

Statistiques à deux variables, cours, terminale, Mathématiques complémentaires

F.Gaudon

13 mai 2023

Table des matières

1	Vocabulaire	2
2	Ajustement d'un nuage de points	3

1 Vocabulaire

Définition :

- Soient x et y deux caractères quantitatifs d'une même population. A chaque individu de la population on associe un couple $(x_i; y_i)$ où x_i et y_i pour $i \in \{1; \dots; n\}$ avec n entier naturel sont les valeurs prises respectivement par x et y . L'ensemble de ces couples constitue une *série statistique à deux variables* x et y .
- Dans un repère $(O; \vec{i}; \vec{j})$, l'ensemble des points M_i de coordonnées $(x_i; y_i)$ est appelé *nuage de points* associé à la série statistique.
- Soit une série statistique à deux variables x et y de moyennes \bar{x} et \bar{y} .

Le point G de coordonnées $(\bar{x}; \bar{y})$ avec

$$\bar{x} = \frac{x_1+x_2+\dots+x_n}{n} \text{ et } \bar{y} = \frac{y_1+y_2+\dots+y_n}{n}$$

est appelé le *point moyen* du nuage de points associé à la série statistique.

Exemple :

Un magasin réalise une étude sur l'influence du prix de vente sur le nombre de machines à laver vendues au cours d'une année. Le tableau suivant donne les résultats de cette étude :

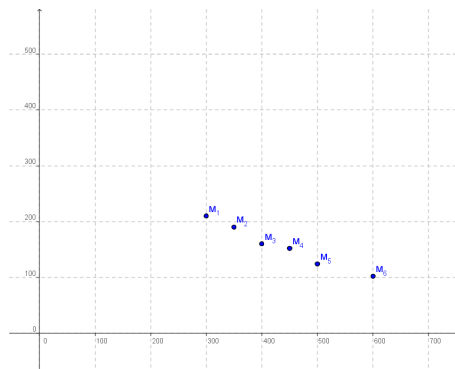
Prix x_i en euros	300	350	400	450	500	600
Nombre de machines vendues	210	190	160	152	124	102

Le nuage de points associé à cette série est constitué des points M_i pour i allant de 1 à 6 dont les coordonnées sont $(300; 210)$, $(350; 190)$, ..., $(600; 102)$.

Le point moyen associé à ce nuage de points est le point G de coordonnées $(\bar{x}; \bar{y})$ données par :

$$\bar{x} = \frac{300+350+\dots+600}{6} = \frac{2285}{6} \approx 433$$

$$\bar{y} = \frac{210+190+\dots+102}{6} = \frac{938}{6} \approx 256$$



2 Ajustement d'un nuage de points

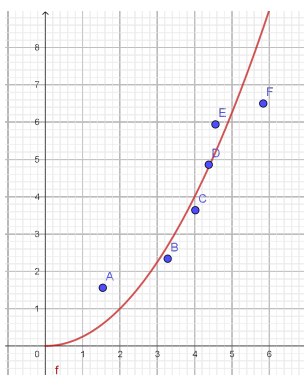
Définition :

- on appelle *ajustement* du nuage de points toute courbe « résumant approximativement » le nuage.
- Toute droite « résumant approximativement » le nuage est appelée *droite d'ajustement affine* du nuage de points.

Remarque :

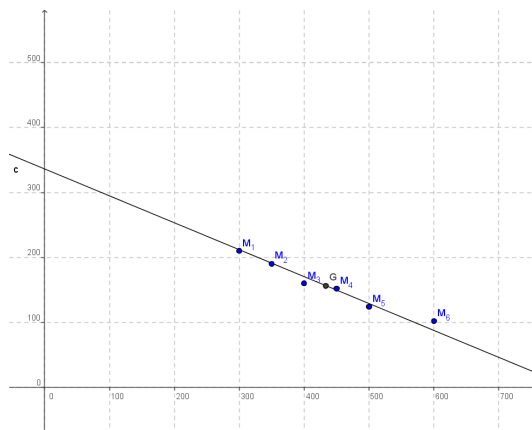
Certains nuages peuvent ne pas sembler être approchables par une quelconque courbe auquel cas les deux variables ne sont pas reliées entre elles.

L'exemple suivant suggère un ajustement par une fonction polynôme du second degré. On parle alors d'ajustement quadratique :



Exemple de méthode approchée :

On trace « au jugé » une droite qui « semble résumer » le nuage de points. C'est une méthode simple mais qui dépend de la droite tracée.



Propriété (rappel) :

Soient A et B de coordonnées $(x_A; y_A)$ et $(x_B; y_B)$ deux points tels que $x_A \neq x_B$. Alors la droite (AB) n'est pas parallèle à l'axe des ordonnées, elle a donc une équation de la forme $y = mx + p$ et on a :

$$m = \frac{y_B - y_A}{x_B - x_A}$$

Exemple de savoir faire [Détermination de l'équation d'une droite dont on connaît les coordonnées de deux points] :

Soit \mathcal{D} la droite passant par les points $M_1(300; 208)$, $M_2(350; 190)$ Son équation est de la forme $y = mx + p$.

$$m = \frac{y_B - y_A}{x_B - x_A} = \frac{208 - 190}{300 - 350} = \frac{18}{-50} = -0.36$$

donc son équation est $y = -0,36x + p$.

