimant par l'appartenance à une modalité :
 $\Omega = \{Homme, Femme\}, \Omega = \{Oui, Non\}, \Omega = \{\text{Noms de villes}\}, \dots$
 - Variable quantitative : s'exprimant par des nombres réels:
par exemple la taille des individus, les résultats d'un examen,..
 - Variables quantitatives discrètes : lorsque Ω est une suite finie ou infinie d'éléments.
 - Variables quantitatives continues : si toutes les valeurs d'un intervalle de \mathbb{R} sont acceptables.

1. MOTS CLÉS

- Dans ce chapitre, on ne s'intéresse qu'au cas où on ne mesure qu'une seule variable sur les individus, comme dans l'exemple des ampoules. On dit alors que l'on fait de la statistique unidimensionnelle.
- Quand on mesure plusieurs variables sur les mêmes individus, on dit que l'on fait de la statistique multidimensionnelle.

2. TABLEAUX STATISTIQUES

Supposons que on veut étudier une variable dans une population ou un échantillon. Alors, une fois les données collectées, la deuxième tâche est de les mettre en ordre dans ce qu'on appelle un tableau statistique ou un tableau de fréquences.

Le tableau statistique ou tableau de fréquences pour un échantillon ou population de taille n est composé de la façon suivante :

- Colonne 1 : la **variable** et ses valeurs classées du plus petit au plus grand :
 x_1, \dots, x_k , pour $x_i \in \Omega$, $i = 1, \dots, k$ (k est le nombre de valeur différentes observés dans le échantillon) .
- Colonne 2 : le nombre d'exemplaires ou **fréquence** de l'échantillon correspondant à chaque valeur de la variable l'absolu f_i .
 - Les fréquences f_i ont les propriétés suivantes :
 - $0 \leq f_i \leq n$, $i = 1, \dots, k$
 - $\sum_{i=1}^k f_i = n$.

2. TABLEAUX STATISTIQUES

- Colonne 3 : la **fréquence cumulée croissante** correspondant à chaque valeur de la variable $F_i \downarrow : F_i \downarrow = \sum_{j=1}^i f_j$.
 - Pour les fréquences cumulées croissantes $F_i \downarrow$ on a $F_k \downarrow = n$.
 - On peut aussi définir la **fréquence cumulée décroissante** correspondant à chaque valeur de la variable $F_i \uparrow : F_i \uparrow = \sum_{j=i+1}^n f_j = n - \sum_{j=0}^i f_j$. Dans ce cas, $F_1 \uparrow = n$.
- Les fréquences cumulées sont calculés uniquement pour les variables quantitatives. Ils n'ont pas du sens pour les variables qualitatives.

2. TABLEAUX STATISTIQUES

- Colonne 4 : la **fréquence relative** h_i . Elle correspond à $h_i = f_i/n$.

- Les fréquences relatives h_i ont les propriétés suivantes :

- $0 \leq h_i \leq 1, i = 1, \dots, k$

- $\sum_{i=1}^k h_i = 1.$

- Colonne 5 : la **fréquence relative cumulée croissante** $H_i \downarrow : H_i \downarrow = \sum_{j=1}^i h_j.$

- Pour les fréquences relatives cumulées croissantes $H_i \downarrow$ on a $H_k \downarrow = 1.$

- On peut aussi définir la fréquence relative cumulée décroissante correspondant à chaque valeur de la variable $H_i \uparrow$:

$$H_i \uparrow = \sum_{j=i+1}^n h_j = 1 - \sum_{j=0}^i h_j. \text{ Dans ce cas, } H_1 \uparrow = 1.$$

- Les fréquences relatives cumulées sont calculés uniquement pour les variables quantitatives. Ils n'ont pas du sense pour les variables qualitatives.

2. TABLEAUX STATISTIQUES

- Dans les tableaux statistiques on peut aussi donner les pourcentages ($h_i \times 100 \%$) et les pourcentages relatifs ($H_i \times 100\%$).

Enfin, le tableau statistique ressemble à ça :

Variable	Fréquences	Fréquences cumulées croissantes	Fréquences cumulées décroissantes	Fréquences relatives	Fréquences relatives cumulées croissantes
X	f_i	$F_i \downarrow$	$F_i \uparrow$	h_i	$H_i \downarrow$
x_1	f_1	f_1	$F_1 \uparrow = n$	h_1	h_1
x_2	f_2	$f_1 + f_2$	$n - f_1$	h_2	$h_1 + h_2$
.
.
.
x_k	f_k	$F_k \downarrow = n$	f_k	h_k	$H_k \downarrow = 1$
Total	n			1	

2. TABLEAUX STATISTIQUES

- EXEMPLE : Pour étudier la latence d'un virus, 90 poussins ont été inoculés. On examina le nombre de jours qui s'étaient écoulés jusqu'à l'apparition des premiers symptômes de la maladie pour chacun d'eux. Les données obtenues étaient les suivantes :

8	10	8	14	16	9	12	13	9	12	12	10	15	8	6
5	9	11	13	5	9	12	13	8	14	8	5	14	6	13
7	8	12	12	8	6	8	9	9	15	8	9	8	13	7
9	12	8	6	9	14	13	8	12	9	11	8	16	10	6
10	13	6	5	14	12	14	6	11	12	10	12	6	7	10
6	15	7	9	5	9	7	10	7	10	8	11	11	14	15

Donner le tableau statistique.

