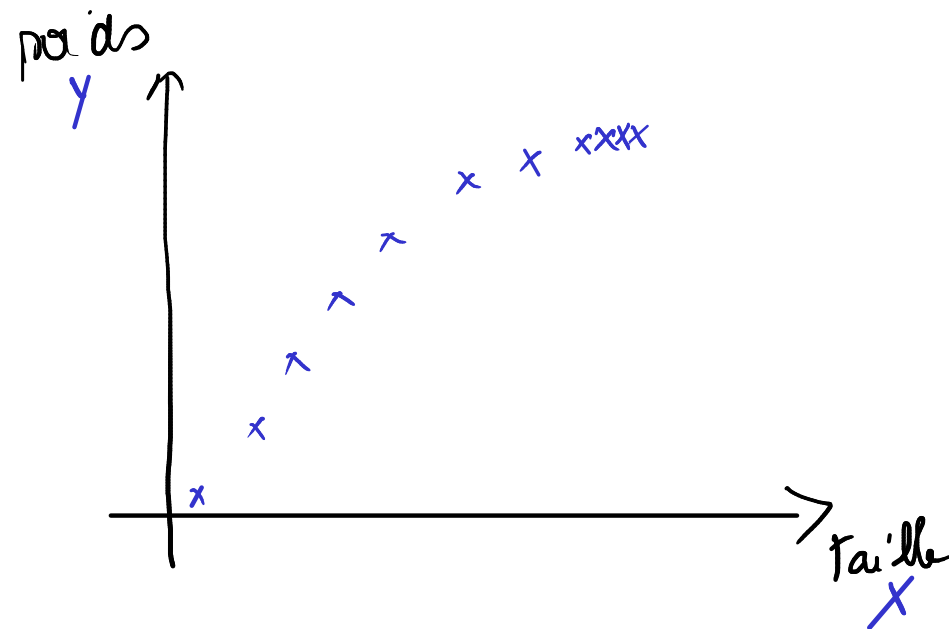


Part III

Régressions linéaires

Carnet de santé:

! on ne regroupe pas par modalité
l'effet de chaque observation est 1



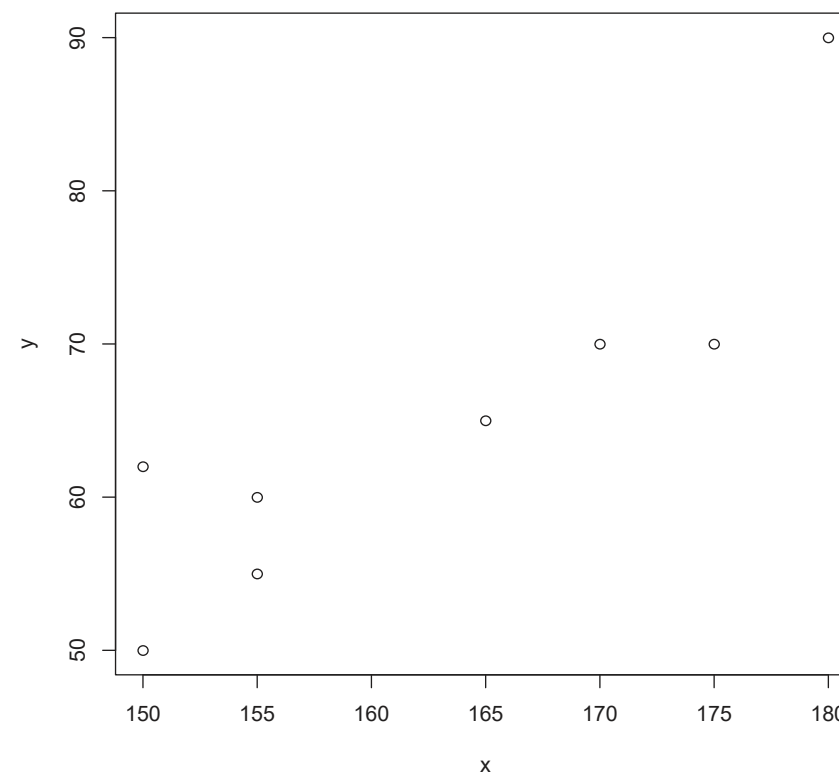
x	x_1	x_2	...	x_n
y	y_1	y_2	...	y_n

4.1 Nuage de points

On relève le couple (taille, poids) de 8 individus. On résume les données dans le tableau suivant.

taille	x	150	155	155	150	165	175	170	180
poids	y	50	55	60	62	65	70	70	90

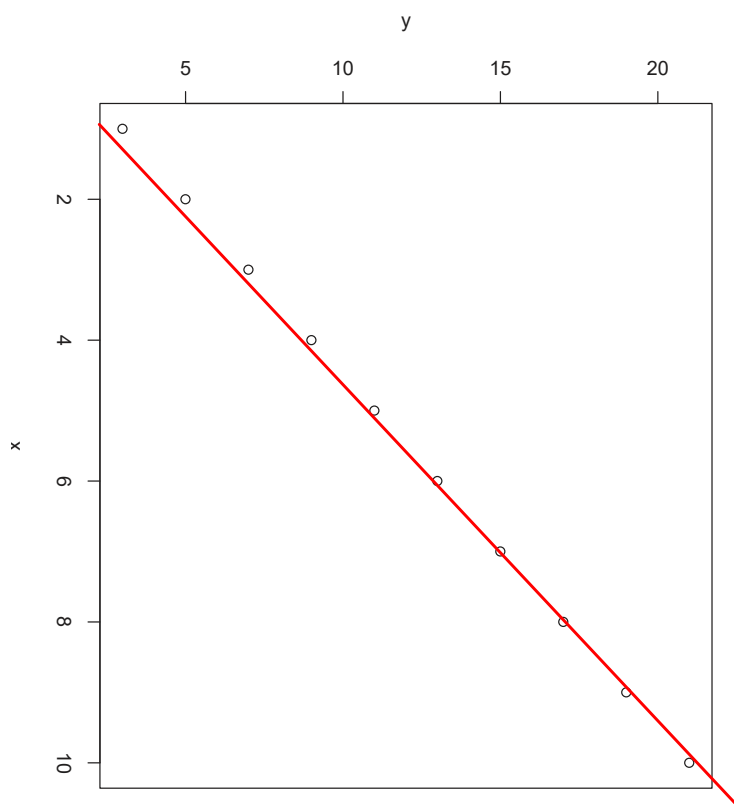
Definition 4.1 Soit une population de N individus. Le graphe des N points (x_i, y_i) est appelé nuage de points de la série.



Nuage de points

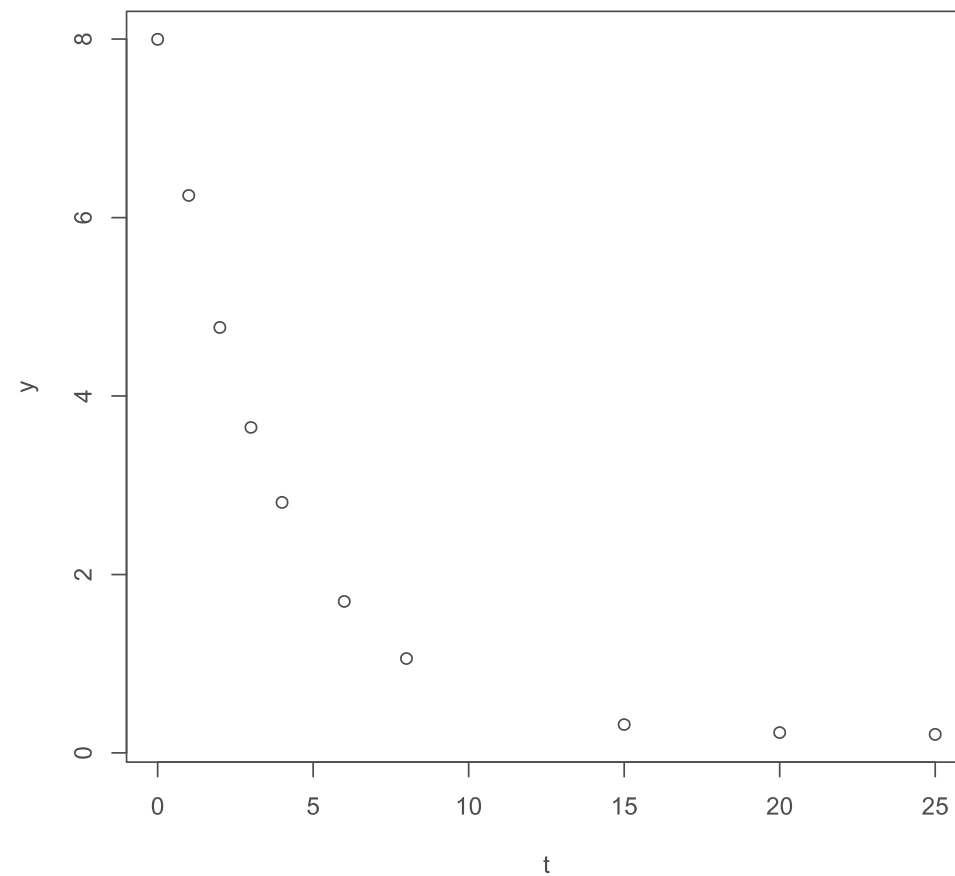
Definition 4.2 Le point ayant pour coordonnées les moyennes (\bar{x}, \bar{y}) est appelé le point moyen.