Or M_1 appartient à \mathcal{D} donc ses coordonnées vérifient l'équation d'où $208 = -0,36 \times 300 + p$ donc $208 = -108 + p$ et $208 + 108 = p$ donc $p = 316$.

L'équation est donc $y = -0,36x + 316$.

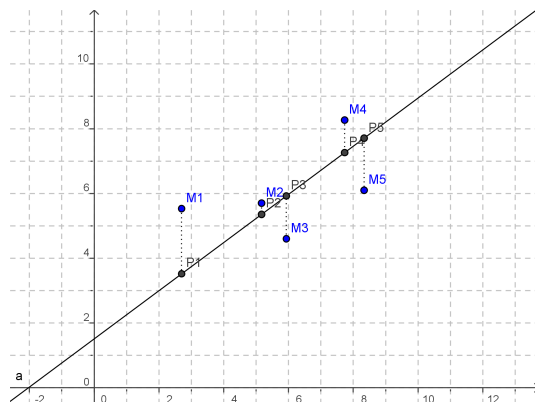
Méthode des moindres carrés :

Avec les notations de la figure ci-dessous, étant donné un nuage de n points M_i , il existe une droite passant par le point moyen G et telle que la somme des carrés des écarts (ou *résidus*) $P_1M_1^2 + P_2M_2^2 + \dots + P_nM_n^2$ soit minimale. Cette droite est appelée *droite de régression de y en x* . On peut montrer que son équation réduite est $y = mx + p$ avec :

$$m = \frac{(x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_p - \bar{x})(y_p - \bar{y})}{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_p - \bar{x})^2}$$

et

$$p = \bar{y} - m\bar{x}$$



Exemple de recherche de l'équation réduite à l'aide des formules :

On reprend l'exemple précédent.

$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
	Total		

D'où $m \approx -0,37$

et $p \approx 315,06$

Exemple de recherche de l'équation à l'aide de la calculatrice :

* TI 82 et plus :

Aller dans le menu **STAT** puis **EDIT**. Entrer les valeurs x_i dans la colonne L_1 et les valeurs y_i dans la colonne L_2 . Quitter (**2nde** **QUIT**) puis menu **STAT** et **CALC**. Choisir **LinReg(ax+b)** puis **2nd L1**, **2nd L2** pour indiquer les deux colonnes à utiliser. Valider ensuite **ENTER**.

* CASIO Graph 25 et plus :

Aller dans le menu **STAT** puis entrer les valeurs x_i dans la colonne 1 et les valeurs y_i dans la colonne 2. Sélectionner ensuite **CALC**. Choisir **SET** et vérifier que la ligne « 2Var XList » est mise à « List1 » et que la ligne « 2Var YList » est mise à « List2 », sinon choisir le menu **LIST** pour indiquer les numéros de liste adaptés. Taper ensuite sur **EXIT** puis choisir **REG** puis **X**.

Obtention de l'équation réduite à l'aide d'un programme :

```
def moindreCarres(L1,L2):  
    mL1=0  
    for k in range(len(L1)):  
        mL1=mL1+L1[k]  
    mL1=mL1/len(L1)  
    mL2=0  
    for k in range(len(L2)):  
        mL2=mL2+L2[k]  
    mL2=mL2/len(L2)  
    numerateur=0  
    denominateur=0  
    for k in range(len(L1)):  
        numerateur=numerateur+(L1[k]-mL1)*(L2[k]-mL2)  
        denominateur=denominateur+(L1[k]-mL1)**2  
    m=numerateur/denominateur  
    p=mL2-a*mL1  
    return m,p
```

Propriété et définition :

On a vu que la droite de régression linéaire de y en x rend minimal le nombre $P_1M_1^2 + \dots + P_nM_n^2$. Le minimum est alors égal à :

$$n(\sigma_y)^2 \left(1 - \left(\frac{\sigma_{xy}}{\sigma_x \sigma_y}\right)^2\right)$$

avec

$$\sigma_x = \frac{1}{n}((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2)$$

$$\sigma_y = \frac{1}{n}((y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_n - \bar{y})^2)$$

$$\sigma_{xy} = \frac{1}{n}((x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}))$$

On appelle *coefficient de corrélation linéaire* le nombre défini par :

$$r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

Plus la valeur absolue de ce nombre est proche de 1, plus la droite de régression linéaire est proche du nuage de points.

Remarque : En pratique, on lira le coefficient de corrélation linéaire sur la calculatrice en même temps que la droite des moindres carrés obtenue : il est noté r .