

TP Statistiques et Séries Chronologiques
Université de Lorraine

**Régression linéaire et théorie des tests
avec R**

Clément Dell'Aiera

La régression linéaire sous R

Ce numéro rappelle les notions nécessaires à l'interprétation d'une sortie R de la fonction lm . Dans toute la suite, *regression* désigne un objet de type lm que l'on a appelé grâce à $lm(Y \sim X_1 + \dots + X_n, data = \dots)$. Si l'on note X_{-i} , le moins signifie que l'on calcule la quantité X sans tenir compte de l'observation i .

Dans le modèle de régression

$$Y = X\beta + \epsilon,$$

la commande `summary(regression)`, où *regression* est un objet de la classe lm , renvoie plusieurs tableaux.

On note $\hat{y}_j = \sum h_{ij}x_j$, i.e.

$$h_{ij} = \frac{1}{n} + \sum \frac{(x_i - \hat{x})(x_j - \hat{x})}{\sum (x_j - \hat{x})^2}.$$

Le premier tableau, *Residuals*, est destiné à donner une idée de la répartition des résidus en affichant les quantiles. Je vous recommande d'afficher tout de même les résidus, et d'observer leur distribution. Mieux, vous pouvez utiliser les résidus studentisés. A priori, bien que tous soient centrés, les résidus n'ont pas même variance (même sous hypothèse d'homoscédasticité!) : $Var[\epsilon_j] = \sigma^2(1 - h_{jj})$. Pour les rendre comparables, on pourrait les réduire, mais si l'on remplace la variance par la variance estimée, celle-ci dépend de l'information contenue dans x_j , ce qui empêche une quantification de l'effet que x_j a sur les coefficients de la régression. Pour palier à ce problème, on introduit

$$\hat{\sigma}_{-j}^2 = \frac{1}{n-3}[(n-2)\hat{\sigma}^2 - \frac{\epsilon_j}{1-h_{jj}}]$$

qui n'est rien d'autre que la variance estimée sur le modèle où l'on a supprimé l'observation x_j . Les résidus studentisés sont définis par $T_j = \frac{\epsilon_j}{\hat{\sigma}_{-j}(1-h_{jj})} \sim T(n-3)$ et suivent une loi de Student à $n-3$ degré de liberté sous des hypothèses raisonnables. Pour détecter une anomalie dans les données, on peut vérifier que les résidus studentisés se répartissent de manière uniforme sur l'intervalle $[-2; 2]$ (sous hypothèse d'homoscédasticité). Repérer des formes suspectes est un moyen facile pour repérer les valeurs aberrantes. On peut par exemple taper : `qqnorm(studres(regression)); qqline(studres(regression))`.

Une autre méthode pour détecter les valeurs aberrantes : utiliser la distance de Cook. Elle est définie par

$$D_j = \frac{\sum_j (\hat{y}_{-i,j} - \hat{y}_j)^2}{2\hat{\sigma}^2},$$

et mesure l'influence d'une observation sur l'ensemble des prévisions (qu'on veut petite!). Encore une règle du pouce : si $D_j > 1$, on enlève l'observation j . Vous pouvez le faire automatiquement en tapant `plot(regression, which = 4)`.

Le deuxième est nommé *Coefficients* :

	Estimate	Std. Error	t-value	$Pr(> t)$
β_j	$\hat{\beta}_j$	$\hat{\sigma}_j$	$\hat{t}_j = \frac{\hat{\beta}_j}{\hat{\sigma}_j}$	p -value

Les trois premières colonnes s'expliquent elles-mêmes, mais à quoi servent les 2 dernières ? A effectuer un test de significativité. Plus précisément, on sait que \hat{t}_j suit une loi de Student à $N - k$ degrés de liberté sous l'hypothèse $H_0 | \beta_j = 0$ contre $H_1 | \beta_j \neq 0$. La quatrième colonne donne donc la valeur de cette statistique, et la dernière sa p -value, définie comme la valeur seuil de confiance α qui fait basculer le test. (Rappelez vous que, mécaniquement, si α diminue assez, on finit par accepter H_0 .) Une règle appliquée par les statisticiens est la suivante :

$p < 0.01$	suspicion très forte contre H_0
$0.01 - 0.05$	suspicion forte contre H_0
$0.05 - 0.1$	suspicion faible contre H_0
> 0.1	peu ou pas de suspicion contre H_0

Donc, si $Pr(> |t|)$ est petit, on rejette l'hypothèse $\beta_j = 0$, ce qui signifie que le coefficient est significatif. R ajoute même des petites étoiles à côté des coefficients les plus significatifs.

