

2. LE COUPLE MOYENNE / ÉCART-TYPE

L'avantage de la moyenne réside dans le fait qu'elle est linéaire : si une classe de 30 élèves a 10 de moyenne et une classe de 20 élèves a 12 de moyenne, on calcule la moyenne des 50 élèves par : $(30 \times 10 + 20 \times 12)/50 = 10,8$. Il n'est pas nécessaire de connaître les différentes notes de chaque classe.

L'inconvénient de la moyenne est qu'elle est sensible aux valeurs extrêmes : dans une entreprise ou l'essentiel des salariés gagne le SMIC et un dirigeant un salaire extravagant, le salaire moyen peut être élevé.

- La *moyenne* de la série est $\bar{x} = \frac{x_1 + \dots + x_N}{N} = \frac{1}{N} \sum_{i=1}^N x_i = \frac{n_1 a_1 + \dots + n_k a_k}{N} = \frac{1}{N} \sum_{j=1}^k n_j a_j$
- La *variance* V de la série est $V = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{1}{N} \left(\sum_{i=1}^N x_i^2 \right) - \bar{x}^2 = \frac{1}{N} \sum_{j=1}^k n_j (a_j - \bar{x})^2$
- L'*écart type* est $\sigma = \sqrt{V}$. (L'écart-type mesure la *dispersion* de la série autour de la moyenne)

Exemple. Dans l'exemple précédent : $\bar{x} = \frac{1}{10}(4 + 7 \times 2 + 9 \times 2 + 11 \times 3 + 15 + 18) = 10,2$.

$$V = \frac{1}{10}((4 - 10,2)^2 + 2(7 - 10,2)^2 + 2(9 - 10,2)^2 + 3(11 - 10,2)^2 + (15 - 10,2)^2 + (18 - 10,2)^2) = 14,76$$

$$\sigma = \sqrt{V} \approx 3,84$$

Si l'on a entré les différents x_i dans la liste L_1 , la calculatrice donne \bar{x} la moyenne, σx l'écart-type (et l'effectif total n , la somme des termes Σ et la somme des carrés Σx^2).

```
EDIT [MODE] TESTS
1:1-Var Stats
2:2-Var Stats
3:Med-Med
4:LinReg(ax+b)
5:QuadReg
6:CubicReg
7↓QuartReg
```

[STAT], Calc (1 :Stat1Var)

```
1-Var Stats L1
```

Stat1Var L_1 ([2nd]+1)

```
1-Var Stats
x=10.2
Σx=102
Σx²=1188
Sx=4.049691346
σx=3.841874542
n=10
```

faire défiler.

3. LE COUPLE MÉDIANE / ÉCART-INTERQUARTILE

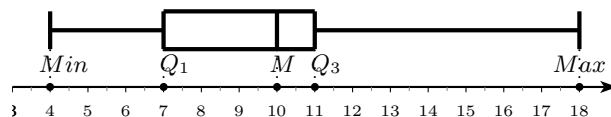
L'inconvénient de ces données est l'absence de linéarité (on ne peut calculer a priori la médiane de deux groupes en connaissant seulement la médiane de chacun des groupes). L'avantage est le peu de sensibilité aux valeurs extrêmes (dans l'exemple de l'entreprise, le salaire médian serait aussi égal au SMIC). On suppose la série des x_i rangée par ordre croissant.

- Une *médiane* M est telle que 50% au moins des caractères sont au dessus de M et 50% en dessous. Par convention : $M = \frac{1}{2}(x_{\frac{N}{2}} + x_{\frac{N}{2}+1})$ si N est pair, $M = x_{\frac{N+1}{2}}$ si N est impair.
- Le *premier quartile* Q_1 est une médiane de la série des x_i où $i > \frac{N+1}{2}$.
- Le *troisième quartile* Q_3 est une médiane de la série des x_i où $i < \frac{N+1}{2}$.
- L'*intervalle interquartile* est $]Q_1, Q_3[$. L'écart interquartile est $Q_3 - Q_1$.