Nuage de point rectiligne.

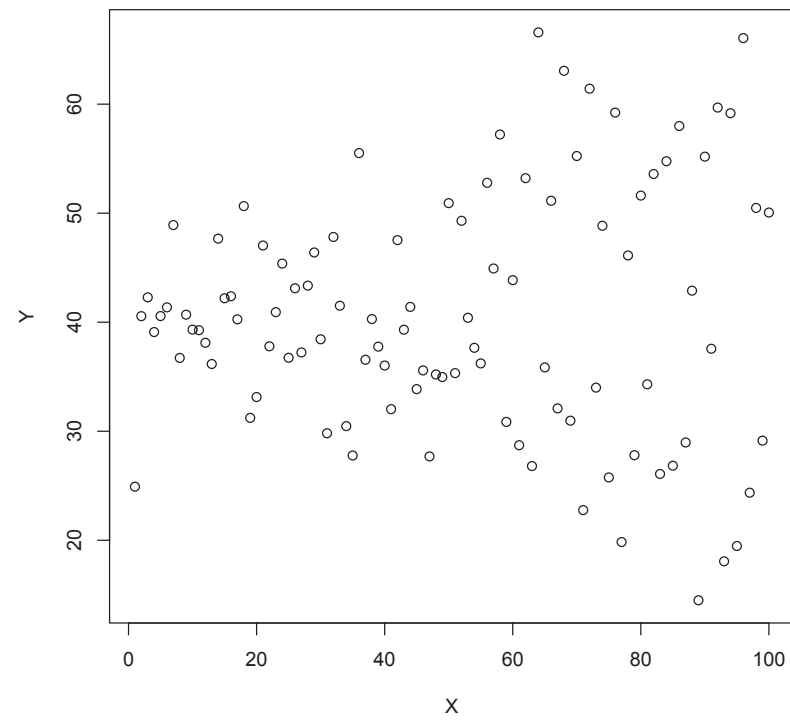


I dée : Cherchons la droite qui se trouve
au plus proche des points !

Nuage ordonné mais non rectiligne



(iii) forme quelconque



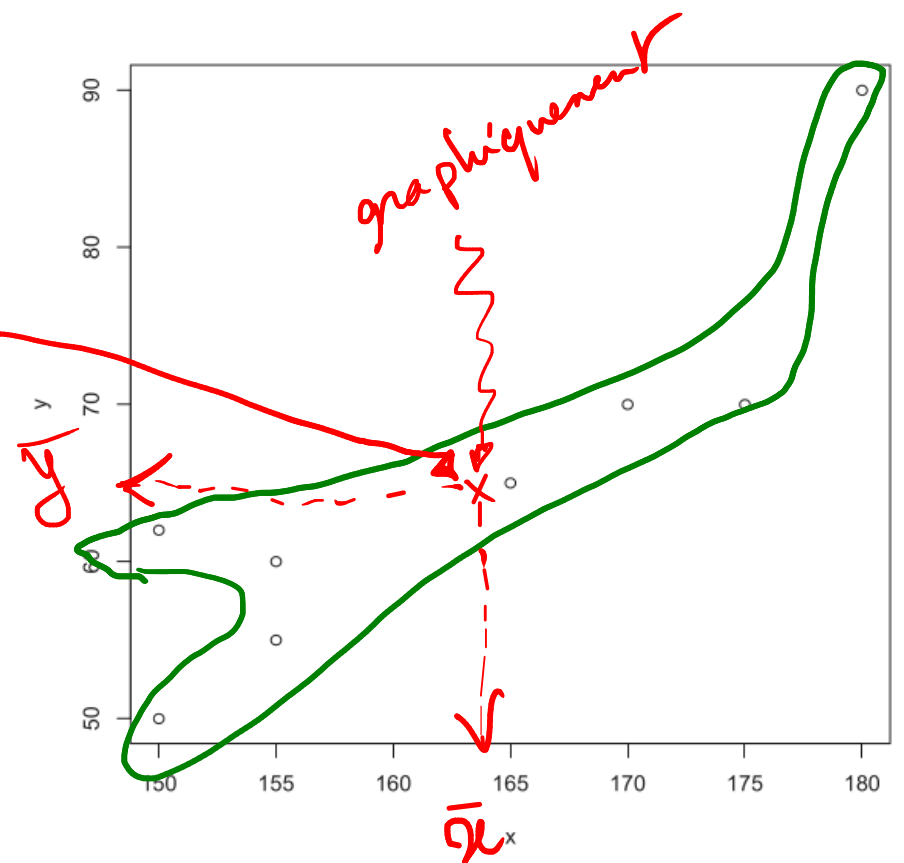
5 Ajustement affine (droite de régression linéaire)

On s'intéresse plus particulièrement au premier cas 4.2.1. Procéder à un ajustement affine revient à chercher une droite D d'équation

$$y = ax + b$$

Centre de gravité (\bar{x}, \bar{y})

Toute droite de régression passe par le centre de gravité.

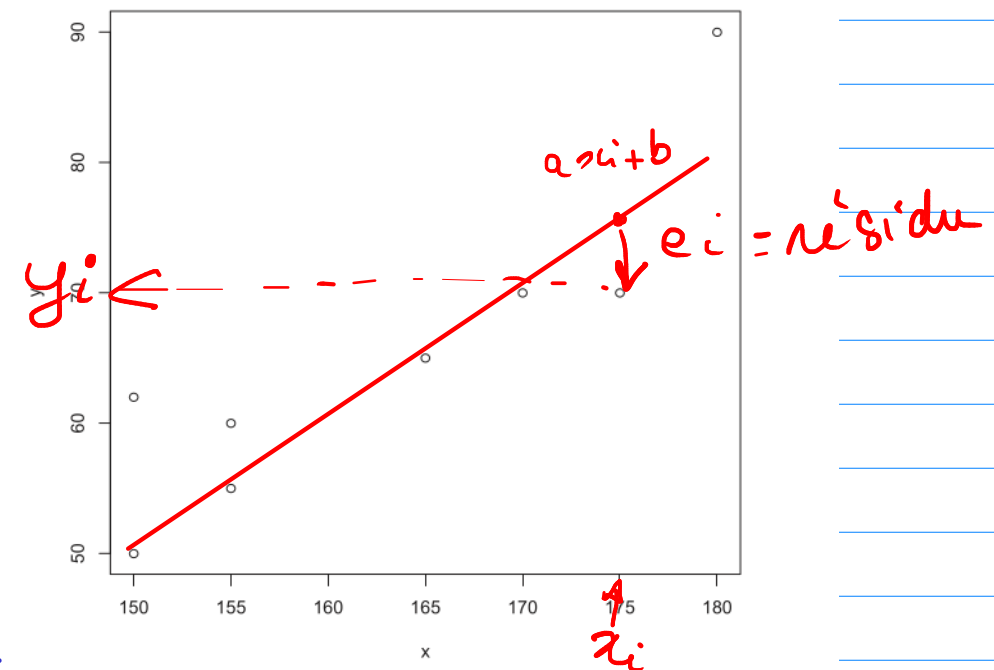


Note: Toute droite de régression passe par le centre de gravité.

$$y_i = ax_i + b + e_i$$

← erreur
entre la droite
et l'observation

$$\begin{aligned} \text{D'où } \bar{y} &= \overline{ax + b + e} \\ &= a\bar{x} + b + \bar{e} \end{aligned}$$



On ne veut pas d'erreur systématique
donc $\bar{e} = 0$. (sinon, on retire \bar{e}
à b pour obtenir un modèle sans erreur systématique).