Reste encore à observer plusieurs indicateurs. L'erreur *Residual standard error* est calculée comme un estimateur de σ sous l'hypothèse de matrice variance-covariance égale à $\sigma^2 I_n$. Il nous reste encore le R^2 et le R^2 ajusté. Rappelons que le coefficient R^2 peut s'interpréter comme le cosinus de l'angle entre le vecteur des observations Y_j et son projeté pour la norme \mathcal{L}^2 sur l'espace linéaire engendré par les observations X_j , soit

$$R^2 = \frac{\sum (\hat{y}_j - \bar{y})^2}{\sum (y_j - \bar{y})^2}.$$

Cet indicateur à des valeurs comprises entre 0 et 1, la proximité avec 1 indiquant une bonne adéquation du modèle aux données. Toutefois, son interprétation est sujette à caution : sa valeur augmente mécaniquement avec l'ajout de variables explicatives. En particulier, pour comparer la qualité de deux modèles au nombre de variables explicatives distinct, on lui préférera le R^2 ajusté, qui prend en compte ce nombre noté k :

$$RR^2 = 1 - (1 - R^2) \frac{n - 1}{n - k - 1}.$$

1 Théorie des tests

1.1 Exemples

1. Une entreprise vend des biens dont elle assure que la durée de vie dépasse 10000 heures. Vous êtes engagés pour vérifier la qualité d'iceux. Vous disposez d'un échantillon de 30 de ces biens. La durée de vie moyenne calculée sur l'échantillon vaut 9900 heures, on suppose que l'écart-type est connu et vaut 120 heures. Pouvez-vous rejeter leur assertion à un niveau de confiance de 0.05%.
Répondez à la même question si l'écart-type n'est plus connu, et que l'écart-type obtenu sur l'échantillon vaut 125 heures.
Calculez la puissance du test, c'est-à-dire la probabilité de l'erreur de seconde espèce.
2. Une firme agroalimentaire assure qu'un cookie qu'elle produit ne contient pas plus de 2 grammes d'un certain composé (graisse, colorant,...). Vous achetez un paquet, contenant 35 cookies, et mesurez une teneur moyenne de 2.1 grammes. En supposant que l'écart-type de l'échantillon est de 0.25, pouvez-vous incriminer la firme à un niveau de confiance de 0.05%. Même question si l'on ne connaît que l'écart-type empirique de 0.3. Calculez la puissance du test.
3. Lors des dernières élections, les médias affirment qu'au moins 60% des citoyens ont voté. Vous interrogez 148 citoyens de façon à obtenir un échantillon représentatif de la population (vous êtes statisticien après tout). Vous obtenez que 85 des personnes interrogées ont voté. A 0.05%, votre test concorde-t-il avec l'affirmation des médias ?

1.2 Analyse of Variance ou procédure ANOVA

1. Une institution de santé publique veut comparer l'effet de trois traitements contre la grippe. Pour cela, 18 hopitaux sont choisis de façon aléatoire, répartis par groupes de 6, chacun appliquant un et un seul des trois traitements. Voici le nombre de personnes guéries au bout d'une semaine de traitement :

Trait. 1	Trait. 2	Trait. 3
22	52	16
42	33	24
44	8	19
52	47	18
45	43	34
37	32	39

- (a) Une méthode pour rentrer les données dans *R* : les taper dans un fichier *.txt* puis utiliser *read.table* pour créer un objet de la classe *data.frame*. (Faites-le)
- (b) Utiliser un test ANOVA pour répondre à la problématique de l'institution.

- (c) Un pays frontalier, lui aussi touché par l'épidémie, décide de répliquer l'expérience avec le même nombre d'hôpitaux et les mêmes traitements. Chaque hôpital est sélectionné de façon aléatoire, et doit appliquer les trois traitements pendant trois semaines, chaque traitement pendant une semaine, l'ordre des traitements étant lui aussi aléatoire. Voici le résultat :

Trait. 1	Trait. 2	Trait. 3
22	52	16
42	33	24
44	8	19
52	47	18
45	43	34
37	32	39

2 Régression sur variables qualitatives

Le but de cet exercice est d'expliquer la concentration en ozone $O3$ en fonction de la température $T12$ et de la direction du vent $vent$ dans la table *ozone.txt*.

1. Télécharger la table, et effectuer des régressions selon les différents modèles.
2. Tester l'égalité des pentes.
3. Tester l'égalité des ordonnées à l'origine.
4. Analyser les résidus.

2.1 ANOVA à 1 facteur

Nous souhaitons modéliser la concentration en ozone en fonction de la direction du vent.

