

# T. D. n° 3

## Régression linéaire simple

### Exercice 1. Les athlètes.

La taille d'un athlète peut jouer un rôle important dans ses résultats en saut en hauteur. Les données utilisées ici présentent donc la taille et la performance de 20 champions du monde.

Observation	Nom	Taille	Performance
1	Jacobs (EU)	1,73	2,32
2	Noji (EU)	1,73	2,31
3	Conway (EU)	1,83	2,40
4	Matei (Roumanie)	1,84	2,40
5	Austin (EU)	1,84	2,40
6	Ottey (Jamaïque)	1,78	2,33
7	Smith (GB)	1,84	2,37
8	Carter (EU)	1,85	2,37
9	McCants (EU)	1,85	2,37
10	Sereda (URSS)	1,86	2,37
11	Grant (GB)	1,85	2,36
12	Paklin (URSS)	1,91	2,41
13	Annys (Belgique)	1,87	2,36
14	Sotomayor (Cuba)	1,96	2,45
15	Sassimovitch (URSS)	1,88	2,36
16	Zhu Jianhua (Chine)	1,94	2,39
17	Brumel (URSS)	1,85	2,28
18	Sjoeberg (Suède)	2,00	2,42
19	Yatchenko (URSS)	1,94	2,35
20	Povarnitsine (URSS)	2,01	2,40

- À partir de l'échantillon proposé, utiliser la méthode des moindres carrés pour estimer les paramètres de la régression linéaire :

$$(\text{Performance}) = \beta_0 + \beta_1 \times (\text{Taille}) + \varepsilon$$

- Compléter le tableau d'analyse de la variance (dit aussi tableau d'ANOVA) :

Source de variation	Degrés de liberté	Somme des carrés	Moyenne des carrés	$F_{obs}$
Régression				
Résiduelle				
Totale				

- Quel pourcentage de la variation totale des performances est expliqué par la variable taille ? Que pensez-vous de ce résultat ? Que faudrait-il faire en tant que chargé de cette étude ?

**Exercice 2. Un exercice pour pratiquer.**

Nous disposons des données suivantes au sujet de deux variables d'intérêt  $X$  et  $Y$  :

$x_i$	7	9	9	10	13	17	19	20	21	25
$y_i$	5	4	6	4	1	2	0	1	1	0

Nous nous référons au modèle linéaire :

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

1. Estimer les paramètres  $\beta_0$  et  $\beta_1$  par la méthode des moindres carrés.
2. Pour chacun de ces deux paramètres, trouver un intervalle de confiance avec un niveau de confiance de 99%.
3. Soit  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  ( $i = 1, \dots, n$ ) où  $\hat{\beta}_0$  et  $\hat{\beta}_1$  sont les estimateurs de  $\beta_0$  et  $\beta_1$  obtenus en 1). Démontrer que nous avons  $\sum \hat{Y}_i = \sum Y_i$ , par deux méthodes (mathématique, et avec R).
4. Donner les intervalles de confiance pour les  $\mu_Y(X)$ .
5. Représenter graphiquement les points  $(x_i, y_i)$ , la droite de régression et l'ensemble des intervalles de confiance pour les  $\mu_Y(X)$ .

**Exercice 3. Exemples.**

Source de variation	Degrés de liberté	Somme des carrés	Moyenne des carrés	$F_{obs}$	$F_t$	$R^2$
Régression	1	501,76	501,76	7,575	4,75	0,387
résiduelle	12	794,90	66,24			
Totale	13	1296,66				

Source de variation	Degrés de liberté	Somme des carrés	Moyenne des carrés	$F_{obs}$	$F_t$	$R^2$
Régression	1	34,186	34,186	43,44	7,71	0,916
résiduelle	4	3,148	0,787			
Totale	5	37,333				

Source de variation	Degrés de liberté	Somme des carrés	Moyenne des carrés	$F_{obs}$	$F_t$	$R^2$
Régression	1	179,76	179,76	27,87	4,08	0,367
résiduelle	48	309,56	6,45			
Totale	49	489,32				

1. En comparant dans les exemples ci-dessus, les liens entre les valeurs  $F_{obs}$ ,  $F_t$  ( $F_t$  est la valeur lue dans la table de Fisher) et  $R^2$ , quelles sont, selon vous, les meilleures régressions ?
2. Avant de calculer le coefficient de détermination  $R^2$ , en n'utilisant que les valeurs  $F_{obs}$  et  $F_t$ , quelle règle pourrions-nous énoncer pour réperer rapidement une bonne analyse de régression ?