$$\text{D'où } \bar{e} = 0 \iff \bar{y} = a\bar{x} + b$$

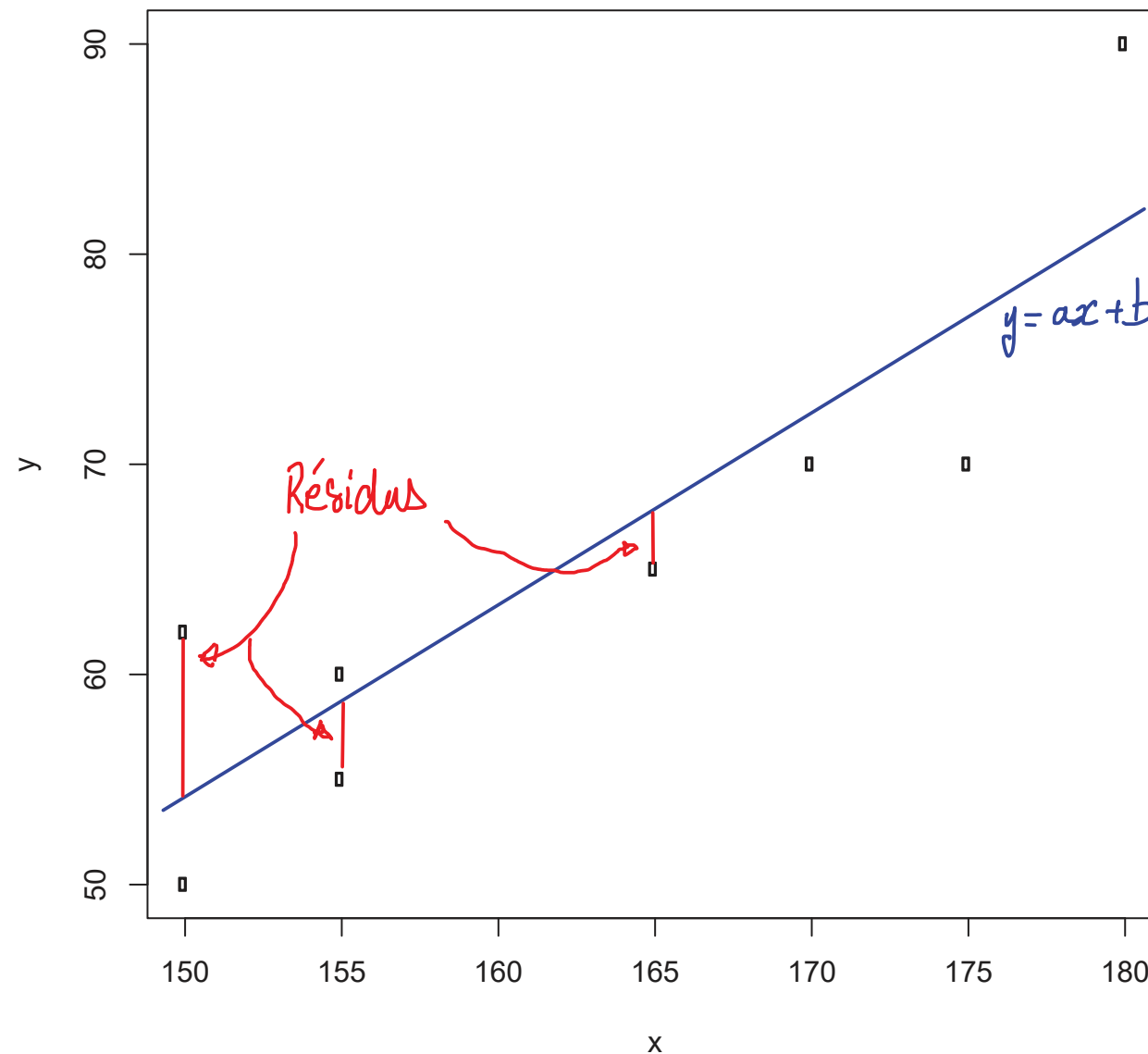
Donc (\bar{x}, \bar{y}) vérifie l'équation de la droite.

Toute droite de régression passe par le centre de gravité.

Pb: Trouver a dans $y = ax + b$!

5.1 La méthode des moindres carrés

L'idée de cette méthode est de chercher la droite $y = ax + b$ qui minimise la somme des carrés des écarts verticaux entre la droite et les points du nuage, les *résidus*.



On cherche à minimiser la distance verticale entre les points et la droite:

$$\sum_i e_i^2$$

↑
résidus

On cherche donc le " a " qui minimise la somme des carrés des résidus $\sum_i e_i^2$

Rappel:

$$\bar{x} = \frac{1}{N} \sum_i x_i$$

$$\bar{y} = \frac{1}{N} \sum_i y_i$$

$$\text{Var } x = \frac{1}{N} \sum_i (x_i - \bar{x})^2$$

$$\text{Var } y \dots$$

$$= m(x^2) - \bar{x}^2$$

$$= \frac{1}{N} \sum_i x_i^2 - \bar{x}^2$$

$$\sigma_x = \sqrt{\text{Var } x}$$

$$\sigma_y = \sqrt{\text{Var } y}$$

$$\text{Cov}(x, y) = \frac{1}{N} \sum_i (x_i - \bar{x})(y_i - \bar{y}) = m(xy) - \bar{x} \bar{y}$$

Note: $\text{Cov}(x, x) = \frac{1}{N} \sum_i (x_i - \bar{x})(x_i - \bar{x})$
 $= \frac{1}{N} \sum_i (x_i - \bar{x})^2 = \text{Var } x$

démontre au S4.

Thm:

$$a = \frac{\text{Cov}(x, y)}{\text{Var } x}$$

et on trouve b en vérifiant $\bar{y} = a\bar{x} + b$.

Exemple:

taille	x	150	155	155	150	165	175	170	180
poids	y	50	55	60	62	65	70	70	90

$$\bar{x} = \frac{150 + 155 + 155 + \dots + 170 + 180}{8} = 162,5$$

$$\bar{y} = \frac{50 + 55 + \dots + 70 + 90}{8} = 65,25$$

$$\begin{aligned} \text{Var } x &= m(x^2) - \bar{x}^2 = \frac{150^2 + 155^2 + \dots + 180^2}{8} - 162,5^2 \\ &= 26\,525 - 26\,406,25 \\ &= 118,75 \end{aligned}$$

$$\text{Var } y = m(y^2) - \bar{y}^2 = 4386,75 - 4257,562 \simeq 129,2$$

Calculons $\text{cov}(x, y)$.

taille	x	150	155	155	150	165	175	170	180
poids	y	50	55	60	62	65	70	70	90

$$\text{Cov}(x, y) = m(x, y) - \bar{x} \bar{y}$$

$$= \frac{50 \times 150 + 55 \times 155 + 60 \times 155 + \dots + 90 \times 180}{8} - 162,5 \times 65,25$$