1. Tracer une boîte à moustaches de la variable $O3$ par rapport aux quatre modalités de la variable $vent$. Le vent semble-t-il avoir une influence sur la concentration en ozone ?
2. On se place dans un modèle d'analyse de la variance à un facteur

$$y_{ij} = \mu + \alpha_j + \epsilon_{ij}$$

- (a) Effectuer la régression linéaire de $O3$ sur $vent$ sous la contrainte $\mu = 0$.
 - (b) Effectuer la régression linéaire de $O3$ sur $vent$ sous la contrainte $\alpha_1 = 0$.
 - (c) Effectuer la régression linéaire de $O3$ sur $vent$ sous la contrainte $\sum n_i \alpha_i = 0$.
 - (d) Effectuer la régression linéaire de $O3$ sur $vent$ sous la contrainte $\sum n_i \alpha_i = 0$.
3. Analyser les résidus afin de constater que l'hypothèse d'homoscédasticité est vérifiée. Pour cela, tracer un boxplot des résidus en fonction de $vent$, les résidus en fonction de $O3$, leurs quantiles théoriques ainsi que la distribution des résidus par modalité de $vent$.

2.2 ANOVA à 2 facteurs

Nous voulons maintenant modéliser la concentration en ozone par le vent et la nébulosité, variable à 2 modalités : SOLEIL et NUAGEUX.

1. Procéder à un examen graphique qui puisse déterminer si l'interaction des facteurs influe sur la variable à expliquer. (voir ce qu'est un *profil*)
2. On suppose la gaussianité des résidus.
 - (a) Tester le modèle avec interaction : **mod1**.
 - (b) Tester le modèle sans interaction : **mod2**.
 - (c) Tester le modèle sans effet du facteur *nebulosité* : **mod3**.
3. Grâce à la commande ANOVA de R, effectuer des analyses de la variance entre les modèles **mod1**, **mod2** et **mod3**.
4. Répondez à la problématique.

3 Problème du voyageur de commerce et algorithme du recuit simulé

D'après le livre de Michel Benaïm et Nicole El Karoui, *Promenade aléatoire, Chaîne de Markov et simulations, martingales et stratégies*, exemple 3.1.8.

La méthode du recuit simulé est un algorithme d'optimisation proche de celui de Metropolis, et consiste à se promener aléatoirement sur l'espace (fini) des états d'un système selon une loi construite de façon à ce que la promenade converge vers un état qui minimise une certaine fonctionnelle.

On se donne une fonction $h :]0; \infty[\rightarrow]0; 1]$ telle que

$$h(x) = xh\left(\frac{1}{x}\right),$$

par exemple $\min(1, x)$ ou bien $\frac{x}{1+x}$. La fonction $V : E \rightarrow \mathbb{R}_+$ est la fonction "coût" à minimiser.

La terminologie "recuit simulé" vient d'une technique métallurgique consistant à faire fondre de façon répétée le métal puis à le faire lentement refroidir pour en améliorer les propriétés. En effet, on va se donner un schéma de décroissance d'une quantité analogue à la température, que nous noterons T_n . Ce schéma est crucial pour que l'algorithme converge vers un minimum global de la fonction V et ne reste pas piégé dans un de ses minima locaux. Vous pourrez choisir l'un des schémas suivants :

— décroissance logarithmique :

$$T_n = \frac{C}{\log(n)}$$

— recuit par palier :

$$T_n = \frac{1}{k} \text{ pour } e^{(k-1)C} \leq n < e^{kC}$$

Voici l'algorithme :

Initialiser X_0 . Choisir le nombre de pas de la marche aléatoire m . Pour n allant de 1 à $m - 1$, répéter :

1. Choisir un voisin y de X_n aléatoirement.
2. Tirer $U \sim \mathcal{U}_{[0,1]}$.
3. Si $U < h(\exp(\frac{1}{T_n})(V(X_n) - V(y)) \frac{N(y)}{N(X_n)})$, accepter $X_{n+1} = y$, sinon refuser i.e. $X_{n+1} = X_n$.

Le problème auquel nous allons appliquer cet algorithme est celui d'un commerçant devant visiter un ensemble fini $E = X_1, \dots, X_N$ de villes, une et une seule fois. Pour minimiser son temps et son argent, il souhaite trouver le chemin l qui minimise la distance, soit, avec nos notations :

$$V(l) = \sum_{j=1}^{N-1} d(X_{l(j)}, X_{l(j+1)}),$$

où l'on voit un chemin comme une permutation de l'ensemble E . Nous travaillerons avec des villes disposées aléatoirement dans le carré $[0; 1] \times [0; 1]$.

1. Créer une fonction qui calcule le coût d'un chemin donné l .
2. Créer une fonction qui affiche un chemin donné l .
3. Choisir une loi de transition sur les chemins, et l'implémenter.
4. Implémenter l'algorithme du recuit simulé sur ce problème, à l'aide d'une fonction si possible. Afficher le chemin obtenu, ainsi que l'évolution de la longueur du chemin en fonction du nombre d'itérations. Qu'en pensez-vous? De combien d'itérations avez-vous besoin pour obtenir un chemin plausible?