Exemple. Dans l'exemple précédent :

4 7 7 9 9 || 11 11 15 18
 $Q_1=7$ $M=\frac{9+11}{2}=10$ $Q_3=11$

Donc $Q_3 - Q_1 = 4$.

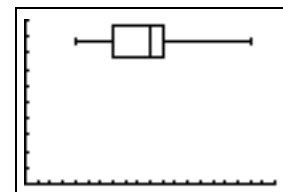


```
1-Var Stats
n=10
minX=4
Q1=7
Med=10
Q3=11
maxX=18
```

après défilement

```
2:Plot1...Off
2:Plot2...Off
3:Plot3...Off
4↓PlotsOff
```

([2nd]+[Y=])



[Graphe] (ajuster la fenêtre)

4. EXEMPLE : STRUCTURE DE L'EMPLOI DU TEMPS PAR GENRE

Le tableau suivant a été publié par l'organisme Eurostat le 8 mars 2006 (journée de la femme). Il présente la structure des emplois du temps des femmes et des hommes âgés de 20 à 74 ans, en heures par jours, selon un sondage réalisé en 2005.

	Emploi, études		Travail domestique		Travail total		
	Femmes	Hommes	Femmes	Hommes	Femmes	Hommes	Différence
Belgique	2.1	3.5	4.5	2.6	6.6	6.1	0.5
Allemagne	2.1	3.6	4.2	2.4	6.3	6.0	
Espagne	2.4	4.7	4.9	1.6	7.3	6.3	
France	2.5	4.1	4.5	2.4	7.0	6.5	
Italie	2.1	4.4	5.3	1.6	7.4	6.0	
Lettonie	3.7	5.2	3.9	1.8	7.6	7.0	
Lituanie	3.7	4.9	4.5	2.2	8.2	7.1	
Hongrie	2.5	3.8	5.0	2.7	7.5	6.5	
Pologne	2.5	4.3	4.8	2.4	7.3	6.7	
Slovénie	3.0	4.1	5.0	2.7	8.0	6.8	
Suède	3.2	4.4	3.7	2.5	6.9	6.9	
Royaume-Uni	2.6	4.3	4.3	2.3	6.9	6.6	
Norvège	2.9	4.3	3.8	2.4	5.7	6.7	
Moyenne	2.7			2.3	7.2		
Variance	0.28			0.13			
Écart-type	0.53			0.36			
Minimum	2.1			1.6			
Quartile 1	2.25			2			
Médiane	2.5			2.4		6.6	
Quartile 3	3.1			2.55			
Maximum	3.7			2.7			

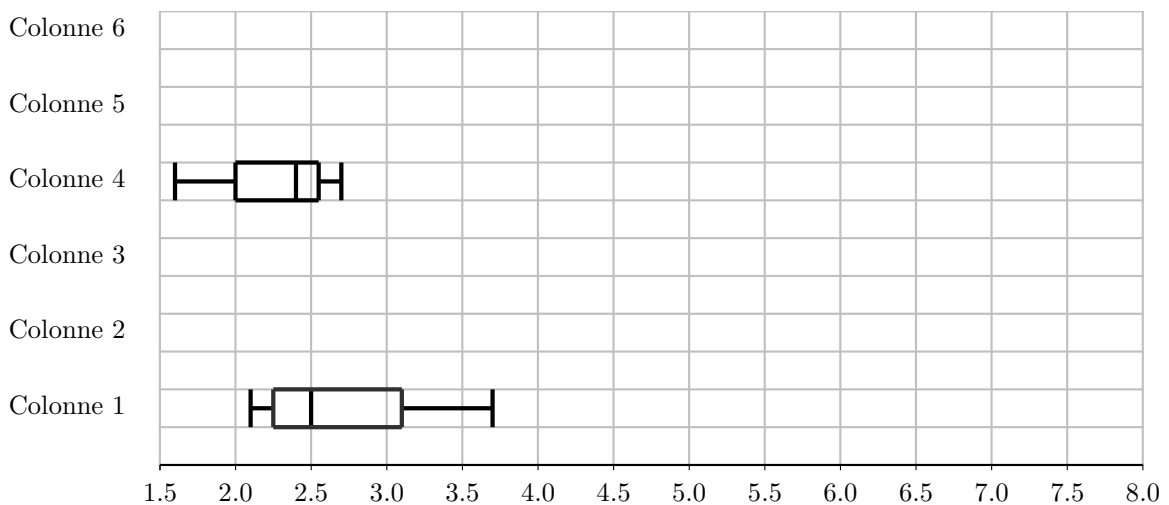
- **Travail payé et études des hommes** : Ranger les termes de la série de la seconde colonne par ordre croissant :

