

La statistique descriptive à deux variables a pour but de mettre en évidence une relation éventuelle qui peut exister entre **deux variables** d'une population, considérées **simultanément**.

## I) Définition

On appelle *série statistique à deux variables*  $X$  et  $Y$ , le relevé simultané de la valeur de deux caractères statistiques  $X$  et  $Y$ . Elle est donc constituée d'une liste de couples de nombres  $(x_i; y_i)$ .

Les données d'une série statistique à deux variables sont la plupart du temps présentées dans un tableau où l'on indique les  $n$  valeurs des variables  $X$  et  $Y$  :

Variable $X$	$x_1$	$x_2$	...	$x_n$
Variable $Y$	$y_1$	$y_2$	...	$y_n$

## II) Nuage de points et point moyen

1) Le plan étant muni d'un repère orthogonal, on peut associer à chaque couple  $(x_i; y_i)$  de la série statistique le point  $M$  de coordonnées  $(x_i; y_i)$ .

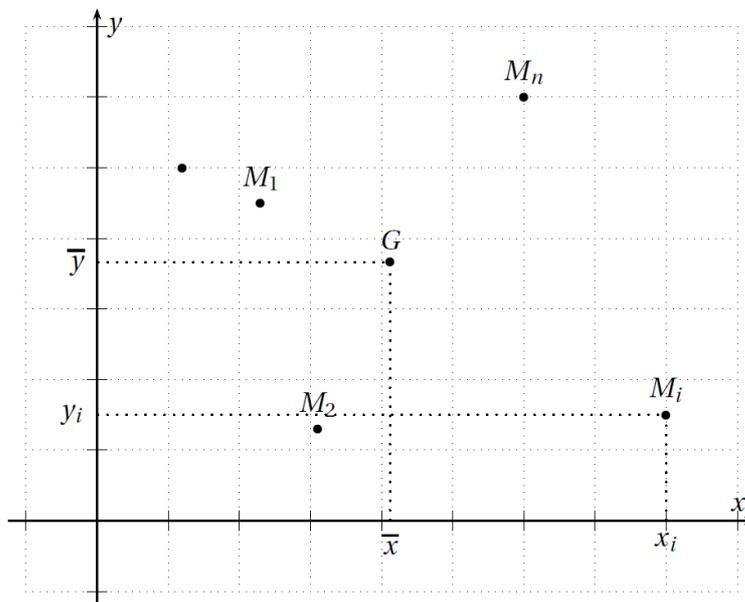
Le graphique ainsi obtenu constitue un *nuage de points*.

2) Le point moyen  $G$  du nuage de points est le point de coordonnées :  $(\bar{x}; \bar{y})$

où :

- l'abscisse est la moyenne de la série  $(x_i)$  :  $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$
- l'ordonnée est la moyenne de la série  $(y_i)$  :  $\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n}$

On dit aussi que c'est le centre de gravité du nuage.



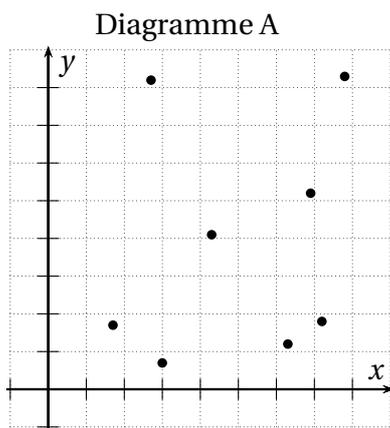
### III) Ajustement

#### 1) Corrélation

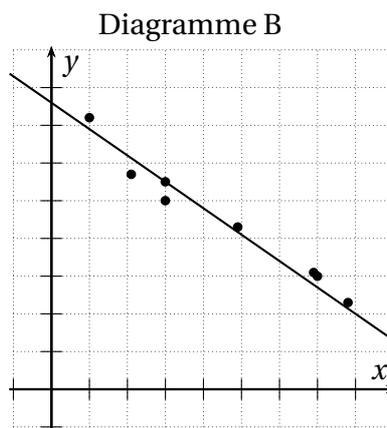
Il y a **corrélation** entre deux variables  $X$  et  $Y$  observées sur les individus d'une même population lorsqu'il y a une **relation** entre  $X$  et  $Y$ .