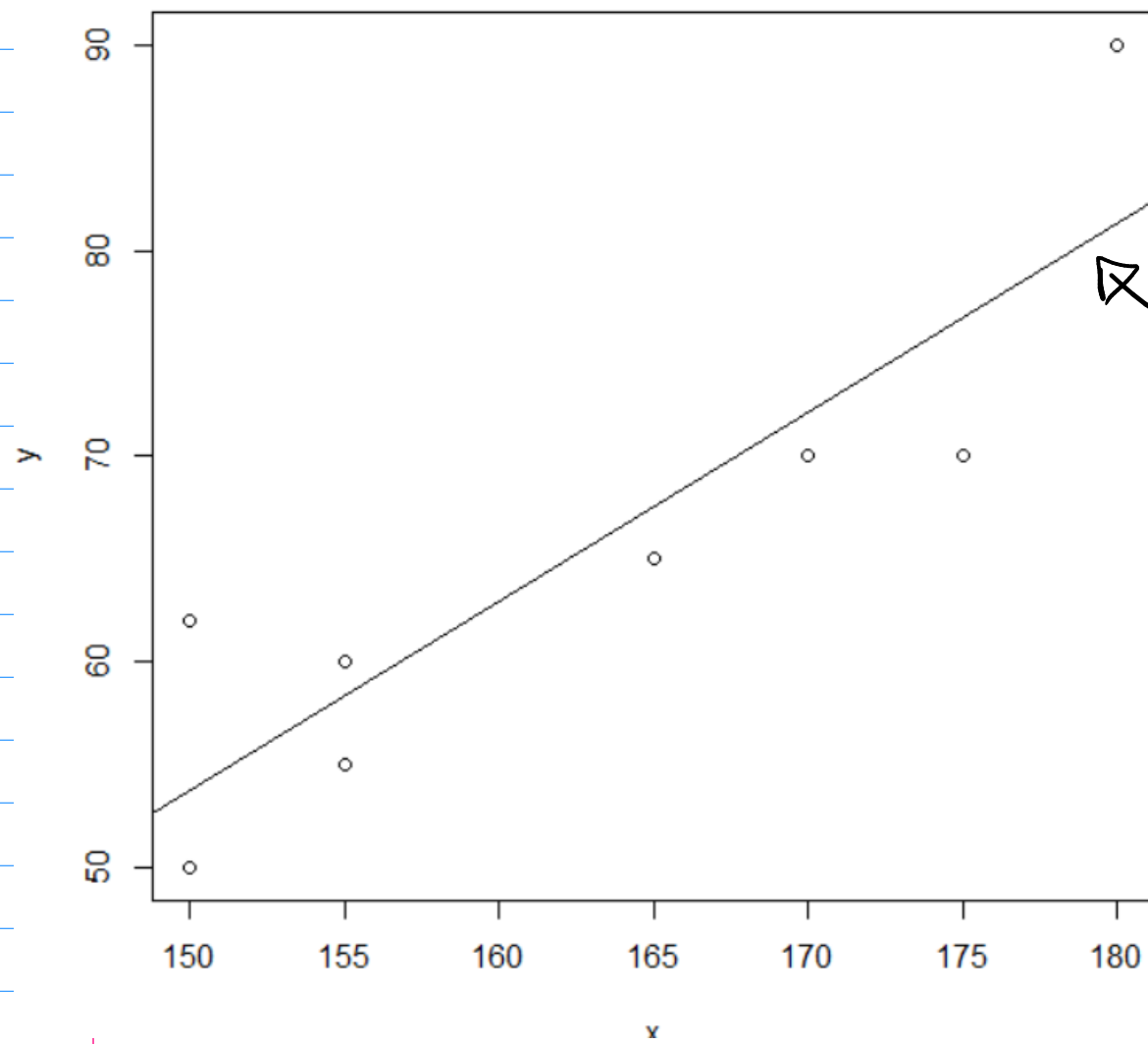
$$= 10712,5 - 10603,125$$

$$= 109,4$$

Équation de la droite des moindres carrés : $y = ax + b$.

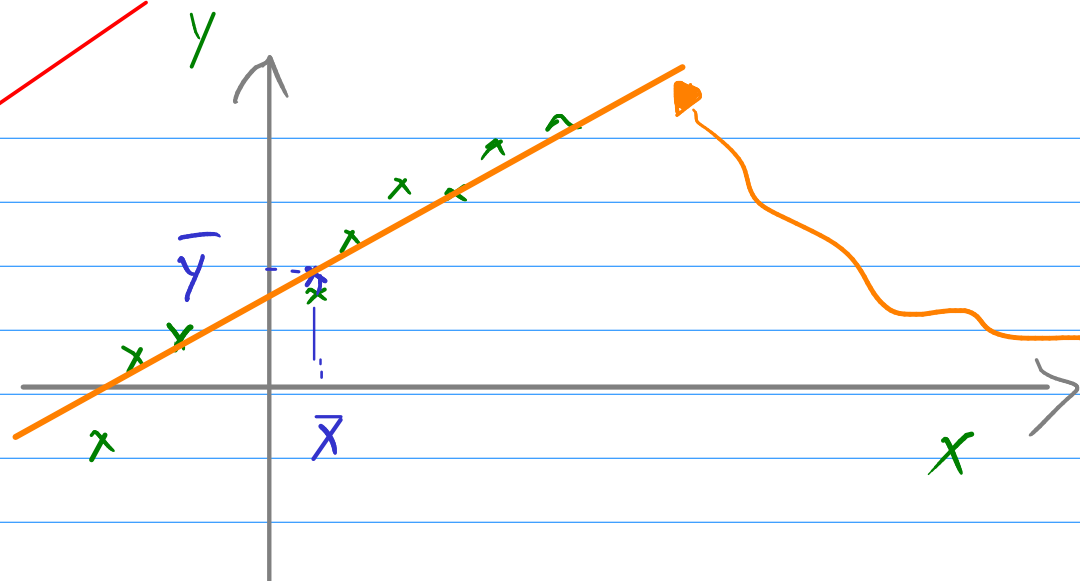
$$a = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{109,4}{118,75} = 0,92$$

$$\text{et } \bar{y} = a\bar{x} + b \Leftrightarrow b = \bar{y} - a\bar{x} = 65,25 - 0,92 \times 162,5 = -84,25$$



$$y = 0,92x - 84,25$$

Rappel



(\bar{x}, \bar{y}) centre de gravité du nuage de points

droite qui passe par le centre de gravité et au plus proche des points

$$y = ax + b$$

$$a = \frac{\text{Cov}(x, y)}{\text{Var } x}$$

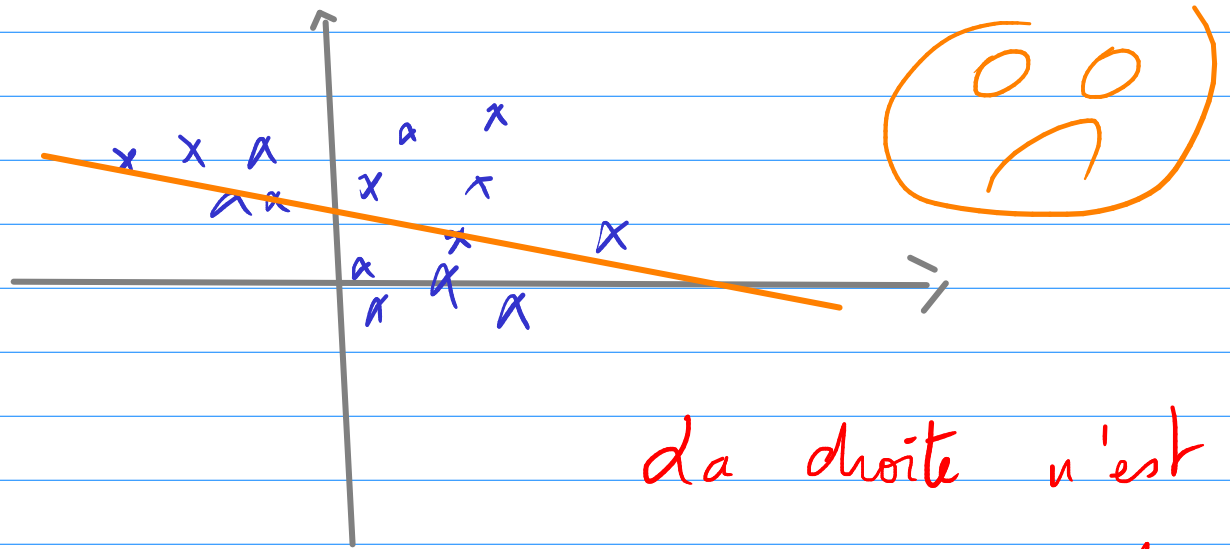
$$\begin{aligned} \text{Cov}(x, y) &= \frac{1}{N} \sum_1^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= m(xy) - \bar{x}\bar{y} \end{aligned}$$

$$\begin{aligned} \text{Var } x &= \frac{1}{N} \sum_1^n (x_i - \bar{x})^2 \\ &= m(x^2) - \bar{x}^2 \end{aligned}$$

$\bar{y} = a\bar{x} + b$ car (\bar{x}, \bar{y}) est sur la droite, donc

$$b = \bar{y} - a\bar{x}$$

Note: Ainsi définies, on peut toujours calculer a et b et ainsi trouver une droite qui passe au plus proche des points -



la droite n'est pourtant pas un modèle satisfaisant.

* Gardons le modèle de droite de régression pour le cas où le nuage est rectiligne -

\Rightarrow C'est le coefficient r_{xy} qui nous l'indique ...

5.2 Coefficient de corrélation linéaire

Notons que la méthode des moindres carrés peut être utilisée pour n'importe quelle série double. On peut tout à fait obtenir une droite de régression dans le cas 4.2.3. Pour s'assurer de façon objective (et non purement visuelle) que l'ajustement est valide, on considère un autre paramètre de la série : le coefficient de corrélation r :

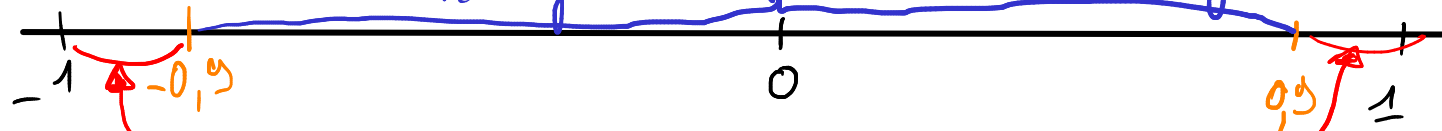
$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$-1 \leq r \leq 1$$

r est du même
signe que a ($y = ax + b$)

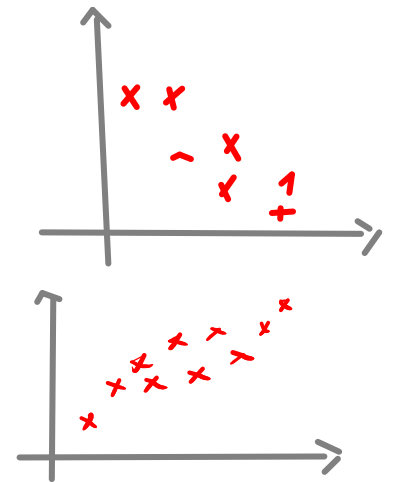
nuage de points non rectiligne

• r_{xy} :



$r < 0$

$r > 0$



régression linéaire
légitime.