≤ ≤ ≤ ≤ ≤ ≤ ≤ ≤ ≤ ≤ ≤

Moyenne =

Variance =

- Remplir les autres colonnes (calculatrice) et compléter le diagramme :



5. STATISTIQUES À DEUX VARIABLES

On considère deux caractères d'une population, dont les valeurs sont x_1, \dots, x_N et y_1, \dots, y_N .

• Dans le plan muni d'un repère orthogonal, l'ensemble des points $M_i(x_i, y_i)$ est le *nuage de points* associé à la série statistique à deux variables.

• Le *point moyen* associé à cette série est le point $G = (\bar{x}, \bar{y})$ où \bar{x} et \bar{y} sont les moyennes respectives des séries x et y .

• Rechercher un *ajustement affine* du nuage revient à déterminer une droite affine d'équation $y = ax + b$ aussi proche que possible des points du nuage.

• Étant donné un ajustement affine $\mathcal{D} : y = ax + b$, $r_i = y_i - (ax_i + b)$ représente, au signe près, l'écart entre le point $M'_i(x_i, ax_i + b)$ de \mathcal{D} et le point du nuage $M_i(x_i, y_i)$. Le nombre r_i est appelé le *résidu* numéro i de l'ajustement.

• On mesure la qualité d'un ajustement affine par la somme des carrés des résidus :

$$S = r_1^2 + r_2^2 + \dots + r_N^2 = \sum_{i=1}^N (y_i - (ax_i + b))^2$$

Plus cette somme est petite, plus l'ajustement est précis.

Exemple. Soit la série double :

i	1	2	3	4
x_i	1	2	4	5
y_i	2	5	3	6

Le nuage des points $M_i(x_i; y_i)$ est représenté :

Le point moyen G a pour coordonnées $(3; 4)$ car :

$\bar{x} = \frac{1}{4}(1 + 2 + 4 + 5) = 3$ et $\bar{y} = \frac{1}{4}(2 + 5 + 3 + 6) = 4$.

On se propose d'ajuster le nuage par la droite d'équation $y = 0,5x + 2,5$.

Les résidus de cet ajustement sont :

$$r_1 = y_1 - (ax_1 + b) = 2 - (0,5 \cdot 1 + 2,5) = 2 - 3 = -1$$

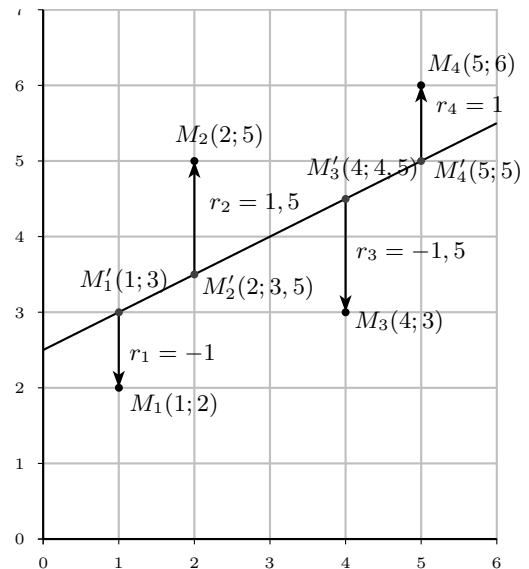
$$r_2 = y_2 - (ax_2 + b) = 5 - (0,5 \cdot 2 + 2,5) = 5 - 3,5 = 1,5$$

