

---

# Statistique descriptive

DUT Informatique, semestre 3

---

Version 1.0



1<sup>er</sup> septembre 2011

Ph. Roux

---

2009-2011

# Table des matières

<b>Table des matières</b>	<b>3</b>
1 Représentation des données statistiques . . . . .	4
2 Paramètres statistiques . . . . .	11
3 Corrélation . . . . .	17
4 Calculs statistiques avec Scilab . . . . .	22
4.1 séries statistiques discrètes . . . . .	22
4.2 séries statistiques regroupées par intervalles . . . . .	26
<b>Index</b>	<b>28</b>

## Avertissement

Pour bien utiliser ce polycopié, il faut le lire au fur et à mesure de l'avancement du cours magistral, et prendre le temps de refaire les exercices types qui y sont proposés.

- Les définitions et théorèmes sont numérotés suivant le même ordre que dans le cours magistral. Ils apparaissent dans un cadre grisé et sont en général suivis de leur démonstration, signalée par une barre dans la marge et un □ à la fin.
- La table des matières et l'index (à la fin du document) permettent de retrouver une notion précise dans ce polycopié.
- Les méthodes et techniques qui seront approfondies en TD sont signalées par un cadre (sans couleurs)
- Des exercices types corrigés, et surtout rédigés comme vous devriez le faire en DS, sont signalés par le symbole :



- Les erreurs et les confusions les plus fréquentes sont signalées dans des cadres rouges avec le symbole :



- Vous êtes libre de réutiliser le contenu de ce document sous les termes de la licence CC-BY-NC-SA



<http://creativecommons.org/licenses/by-nc-sa/3.0/deed.fr>

## 1 Représentation des données statistiques

La statistique est la partie des mathématiques qui s'intéresse à la détermination des caractéristiques d'un ensemble de données (généralement très important). Ici nous ne nous intéresserons pas aux problèmes posés par la collecte d'un grand nombre de données, mais seulement aux deux aspects suivants :

- le traitement des données collectées qui est appelé *la statistique descriptive*
- l'interprétation des données qui est appelée *l'inférence statistique*

Il existe deux types de données statistiques :

- les caractères qualitatifs (qui ne peuvent pas être mesurées) qui constitue une nomenclature, par exemple le sexe (Masculin/Féminin), les couleurs . . .
- les caractères quantitatifs (qui peuvent être mesurées) et que l'on peut représenter par un nombre réel, par exemple taille, le poids, . . .

Dans ce cours nous nous intéresserons aux caractères quantitatifs, qu'on va représenter par des variables statistiques.

**Définition 1.1 (série statistique)** Soit  $P$  un ensemble, appelé « population », alors un variable statistique est une application :

$$\begin{aligned} X : P &\longrightarrow \mathbb{R} \\ i &\longmapsto X(i) \end{aligned}$$

$X(P)$  est appelé univers image de  $X$  et on dira que :

- $X$  est discrète si  $X(P)$  est un ensemble discret (fini ou dénombrable)
- $X$  est continue si  $X(P)$  est un intervalle de  $\mathbb{R}$

$N = \text{Card}(P)$  sera appelé l'effectif total

Dans la pratique on s'intéresse surtout à l'univers image  $X(P)$ <sup>1</sup> et on identifie souvent une variable statistique avec la liste des valeurs  $X(i)$  prises par la variable. Dans ce cas on parle en général de *série statistique*.

**1.1 Exemples de séries statistiques** : on considère la population  $P$  composée des étudiants inscrits au semestre 1 du DUT INFORMATIQUE à Lannion (on admettra qu'on a un effectif total  $N = 100$ ) et nous définissons deux variables statistiques :

- $X$  = « nombre d'absence injustifiées (en maths) au semestre 1 »
- $Y$  = « taille (exprimée en cm) de chaque individu »

ce sont bien des caractères quantitatifs qui peuvent donc être décrits par une série statistique, qui correspondent aux listes de nombres :

$X = (2, 4, 2, 1, 0, 0, 3, 1, 3, 0, 2, 3, 0, 3, 1, 0, 2, 0, 2, 1, 0, 3, 1, 2, 2, 3, 1, 3, 2, 2, 0, 2, 0, 3, 2, 1, 2, 2, 3, 2, 0, 2, 0, 2, 0, 2, 1, 2, 0, 0, 3, 1, 0, 2, 0, 2, 1, 2, 2, 2, 1, 2, 0, 2, 1, 0, 0, 4, 1, 0, 3, 0, 2, 0, 0, 3, 1, 3, 0, 0, 1, 2, 1, 2, 6, 2, 2, 1, 1, 2, 3, 1, 0, 1, 1, 1, 3, 1, 3, 2)$

et

$Y = (188, 169, 181, 180, 159, 180, 164, 184, 177, 159, 187, 174, 170, 173, 175, 150, 171, 166, 173, 182, 167, 173, 174, 175, 166, 173, 180, 188, 165, 173, 178, 182, 175,$

<sup>1</sup> pour des raisons de confidentialité on ne peut pas en général identifier chaque individu de la population  $P$

186, 162, 193, 173, 184, 184, 184, 182, 187, 193, 161, 193, 169, 163, 179, 179, 184, 182, 182, 172, 175, 193, 170, 176, 166, 177, 163, 191, 171, 189, 183, 178, 197, 166, 187, 180, 172, 181, 167, 177, 177, 186, 174, 168, 182, 183, 182, 170, 186, 193, 167, 184, 159, 169, 192, 172, 167, 178, 176, 177, 175, 172, 175, 182, 176, 179, 188)  
 c'est ce qu'on appelle des données brutes.

Évidement il n'est pas très parlant de représenter une série statistiques sous forme de données brutes c'est pourquoi on va utiliser diverses représentations permettant de faire ressortir la répartition des données.

**Définition 1.2 (modalités)**

Soit  $X$  une variable statistique alors on appelle modalités de  $X$  :

- l'ensemble des valeurs différentes  $\{x_1; x_2; \dots; x_n\}$  prises par  $X$  si  $X$  est discrète
- un ensemble d'intervalles  $\{[x_1; x_2[; [x_2; x_3[; \dots; [x_{n-1}; x_n]$  formant une partition de  $X(P)$  si  $X$  est continue

un tableau statistique est un tableau regroupant les caractéristiques d'une série statistiques par modalités :

série discrète

série continue

modalités	$x_1$	$x_2$	$x_3$	...
effectif $n_i$				
fréquence $f_i$				

modalités	$[x_1; x_2[$	$[x_2; x_3[$	...
effectif $n_i$			
fréquence $f_i$			

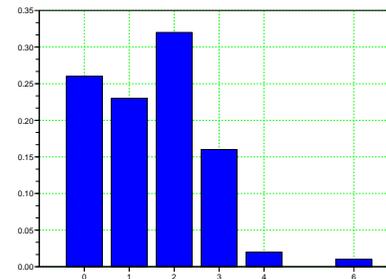
Si  $N$  est l'effectif total alors  $f_i = \frac{n_i}{N} \forall i = 1, \dots, n$ .

On représentera graphiquement les données d'une série discrète (resp. continue) en dessinant des colonnes de hauteurs (resp. surface) proportionnelle à l'effectif de chaque modalité, c'est ce qu'on appelle un diagramme en bâtons (resp. un histogramme).