*Remarque.* L'existence d'une corrélation entre deux variables peut être décelée **dans un premier temps** à l'aide d'un nuage de points.

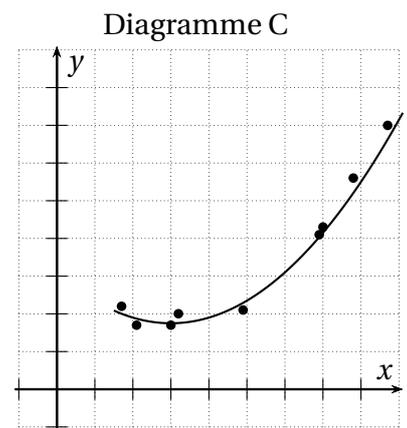
*Exemples.* Considérons les diagrammes suivants :



Absence de corrélation entre  $X$  et  $Y$



Corrélation entre  $X$  et  $Y$



Corrélation entre  $X$  et  $Y$

#### 2) Ajustement de $y$ en $x$

Lorsque les valeurs de  $x$  sont connues, effectuer un ajustement de  $y$  en  $x$  d'un nuage de points consiste à trouver une fonction dont la courbe représentative d'équation  $y = f(x)$  est la plus « proche » du nuage.

*Remarques.*



- Un ajustement permet de faire des estimations : *interpolation* (dans l'intervalle d'étude) et *extrapolation* (en dehors). Extrapolation = prévision.
- Lorsque les points du nuage sont *presque alignés*, comme pour le diagramme B, on recherche une droite qui passe le plus près possibles des points. On effectue alors un **ajustement affine**.
- On verra qu'il existe des ajustements qui ne sont pas affines, comme sur le diagramme C.

### 3) Ajustement affine

Une droite d'ajustement affine est une droite qui passe **au plus près** du nuage de points.

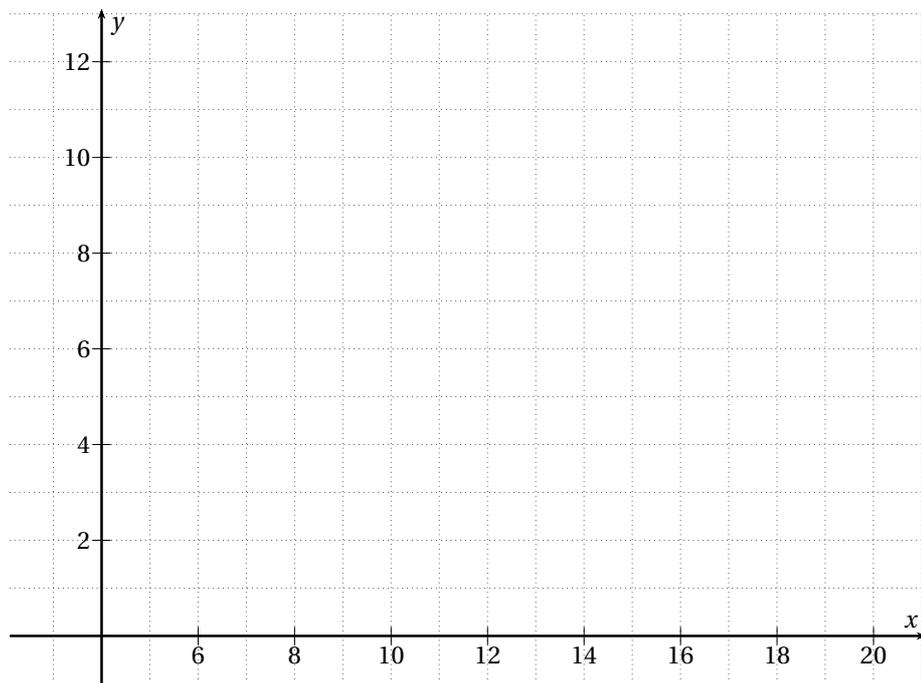
#### Ajustement au jugé

Le gérant d'un hypermarché, disposant d'un potentiel de vingt-huit caisses enregistreuses, a fait réaliser une étude statistique sur le temps moyen (en minutes) d'attente d'un client à la caisse.

