

5UAA1 : Statistique 2 variables

1 Rappel vocabulaire statistique à une variable

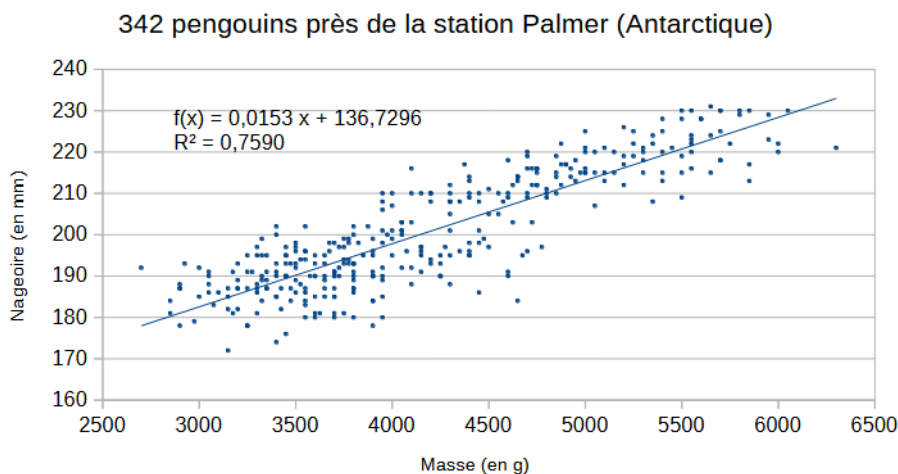
1.1 Population, individus, caractère

Une étude statistique porte sur un ensemble d'**individus** (parfois un échantillon de cet ensemble) appelé **population**.

Elle collecte un (statistique à une variable) ou plusieurs (statistique à plusieurs variables) caractères (propriétés, caractéristiques) que chaque individu possède.

Considérons la taille des nageoires et la masse de 342 pingouins près de la station Palmer, en Antarctique¹.

Traçons le **nuage de points** correspondant.



On note $n = 342$ le nombre d'individus, c'est à dire l'**effectif total** de la population observée.

1.2 Mesures de position, de dispersion

Nous avons défini en quatrième des grandeurs qui permettent de décrire les séries de données collectées.

Certaines donnent des informations sur la position : médiane, moyenne

D'autres donnent des informations sur la dispersion : écart type, écart interquartile

La moyenne des masses des 342 pingouins : $\bar{x} = \frac{x_1 + x_2 + \dots + x_{342}}{342} = 4201,8 \text{ g}$

La variance des masses des 342 pingouins :

$$\text{Var}(x) = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_{342} - \bar{x})^2}{342} = 641\,250 \text{ g}^2$$

L'écart type des masses des 342 pingouins : $\sigma(x) = \sqrt{\text{Var}(x)} = \sqrt{641\,250} = 800,8 \text{ g}$

Les masses se situent essentiellement dans l'intervalle $[\bar{x} - 2 * \sigma; \bar{x} + 2 * \sigma] = [2600 \text{ g}; 5803 \text{ g}]$

La moyenne des longueurs des nageoires des 342 pingouins : $\bar{y} = \frac{y_1 + y_2 + \dots + y_{342}}{342} = 200,9 \text{ mm}$

1. Horst AM, Hill AP, Gorman KB (2020). palmerpinguins : Palmer Archipelago (Antarctica) penguin data. R package version 0.1.0. <https://allisonhorst.github.io/palmerpinguins/>. doi : 10.5281/zenodo.3960218. R for datascience <https://r4ds.hadley.nz/data-visualize>

La variance des longueurs des nageoires des 342 pingouins :

$$\text{Var}(y) = \frac{(y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_{342} - \bar{y})^2}{342} = 197,15 \text{ mm}^2$$

L'écart type des longueurs des nageoires des 342 pingouins : $\sigma(y) = \sqrt{\text{Var}(y)} = \sqrt{197,2} = 14,04 \text{ mm}$

La plus grande partie des longueurs se situent dans l'intervalle $[\bar{y} - 2 * \sigma; \bar{y} + 2 * \sigma] = [173 \text{ mm}; 229 \text{ mm}]$

1.3 Définition de la variance et de l'écart type

La variance $\text{Var}(x)$ des x_i est la moyenne des carrés des écarts $x_i - \bar{x}$ à la moyenne \bar{x} .