3. TABLEAUX STATISTIQUES : DONNÉES ASSOCIÉES

Données associées : Si le nombre de variables est trop grand, ou on ne connaît pas les valeurs spécifiques des variables on va associer les données.

La méthode :

- Le nombre des classes k est calculé à l'aide de la règle de Sturges :

$$k = 1 + 3.322 \log_{10} n \text{ (approximative)}$$

- Les classes sont définies par des intervalles : $]l_1, l_2],]l_2, l_3], \dots,]l_k, l_{k+1}]$.
Chaque donnée est classée sur une seule classe. Il est recommandable que toutes les classes ont la même taille, mais pas obligatoire.

- La marque de la classe : le point intermédiaire: $x_i = \frac{l_{i+1} + l_i}{2}, i = 1, \dots, k$

- La taille ou longueur de la classe : $a_i = l_{i+1} - l_i, i = 1, \dots, k$

3. TABLEAUX STATISTIQUES : DONNÉES ASSOCIÉES

Pour le exemple d'avant, on va :

1. Définir la plage de la variable : $16-5 = 11$.
2. Obtenir la valeur $k = 1 + 3.322 \log_{10} 90 = 7.49$. Alors, $k = 7$.
3. La talle de chaque classe va être environ de $plage/k = 11/7 = 1.57$ unites.
4. Avec ces données on décide de prendre de classes de talle $a_i = 2$. Et en faire 6 classes.

Nombre de jours	x_i	a_i	f_i	$F_i \downarrow$	h_i	$H_i \downarrow$
[4.5, 6.5)	5.5	2	14	14	0.1556	0.1556
[6.5, 8.5)	7.5	2	20	34	0.2222	0.3778
[8.5, 10.5)	9.5	2	20	54	0.2222	0.6
[10.5, 12.5)	11.5	2	16	70	0.1778	0.7778
[12.5, 14.5)	13.5	2	14	84	0.1556	0.9334
[14.5, 16.5)	15.5	2	6	90	0.0666	1
Total			90		1	

4. REPRÉSENTATIONS GRAPHIQUES

- **Variables qualitatives ou variables quantitatives discrètes :**
 - La différence fondamentale entre les représentations pour des variables qualitatives et quantitatives tient au fait qu'il existe un ordre naturel sur les modalités (qui sont des nombres réels) pour les variables quantitatives, alors qu'aucun ordre n'est prédéfini pour les variables qualitatives.
- **Variables qualitatives ou variables quantitatives discrètes :**
 - Diagrammes en bâtons
 - Diagrammes sectoriels (ou en camemberts)
- **Variables quantitatives continus :**
 - Histogramme
 - Histogramme cumulé

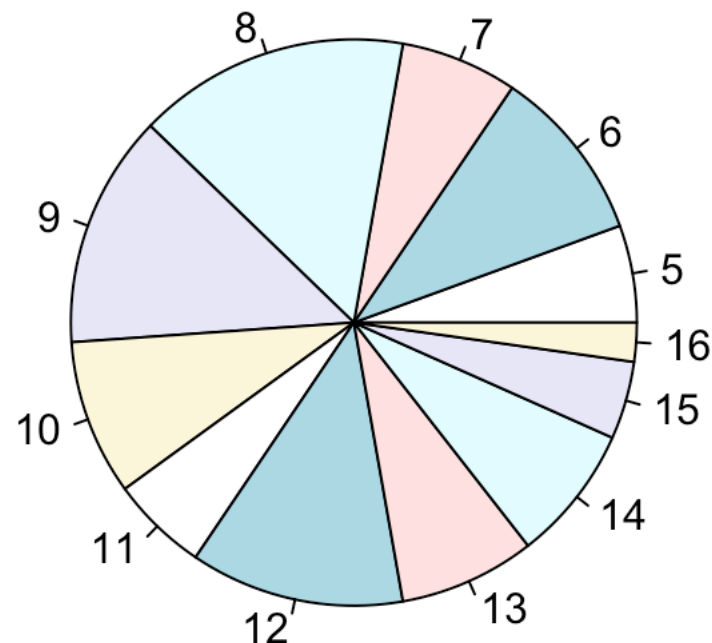
5. R. GRAPHIQUES : QUALITATIVES ET QUANTITATIVES DISCRETS

Variables qualitatives ou variables quantitatives discrètes :

- **Diagrammes sectoriels (ou en camemberts) :** à chaque modalité correspond un secteur de disque dont l'aire est proportionnelle à la fréquence relative de la modalité.