 **1.2 tableau statistique d'une série statistique** Pour faire le tableau statistique de la série  $X$ , dont les modalités sont  $\{0; 1; 2; 3; 4; 6\}$ , il faudra à partir des données brutes construire le tableau :

$X_i =$	0	1	2	3	4	6	→ modalités
$n_i$	26	23	32	16	2	1	→ $\sum n_i = N = 100$
$f_i$	0.26	0.23	0.32	0.16	0.02	0.01	→ $\sum f_i = 1$

On pourra aussi représenter graphiquement la série statistique à l'aide d'un diagramme en bâtons :

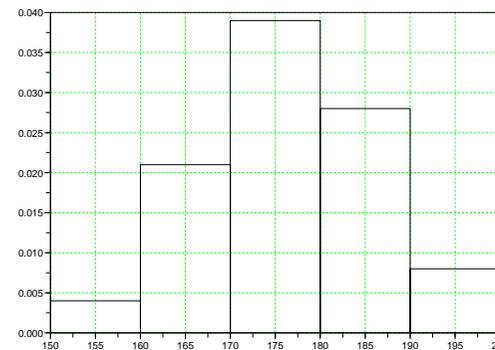


ou chaque colonne à une hauteur proportionnelle à l'effectif de la modalité.

Pour la série  $Y$ , même si les tailles ne prennent que des valeurs entières, il n'est pas opportun de la représenter comme une série discrète puisque dans ce cas on se retrouverait avec 34 modalités différentes ! C'est pour cette raison qu'on va considérer cette série comme une série continue prenant ses valeurs dans l'intervalle  $[150; 200]$ . Si on découpe cet intervalle en modalités de même taille on obtient :

$Y_i =$	$[150, 160[$	$[160, 170[$	$[170, 180[$	$[180, 190[$	$[190, 200[$
$n_i$	4	21	39	28	8
$f_i$	0.04	0.21	0.39	0.28	0.08

ce qui donne :



**Proposition 1.3 (effectifs corrigés)** Soit  $X$  une variable statistique continue alors ayant des modalités  $[x_1; x_2[ [x_2; x_3[ \dots [x_{n-1}; x_n[$  d'amplitude différentes alors pour chaque modalité on appelle « effectif corrigé » ou « hauteur corrigée » la hauteur qu'il faut donner dans l'histogramme à la colonne pour que sa surface soit proportionnelle à son effectif réel.

Dans un histogramme normalisé on calcule les hauteurs corrigées par la formule :

$$h_i = \frac{f_i}{|x_{i+1} - x_i|} = \frac{n_i}{N|x_{i+1} - x_i|}$$

de telle sorte que la surface totale des colonnes soit égale à 1.

**Preuve :** chaque modalité  $[x_i, x_{i+1}[$  à pour largeur  $|x_{i+1} - x_i|$  et pour hauteur  $h_i$  sa surface est donc :

$$h_i \times |x_{i+1} - x_i| = f_i = \frac{n_i}{N}$$

est bien proportionnel à l'effectif réel, de plus il est évident que la somme des aires est alors  $\sum_i f_i = 1$ .  $\square$

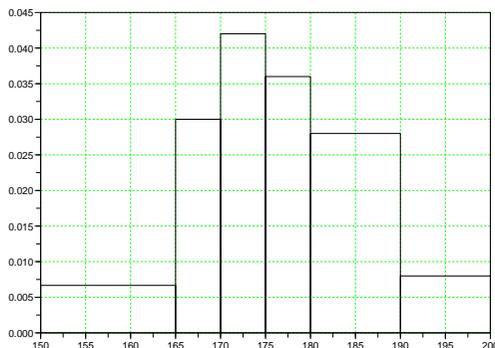
 **1.3** Reprenons le cas de la variable  $Y$  que l'on étudie avec les modalités :

$[150, 165[$ ,  $[165, 170[$ ,  $[170, 175[$ ,  $[175, 180[$ ,  $[180, 190[$ ,  $[190, 200[$

Après avoir calculé les effectifs de chaque modalité, il faut cette fois calculer les hauteurs corrigées :

$Y_i =$	$[150, 165[$	$[165, 170[$	$[170, 175[$	$[175, 180[$	$[180, 190[$	$[190, 200[$
$f_i$	0.1	0.15	0.21	0.18	0.28	0.08
$ x_{i+1} - x_i $	15	5	5	5	10	10
$h_i$	0.0066667	0.03	0.042	0.036	0.028	0.008

pour pouvoir dessiner l'histogramme



Un autre représentation permettant d'identifier la nature de la répartition des données statistiques est donnée par la fonction de répartition, qui est étroitement associée à la notion de cumul des fréquences.

**Définition 1.4 (Fonction de répartition)** Soit  $X$  une série statistique alors la fonction de répartition de  $X$

$$F: \mathbb{R} \rightarrow [0; 1]$$

$$x \mapsto \sum_{\{i|x_i \leq x\}} f_i$$

On retiendra que la fonction de répartition vérifie les conditions suivantes :

**Proposition 1.5** Soit  $F$  la fonction de répartition d'une série statistique  $X$  alors

- $F$  est croissante
- $\lim_{t \rightarrow -\infty} F(t) = 0$
- $\lim_{t \rightarrow +\infty} F(t) = 1$
- si  $X$  est une variable discrète alors  $F$  est constante par morceau.
- si  $X$  est une variable continue alors  $F$  on la représentera par une fonction continue et affine par morceau.

**Preuve :**

- si  $s \leq t$  alors

$$F(s) = \sum_{\{i|x_i \leq s\}} f_i \leq \sum_{\{i|x_i \leq t\}} f_i = F(t)$$

- si  $s < x_1$  (la première modalité) alors

$$F(s) = \sum_{\{i|x_i \leq s\}} f_i = 0 \implies \lim_{s \rightarrow -\infty} F(s) = 0$$

- inversement si  $t > x_n$  (dernière modalité) alors

$$F(t) = \sum_{\{i|x_i \leq t\}} f_i = f_1 + \dots + f_n = 1 \implies \lim_{t \rightarrow +\infty} F(t) = 1$$

Si  $X$  est une variable discrète, l'effectif d'une modalité  $] -\infty, x]$  ne va varier que lorsque  $x$  prend pour valeur une des modalités  $x_i$ . Il est donc normal que le graphe de  $F$  soit constant par morceau. Par contre la représentation du graphe de  $F$  par une fonction affine par morceau dans le cas où  $X$  est une variable continue, est une approximation qui se justifiera d'un point de vue pratique (voir les exemples).  $\square$

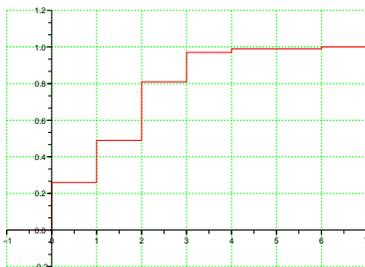
 **1.4** Pour construire la fonction de répartition de la variable  $X$  il faut donc faire le cumul des fréquences. Pour cela repartons du tableau statistique et ajoutons y une ligne pour calculer le cumul des fréquences :

$X_i =$	0	1	2	3	4	6
$f_i$	0.26	0.23	0.32	0.16	0.02	0.01
cumul	0.26	0.49	0.81	0.97	0.99	1

Cette dernière ligne donne la valeur de la fonction de répartition sur l'intervalle qui démarre à la modalité correspondante :