On note  $x_i$  le nombre de caisses ouvertes,  $y_i$  le temps moyen correspondant en minutes.

$x_i$	4	5	6	7	8	9	10	11	12	13
$y_i$	12,25	12	11,5	11,75	10	10	9,75	9	8,25	8

1. (a) Construire dans le repère de la figure ci-dessous le nuage de points associé à la série.



- (b) Calculer les coordonnées du point moyen  $G$  du nuage. Placer  $G$  dans le repère.

2. Peut-on penser à un ajustement affine? Justifier.

3. On propose comme droite d'ajustement affine la droite  $\mathcal{D}$  qui a pour équation  $y = -0,5x + 14,5$ .

(a) Tracer la droite  $\mathcal{D}$  sur le graphique.

(b) Le point  $G$  appartient-il à la droite  $\mathcal{D}$ ? Justifier par un calcul.

(c) En utilisant cette droite :

- déterminer le temps moyen d'attente d'un client à la caisse lorsque 20 caisses sont ouvertes
- déterminer le nombre de caisses à ouvrir pour que le temps moyen d'attente d'un client à la caisse soit de trois minutes



- déterminer le temps moyen pour 7 caisses ouvertes (valeur existante). Que constatez-vous ? L'écart entre la valeur observée et la valeur estimée est appelé **résidu**.

## IV) Ajustement par la méthode des moindres carrés

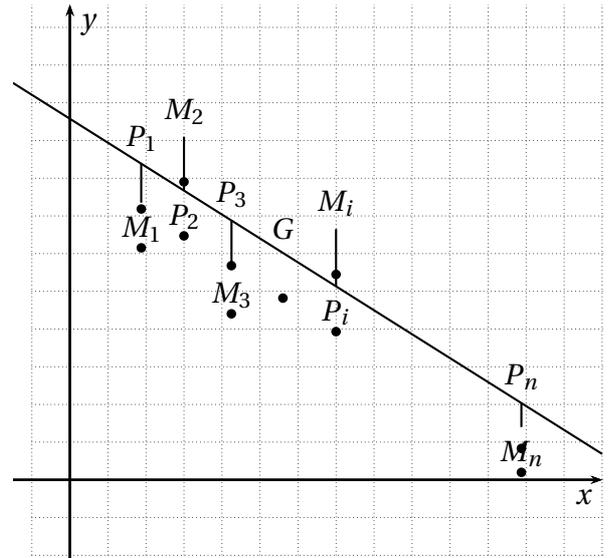
On connaît les valeurs  $x_i$ , on cherche à obtenir une droite d'ajustement dont les valeurs  $y$  sont les plus proches possibles des  $y_i$  « verticalement ».

Les points  $M_1, M_2, \dots, M_n$  sont les points du nuage.

Les points  $P_1, P_2, \dots, P_n$  sont les points d'une droite  $\mathcal{D}$  de mêmes abscisses que, respectivement,  $M_1, M_2, \dots, M_n$ , d'équation  $y = ax + b$  qui est telle que la somme  $S = M_1P_1^2 + M_2P_2^2 + \dots + M_nP_n^2$  soit minimale.

On admet qu'une telle droite existe toujours et on dit que cette droite réalise un ajustement affine du nuage de  $y$  en  $x$  par la méthode dite *des moindres carrés*.

**Elle passe toujours par le point moyen  $G$  du nuage.**



### 1) Définition :

La droite  $\mathcal{D}$  d'équation  $y = ax + b$  obtenue par la méthode des moindres carrés est appelée *droite de régression de  $y$  en  $x$* .

