

Feuille de Travaux Dirigés n° 1

Régression linéaire simple avec R

Les exercices **1**, **4** et **11** proviennent du livre de G. Baillargeon, *Probabilités, Statistique et Techniques de régression*, aux éditions SMG, 1995.

Les exercices **2** et **3** proviennent du livre de Frontier, Davoult, Gentilhomme, Lagadeuc, *Statistique pour les sciences de la vie et de l'environnement*, aux éditions Dunod, 2001.

Les exercices **7**, **8**, **9** et **10** proviennent du livre de Y. Dodge, *Analyse de régression appliquée*, aux éditions Dunod, 1999.

Exercice I.1. Cet exercice doit se traiter en grande partie avec R.

Nous donnons les couples d'observations suivants :

x_i	18	7	14	31	21	5	11	16	26	29
y_i	55	17	36	85	62	18	33	41	63	87

1. La première étape est d'obtenir les données. Pour cela, vous pouvez les télécharger sur mon site, puis les enregistrer sur le bureau du poste. Par exemple, depuis le bureau de mon ordinateur portable, les lignes de commande à taper sous R sont les suivantes :

```
> setwd("C:\\Documents and Settings\\Bertrand\\Bureau")
> Exo1<-read.csv(file.choose())
```

2. Tracer le diagramme de dispersion des couples $(x_i; y_i)$. À la vue de ce diagramme, pouvons-nous soupçonner une liaison linéaire entre ces deux variables ?
3. Déterminer pour ces observations la droite des moindres carrés, c'est-à-dire donner les coefficients de la droite des MC.
4. Donner les ordonnées des y_i calculés par la droite des moindres carrés correspondant aux différentes valeurs des x_i .
5. Tracer ensuite la droite sur le même graphique.
6. Quelle est une estimation plausible de Y à $x_i = 21$?
7. Quel est l'écart entre la valeur observée de Y à $x_i = 21$ et la valeur estimée avec la droite des moindres carrés ? Comment appelons-nous cet écart ?
8. Est-ce que la droite des moindres carrés obtenue à la question 3 passe par le point $(\bar{x}; \bar{y})$? Pouvons-nous généraliser cette conclusion à n'importe laquelle droite de régression ?

Remarque : Voici quelques lignes de commande qui pourront vous aider à répondre aux questions. À vous de savoir à quoi elles répondent.

```
> setwd("C:\\Documents and Settings\\Bertrand\\Bureau")
> Exo1<-read.csv(file.choose())
> Exo1
```

```

      x_i y_i
1    18 55
2     7 17
3    14 36
4    31 85
5    21 62
6     5 18
7    11 33
8    16 41
9    26 63
10   29 87
> str(Exo1)
'data.frame':  10 obs. of  2 variables:
 $ x_i: int  18 7 14 31 21 5 11 16 26 29
 $ y_i: int  55 17 36 85 62 18 33 41 63 87
> plot(Exo1)
> Droite<-lm(y_i~x_i,data=Exo1)
> coef(Droite)
(Intercept)          x_i
  1.021341    2.734756
> fitted(Droite)
      1      2      3      4      5      6      7
50.24695 20.16463 39.30793 85.79878 58.45122 14.69512 31.10366
      8      9     10
 44.77744 72.12500 80.32927
> abline(coef(Droite),col="red")
> residuals(Droite)
      1      2      3      4      5      6
4.7530488 -3.1646341 -3.3079268 -0.7987805  3.5487805  3.3048780
      7      8      9     10
 1.8963415 -3.7774390 -9.1250000  6.6707317
> residuals(Droite)[5]
      5
3.548780

```

Exercice I.2.

On étudie l'influence d'un antibiotique sur une culture bactérienne. On répartit dans 10 tubes des volumes égaux de culture additionnés d'une quantité X d'antibiotique, et on mesure, après incubation, la densité optique D . Les résultats sont les suivants.