$I$	$] -\infty, 0[$	$[0, 1[$	$[1, 2[$	$[2, 3[$	$[3, 4[$	$[4, 6[$	$[6, +\infty[$
$F(t)$	0	0.26	0.49	0.81	0.97	0.99	1

ce qui donne :



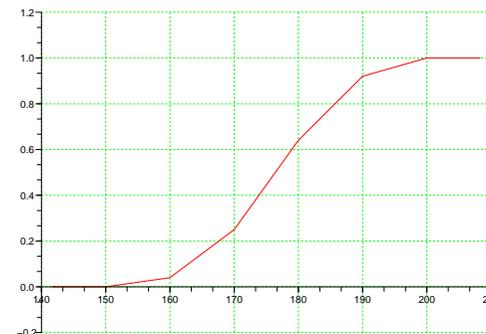
de même pour  $Y$  on repart du tableau statistique :

$Y_i =$	$[150, 160[$	$[160, 170[$	$[170, 180[$	$[180, 190[$	$[190, 200[$
$f_i$	0.04	0.21	0.39	0.28	0.08
cumul	0.04	0.25	0.64	0.92	1

permet de calculer les valeurs de la fonction de répartition :

$t$	$] -\infty, 150[$	160	170	180	190	200	$[200, +\infty[$
$F(t)$	0	0.04	0.25	0.64	0.92	1	1

par contre on doit représenter la fonction de répartition comme une fonction continue affine par morceau (*i.e.* son graphe est une suite de segments de droites) entre chaque borne des modalités :



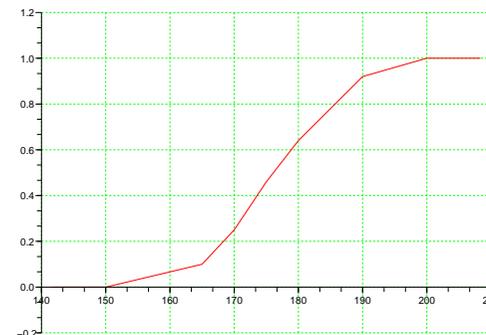
ce procédé permet d'obtenir des résultats numériquement cohérents quand on change les modalités lors de l'étude d'une série statistique. Reprenons le cas de  $Y$  avec des modalités de taille variable, ici on aura pas de problème d'effectifs corrigés, le tableau des fréquences devient :

$Y_i =$	$[150, 165[$	$[165, 170[$	$[170, 175[$	$[175, 180[$	$[180, 190[$	$[190, 200[$
$f_i$	0.1	0.15	0.21	0.18	0.28	0.08
cumul	0.1	0.25	0.46	0.64	0.92	1

ce qui donne cette fois pour la fonction de répartition :

$t$	$] -\infty, 150[$	165	170	175	180	190	200	$[200, +\infty[$
$F(t)$	0	0.1	0.25	0.46	0.64	0.92	1	1

et pour le graphe :



## 2 Paramètres statistiques

L'analyse des données statistiques doit permettre d'isoler certaines valeurs remarquables définissant la manière dont sont réparties les données statistiques. Le plus connu est la moyenne empirique, mais il y a aussi la variance, la corrélation ...

⚠ Dans le cas d'une variable continue, il n'est pas possible de calculer la valeur exacte des paramètres statistiques, puisqu'on ne connaît plus exactement les valeurs valeurs mais seulement leur nombre dans un intervalle donné! Pour chaque indicateur statistique on aura donc deux définitions différentes : une dans le cas des séries discrètes et l'autre dans le cas des séries continues regroupées par intervalles.

La moyenne permet de cerner facilement la valeur centrale de la série statistique.

### Définition 2.1 (Moyenne empirique)

Soit  $X$  une série statistique alors la moyenne de  $X$ , notée  $\bar{X}$ , est la valeur :

$$\bar{X} = \underbrace{\frac{1}{N} \sum_{k \in P} X(k)}_{\text{données brutes}} = \underbrace{\frac{1}{N} \sum_{i=1}^n n_i x_i}_{\text{effectifs par modalités}} = \underbrace{\sum_{i=1}^n f_i x_i}_{\text{fréquences par modalités}}$$

Dans le cas d'une série continue regroupée par intervalle on a la définition :

### Définition 2.2 (Moyenne approchée)

Soit  $X$  une variable continue regroupée suivant les modalités  $[x_1; x_2[ [x_2; x_3[ \dots [x_{n-1}; x_n]$ . Si on pose  $x_c = \frac{x_{i+1} + x_i}{2}$  le milieu de chaque modalité alors la moyenne approchée est définie par

$$\text{moyenne approchée} = \sum_{i=1}^n f_i \times x_c$$

✎ 2.1 Pour une variable discrète  $X$  on peut calculer la moyenne directement à partir du tableau statistique :

$X_i =$	0	1	2	3	4	6
$n_i$	26	23	32	16	2	1
$f_i$	0.26	0.23	0.32	0.16	0.02	0.01

en faisant

$$\bar{X} = 0.26 \times 0 + 0.23 \times 1 + 0.32 \times 2 + 0.16 \times 3 + 0.02 \times 4 + 0.01 \times 6 = 1.49$$

de même pour la variable  $Y$  on trouvera  $\bar{Y} = 176.71$ . Mais dans le cas d'une variable continue comme  $Y$ , si on a pas accès aux données brutes on peut quand même calculer une valeur approchée de la moyenne à partir du tableau des fréquences :

$Y_i =$	[150, 160[	[160, 170[	[170, 180[	[180, 190[	[190, 200[
$f_i$	0.04	0.21	0.39	0.28	0.08

$$\begin{aligned} \text{moyenne approchée} &= 0.04 \times 155 + 0.21 \times 165 + 0.39 \times 175 + 0.28 \times 185 + 0.08 \times 195 \\ &= 176.5 \approx 176.71 = \bar{Y} \end{aligned}$$

la valeur est proche de la moyenne calculée avec les données brutes. Si on change les modalités on va trouver un résultat proche :

$Y_i =$	[150, 165[	[165, 170[	[170, 175[	[175, 180[	[180, 190[	[190, 200[
$f_i$	0.1	0.15	0.21	0.18	0.28	0.08

cette fois on trouve :

$$\begin{aligned} \text{moyenne approchée} &= 0.1 \times 157.5 + 0.15 \times 167.5 + 0.21 \times 172.5 + 0.18 \times 177.5 + 0.28 \times 185 + 0.08 \times 195 \\ &= 176.45 \approx 176.71 = \bar{Y} \end{aligned}$$

On retiendra aussi que si on manipule plusieurs séries statistiques sur la même population alors la moyenne est linéaire :

### Proposition 2.3 (linéarité de la moyenne)

Soit  $X, Y$  deux séries statistiques (sur une même population) et  $a, b \in \mathbb{R}$  alors

$$\overline{aX + bY} = a\bar{X} + b\bar{Y}$$

en particulier  $\overline{aX + b} = a\bar{X} + b$ .

La moyenne n'est pas suffisante pour comprendre la manière dont sont réparties les valeurs d'une série statistique  $X$ . On a besoin d'avoir une indication sur la dispersion des données. Le meilleur de ces indicateurs est l'écart-type qui est défini comme suit.