Dans l'exemple:

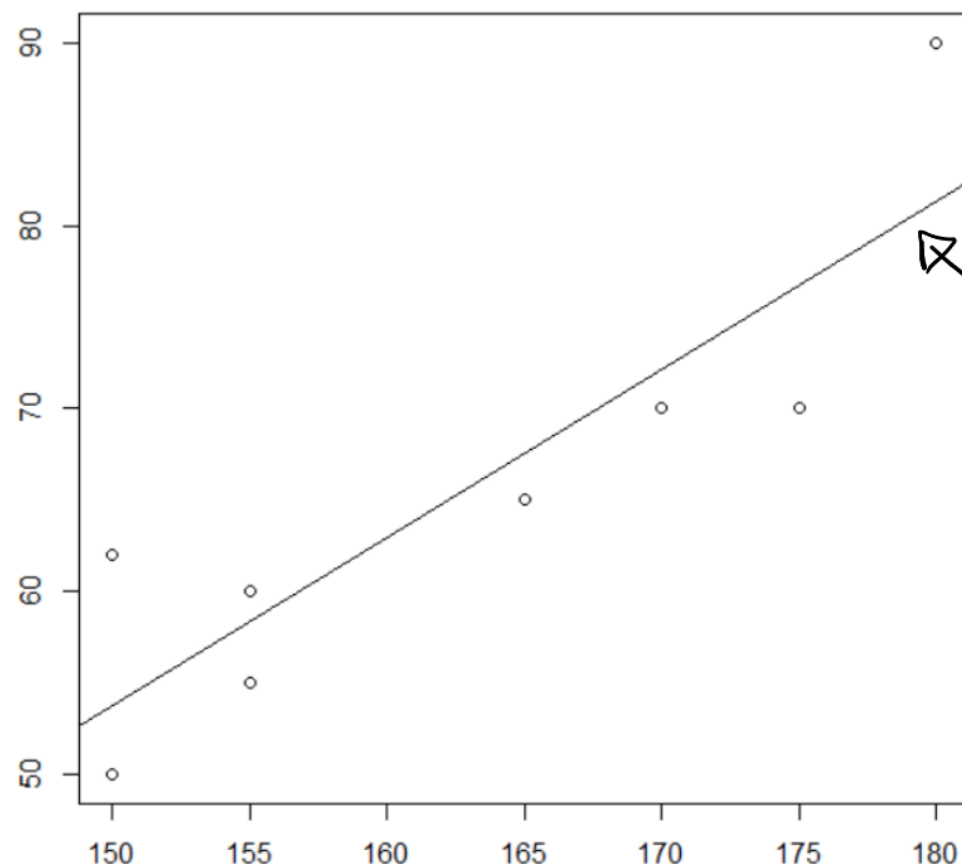
taille	x	150	155	155	150	165	175	170	180
poids	y	50	55	60	62	65	70	70	90

$$\text{Cov}(x, y) = 109,4$$

$$\text{Var } x = 118,75$$

$$\text{Var } y = 129,2$$

$$r_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var } x} \sqrt{\text{Var } y}} = \frac{109,4}{\sqrt{118,75} \sqrt{129,2}} = 0,88$$



$$y = 0,92x - 84,25$$

$$-0,9 < r_{xy} < 0,9$$

La droite n'est pas
très proche des points.

On devrait éviter d'utiliser ce
modèle.

$$R^2 = r_{xy}^2 = 0,88^2 = 0,78.$$

78 % de la variance des observation est expliquée par la droite.

5.3 Coefficient de détermination R^2

La variance est une bonne mesure de l'hétérogénéité d'une série (contrairement à la moyenne qui considère tous les individus comme semblables). La variance de la série x se décompose comme la variance expliquée par la droite de régression plus celle de l'erreur (résidus) :

$$Var(y) = Var(ax + b) + Var(e).$$

Le coefficient de détermination est le rapport de variance de y expliquée par la régression :

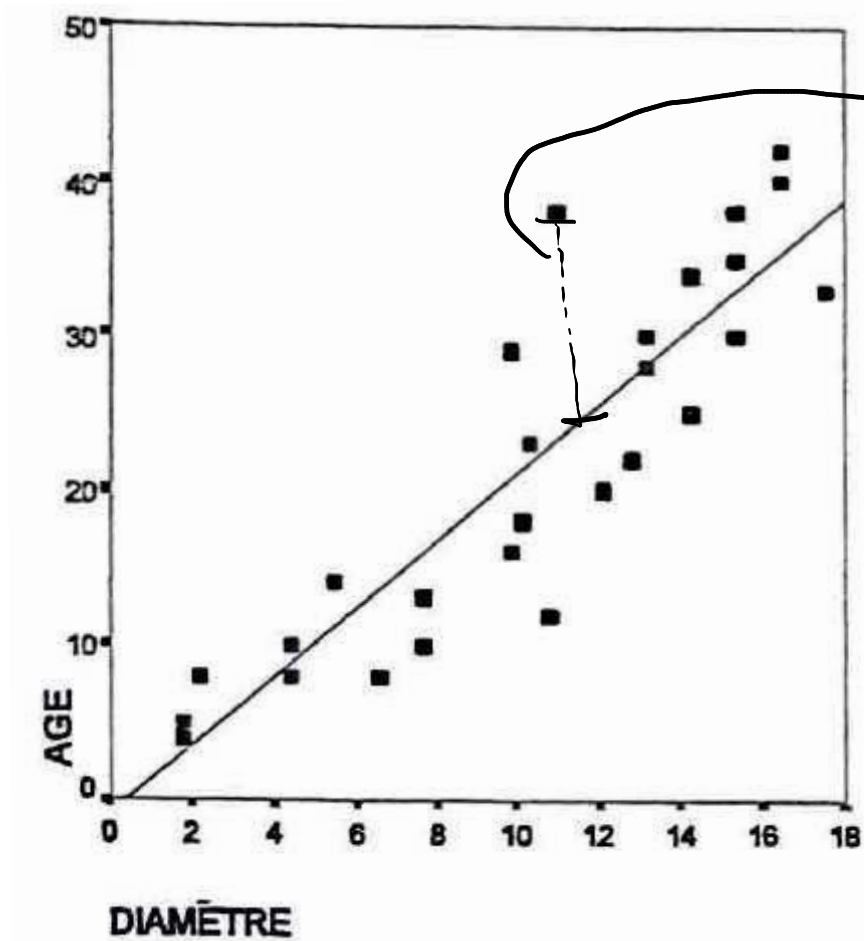
$$R^2 = \frac{Var(ax + b)}{Var(y)}. \quad (5.1)$$

Il se trouve que R^2 est le carré du coefficient de corrélation linéaire r_{xy} .