Tableau statistique

Nombre de jours	f_i	$F_i \downarrow$	$F_i \uparrow$	h_i
5	5	5	90	0.0556
6	9	14	85	0.1
7	6	20	76	0.0667
8	14	34	70	0.1555
9	12	46	56	0.1333
10	8	54	44	0.0889
11	5	59	36	0.0556
12	11	70	31	0.1222
13	7	77	20	0.0778
14	7	84	13	0.0778
15	4	88	6	0.0444
16	2	90	2	0.0222
Total	90			1



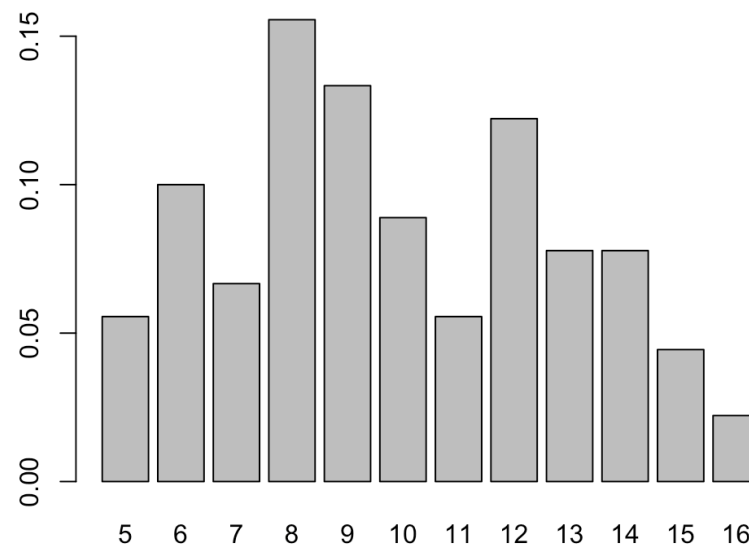
5. R. GRAPHIQUES : QUALITATIVES ET QUANTITATIVES DISCRETS

Variables qualitatives ou variables quantitatives discrètes :

- **Diagrammes en bâtons:** chaque modalité correspond un rectangle vertical dont la hauteur est proportionnelle à la fréquence relative de la modalité.

Tableau statistique :

Nombre de jours	f_i	$F_i \downarrow$	$F_i \uparrow$	h_i
5	5	5	90	0.0556
6	9	14	85	0.1
7	6	20	76	0.0667
8	14	34	70	0.1555
9	12	46	56	0.1333
10	8	54	44	0.0889
11	5	59	36	0.0556
12	11	70	31	0.1222
13	7	77	20	0.0778
14	7	84	13	0.0778
15	4	88	6	0.0444
16	2	90	2	0.0222
Total	90			1



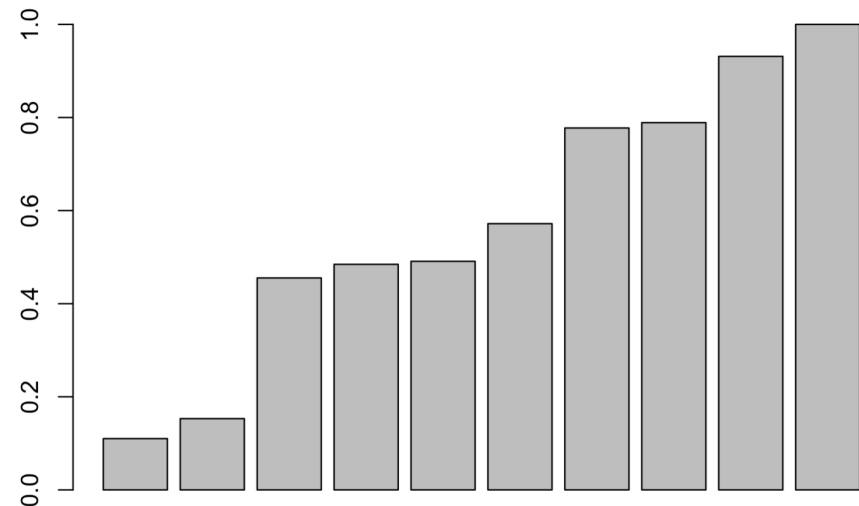
5. R. GRAPHIQUES : QUALITATIVES ET QUANTITATIVES DISCRETS

Variables quantitatives discrètes :

- **Diagramme en bâtons cumulé** : chaque modalité correspond un rectangle vertical dont la hauteur est proportionnelle à la fréquence relative **cumulé croissante** de la modalité.