**Définition 2.4 (variance et écart-type)** Soit  $X$  une série statistique de moyenne  $\mu = \bar{X}$  alors on appelle variance de  $X$  la valeur

$$\text{Var}(X) = \overline{(X - \mu)^2}$$

et l'écart-type de  $X$  est  $\sigma_X = \sqrt{\text{Var}(X)}$ .

On remarquera qu'on a toujours  $\text{Var}(X) \geq 0$  puisque c'est la moyenne de nombres positifs et l'écart-type est donc toujours bien défini!

⚠ L'écart-type  $\sigma$  donne une estimation de l'intervalle  $[\bar{X} - \sigma; \bar{X} + \sigma]$  centré sur la moyenne où sont concentrées les 2/3 des données. il permet donc de comprendre à quel point les données sont concentrées (ou pas) autour de la moyenne.

Pour une variable continue si on a pas accès aux données brutes on pourra aussi calculer une variance et un écart-type approché (comme pour la moyenne approché).

Dans la pratique on calculera la variance d'une série via le théorème suivant.

**Théorème 2.5 (Koening)**

Soit  $X$  une variable statistique alors

$$\text{Var}(X) = \overline{X^2} - \bar{X}^2$$

**Preuve :** On a que  $(X - \mu)^2 = X^2 - 2\mu X + \mu^2$  donc en utilisant la linéarité de la moyenne on obtient

$$\begin{aligned} \text{Var}(X) &= \overline{(X - \mu)^2} \\ &= \overline{X^2 - 2\mu X + \mu^2} \\ &= \overline{X^2} - 2\mu \bar{X} + \mu^2 \\ &= \overline{X^2} - 2\bar{X} \times \bar{X} + \bar{X}^2 \quad \text{car } \mu = \bar{X} \\ &= \overline{X^2} - \bar{X}^2 \end{aligned}$$

□

⚠ Dans la démonstration de la formule de Koening si on remplace  $\mu$  par  $\bar{X}$  il faut bien faire attention que :

$$\overline{\bar{X}^2} = \bar{X}^2 \neq \overline{X^2} \quad \text{et} \quad \overline{\bar{X}X} = \bar{X} \bar{X} = \bar{X}^2 \neq \overline{X^2}$$

📎 **2.2** Calculons la variance de  $X$  à partir du tableau statistique :

$X_i =$	0	1	2	3	4	6
$n_i$	26	23	32	16	2	1
$f_i$	0.26	0.23	0.32	0.16	0.02	0.01

on a déjà obtenu que  $\bar{X} = 1.49$  de même on obtient

$$\begin{aligned} \overline{X^2} &= 0.26 \times 0 + 0.23 \times 1 + 0.32 \times 4 + 0.16 \times 9 + 0.02 \times 16 + 0.01 \times 36 \\ &= 3.63 \end{aligned}$$

donc en utilisant le théorème de Koening on en tire pour la variance  $\text{Var}(X) = 3.63 - 1.49^2 = 3.63 - 2.2201 = 1.4099$ . On aurait trouvé le même résultat en calculant :

$$\begin{aligned} \text{Var}(X) &= 0.26 \times (0 - 1.49) + 0.23 \times (1 - 1.49) + 0.32 \times (2 - 1.49) \\ &\quad + 0.16 \times (3 - 1.49) + 0.02 \times (4 - 1.49) + 0.01 \times (6 - 1.49) \\ &= 1.4099 \end{aligned}$$

dans tous les cas on trouve l'écart-type  $\sigma_X = \sqrt{1.4099} = 1.1873921$ .

Pour les cas où on a pas directement accès aux données brutes (en particulier pour les variables statistiques continues) il existe d'autres indicateurs de centralité et de dispersion qui peuvent être calculés à partir de l'histogramme ou du graphe de la fonction de répartition .

**Définition 2.6 (Mode et médiane)**

Soit  $X$  une série statistique alors

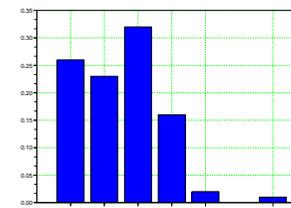
- Le mode est la modalité la plus représentée (effectif maximum)
- la médiane  $m$  est la valeur qui permet de couper la population en 2 part égales :  $\{i|X(i) \geq m\}$  et  $\{i|X(i) \leq m\}$ 
  - si l'effectif total  $N = 2p + 1$  est impair alors la médiane est la  $p + 1^{\text{ème}}$  valeur
  - si l'effectif total  $N = 2p$  est pair alors la médiane est la moyenne des  $p^{\text{ème}}$  et  $p + 1^{\text{ème}}$  valeurs

On peut aussi caractériser la médiane à partir de la fonction de répartition, ce qui permet d'en calculer une valeur approchée si on a pas accès aux données brutes.

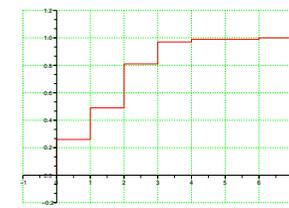
**Proposition 2.7 (médiane et fonction de répartition)** Soit  $X$  une série statistique alors et  $F$  la fonction de répartition associée comme étant la solution de  $F(x) = 0.5$  que l'on pourra déterminer graphiquement pour une variable statistique continue. Pour une variable statistique discrète la médiane il y a 2 cas :

- la première valeur telle que  $F(x) > 0.5$  (si  $F(x) = 0.5$  n'a pas de solution)
- la moyenne des valeurs telle que  $F(x) = 0.5$  (si  $F(x) = 0.5$  a plusieurs solutions)

📎 **2.3** On reprend la variable  $X$  (nombre d'absences injustifiées) la moyenne est  $\bar{X} = 1.49$  le mode et la médiane se trouvent facilement sur l'histogramme et le graphe de la fonction de répartition



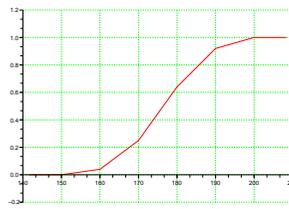
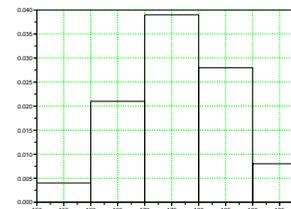
mode 2



médiane 2

on peut remarquer qu'ici il n'y a pas de solution à l'équation  $F(x) = 0.5$  et que la première valeur telle que  $F(x) > 0.5$  est bien  $x = 2$  d'après le tableau :

Pour la série continue  $Y$  on a pour moyenne brute  $\bar{Y} = 176.71$ .