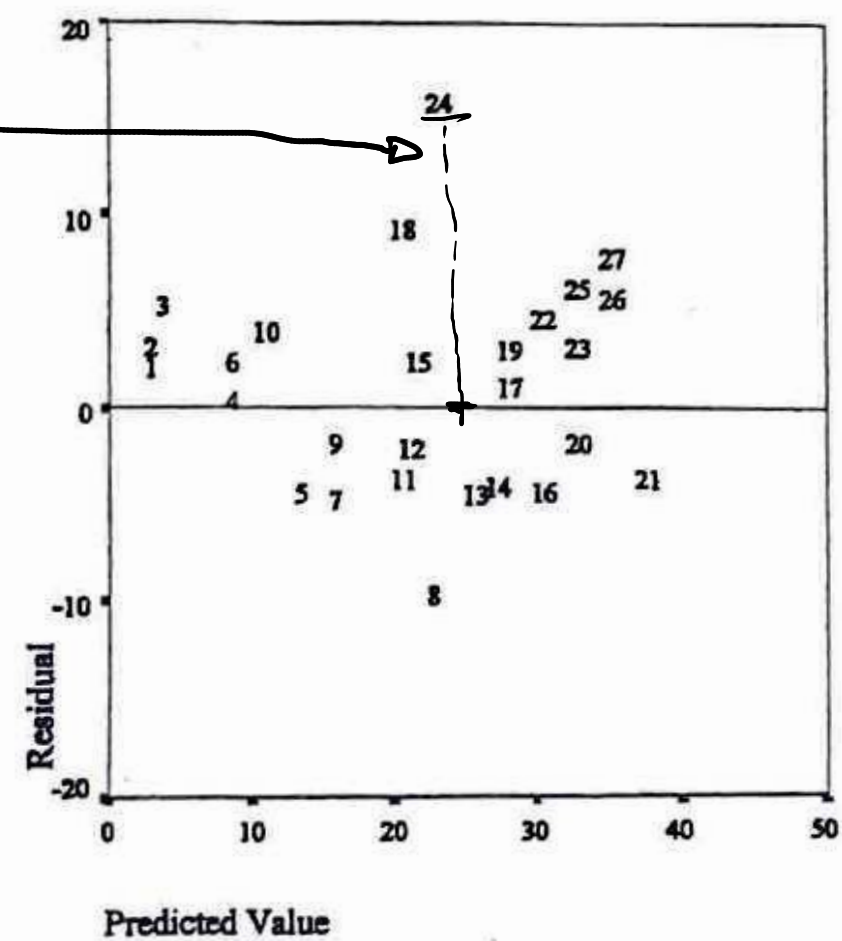
$$r_{xy} > 0,9 \text{ ou } r_{xy} < -0,9 \Leftrightarrow R^2 > 0,8 \Rightarrow \text{80\% de variance expliquée.}$$

Sur logiciel : Age et diamètre de 27 arbres .

Nuage de points



Résidus



Résumé de la régression

Allez chercher les infos importantes...

	Somme des carrés	ddl	Carré moyen	F	Signification
Régression	2905,549	1	2905,55	93,44	,000
Résidu	777,414	25	31,097		
Total	3682,963	26			

	Coefficients non standardisés		Coefficients standardisés	t	Signification
	B	Erreur standard	Bêta		
(constante)	-,974	2,604		-,374	,711
DIAMETRE	2,206	,228	,888	9,67	,000

coefficients

$$\text{Age} = 2,206 \text{ diamètre} - 0,974$$

$$r_{xy} = 0,888$$

$$|r_{xy}| < 0,9$$

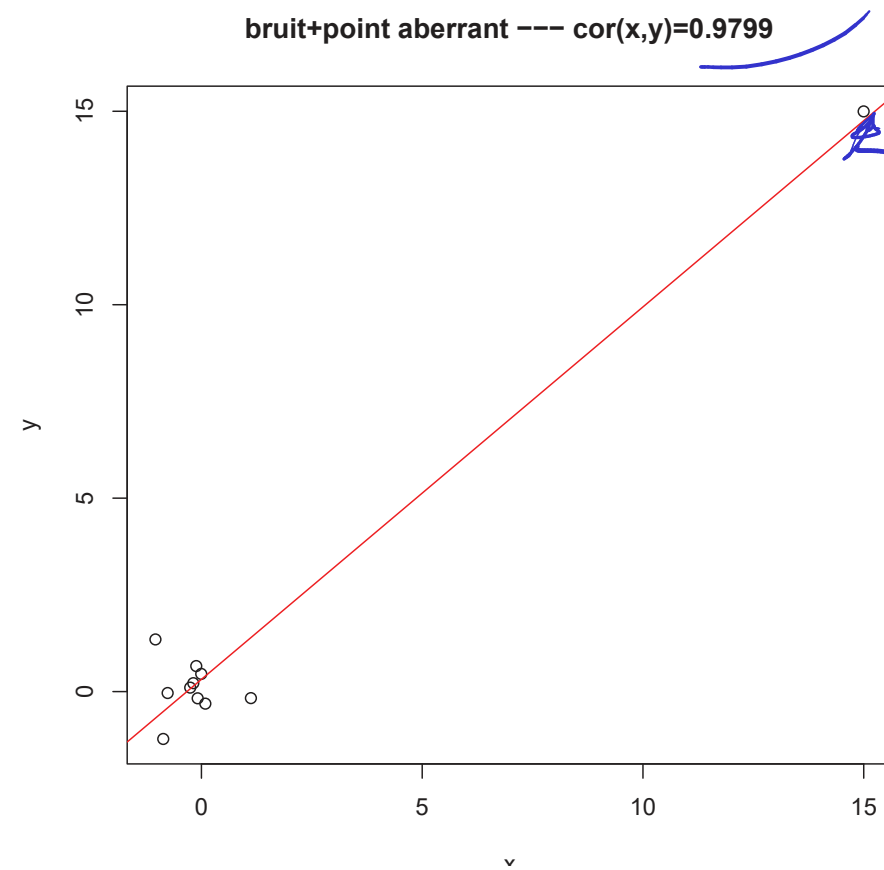
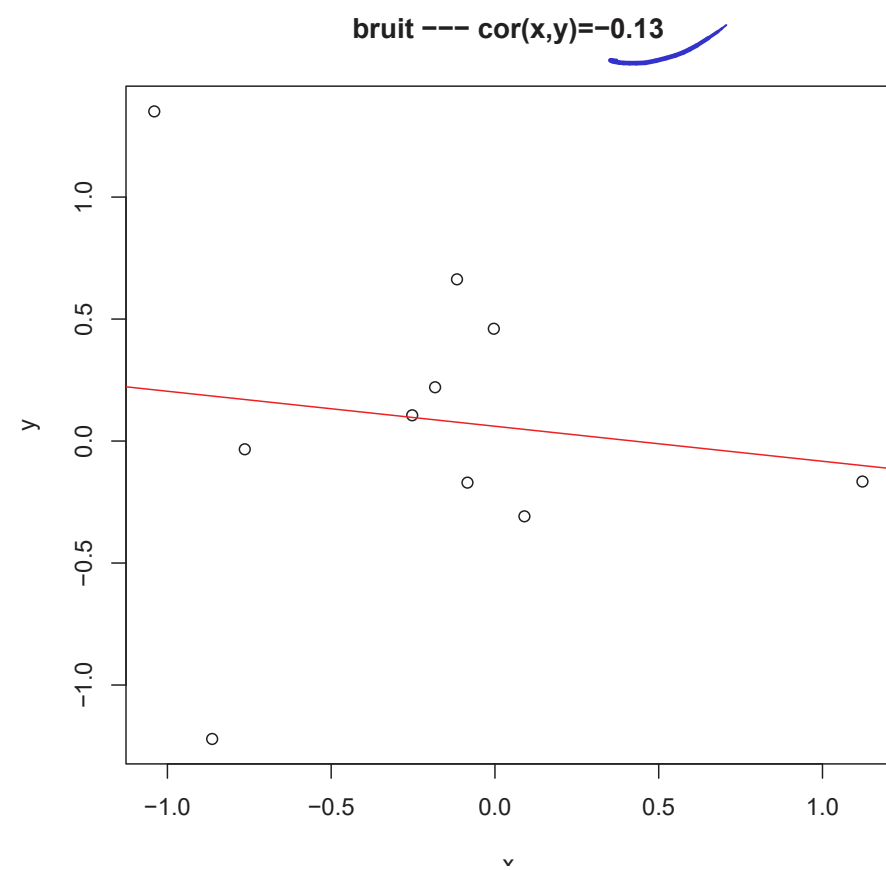
ça donne une tendance mais pas une prévision fiable.

6 Discussions

6.1 Manipulation du coefficient r

Le coefficient de corrélation linéaire r ou le coefficient de détermination R^2 mesurent le caractère rectiligne du nuage. Il ne suffira pas à décrire, seul, la pertinence d'un ajustement.