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

L'écart type (de Pearson) des x_i , noté $\sigma(x)$ en quatrième, est la racine carrée de la variance des données observées : $\sigma(x) = \sqrt{\text{Var}(x)}$.

1.4 Une formule très utile de la variance

La variance est égale à **la moyenne des carrés moins le carré de la moyenne** : $\text{Var}(x) = \overline{x^2} - \bar{x}^2$

Démonstration :

$$\text{Var}(x) = \frac{1}{n} \sum_i (x_i - \bar{x})^2 = \frac{1}{n} \sum_i x_i^2 - 2\bar{x} \cdot \frac{1}{n} \sum_i x_i + \frac{1}{n} \sum_i \bar{x}^2 = \overline{x^2} - 2\bar{x}^2 + \frac{1}{n} \cdot n \cdot \bar{x}^2 = \overline{x^2} - \bar{x}^2$$

2 Statistique à deux variables

2.1 Objectif

Observer le nuage constitué de n points de coordonnées $(x_i; y_i)$. Si une relation linéaire semble exister (les points ont l'air de se répartir sur une droite), trouver l'équation $y = ax + b$ de la droite qui approche ou « **ajuste** » au mieux ce nuage **au sens des moindres carrés**.

Remarque : un document annexe, disponible sous moodle, donne des détails pour justifier les formules données ici. Il permet également de calculer une incertitude expérimentale sur la pente a de la droite trouvée sous l'hypothèse qu'il existe vraiment une relation linéaire entre x et y . Le calcul de cette incertitude est par exemple très utile en sciences quand on effectue plusieurs mesures $(x_i; y_i)$ d'un phénomène physique (a est alors une mesure d'une grandeur que l'on cherche à déterminer, et on veut connaître la précision de cette mesure).

Dans l'exemple des 342 pingouins, on cherche juste à déterminer une tendance.

2.2 Méthode des moindres carrés

Soit $y = ax + b$ l'équation de la droite cherchée.

Il faut trouver sa pente a et son ordonnée à l'origine b afin de **minimiser la somme des carrés des écarts verticaux** :

$$\sum_{i=1}^n (y_i - (ax_i + b))^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{en notant : } \boxed{\hat{y}_i = ax_i + b}$$

$e_i = y_i - \hat{y}_i$ est l'écart vertical entre le point $(x_i; y_i)$ du nuage et le point correspondant $(x_i; \hat{y}_i)$ sur la droite de régression (d'abscisse aussi égale à x_i).

L'écart vertical e_i est souvent appelé écart résiduel ou résidu.

e_i est l'écart entre la valeur observée y_i et sa valeur estimée $\hat{y}_i = ax_i + b$ connaissant x_i

La méthode des moindres carrés permet de minimiser $\sum_{i=1}^n e_i^2$

2.3 Calculs de a et b ; Covariance de x et y ; Variables centrées

Un document annexe sur <https://ophysis.net> donne les justifications pour ceux qui sont intéressés.

$$a = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{x^2 - \bar{x} \cdot \bar{x}} \quad \text{et} \quad \boxed{b = \bar{y} - a\bar{x}}$$

On reconnaît au dénominateur de a la variance $\text{Var}(x) = \overline{x^2} - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ des x_i .

On définit la covariance $\text{Cov}(x; y)$ des x_i et y_i par $\boxed{\text{Cov}(x; y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}$

Ainsi, par définition, la covariance $\text{Cov}(x; y)$ des x_i et y_i est la moyenne du produit des **variables centrées** $\boxed{X_i = x_i - \bar{x}}$ et $\boxed{Y_i = y_i - \bar{y}}$.

Démontrer en exercice que $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \overline{xy} - \bar{x} \cdot \bar{y}$

La covariance est égale à **la moyenne des produits moins le produit des moyennes**.