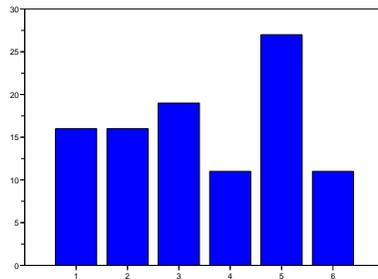


mode 175 ∈ [170, 180[                      médiane 177

alors que graphiquement on trouve que la médiane vaut ≈ 176.4.  
 La médiane est en général un bon indicateur de centralité mais sa détermination se fait graphiquement (dont de manière imprécise), au contraire le mode est plus facile à déterminer mais peut être très éloigné de la moyenne pour des séries fortement dispersées. Par exemple pour une série Z de résultats de N = 100 lancés d'un dé à 6 faces on obtient :

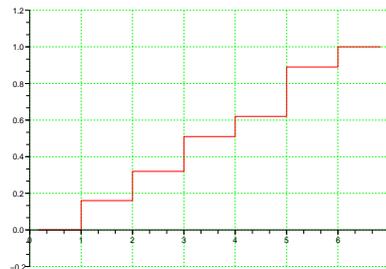
$Z_i =$	1	2	3	4	5	6
$n_i$	16	16	19	11	27	11
$f_i$	0.16	0.16	0.19	0.11	0.27	0.11

effectif total N = 100



le mode sera donc 5 assez éloigné de la moyenne empirique  $\bar{Z} = 3.5$ . Par contre le calcul de la médiane donne 3 d'après la fonction de répartition :

I	] -∞, 1[	[1, 2[	[2, 3[	[3, 4[	[4, 5[	[5, 6[	[6, +∞[
F(t)	0	0.16	0.32	0.51	0.62	0.89	1



**Définition 2.8 (quartiles)**

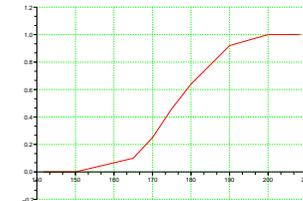
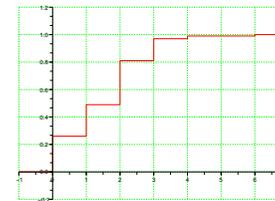
Soit X une série statistique alors on définit les quartiles  $Q_1, Q_3$  à partir de la fonction de répartition de la manière suivante :

- $Q_i$  est la solution de  $F(Q_1) = 0.25$  et  $F(Q_3) = 0.75$  ou à défaut (pour une variable discrète) :
  - $Q_i = x$  la première valeur telle que  $F(x) > 0.25 \times i$  (si  $F(x) = 0.25 \times i$  n'a pas de solution)
  - la moyenne des valeurs telle que  $F(x) = 0.25 \times i$  (si  $F(x) = 0.25 \times i$  a plusieurs solutions)

L'inter-quartile est la valeur  $Q_3 - Q_1$

La définition des quartiles est très similaire à celle de la médiane (d'ailleurs la médiane peut être vue comme le quartile  $Q_2$ ). L'inter-quartile est un indicateur de dispersion, il donne la largeur d'un intervalle  $[Q_1; Q_3]$  contenant plus de 50% des données. De ce fait en général l'inter-quartile est ≈ 1.3333 $\sigma_X$ .

2.4 on peut reprendre le calcul des inter-quartile pour les séries X et Y :



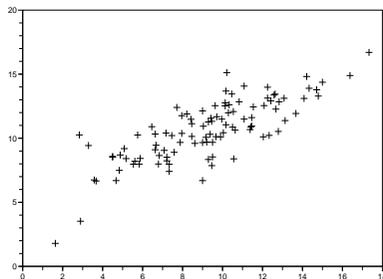
- écart-type  $\sigma_X = 1.1873921$
- inter-quartile  $q_3 - q_1 = 2 - 0 = 2$
- écart-type  $\sigma_Y = 9.1129523$
- inter-quartile  $q_3 - q_1 = 12.5$

Pour une variable statistique continue la médianes et les quartiles sont calculées graphiquement à partir du graphe de la fonction de répartition (ou plutôt de son approximation affine).

### 3 Corrélation

En statistique il est souvent important de rechercher s'il existe un lien entre deux variables  $X$  et  $Y$ , lien qui, dans l'idéal, pourrait s'exprimer par une équation comme  $Y = aX + b$ . Prenons un exemple simple : on considère la population des étudiants inscrit au S1 du DUT INFORMATIQUE de Lannion en 2008 ( $N = 103$  étudiants) et les deux variables aléatoires associées à cette population

$M$  = « note de maths au S1 » et  $U$  = « moyenne générale au S1 »  
 On voudrait quantifier le lien entre ces deux variables statistiques. Si on place les points de coordonnées  $(M_i, U_i)$  sur un graphe on obtient le "nuage" suivant :



ces points semblent grossièrement alignés sur une droite, cela signifie qu'on peut quantifier le lien entre  $M$  et  $U$  par une équation de la forme  $U = aM + b$ . Le but est donc de trouver les coefficients  $a$  et  $b$  de telle sorte que la droite  $y = ax + b$  passe au plus près du maximum de points comme sur la figure FIG.1. C'est ce qu'on appelle faire une *régression linéaire*.

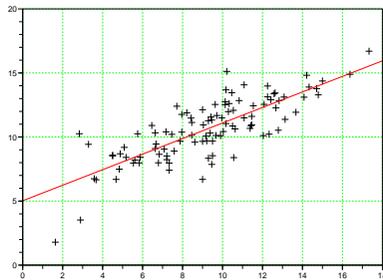


FIGURE 1 – droite de régression

Cette méthode empirique n'est pas satisfaisante, pour obtenir une méthode plus consistante nous allons devoir définir une nouvelle quantité statistique : la

corrélation.

**Définition 3.1 (corrélation)** Soient  $X$  et  $Y$  deux séries statistiques d'effectif  $N$  alors la corrélation de  $X$  et  $Y$  est la valeur

$$\sigma_{X,Y} = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$

Cette définition est une généralisation de la définition de la variance, d'ailleurs on peut remarquer que  $\sigma_{X,X} = \text{Var}(X) = \sigma_X^2$ . Et comme pour la variance on a une formule de Koenig pour simplifier le calcul de la covariance :

**Théorème 3.2 (Koenig)** Soient  $X$  et  $Y$  deux séries statistiques de moyennes  $\bar{X}$  et  $\bar{Y}$  alors

$$\sigma_{X,Y} = \overline{XY} - \bar{X}\bar{Y}$$

**Preuve :** La démonstration de cette formule est très similaire à c'autre formule de Koenig exprimant la variance d'une série statistique en fonction de la moyenne des carrés :

$$\begin{aligned} \sigma_{X,Y} &= \overline{(X - \bar{X})(Y - \bar{Y})} \\ &= \overline{XY - \bar{X}Y - X\bar{Y} + \bar{X}\bar{Y}} \\ &= \overline{XY} - \bar{X}\bar{Y} - \bar{X}\bar{Y} + \bar{X}\bar{Y} \\ &= \overline{XY} - \bar{X}\bar{Y} \end{aligned}$$

□

On a aussi l'inégalité suivante qui va prendre tout son sens plus loin.