X	0,2	0,2	0,4	0,4	0,6	0,6	0,8	0,8	1,0	1,0
D	19	21	35	38	64	66	115	130	200	210

- a) Un ajustement linéaire semble-t-il justifié? Pour répondre à cette question, utiliser \mathbf{R} . Que devez-vous calculer comme coefficient avec \mathbf{R} ?

- b) En transformant une des deux variables avec une fonction adaptée, déterminer une équation de régression en précisant quelles sont la variable explicative et la variable expliquée ?
- c) Donner à l'aide de \mathbf{R} , une prévision de D pour une quantité d'antibiotique $X = 0,5$. Calculer, toujours à l'aide de \mathbf{R} , l'intervalle de sécurité à 95% de cette prévision.

Exercice I.3.

On mesure le poids frais et le poids sec de 20 prélèvements de plancton. Les résultats sont les suivants (exprimés en g par 10 m³ d'eau de mer)

poids frais	20,4	28,4	48,7	28,8	32,9	85,2	32,2	27,8	27,0	36,7
poids sec	3,6	3,4	5,6	4,1	3,3	9,3	3,7	3,2	2,9	4,5
poids frais	20,4	24,3	24,3	18,0	31,7	25,7	41,2	53,0	61,0	61,2
sec	2,6	2,8	3,1	2,6	4,4	2,8	4,6	6,0	7,2	6,3

- a) Calculer à l'aide de \mathbf{R} le coefficient de corrélation linéaire entre le poids frais et le poids sec. Est-il significatif et à quel seuil ? Repérer un *outsider* parmi les couples de valeurs ; l'éliminer et reprendre la question. Pour cette dernière partie de question, vous devez appliquer la procédure étudiée en cours.
- b) La teneur en eau de chaque prélèvement planctonique est estimée par la différence entre poids frais et poids sec. Estimer sa variance à l'aide de \mathbf{R} .
- c) Y a-t-il un sens à calculer le coefficient de corrélation entre le poids frais et la teneur en eau ainsi estimée, et pourquoi ?
- d) Donner, à l'aide de \mathbf{R} , la droite permettant de connaître approximativement le poids sec après une mesure de poids frais. Quelle est, dans ces conditions, la proportion de variance du poids sec expliquée par la régression ?
- e) Soit un poids frais de 40 grammes. Calculer, à l'aide de \mathbf{R} , la valeur la plus probable du poids sec, et son intervalle de sécurité à 95%.

Exercice I.4.

Cet exercice doit se traiter en grande partie avec \mathbf{R} .

La société de Transport Bertrand veut établir une politique d'entretien des camions de sa flotte. Tous sont de même modèle et utilisés à des transports semblables. La direction de la société est d'avis qu'une liaison statistique entre le coût direct de déplacements (*cents* par *km*) et l'espace de temps écoulé depuis la dernière inspection de ce camion serait utile. Nous avons donc recueilli un certain nombre de données sur ces deux variables. Nous souhaitons utiliser la régression linéaire comme modélisation statistique.

Coût direct	10	18	24	22	27	13	10	24	25	8	16
Nombre de mois	3	7	10	9	11	6	5	8	7	4	6
Coût direct	20	28	22	19	18	26	14	20	26	30	12
Nombre de mois	9	12	8	10	9	11	6	8	10	12	5

1. Quelle variable devrions-nous identifier variable dépendante (Y) et laquelle devrions-nous identifier variable explicative (X) ?

2. Tracer le diagramme de dispersion de ces observations. Est-ce que le nuage de points suggère une forme de liaison particulière ?
3. Calculer l'équation de la droite des moindres carrés.
4. Avec l'équation de la droite des moindres carrés, quelle est l'estimation la plus plausible du coût direct de déplacement pour des camions dont la dernière inspection remonte à 6 mois ?
5. D'après les résultats de cette étude, un délai supplémentaire d'un mois pour l'inspection d'un camion occasionnera-t-il une augmentation ou une diminution du coût direct ? Quelle sera vraisemblablement la valeur de cette variation de coût ?
6. Déterminer la variation totale dans le coût direct de déplacement.
7. L'équation de la droite des moindres carrés pour les données de la société est : $\hat{y}_i = 1,54941 + 2,26087 x_i$. Calculer la variation qui est expliquée par la droite des moindres carrés.
8. Quelle est la variation résiduelle ?
9. Calculer le coefficient R^2 et interpréter le résultat.