Tableau statistique

Nombre de jours	f_i	$F_i \downarrow$	$F_i \uparrow$	h_i
5	5	5	90	0.0556
6	9	14	85	0.1
7	6	20	76	0.0667
8	14	34	70	0.1555
9	12	46	56	0.1333
10	8	54	44	0.0889
11	5	59	36	0.0556
12	11	70	31	0.1222
13	7	77	20	0.0778
14	7	84	13	0.0778
15	4	88	6	0.0444
16	2	90	2	0.0222
Total	90			1



6. R. GRAPHIQUES : QUANTITATIVES CONTINUES

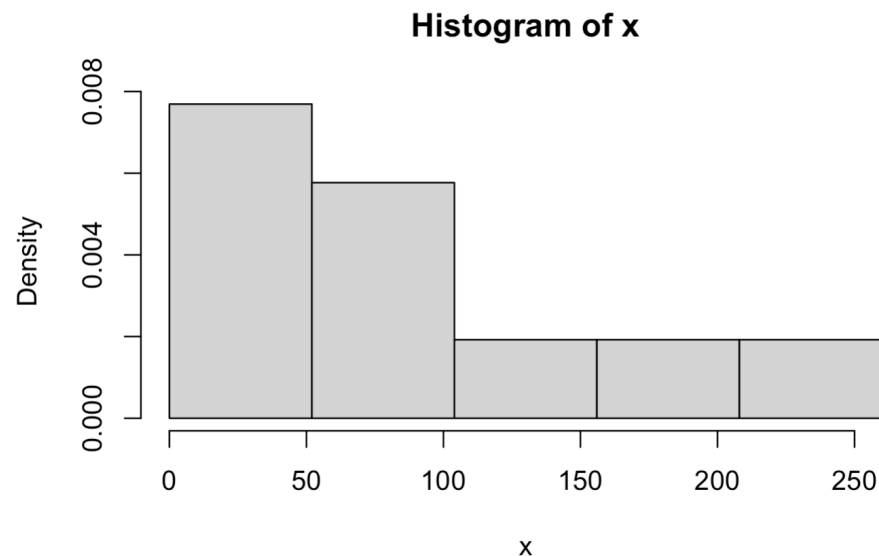
Variables quantitatives continues :

- Quand la variable étudiée est continue, les représentations du type diagramme en bâtons sont sans intérêt, car les données de x sont en général toutes distinctes, donc les effectifs tous égaux à 1.
- On considère deux types de représentations graphiques :
 - L'Histogramme
 - La fonction de répartition empirique.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : HISTOGRAMME

Histogramme : L'histogramme est la figure constituée de rectangles dont les bases sont les classes et dont les aires sont égales aux fréquences de ces classes.

- La représentation par histogramme consiste à regrouper les observations « proches » en classes.
- L'histogramme fournit bien une visualisation de la répartition des données.
- L'histogramme nous fourni une estimation de la densité des observations.



Histogramme : Comment generer un histogramme :

Considérons une ensemble d'observations x_1, \dots, x_n ordonnées de façon croissante.

- Étape 1 : Fixer une borne inférieur $a_0 < x_1$ et une borne supérieur $a_n > x_n$.
- Étape 2 : Partitioner l'intervalle des donnes $[a_0, a_n]$ en k intervalles $[a_{i-1}, a_i]$ appelés classes. La longueur de la classe i est $l_i = a_i - a_{i-1}$.
 - Si toutes les classes sont de la même longueur \longrightarrow histogramme à pas fixe
 - Si les classes ont de longueurs différentes \longrightarrow histogramme à pas variable
- Étape 3 : On calcule les fréquences f_i et fréquences relatives $h_i = f_i/n$ de chaque classe.
- Étape 4: Designer la figure : sur l'axe x on a les bases des rectangles, donnés par les longueurs classes. L'hauteur de la classe i est donne par h_i/l_i .

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : HISTOGRAMME

Remarque : L'histogramme dépend de plusieurs paramètres : les bornes inférieure et supérieure a_0 et a_n , le nombre k et la largeur h_i des classes.

- CONSEILS :

- Nombre des classes k : recommandé avoir entre 5 et 20 classes. Choisir à l'aide de la règle de Sturges.
- Borne sup a_n et inf a_0 : Doivent respecter une certaine homogénéité des largeurs de classes. Un choix fréquent est $a_0 = x_1 - 0.025(x_n - x_1)$ et $a_n = x_n + 0.025(x_n - x_1)$.
- Longueur des classes : Le choix le plus fréquent est l'histogramme à pas fixe; toutes les classes même largeur $l = (a_n - a_0)/k$.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : HISTOGRAMME

