

# TD D'INFORMATIQUE : RÉGRESSION LINÉAIRE

On souhaite calculer la pente et l'ordonnée à l'origine d'une droite de régression linéaire, pour des données de mesures expérimentales et valider l'hypothèse d'une relation linéaire entre ces données :

On envisage ainsi ces données comme deux variables aléatoires  $X_i$  et  $Y_i$  correspondant aux  $N$  mesures indexées par  $i$  de 0 à  $N-1$ . On veut également pondérer les différentes mesures en fonction de l'incertitude de chacune des mesures sur la variable  $Y$  mesurée. La question que l'on pose est donc de déterminer le degré de corrélation entre ces deux variables aléatoires.

## PRÉSENTATION DU PROBLÈME [bien lire cette partie]

Les valeurs expérimentales nous sont données sous la forme de deux listes  $X$  et  $Y$  de  $N$  éléments, ainsi que la listes  $Sy$  [ $Sy$  pour sigma\_y] des  $N$  incertitudes respectives aux  $N$  valeurs de  $Y$  mesurées.

**On veut évaluer la relation linéaire entre les deux variables aléatoires :  $Y = a.X + b$**

Soit  $y_i^{th} = ax_i + b \quad \forall i$  : la valeur théorique attendue pour  $Y_i$  en  $X_i$ , et soit  $y_i^{obs}$  la valeur réellement mesurée. On peut quantifier l'écart global de nos mesures à la théorie en calculant l'écart-quadratique moyen entre la théorie et l'expérience (divisé par l'incertitude) :

$$S = \sum_0^{N-1} \left( \frac{y_i^{th} - y_i^{obs}}{\sigma_i} \right)^2 \quad \text{soit} \quad S(a,b) = \sum_0^{N-1} \frac{(ax_i + b - y_i^{obs})^2}{\sigma_i^2}$$

$S(a, b)$  est l'écart quadratique cumulé entre la théorie et l'expérience, et pondéré par l'inverse de l'incertitude de mesure sur  $Y_i$  : plus on est confiant sur un point de mesure, plus ce point pèsera sur  $S$  et inversement. **C'est la méthode dite du  $\chi^2$  de Pearson !**

Pour trouver la meilleure droite passant à travers les données, on cherche donc à minimiser cet écart : il faut trouver  $a$  et  $b$  qui rendent  $S(a, b)$  minimale. On écrit :

$$\frac{dS(a,b)}{da} = 0 \quad \text{qui amène :} \quad \sum_0^{N-1} \frac{x_i (ax_i + b - y_i^{obs})}{\sigma_i^2} = 0 \quad (\text{eq}^\circ \text{ a})$$

$$\frac{dS(a,b)}{db} = 0 \quad \text{qui amène :} \quad \sum_0^{N-1} \frac{(ax_i + b - y_i^{obs})}{\sigma_i^2} = 0 \quad (\text{eq}^\circ \text{ b})$$

On peut alléger les notations en introduisant l'opérateur de moyenne  $\langle . \rangle$  qui donne la moyenne sur l'ensemble des valeurs prises par l'indice  $i$  d'une variable aléatoire :

$$\langle f \rangle = \sum_0^{N-1} \frac{f_i}{\sigma_i^2} / \sum_0^{N-1} \frac{1}{\sigma_i^2} \quad (\text{eq}^\circ \text{ c}) \quad \text{Cette moyenne étant pondérée par l'inverse de la variance.}$$

En pratique on écrira une fonction en Python qui lit une liste de mesures et une liste d'incertitudes et qui renvoie la valeur moyenne [cf fonction `bracket(x)` au verso].