$$r_3 = y_3 - (ax_3 + b) = 3 - (0,5 \cdot 4 + 2,5) = 3 - 4,5 = -1,5$$

$$r_4 = y_4 - (ax_4 + b) = 6 - (0,5 \cdot 5 + 2,5) = 6 - 4,5 = 1,5$$

La somme des carrés des résidus est :

$$S = r_1^2 + r_2^2 + r_3^2 + r_4^2 = 6,5$$



Théorème. Soit x_1, \dots, x_N et y_1, \dots, y_N une série double. Si la variance de la série des x_i est non nulle, il existe un ajustement affine qui rend minimale la somme des carrés des résidus. Cet ajustement est la *droite des moindres carrés* (ou droite de régression linéaire) d'équation :

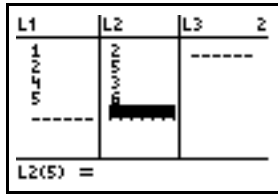
$$y = ax + b \text{ avec } a = \frac{\text{cov}(x, y)}{V(x)} = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2} \text{ et } b = \bar{y} - a\bar{x}$$

Remarque. La droite des moindres carrés passe toujours par le point moyen G . Elle est le meilleur ajustement affine au sens où elle rend minimale la somme des carrés des résidus.

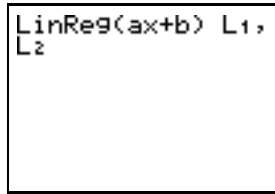
Exemple. Dans l'exemple précédent, déterminer l'équation de la droite des moindres carrés et montrer que la somme des carrés des résidus est inférieure à celle de l'ajustement proposé.

6. UTILISER LA CALCULATRICE

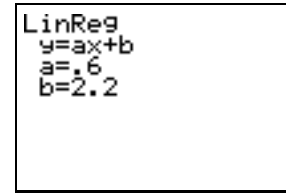
- Équation de la droite des moindres carrés :



$x \rightarrow L_1, y \rightarrow L_2$

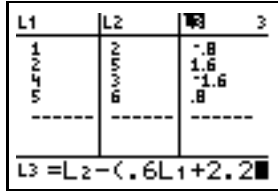


([STAT]+Calc) LinReg(ax+b) L1, L2

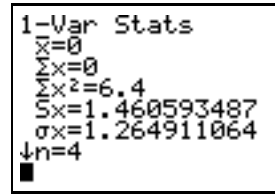


Résultat

- Somme des carrés des résidus.



Se placer sur L3+[Entrer]



Stat1Var L3 voir Σx^2

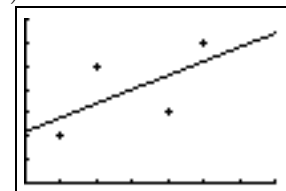
- Graphe (au préalable, ajuster la fenêtre et entrer $Y_1 = 0.6X + 2.2$)



[2nde]+[GrapheStat]



Paramétrage



[Graphe]

7. EXEMPLE : INSCRITS AU PÔLE EMPLOI DEPUIS SEPTEMBRE 2008

D'après l'INSEE, voici le nombre d'inscrits au pôle emploi (en milliers) entre septembre 2008 et août 2009 :

Mois	09-08	10-08	11-08	12-08	01-09	02-09	03-09	04-09	05-09	06-09	07-09	08-09
Mois $n^o x_i$	0	1	2	3	4	5	6	7	8	9	10	11
Inscrits y_i	3299	3348	3387	3438	3524	3605	3688	3786	3843	3851	3888	3923
Résidu r_i												

Les coordonnées du point moyen sont $(\bar{x}, \bar{y}) = \dots\dots\dots$

La droite d'ajustement des moindres carrés a pour équation : $\dots\dots\dots$

La somme des carrés des résidus r_i est $\sum_{i=0}^{11} r_i^2 = \dots\dots\dots$

Combien de milliers d'inscrits peut-on prévoir pour septembre 2009 ? $\dots\dots\dots$