- **Exemple:** durée de vie des ampoules :

$$x = 91.6; 35.7; 251.3; 24.3; 5.4; 67.3; 170.9; 9.5; 118.4; 57.1$$

1. Mettre les variables dans l'ordre croissante :

$$x^* = 5.4; 9.5; 24.3; 35.7; 57.1; 67.3; 91.6; 118.4; 170.9; 251.3$$

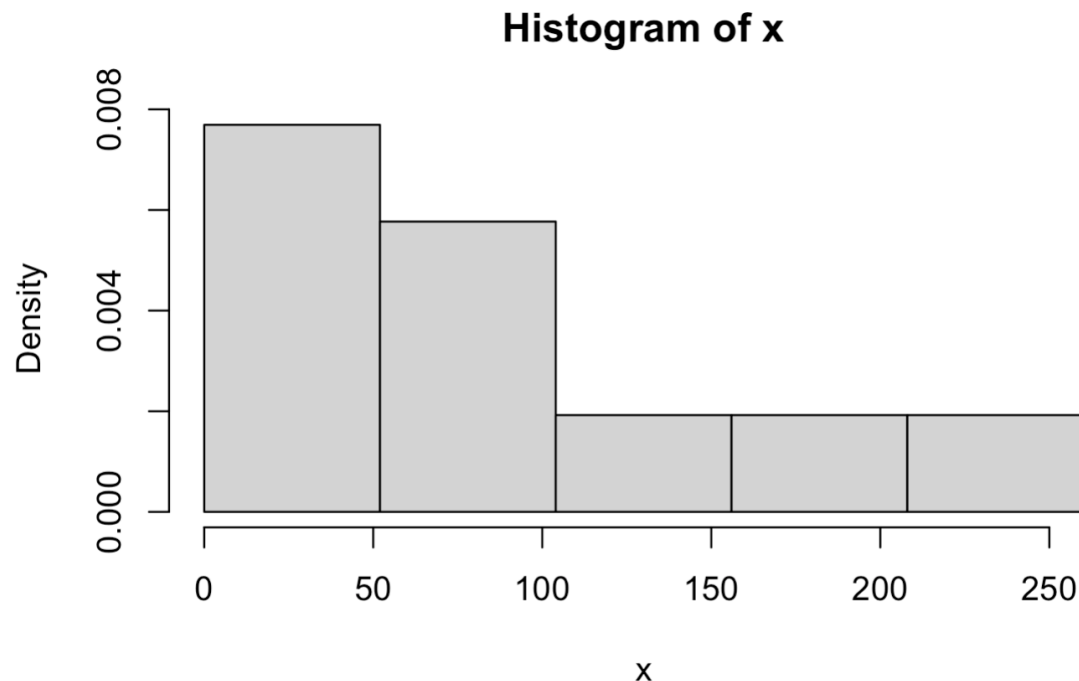
2. $n = 10$, $k = 5$, $a_0 = -0.74 \approx 0$, $a_6 = 257.4 \approx 260$, $l = 52 \forall i = 1, \dots, 5$

3.

Classes $]a_{i-1}, a_i]$	$]0, 52]$	$]52, 104]$	$]104, 156]$	$]156, 208]$	$]208, 260]$
Fréquences f_i	4	3	1	1	1
Fréquences relatives h_i	40%	30%	10%	10%	10%
Hauteurs - densité h_i/l	0.0077	0.0058	0.0019	0.0019	0.0019

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : HISTOGRAMME

- **Exemple:** L'histogramme
- La conclusion : La plus part des observations sont concentrés sur des petites valeurs de la durée de vie des ampoules. Et plus la durée de vie grandit, moins il y a d'observations.



6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : HISTOGRAMME

Remarques : **Approximation de la densité** :

- Notons \bar{f} la fonction en escalier constante sur les classes et qui est constante sur les classes et qui vaut h_i/l_i sur la classe $]a_{i-1}, a_i]$.

- L'aire du i ème rectangle est la fréquence relative de la classe i :

$$\text{aire} = \text{hauteur} \times l = \frac{h_i}{l} \times l = h_i = f_i/n = \int_{a_{i-1}}^{a_i} \bar{f}(x)dx$$

- La fréquence est défini aussi comme le pourcentage d'observations appartenant à la classe i . Donc c'est une estimation naturelle de la probabilité qu'une observation appartienne à cette classe. Cette probabilité est :