**Proposition 3.3 (corrélation)** Soient  $X$  et  $Y$  deux séries statistiques d'écart-type  $\sigma_X$  et  $\sigma_Y$  alors

$$|\sigma_{X,Y}| \leq \sigma_X \sigma_Y$$

**Preuve :** On pose  $S = (X - \bar{X})$  et  $T = (Y - \bar{Y})$  il est alors évident que

$$\overline{S^2} = \overline{(X - \bar{X})^2} = \text{Var}(X), \quad \overline{T^2} = \overline{(Y - \bar{Y})^2} = \text{Var}(Y), \quad \overline{ST} = \sigma_{X,Y}$$

Maintenant on décide de calculer :

$$\overline{(S + \lambda T)^2} = \overline{S^2 + 2\lambda ST + \lambda^2 T^2} = \overline{S^2} + 2\lambda \overline{ST} + \lambda^2 \overline{T^2} = \sigma_X^2 - 2\lambda \sigma_{X,Y} + \lambda^2 \sigma_Y^2$$

c'est donc un trinôme du second degré par rapport à la variable  $\lambda$ , mais ce polynôme est toujours positif car  $\overline{(S + \lambda T)^2} \geq 0$ , donc son discriminant est négatif ou nul :

$$\sigma_X^2 - 2\lambda \sigma_{X,Y} + \lambda^2 \sigma_Y^2 \geq 0 \implies \sigma_{X,Y}^2 - \sigma_X^2 \sigma_Y^2 \leq 0 \implies |\sigma_{X,Y}| \leq \sigma_X \sigma_Y$$

□

Maintenant pour calculer l'équation de la "meilleure" droite passant parmi les points du nuage, on va pouvoir utiliser une méthode découverte par Gauss au début du XIX<sup>ème</sup> siècle : la méthode des moindres carrés. Cette méthode consiste à minimiser l'erreur globale qu'on comment en écrivant que  $Y = aX + b$  :

**Théorème 3.4 (méthode des moindres carrés)** Soient  $X$  et  $Y$  deux séries statistiques alors la méthode des moindres carrés consiste à chercher la "meilleure" droite, d'équation  $y = ax + b$ , passant par le nuage de points  $(X, Y)$  comme étant la droite qui minimise la somme des carrés des écarts entre les points  $(X(i), Y(i))$  et  $(X(i), aX(i) + b)$  c'est à dire on cherche  $a$  et  $b$  qui minimisent la fonction

$$\phi(a, b) = \sum_{i=1}^N (Y(i) - aX(i) - b)^2$$

En utilisant nos connaissances en analyse, on peut maintenant trouver le minimum de la fonction  $\phi(a, b)$ .

**Théorème 3.5 (ajustement linéaire)** Soient  $X$  et  $Y$  deux séries statistiques alors la droite d'ajustement linéaire de  $Y$  en  $X$  a pour équation  $y = ax + b$  avec :

$$a = \frac{\sigma_{X,Y}}{\text{Var}(X)} \quad \text{et} \quad b = \bar{Y} - a\bar{X}$$

Enfin on peut quantifier la qualité de l'approximation  $Y = aX + b$  en utilisant le coefficient de corrélation.

**Théorème 3.6 (coefficient de corrélation)** Soient  $X$  et  $Y$  deux séries statistiques, on suppose que la droite d'ajustement linéaire de  $Y$  en  $X$  a pour équation  $y = ax + b$  et on définit son coefficient de corrélation par :

$$\rho = \frac{\sigma_{X,Y}^2}{\text{Var}(X)\text{Var}(Y)} \in [0, 1]$$

alors le nuage de points  $(X, Y)$  est d'autant plus proche de la droite d'ajustement linéaire de  $Y$  en  $X$  que  $\rho$  est proche de 1. En particulier tous les points du nuage sont sur la droite  $y = ax + b$  si et seulement si  $\rho = 1$ .

La corrélation sera dite forte (resp. faible) si  $\rho \geq \frac{\sqrt{3}}{2} \approx 0.866$  (resp.  $\rho < \frac{\sqrt{3}}{2} \approx 0.866$ )

**Preuve :** le but est de trouver 2 équations qui permettent de trouver les  $a$  et  $b$  qui minimise la fonction  $\phi(a, b)$ . Pour cela on va minimiser  $\phi$  d'abord par rapport à  $b$  et ensuite par rapport à  $a$  :

- recherche du minimum par rapport à  $b$

$$\begin{aligned} \phi(a, b) &= \sum_{i=1}^N (Y(i) - (aX(i) + b))^2 \\ &= \sum_{i=1}^N ((Y(i) - aX(i)) - b)^2 \\ &= \sum_{i=1}^N (Y(i) - aX(i))^2 - 2b(Y(i) - aX(i)) + b^2 \\ &= N \left( b^2 - 2b\overline{Y - aX} + \overline{(Y - aX)^2} \right) \end{aligned}$$

ce trinôme atteint son minimum quand  $b = \overline{Y - aX} = \bar{Y} - a\bar{X}$ , on peut remarquer que cette équation signifie que la "meilleure" droite passe par le point "moyen"  $(\bar{X}, \bar{Y})$  puisque  $\bar{Y} = a\bar{X} + b$ .

- recherche du minimum par rapport à  $a$

$$\begin{aligned} \phi(a, b) &= \sum_{i=1}^N (Y(i) - (aX(i) + b))^2 \\ &= \sum_{i=1}^N ((Y(i) - aX(i)) - (\bar{Y} - a\bar{X}))^2 \\ &= \sum_{i=1}^N ((Y(i) - \bar{Y}) - a(X(i) - \bar{X}))^2 \\ &= \sum_{i=1}^N (Y(i) - \bar{Y})^2 - 2a(Y(i) - \bar{Y})(X(i) - \bar{X}) + a^2(X(i) - \bar{X})^2 \\ &= N \left( \overline{(Y - \bar{Y})^2} - 2a\overline{(Y - \bar{Y})(X - \bar{X})} + \overline{(Y - \bar{Y})^2 X^2} \right) \\ &= N (a^2 \text{Var}(X) - 2a\sigma_{X,Y} + \text{Var}(Y)) \end{aligned}$$

ce trinôme atteint son minimum quand  $a = \frac{\sigma_{X,Y}}{\text{Var}(X)}$

- Pour finir calculons la somme des écarts pour la "meilleure" droite :

$$\begin{aligned} \phi(a, b) &= N \left( \frac{\sigma_{X,Y}^2}{\text{Var}(X)^2} \text{Var}(X) - 2 \frac{\sigma_{X,Y}^2}{\text{Var}(X)} - \text{Var}(Y) \right) \\ &= N \left( \text{Var}(Y) - \frac{\sigma_{X,Y}^2}{\text{Var}(X)} \right) = N \text{Var}(Y) \left( 1 - \frac{\sigma_{X,Y}^2}{\text{Var}(X)\text{Var}(Y)} \right) \\ &= N \text{Var}(Y) (1 - \rho) \end{aligned}$$

$\phi(a, b)$  étant toujours positive (par définition) on doit avoir que  $\rho \leq 1$ , ce qui est bien vérifié d'après la proposition 3.3. Maintenant on en déduit facilement que  $\rho = 1 \iff \phi(a, b) = 0$  ce qui signifie bien que tous les écarts sont nuls si  $\rho = 1$ .

□

 **3.1** Reprenons l'exemple des notes du premier semestre :

$M = \ll \text{note de maths au S1} \gg$  et  $U = \ll \text{moyenne générale au S1} \gg$

et appliquons la méthode des moindres carrés. On calcule d'abord les différents paramètres :

- $\bar{M} = 9.2500971$   $\text{Var}(M) = 9.8811854$  et  $\sigma_M = 3.1434353$
- $\bar{U} = 10.636117$   $\text{Var}(U) = 5.7622801$  et  $\sigma_U = 2.400475$
- $\overline{MU} = 104.37805$

On en déduit l'ajustement de  $U$  par rapport à  $M$

- $a = \frac{\sigma_{M,U}}{\text{Var}(U)} = \frac{5.9929373}{3.1434353^2} = 0.6064998$
- $b = \bar{U} - a\bar{M} = 10.636117 - 0.6064998 \times 9.2500971 = 5.0259342$
- $\rho = \frac{\sigma_{M,U}^2}{\text{Var}(M)\text{Var}(U)} = \frac{5.9929373^2}{9.8811854 \times 5.7622801} = 0.6307773$  (corrélation faible)

ce qui donne :

$$U = 0.6064998M + 5.0259342$$

qui correspond bien à ce qu'on obtient graphiquement (voir droite rouge sur les figures FIG.1 et FIG.2).