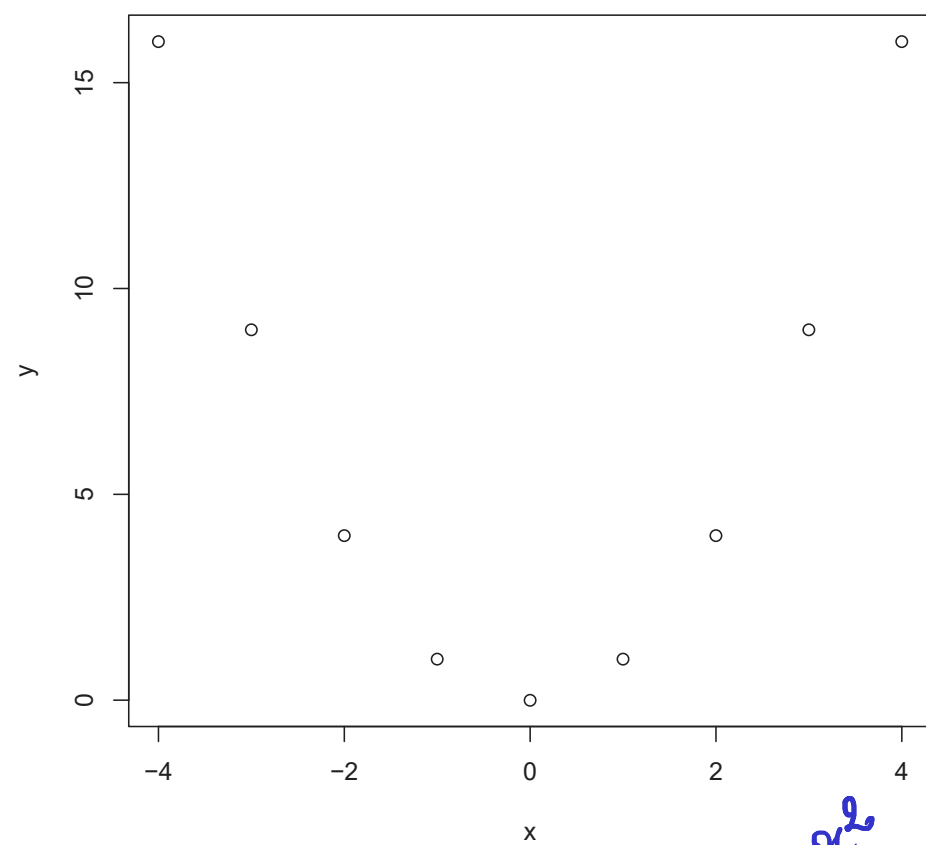
Il faudra, en particulier, faire attention aux points aberrants.



point aberrant

Toujours tracer le nuage de points.

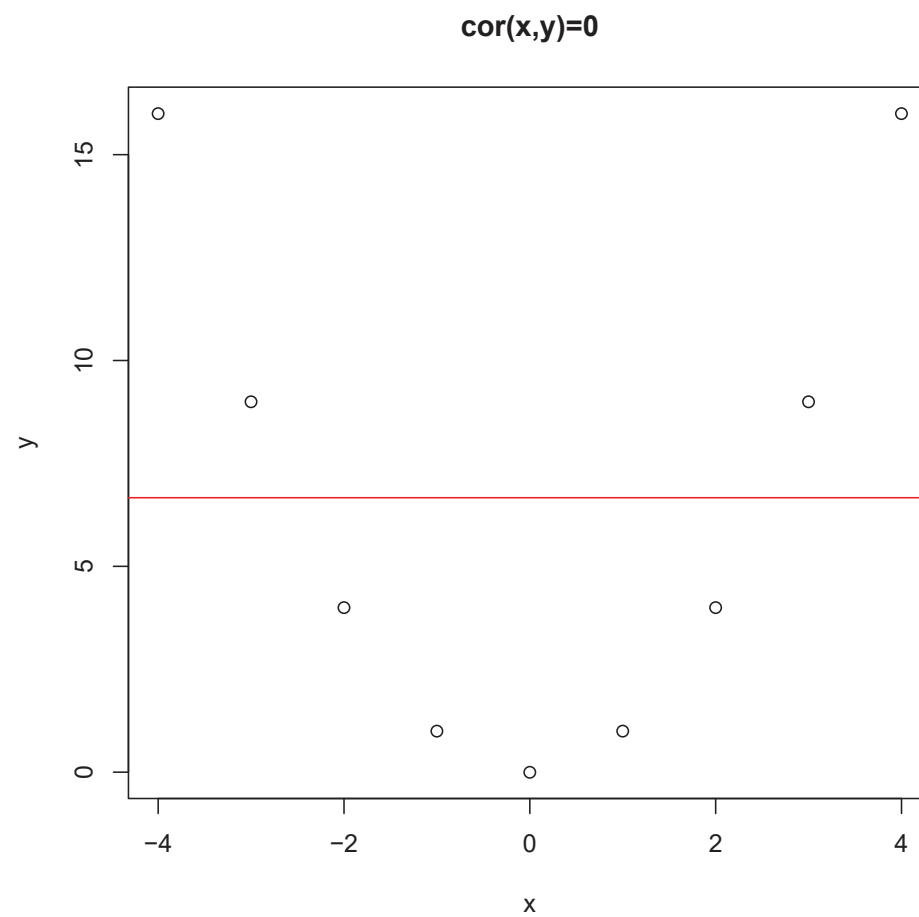
Notons également qu'ils passent pleinement à côté d'une forme en cloche. Bien qu'il y ait une liaison entre les variables, la nature de cette liaison n'est pas linéaire.



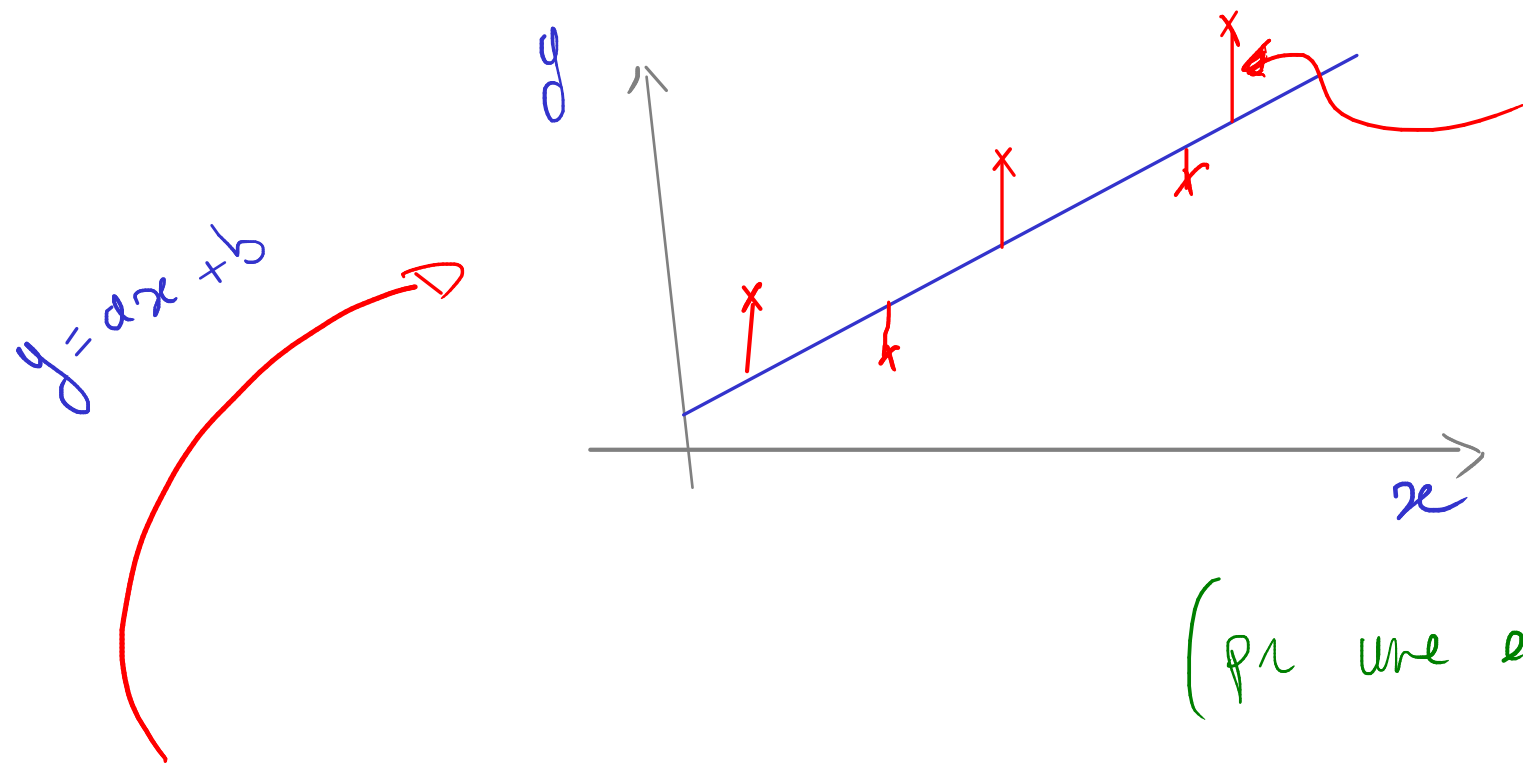
$$y = x^2$$

↳ forte corrélation mais non linéaire.

$$\Rightarrow r_{xy} \approx 0$$



changement
de variable ?



on a minimisé ces distances
distances verticales.