Exercice I.5. Cet exercice doit se traiter en grande partie avec R.

Une étudiante en sociologie veut analyser, dans le cadre d'un projet de fin de session, s'il existe une relation linéaire entre la densité de population dans les régions métropolitaines et le taux de criminalité correspondant dans ces régions.

Le taux de criminalité (Y) est indiqué en nombre de crimes par 10 000 habitants et la densité de population (X) est mesurée en milliers d'habitants par km^2 .

Région	1	2	3	4	5	6	7	8	9	10	11	12
x_i	7,7	5,8	11,5	2,1	3,7	3,6	7,5	4,2	3,8	10,3	8,6	7,2
y_i	12	9	15	4	4	2	10	3	5	11	10	11

1. Tracer le diagramme de dispersion de ces observations.
2. Calculer les coefficients de la droite des moindres carrés.
3. À quelle augmentation du taux de criminalité pouvons-nous nous attendre pour une variation unitaire (ici 1 000 habitants par km^2) de la densité de population ?
4. Estimer le taux de criminalité le plus plausible pour une densité de population de 7 500 habitants par km^2 .
5. À l'aide des calculs préliminaires, calculer la variation totale du taux de criminalité.
6. Calculer la variation qui est expliquée par la droite des moindres carrés.
7. Quelle proportion de la variation totale est expliquée par la droite des moindres carrés ?

Exercice I.6. Cet exercice doit se traiter en grande partie avec R.

Un étudiant en techniques forestières veut utiliser la régression linéaire pour estimer le volume en bois utilisable d'un arbre debout en fonction de l'aire du tronc mesuré à 25 cm du sol. Il a choisi au hasard 10 arbres et a mesuré, à la base, l'aire correspondante (en cm^2). Il a par la suite enregistré, une fois l'arbre coupé, le volume correspondant en m^3 .

Vol.	0,152	0,284	0,187	0,350	0,416	0,230	0,242	0,276	0,383	0,140
Aire	297	595	372	687	790	520	473	585	762	232

- Déterminer les coefficients de la droite des moindres carrés.
- Son professeur lui mentionne qu'il peut, à l'oeil, évaluer avec une assez bonne précision le volume d'un arbre. L'étudiant un peu perplexe lui lance un défi : « Je gage 1 euro que je fais mieux que vous avec le modèle des moindres carrés. »
« D'accord. »
Ayant justement un arbre tout près, le professeur lui dit, après une expertise de quelques minutes que cet arbre a un volume de $0,22 m^3$. Sans plus tarder, l'étudiant mesure l'aire de la base de l'arbre et obtient $465 cm^2$. Calculer avec la droite des moindres carrés, l'estimation la plus plausible du volume de l'arbre.
- L'étudiant s'acharne par la suite à couper l'arbre et le volume correspondant est $0,24 m^3$. Celui qui a le plus faible écart de prévision empoche le pari. Lequel s'est enrichi de 1 euro ?
- Est ce que le volume moyen des arbres échantillonnés aurait donné une estimation aussi bonne que la droite des moindres carrés pour cet arbre ?

Exercice I.7. 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

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

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

2. 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				

3. 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 I.8. 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 β_0 et β_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 β_0 et β_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 I.9. Trois 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	1 296,66				

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 repérer rapidement une bonne analyse de régression ?

Exercice I.10. 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 β_0 et β_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 I.11. 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 augmente 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 $s(\hat{\beta}_1)$ et tester les deux hypothèses suivantes avec la loi de Student.
$$\mathcal{H}_0 : \beta_1 = 0 \quad \text{contre} \quad \mathcal{H}_1 : \beta_1 \neq 0.$$
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. ?