*Remarque*

- La droite de régression de  $y$  en  $x$  minimise la somme des carrés des distances en ordonnée

### 2) Corrélation :



Le nombre qui décrit la **validité** de la droite d'ajustement et qui mesure le **degré de dépendance** linéaire entre les variables  $x$  et  $y$  est le coefficient de corrélation (de Pearson), noté  $r$ .

### 3) Propriétés :



**$-1 \leq r \leq 1$   $r$  est toujours compris entre  $-1$  et  $+1$ .**

Si  $r > 0$  : entre  $x$  et  $y$  il y a une **corrélation positive** (dépendance linéaire positive).

Si  $r < 0$  : entre  $x$  et  $y$  il y a une **corrélation négative**.

Si  $r = 0$  ou voisin de  $0$  : il n'y a pas de dépendance linéaire entre  $x$  et  $y$ .

Si  $r = 1$  ou  $r = -1$  : les points du nuage sont rigoureusement alignés ; la dépendance linéaire est parfaite.

Si  $r \geq 0,95$  alors la corrélation linéaire entre  $X$  et  $Y$  est forte. Dans ce cas un ajustement affine est justifié (Les points du nuage sont dans une situation proche de l'alignement).

Une forte corrélation entre deux grandeurs  $x$  et  $y$  ne signifie pas nécessairement qu'il y a un **lien de causalité** entre ces grandeurs. Par exemple, il est possible que les deux grandeurs soient des **effets d'une même cause**. Par exemple, on peut constater une forte corrélation entre les notes en latin et en mathématiques dans un groupe d'étudiants, ce qui ne veut pas dire pour autant que la bonne note obtenue en latin favorise une bonne note en mathématiques ou vice-versa. Cf fiche [Corrélation et causalité.pdf](#).



*Remarque.* On détermine  $a$ ,  $b$  et  $r$  avec la calculatrice ou le tableur (voir page suivante).

**Exemple.** On reprend l'exemple du gérant d'hypermarché.

1. Avec la calculatrice, déterminer l'équation de la droite d'ajustement  $\Delta$  obtenue par la méthode des moindres carrés. On arrondira  $a$  et  $b$  à  $10^{-2}$ .
2. Déterminer le coefficient de corrélation  $r$ , arrondi à  $10^{-2}$ . Que peut-on penser de la qualité de cet ajustement ?

## V) Utilisation de la calculatrice

On peut retrouver tous ces paramètres statistiques en utilisant les listes d'une calculatrice.

**Calculatrice Casio**

- Effacer les données précédentes : Dans le menu **STAT, DELA (F6), YES** et **EXE**.
- Saisir les valeurs de  $X$  dans la 1<sup>ère</sup> colonne (**LIST1**) et celles de  $Y$  dans la 2<sup>ème</sup> colonne (**LIST2**).

--	--	--

**Utilisation du graphique :**

- Pour paramétrer le graphique : **GRPH (F1), SET (F6)**. Choisir **Scatter** pour Graph Type, **List1** pour XList et **List2** pour YList et **EXE**.
- Pour afficher le nuage : **GPH1 (F1)**.
- Pour afficher les coordonnées du point G : **CALC (F2), 2VAR (F2)**. Faire défiler les valeurs.
- Pour afficher les résultats  $a$ ,  $b$  et  $r$ , de la « régression linéaire » (ajustement affine) : **CALC (F2), REG (F3) X (F1)** et **ax+b (F1)**.
- Pour dessiner la droite d'ajustement : **COPY (F6)**, puis dans le menu **GRAPH**, la dessiner.

--	--	--

**Calcul des paramètres :**

- Paramétrer les listes : **CALC (F2)** et **SET (F6)**. Choisir **List1** pour 2Var XList et **List2** pour 2Var YList et **List2** pour YList, puis **EXE**.
- Pour afficher les coordonnées du point G : **CALC (F1)** et **2VAR (F1)**. Faire défiler avec la touche  $\downarrow$ .
- Pour afficher les résultats  $a$ ,  $b$  et  $r$ , de la « régression linéaire » (ajustement affine) : **REG (F3)** et **X (F1)**.