⇒ Pas d'erreur sur x

Pas de souci si y
dépend de x

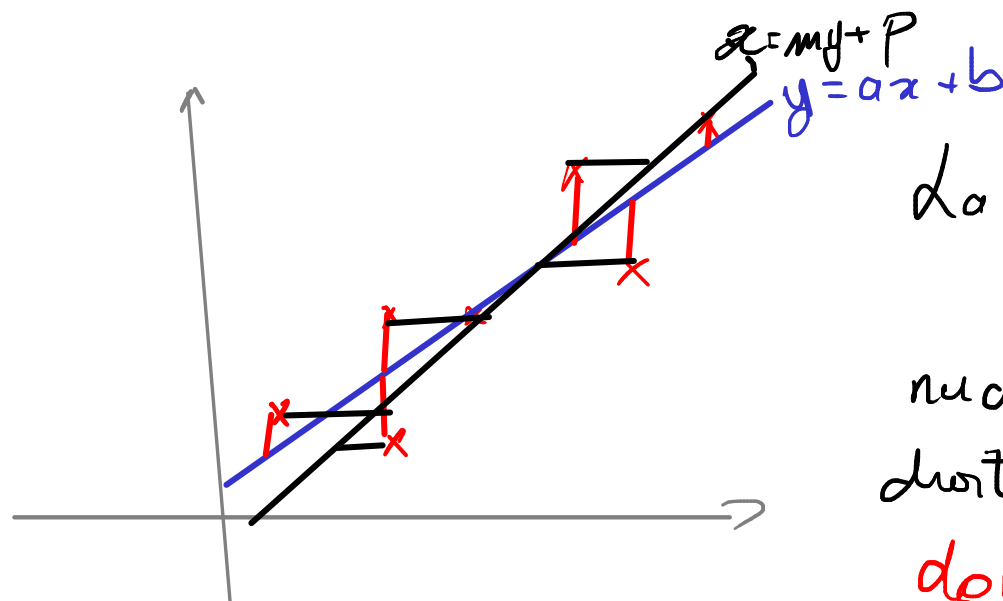
(par une évolution temporelle par exemple)

6.2 Régression $x = my + p$

x chiffre d'affaires
 y budget pub

y dépend de x
et x dépend de y .

il y a de l'erreur verticale (y)
et horizontale (x)



Laquelle choisir ?

Si R^2 est proche de 1, le
nuage est rectiligne et les
droites sont proches des points
done les droites sont très proches.

Faux problème.

6.3 Changement de variable

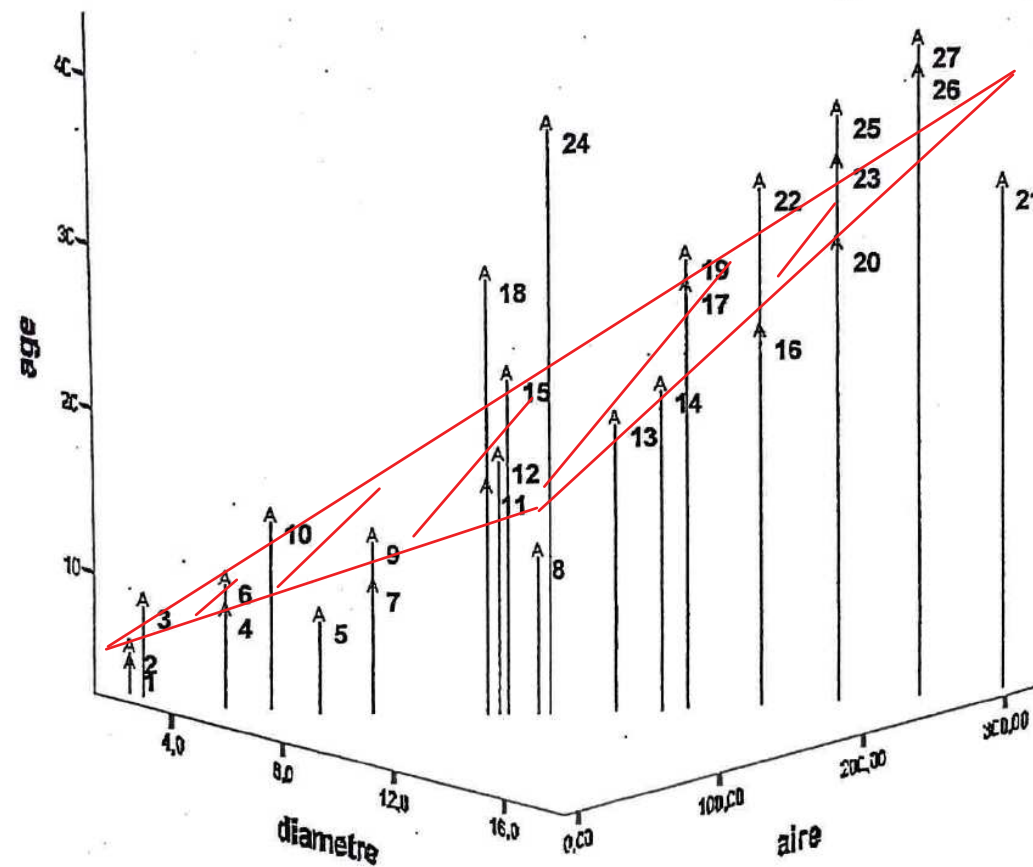
6.4 Régressions linéaire multiple

On peut souhaiter expliquer une variable réponse y par plusieurs variables explicatives x_1, \dots, x_d en utilisant une équation cartésienne :

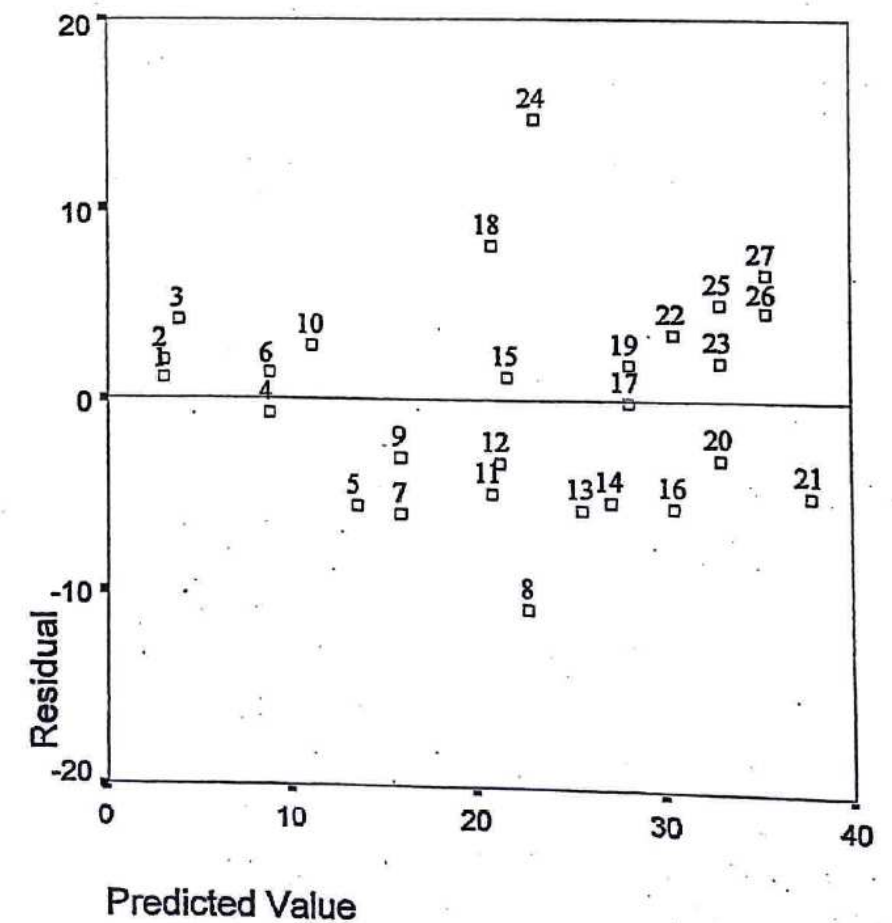
$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d,$$

Reprenons notre exemple d'âge des arbres, et intéressons nous au diamètre et à la surface de la section (les arbres ne sont pas des cylindres parfaits, le diamètre n'explique pas totalement la surface). On obtient les graphes suivants

Nuage de points



Résidus



Résumé de la régression

	Coefficients non standardisés			
	B	Erreur standard		
(constante)	2,720	4,317	R	R-deux
DIAMETRE	1,161	1,001		
AIRE	,055	,051		
			,894 ^a	,799