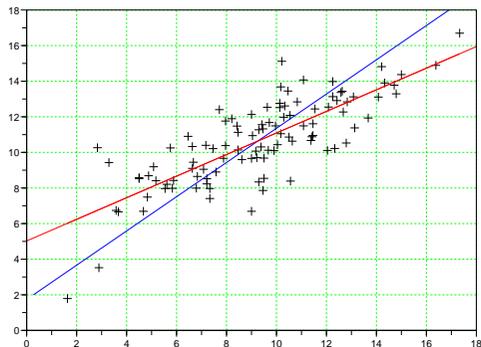


FIGURE 2 – comparaison des ajustements de  $U$  par rapport à  $M$  (en rouge) et de  $M$  par rapport à  $n$  (en bleu) avec une corrélation faible  $\rho = 0.6307773$

On aurait pu de la même manière chercher à faire l'ajustement de  $M$  par rapport à  $U$ . Il suffit pour cela d'échanger les rôles de  $U$  et  $M$ . On peut remarquer que ce changement ne modifie pas la valeur de la corrélation (car  $\sigma_{X,Y} = \sigma_{Y,X}$ ) ni du coefficient de corrélation  $\rho = 0.6307773$  (corrélation faible), on peut donc reprendre ces valeurs pour calculer les nouveaux coefficients de la droite d'ajustement de  $M$  par rapport à  $U$  :

- $a = \frac{\sigma_{U,M}}{\text{Var}(U)} = \frac{5.9929373}{2.400475^2}$
- $b = \bar{M} - a\bar{U} = 9.2500971 - 1.0400288 \times 10.636117 = -1.8117705$

on trouverait donc :

$$M = 1.0400288U - 1.8117705$$

⚠ on remarquera que l'ajustement de  $X$  par rapport à  $Y$  n'est pas le même que celui trouvé en ajustant  $Y$  par rapport à  $X$ , on peut s'en convaincre facilement sur la figure FIG.2! Ceci est lié au fait que la corrélation soit faible. Les deux méthodes donnent la même équation si et seulement si  $\rho = 1$ .

## 4 Calculs statistiques avec Scilab

À chaque TP vous trouverez un fichier \*.sce contenant des séries statistiques sous forme de d'une matrice à une ligne ou 1 colonne :

```
X=[11;3;1;10;3;1;13;2;1;10;2;4;4;9;4;10; ...]
```

Il vous suffira de charger cette matrice dans l'environnement scilab pour pouvoir commencer l'étude de la série statistique.

### 4.1 séries statistiques discrètes

- L'effectif total  $N$  de la variable statistique  $X$  (le nombre de valeurs valeurs) s'obtient en récupérant la taille de la matrice :

```
-->X
X =
 11.
  3.
  1.
 10.
  3.
  1.
[More (y or n) ?]

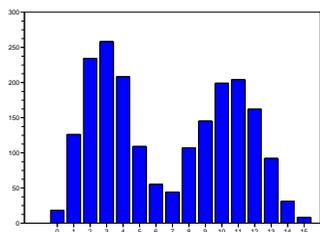
-->size(X)
ans =
 2000.    1.

-->N=length(X)
N =
 2000.
```

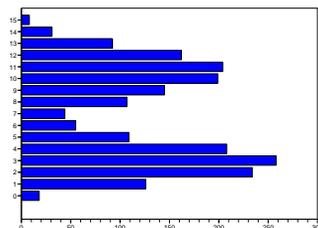
- Scilab permet de calculer facilement les effectifs  $n_i$  et les fréquences  $f_i$  des modalités de la variable  $X$  avec la commande tabul (il y a aussi nfreq mais moins pratique) :

ordre croissant	valeurs	effectifs	fréquences
-->m=tabul(X, 'i')	-->x=m(:,1)	-->n=m(:,2)	-->f=n/N
m =	x =	n =	f =
0.	18.	0.	18.
1.	126.	1.	126.
2.	234.	2.	234.
3.	258.	3.	258.
4.	208.	4.	208.
5.	109.	5.	109.
6.	55.	6.	55.
7.	44.	7.	44.
8.	107.	8.	107.
9.	145.	9.	145.
10.	199.	10.	199.
11.	204.	11.	204.
12.	162.	12.	162.
13.	92.	13.	92.
14.	31.	14.	31.
15.	8.	15.	8.

- On peut ensuite tracer différents types d'**histogrammes** de la variable  $X$  avec les commandes `bar`, `barh` pour les diagrammes en bâton :

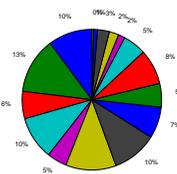


`bar(x,f)`

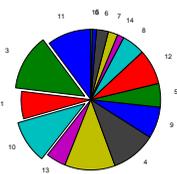


`barh(x,f)`

et `pie` pour les « camemberts » :



`pie(f)`



`pie(f, bool2s((x==3) | (x==10)), string(x))`

- Scilab permet de calculer facilement les principaux **estimateurs de centralité** avec `mean` et `median` ou par calcul :

<code>--&gt;median(X)//médiane</code>	<code>--&gt;mean(X)//moyenne</code>	<code>--&gt;sum(X)/N</code>
ans =	ans =	ans =
6.	6.736	6.736

- On peut facilement calculer les **estimateurs de dispersion** définis en cours :

<code>--&gt;//calcul de la variance</code>	<code>--&gt;//par la formule de Koenig</code>
<code>--&gt;sum((X-mean(X)).^2)/N</code>	<code>--&gt;mean(X.^2)-mean(X)^2</code>
ans =	ans =
16.377304	16.377304
	<code>--&gt;sqrt(ans)//écart-type</code>
	ans =
	4.0468882

⚠ Attention scilab possède aussi une fonction `variance` et une fonction `écart-type` (`stdev` mais qui utilisent **une définition légèrement différente** de celle vue en cours.

<code>--&gt;//estimateurs non-biaisés</code>	<code>--&gt;// vérification</code>
<code>--&gt;variance(X)</code>	<code>--&gt;sum((X-mean(X)).^2)/(N-1)</code>
ans =	ans =
16.385497	16.385497
<code>--&gt;stdev(X)//écart-type</code>	<code>--&gt;sqrt(ans)</code>
ans =	ans =
4.0479003	4.0479003

En fait ce sont les **estimateurs non-biaisés** qui sont calculés par Scilab.

- On peut enfin calculer les **quartiles** et les **déciles** :

<code>--&gt;//quartiles</code>	<code>--&gt;p=perctl(X,90)//9ième décile</code>
<code>--&gt;Q=quart(X)</code>	p =
Q =	12. 822.
3.	
6.	
10.	<code>--&gt;p=perctl(X,[10:10:100])</code>
	p =
<code>--&gt;//inter-quartile</code>	
<code>--&gt;iqr(X)</code>	2. 1466.
ans =	3. 1817.
	3. 298.
	4. 483.
	6. 371.
<code>--&gt;//vérification</code>	9. 1397.
<code>--&gt;Q(3)-Q(1)</code>	10. 968.
ans =	11. 1152.
	12. 822.
	15. 229.