**Pour obtenir des estimations :**

- Pour estimer la valeur de  $Y$  pour la valeur  $X = x_0$  : Dans le menu **RUN**,  $x_0$ , **OPTN STAT (F5) ŷ (F2)** et **EXE**.
- Pour estimer la valeur de  $X$  pour la valeur  $Y = y_0$  : Dans le menu **RUN**,  $y_0$ , **OPTN STAT (F5) x̂ (F1)** et **EXE**.

**Calculatrice TI**

- Effacer les données précédentes : **EDIT, 4:EffacerList** **entrer**, **L1,L2 (2nde 1)**, **2nde 2)** et **entrer**.
- Dans le menu **STATS 1:EDIT** **entrer**, saisir les valeurs de  $X$  dans **L1**, et celles de  $Y$  dans **L2** et quitter (**2nde Mode**).

--	--	--	--

**Pour utiliser le graphique :**

- Pour paramétrer le graphique : **graph stats (2nde fx)**, **Graph1** **entrer**.
- Mettre **On** en surbrillance, icône nuage en surbrillance, **L1** pour Xlist et **L2** pour Ylist, etc...
- Pour afficher le nuage : **Graph**. Modifier la fenêtre graphique **fenetre** en prenant en compte les valeurs des  $x_i$  et  $y_i$ .

--	--	--	--

**Pour obtenir les paramètres :**

- Paramétrer les listes :
- Pour afficher les coordonnées du point G : **CALC, 2:STATS 2-Var** **entrer** et **L1,L2 (2nde 1)**, **2nde 2)**. Faire défiler les valeurs.
- Pour afficher les résultats  $a$ ,  $b$  de la « régression linéaire » (ajustement affine) : **CALC, 4:RégLin(ax+b)** **entrer**, **L1,L2** et **entrer**.
- Pour afficher le coefficient de corrélation  $r$  : **VAR, 5:Statistiques** **entrer**, mettre **EQ** en surbrillance **7:r** **entrer** **entrer**.

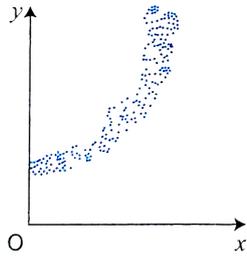
**Pour obtenir des estimations :**

- Pour estimer la valeur de  $Y$  pour la valeur  $X = x_0$  : **def tbl (2nde fenetre)**.
- Compléter avec la valeur de  $x_0$ , **table (2nde Graph)**. Affichage sur la 1<sup>ère</sup> ligne de la valeur de  $y_0$ .

## VI) Changement de variable

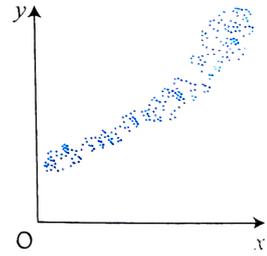
Dans certains cas, les points du nuage semblent se répartir autour d'une courbe autre qu'une droite. Il est parfois possible de se ramener à un ajustement affine à l'aide d'un changement de variable.

Dans un tel cas, on peut penser qu'il existe une relation du type  $y = ax^2 + b$  entre  $x$  et  $y$ .



En posant  $u = x^2$ , on se ramène à  $y = au + b$  et on peut déterminer  $a$  et  $b$  avec la méthode des moindres carrés.

Dans un tel cas, on peut penser qu'il existe une relation du type  $y = ae^x + b$  entre  $x$  et  $y$ .



En posant  $u = e^x$ , on se ramène à  $y = au + b$  et on peut déterminer  $a$  et  $b$  avec la méthode des moindres carrés.

Méthode :

Exercice résolu n° 6 page 261 du manuel numérique Hyperbole ed. Nathan