On peut ainsi reformuler nos équations (eq° a) et (eq° b) en développant et en divisant par la somme des inverses des variances. Soit :

$$\langle x_i(ax_i + b - y_i^{obs}) \rangle = 0 \Rightarrow a \langle x_i^2 \rangle + b \langle x_i \rangle - \langle x_i y_i^{obs} \rangle = 0 \quad (\text{eq}^\circ \text{a}')$$

$$\langle ax_i + b - y_i^{obs} \rangle = 0 \Rightarrow a \langle x_i \rangle + b - \langle y_i^{obs} \rangle = 0 \quad (\text{eq}^\circ \text{b}')$$

Les inconnues étant a et b c'est un simple système de deux équations à deux inconnues. Soit :

$$a = \frac{\langle x_i y_i^{obs} \rangle - \langle x_i \rangle \langle y_i^{obs} \rangle}{\langle x_i^2 \rangle - \langle x_i \rangle^2} \quad \text{et} \quad b = \langle y_i^{obs} \rangle - a \langle x_i \rangle$$

On note que pour calculer les termes  $\langle x_i y_i^{obs} \rangle$  ou encore  $\langle x_i^2 \rangle = \langle x_i \cdot x_i \rangle$ , il faut construire au préalable la liste des N produits  $X_i \cdot Y_i$  ou  $X_i \cdot X_i$ . Il faudra donc créer une fonction qui prend en argument deux listes et renvoie la liste des N produits deux à deux [cf fonction ProdXY(x) au verso].

### **Calcul du coefficient de corrélation :**

(En général, les calculatrices indiquent R au carré.)

$$R = \frac{COV(x_i, y_i^{obs})}{\sigma_x \sigma_y} \quad (\text{eq}^\circ \text{d})$$

On définit pour cela :

- la covariance des deux variables aléatoires :  $COV(x_i, y_i^{obs}) = \frac{\langle (x_i - \mu_x)(y_i^{obs} - \mu_y) \rangle}{N}$

=> Celle-ci nécessite une fonction moyenne qui calcule la moyenne d'une liste [ à utiliser sur X et Y ].

- Les écarts types des listes X et Y :  $\sigma_x$  et  $\sigma_y$

### **Comparaison au résultat obtenu avec votre machine à calculer :**

Nous testerons nos résultats sur les listes de données suivantes :

##données

x=[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

y=[0.8, 2.4, 2.0, 2.2, 2.7, 3.5, 3.0, 4.0, 5.0, 4.8]

#incertitude sur Yi

Sy=[0.5, 1.0, 0.2, 0.5, 0.5, 0.5, 0.5, 1.0, 0.2, 0.5]

Rq : A priori votre calculatrice ne prendra pas en compte les incertitudes de mesure. Toutes les mesures ont alors le même poids qu'elles soient fiables ou non.

Pour comparer à votre machine, il suffit de poser  $Sy[i] = 1$  pour toutes les mesures i. Il est bien sûr plus fiable de tenir compte des incertitudes de mesures.

\*\*\* Y = aX + b \*\*\*

a : 0.48      b : 0.50      R2 : 0.89366      R : 0.94534

# IMPLÉMENTATION EN PYTHON

## Listes de fonctions à créer :

**Dans chaque cas on veillera à minimiser les calculs redondants donc inutiles.**

Pb : les fonctions ci-dessous ne prennent pas en entrée la/les moyenne(s) des listes.

```
def Mean(x):          #moyenne d'une variable aléatoire
                    #retourne la moyenne d'une liste x    ex: mean([1, 2, 3]) retourne le réel 2.0
```

```
def Sigma(x):        #écart-type d'une variable aléatoire
                    #retourne l'écart type d'une liste x (soit un réel)
```

```
def Cov(x, y):      #Covariance des variables aléatoires X et Y
                    #retourne la covariance de deux listes x et y (soit un réel)
```

```
def ProdXY(x, y):   #Construction de la liste produit  $X_i \cdot Y_i$ 
                    #exemple ProdXY([1, 2, 3], [4, 5, 6]) retourne [4, 10, 18] donc une liste !
```

```
def Bracket(x, Sy=[1.]*len(x)):  # cf (eq° c)
                                    # Sy=[1.]*len(x) permet de mettre par défaut toutes les
                                    # incertitudes à 1.0
```

#Attention par définition de notre opérateur de moyenne, Bracket prend en entrée :

#une liste x quelconque de taille N et les N incertitudes Sy d'indices i associées à Y.

# Rq : Sy=[1.]\*len(x) ne marche pas si x n'est pas déjà en mémoire ...

```
def Regression(x, y, Sy=[1.]*len(y)):  ## Calcul de la régression
```

```
    # Calcul de a
    # Calcul de b
    # Calcul de R
    # Affichage
```