$$P(a_{i-1} < X \leq a_i) = F(a_i) - F(a_{i-1}) = \int_{a_{i-1}}^{a_i} f(x)dx$$

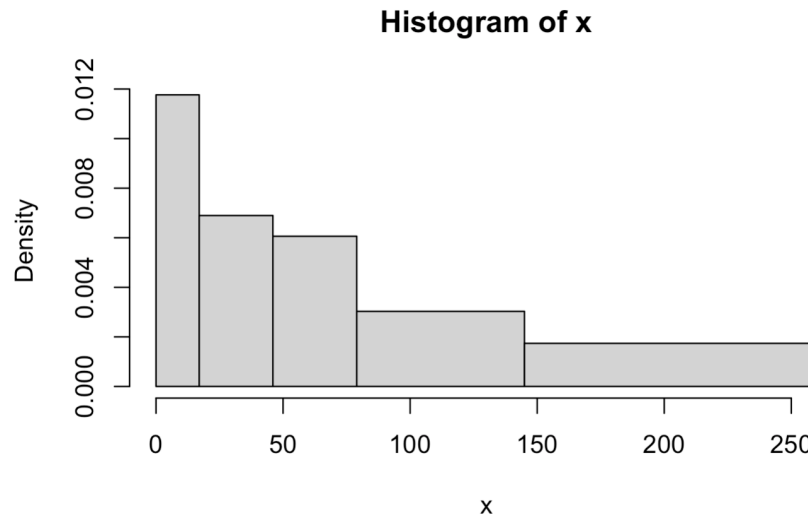
ALORS, aire rectangle = $P(a_{i-1} < X \leq a_i)$.

Remarques : Approximation de la densité II

- On en déduit que l'histogramme fournit une estimation de la densité des observations. L'estimation de la densité en un point x , $\hat{f}(x)$ est égale à la hauteur $\bar{f}(x)$ du rectangle correspondant à la classe à laquelle x appartient.
- La densité peut être interprétée aussi comment la probabilité qu'une voiture sélectionnée à partir des données se trouve dans cette fourchette de prix:
$$\text{hauteur} = h_i/l \text{ et } h_i \text{ est une approximatif de la probabilité .}$$
- L'allure de l'histogramme permettra donc de proposer des modèles probabilistes vraisemblables pour la loi de X en comparant la forme de \hat{f} à celle de densités de lois de probabilité usuelles.

Remarques : Pas fin vs. pas variable

- L'inconvénient d'un histogramme à pas fixe est que certaines classes peuvent être très chargées et d'autres pratiquement vides. Par exemple ici, la classe 1 contient plus d'observations à elle seule que les classes 3, 4 et 5 réunies.
- On pourrait avoir un histogrammes du même effectif (observations): $f_i = f$ pour tout $i = 1, \dots, k$ et ajouter les longueurs des classes selon le nombre des effectifs.

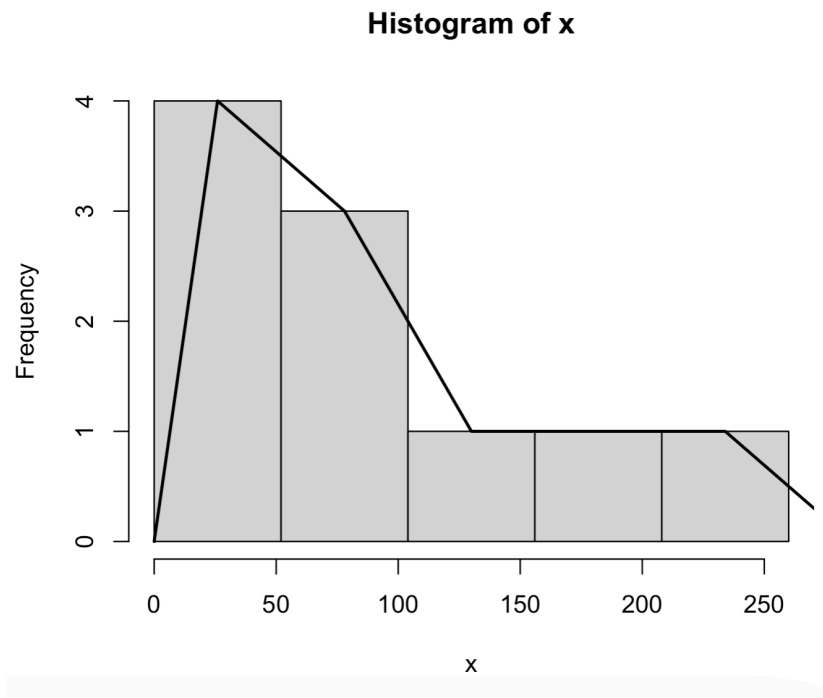


Remarques : Pas fin vs. pas variable II

- On voit que des histogrammes distincts sur les mêmes données peuvent être sensiblement différents.
- Donc il faudra se méfier des histogrammes si on veut estimer la densité des observations. On se contentera de dire que l'histogramme donne une allure générale de cette densité.