2.4 Autres formules

$$a = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{x^2 - \bar{x} \cdot \bar{x}} = \frac{\text{Cov}(x; y)}{\text{Var}(x)} = \frac{\sum_{i=1}^n X_i \cdot Y_i}{\sum_{i=1}^n X_i^2}$$

La dernière égalité s'obtient en multipliant par n le numérateur et le dénominateur.

L'équation de la droite est $\boxed{y - \bar{y} = \frac{\text{Cov}(x; y)}{\text{Var}(x)} \cdot (x - \bar{x})}$ ou encore $\boxed{y = a \cdot x + b}$ avec $\boxed{b = \bar{y} - a\bar{x}}$

En effet, le point moyen (\bar{x}, \bar{y}) est sur la droite de régression linéaire : $y = a \cdot x + b$.

2.5 Moyenne et variance observées des résidus; Coefficient de corrélation

Connaissant x_i , on peut estimer y_i en calculant $\hat{y}_i = ax_i + b = \bar{y} + \frac{\text{Cov}(x; y)}{\text{Var}(x)} \cdot (x_i - \bar{x})$.

La moyenne observée des résidus $e_i = y_i - (a \cdot x_i + b)$ est nulle : $\bar{e} = 0$.

La variance observée des e_i est $\text{Var}(e) = \overline{e^2} = \text{Var}(y) \cdot (1 - R(x; y)^2)$ avec : $\boxed{R(x; y) = \frac{\text{Cov}(x; y)}{\sigma(x) \cdot \sigma(y)}}$

$R(x; y)$ est appelé **coefficient de corrélation** entre x et y .

Comme $\text{Var}(e)$ et $\text{Var}(y)$ sont positifs, on en déduit que $1 - R(x; y)^2 \geq 0$, soit $\boxed{-1 \leq R(x; y) \leq 1}$

De plus, si $R(x; y)^2 = 1$ (donc $R(x; y) = \pm 1$), alors $e_i = 0$ (puisque $\overline{e^2} = 0$) et donc $y_i = \hat{y}_i$ pour les n points qui sont donc parfaitement alignés. La connaissance de x_i permet de déterminer la valeur de y_i grâce à la droite de régression.

2.6 Interprétation du coefficient de corrélation $R(x; y)$

On remarque que :

$$a = \frac{\text{Cov}(x; y)}{\text{Var}(x)} = \frac{\text{Cov}(x; y)}{\sigma(x) \sigma(x)} \cdot \frac{\sigma(y)}{\sigma(y)} = \frac{\text{Cov}(x; y)}{\sigma(x) \sigma(y)} \cdot \frac{\sigma(y)}{\sigma(x)} = R(x; y) \cdot \frac{\sigma(y)}{\sigma(x)}$$

Ainsi, si $R(x; y) = 1$, la pente de la droite de régression est le rapport $\frac{s_y}{s_x}$, ce qui est cohérent avec la définition de la pente (rapport des variations en y sur les variations en x).

Si $R(x; y) = -1$, la pente de la droite de régression est le rapport $-\frac{\sigma(y)}{\sigma(x)} < 0$.

De façon générale, si $R(x; y) < 0$, alors $a < 0$: quand x_i augmente, y_i diminue.

Si $R(x; y) > 0$, alors $a > 0$: quand x_i augmente, y_i aussi.

3 Retour en Antarctique

On trouve que $\text{Cov}(x; y) = 9795,7 \text{ mm} \cdot \text{g}$ et donc $a = \frac{\text{Cov}(x; y)}{\text{Var}(x)} = \frac{9795,7}{641250} = 0,0153 \text{ mm/g}$

Et donc, $b = \bar{y} - a\bar{x} = 200,9 - 0,0153 \cdot 4201,8 = 136,6 \text{ mm}$

et $R(x; y) = \frac{\text{Cov}(x; y)}{\sigma(x) \cdot \sigma(y)} = \frac{9795,7}{800,8 \cdot 14,04} = 0,871$ ce qui donne $R(x; y)^2 = 0,759$

On retrouve ainsi les valeurs calculées par le tableur **Calc** de **libre office** lors du calcul de la courbe de tendance linéaire (voir le graphique en première page).

En moyenne, la longueur de la nageoire des pingouins augmente d'environ 1,5 mm par kilogramme supplémentaire.