**Exercice 4. Calories.**

Soient les données présentées dans le tableau ci-dessous. Il s'agit du nombre de calories consommées par jour et du pourcentage de population agricole dans 11 pays.

Observation $i$	Pays	% Population agricole	Calories par jour et par personne
1	Suisse	4,0	3 432
2	France	5,7	3 273
3	Suède	4,9	3 049
4	USA	3,0	3 642
5	Ex-URSS	14,8	3 394
6	Chine	69,6	2 628
7	Inde	63,8	2 204
8	Brésil	26,2	2 643
9	Pérou	38,3	2 192
10	Algérie	24,7	2 687
11	Ex-Zaire	65,7	2 159

1. Représenter graphiquement  $Y$  en fonction de  $X$ .
2. Estimer les paramètres  $\beta_0$  et  $\beta_1$  du modèle :
 
$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$
3. Construire le tableau d'analyse de la variance correspondant à cette régression.
4. Construire un intervalle de confiance à 95% autour de la droite de régression.
5. Représenter sur le graphique de la question (1) la droite de régression et l'intervalle de confiance calculé à la question (4).

**Exercice 5. Composant électronique.**

Un certain composant électronique est fabriqué une fois par mois par l'entreprise Micro-Systèmes. La quantité fabriquée varie avec la demande du marché. Dans le but de planifier la production et d'établir certaines normes sur le nombre d'hommes-minutes exigés pour la production de différents lots de ce composant électronique, le responsable de la production a relevé l'information suivante pour 15 cédules de production. Le nombre d'hommes-minutes est identifié par  $Y$  et la quantité fabriquée par  $X$ .

$y_i$	150	192	264	371	300	358	192	134	242	238	226	302	340	182	169
$x_i$	35	42	64	88	70	85	40	30	55	60	51	72	80	44	39

1. Quelle serait la première étape à franchir avant d'aborder tout calcul préliminaire ?
2. Le responsable de la production envisage d'utiliser le modèle linéaire simple comme modèle prévisionnel. Spécifier ce modèle et bien identifier chacune des composantes du modèle dans le contexte de ce problème.
3. Déterminer l'équation de régression.

4. D'après l'équation de régression, si le nombre d'unités à fabriquer augment de 10, quelle sera l'augmentation correspondante du nombre moyen d'hommes-minutes requis ?
5. En l'absence de l'information que nous donne la quantité à fabriquer, quelle serait une bonne estimation du nombre d'hommes-minutes requis ?
6. Quelle correction peut-il apporter à son estimation du nombre moyen d'hommes-minutes requis en introduisant la connaissance de  $X$  dans son analyse ?
7. Donner la valeur de  $s(\beta_1)$  et tester les hypothèses suivantes avec la loi de Student.

$$\mathcal{H}_0 : \beta_1 = 0, \quad \mathcal{H}_1 : \beta_1 \neq 0.$$

8. Donner la variation qui est expliquée par la droite de régression et la variation qui est inexpliquée par la droite.
9. Déterminer le pourcentage de variation qui est expliqué par la droite de régression.
10. Donner une estimation du nombre moyen d'hommes-minutes requis pour :  
 $x_h = 42$ ;  $x_h = 57$ ;  $x_h = 72$ .
11. Pour quelle quantité  $X_n$ , l'estimation du nombre moyen d'hommes-minutes requis serait-elle la plus précise ?
12. Entre quelles valeurs peut se situer le vrai nombre moyen d'hommes-minutes requis pour les lots dont la quantité a été déterminée à la question 11. ? Utiliser un niveau de confiance de 95%.
13. Quelle est la marge d'erreur dans l'estimation effectuée en 12. ?