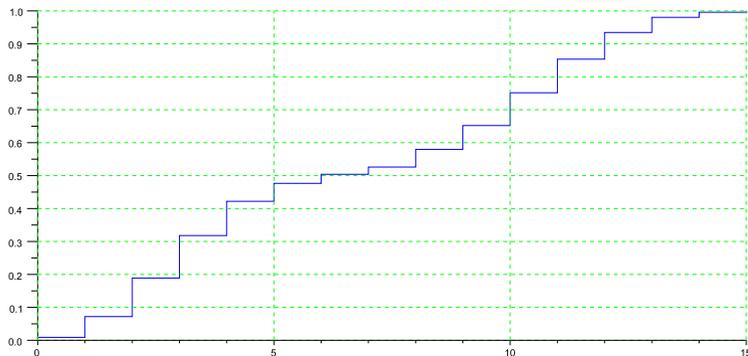
⚠ la deuxième colonne du résultat renvoyé par `perctl` correspond à la position du **décile** (la valeur dans la première colonne) dans la série statistique (matrice  $X$ ).

- Le calcul des **effectifs cumulés** peut se faire à partir de la matrice contenant les effectifs avec la fonction `cumsum` :

```

cumul des effectifs      | cummul fréquence
-->F=cumsum(n)/N        | -->F=cumsum(f)
F =                     | F =
0.009                   | 0.009
0.072                   | 0.072
0.189                   | 0.189
0.318                   | 0.318
0.422                   | 0.422
0.4765                  | 0.4765
0.504                   | 0.504
0.526                   | 0.526
0.5795                  | 0.5795
0.652                   | 0.652
0.7515                  | 0.7515
0.8535                  | 0.8535
0.9345                  | 0.9345
0.9805                  | 0.9805
0.996                   | 0.996
1.                      | 1.
    
```

- on peut ensuite tracer le **graphe de la fonction de répartition d'une série statistique discrète** avec la commande `plot2d2`



```
plot2d2(x,F,2),xgrid(3)
```

⚠ Attention, cette commande ne dessine le graphe de la fonction de répartition  $F$  que sur l'intervalle  $[\min_i(x_i), \max_i(x_i)]$  ce graphe doit être prolongé par :  
 - 0 pour  $x < \min_i(x_i)$   
 - 1 pour  $x > \max_i(x_i)$

## 4.2 séries statistiques regroupées par intervalles

Le cas des séries statistiques regroupées par intervalles est plus compliqué :

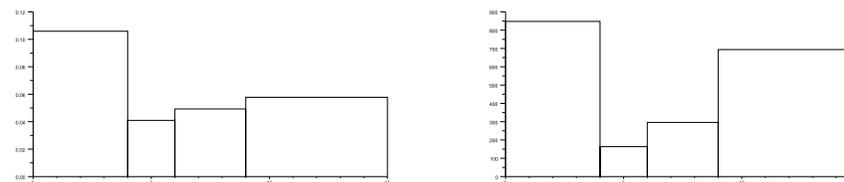
- On peut reprendre l'étude de la série  $X$  en groupant les valeurs en 4 intervalles :  $[0, 4]$   $]4, 7]$   $]7, 11]$  et  $]11, 15]$ . Pour cela il faut définir une matrice  $I$  contenant les bornes des intervalles

```

intervalles             | les centres des intervalles
-->I=[0,4,6,9,15]      | -->xc=[2 5 7.5 12]
I =                    | xc =
0.   4.   6.   9.  15. | 2.   5.   7.5  12.
    
```

⚠ attention il y a une entrée de plus dans la matrice  $I$  que le nombre d'intervalles !

- on peut faire l'histogramme en regroupant les valeurs par classes avec `histplot`. Dans l'histogramme obtenu **la surface de chaque colonne est proportionnelle à l'effectif par défaut** et l'échelle de l'axe des ordonnées est telle que la surface des colonnes vaut 1. Si on veut avoir les effectifs réels en ordonnée (au lieu des effectifs corrigés) il faut rajouter l'option `normalization=%f` :



```
histplot([0,4,7,11,15],X)  histplot([0,4,7,11,15],X,normalization=%f)
```

- si on veut récupérer les valeurs exactes des effectifs et des fréquences il faut utiliser `dsearch` :

```

-->[ind,n]=dsearch(X,I);n      | -->f=n/N
n =                             | f =
848.   164.   296.   694.    | 0.4235764  0.0819181  0.1478521  0.3466533
    
```

on peut vérifier que les effectifs corrigés sont bien obtenus en divisant les effectifs réels par la longueur de l'intervalle correspondant :

```

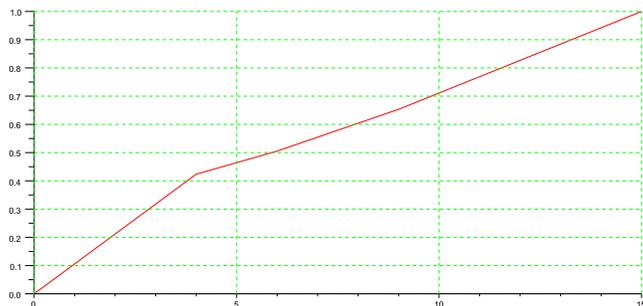
longueur des intervalles      | fréquences corrigées
-->l=[4 2 3 6]               | -->h=f./l
l =                            | h =
4.   2.   3.   6.            | 0.1058941  0.0409590  0.0492840  0.0577756
    
```

- on peut ensuite calculer les fréquences cumulées et tracer le **graphe de la fonction de répartition d'une série statistique continue** avec les commandes cumsum et plot2d :

```
-->F=cumsum(f)
F =

    0.4235764    0.5054945    0.6533467    1.

-->plot2d(I,[0 F],5),xgrid(3)
```



ce qui permet de déterminer graphiquement médiane et quartiles.

⚠ La matrice  $I$  définissant les bord intervalles possède une case de plus que la matrice  $F$  (pour  $k$  intervalles il y a  $k + 1$  bornes). Il faut donc rajouter une valeur dans  $F$  : c'est le 0 dans  $[0 F]$  qui la valeur de départ (donc à gauche) des fréquences cumulées.

- on peut aussi calculer les valeurs approchées des différents paramètres statistiques en remplaçant les modalités discrètes  $x$  par les centres des intervalles  $xc$ . On peut comparer les valeurs trouvées aux valeurs exactes :

moyenne approchée	variance approchée	écart-type approché
-->M=sum(xc.*f)	-->sum(f.*(xc-M).^2)	-->sqrt(ans)
M =	ans =	ans =
6.5254745	19.395205	4.4039988

## Index

- ajustement linéaire, 19
- caractère qualitatifs, 4
- caractère quantitatifs, 4
- coefficient de corrélation, 19
- corrélation, 18
- diagramme en bâtons, 5
- droites d'ajustement linéaire, 21
- écart-type, 12
- effectif corrigé, 7
- fonction de répartition, 8
- histogramme, 5
- inter-quartile, 16
- Koening, 18
- linéarité de la moyenne, 12
- médiane, 14
- méthode des moindres carrés, 19
- modalité, 5
- mode, 14
- moyenne approchée, 11
- moyenne empirique, 11
- population, 4
- quartile, 16
- regression linéaire, 17
- série statistique, 4
- tableau statistique, 5
- théorème
  - Koening, 13
- variance statistique, 12