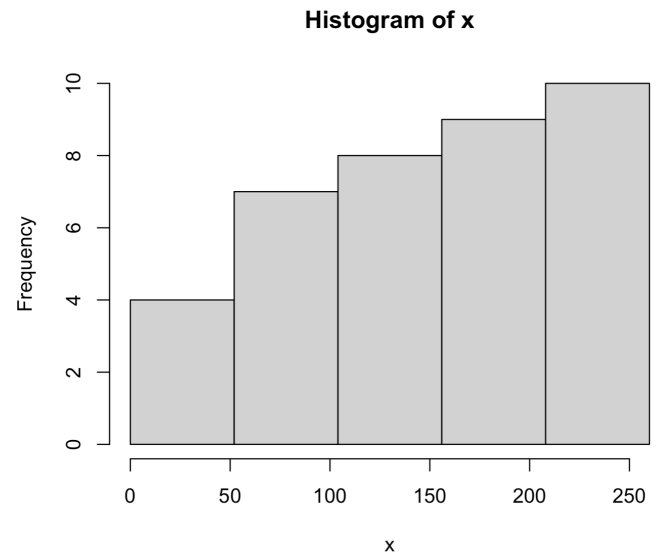
Éléments additionnelles :

- **Le polygone des fréquences** : s'obtiens en enchaînant les points intermédiaires des points supérieurs des rectangles de l'histogramme.



Éléments additionnelles :

- **Histogramme cumulé** : Remarque : Si au lieu des fréquences f_i , on considère les effectifs cumulés $m_i = \sum_{j=1}^i h_j$ on construit un histogramme cumulé, qui fournit une estimation de la fonction de répartition de la variable étudiée.



6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

Fonction de répartition empirique, FdRE, F_n associée à un échantillon x_1, \dots, x_n est la fonction définie par :

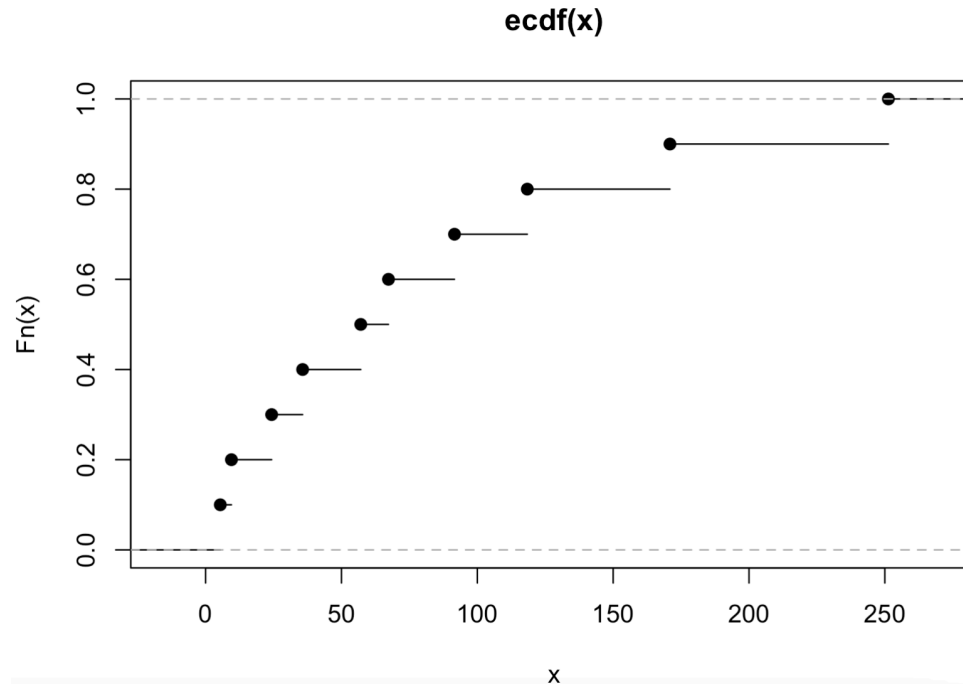
$\forall x \in \mathbb{R}, F_n(x) = \text{pourcentage d'observations inférieures à } x.$

$$\forall x \in \mathbb{R}, F_n(x) = \frac{1}{n} \sum_{j=1}^n 1_{(x_j \leq x)}.$$

- La fonction de répartition de X , $F(x) = P(x \leq X)$, donne la probabilité qu'une observation soit inférieure à x , tandis que $F_n(x)$ est le pourcentage d'observations inférieures à x . Par conséquent, $F_n(x)$ est une estimation de $F(x)$.
- $F_n(x)$ est une fonction en escalier qui fait des sauts de hauteur $1/n$ en chaque point de l'échantillon.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

Pour l'exemple des ampoules :



6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

- **L'utilité de la fonction de répartition empirique :**
 - Nous aider à déterminer un modèle probabiliste acceptable pour les observations de notre base des données.
 - Si nous essayons de tracer la FdRE et le comparer aux fonctions de répartition connues, nous n'allons pas arriver. Car, toutes les fonctions de répartition vont être très proches à la réponse, à vue d'œil. Nous allons le faire, alors, avec le graphe des probabilités.
- **Proposition :** Soit F la fonction de répartition d'une loi de probabilité, dépendant d'un paramètre inconnu λ . S'il existe des fonctions h , g , α et β telles que
$$\forall x \in \mathbb{R}, h[F(x)] = \alpha(\lambda)g(x) + \beta(\lambda),$$

alors le nuage des points $\{(g(x_i), h(1/i))\}_{i=1, \dots, n}$ est le **graphe de probabilités (qqplot en anglais)** pour la loi de fonction de répartition F .

Si les points du nuage sont approximativement alignés, on admettra que F est une fonction de répartition plausible pour les observations.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

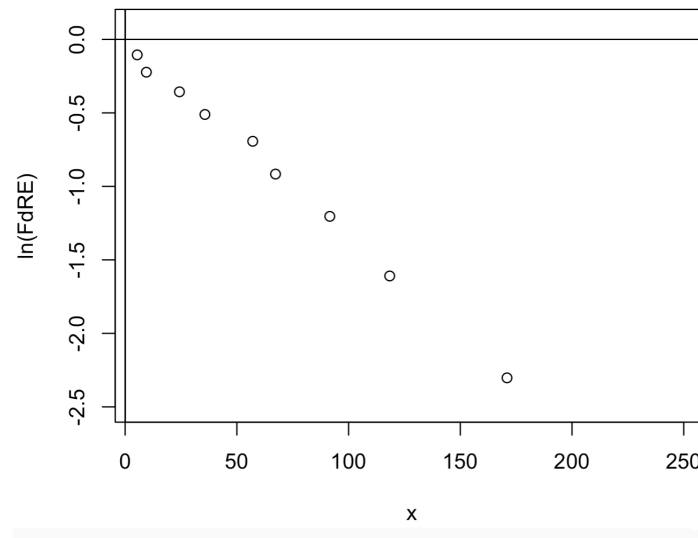
- **Dans la pratique :**

1. Appliquer une transformation à la fonction de répartition empirique qui permette de reconnaître visuellement une caractéristique d'une loi de probabilité.
2. Faire un **graphe de probabilités** (Q-Q plot).
3. Si les points sont approximativement alignés les observations proviennent d'une loi de probabilité bien précise.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

Exemple 1 : Voir si les données des durées de vie des ampoules suivent une loi exponentielle : La loi exponentielle $F(x) = 1 - e^{-\lambda x}$.

1. Transformation à appliquer : $\ln(1 - FdRE)$, et $\ln(1 - F(x)) = -\lambda x$.
2. Faire le graphique de probabilités qqplot : $\{(x_i, \ln(1 - i/n))\}_{i=1, \dots, n}$

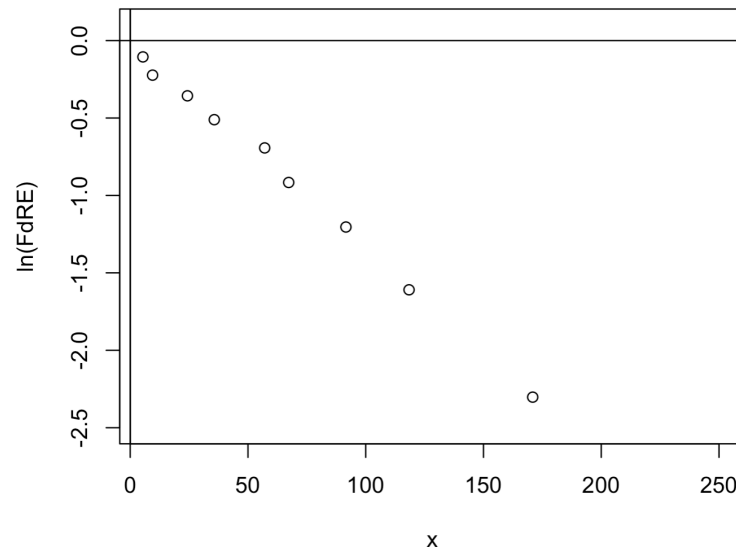


Si ces points sont approximativement alignés sur une droite de pente négative et passant par l'origine, on pourra considérer que la loi exponentielle est un modèle probabiliste vraisemblable pour ces observations. La pente de la droite fournit alors une estimation graphique de λ . Inversement, si ce n'est pas le cas, il est probable que les observations ne soient pas issues d'une loi exponentielle.

6. R. GRAPHIQUES : QUANTITATIVES CONTINUES : FONCTION DE REPARTITION EMPIRIQUE

Exemple 1 : Voir si les données des durées de vie des ampoules suivent une loi exponentielle : La loi exponentielle $F(x) = 1 - e^{-\lambda x}$.

3. Conclusion :



- Si ces points sont approximativement alignés sur une droite de pente négative et passant par l'origine, \implies on pourra considérer que la loi exponentielle est un modèle probabiliste vraisemblable pour ces observations. La pente de la droite fournit alors une estimation graphique de λ .
- Si ce n'est pas le cas, \implies il est probable que les observations ne soient pas issues d'une loi exponentielle.