

# TP Statistiques et Séries Chronologiques

## Université de Lorraine

### Partiel 2015

Clément Dell'Aiera

#### 1 Estimation de la fonction de répartition

On se donne un  $n$ -échantillon  $X_1, \dots, X_n$  i.i.d suivant une loi donnée par la même fonction de répartition (f.d.r)  $F$  sur  $\mathbb{R}$ .  $\mathcal{F}$  dénote l'ensemble des fonctions de répartition sur  $\mathbb{R}$ .

1. Décrire l'expérience statistique.
2. Le modèle est-il dominé ?
3. On veut estimer  $F(x) = \mathbb{P}(X \leq x)$ .
  - (a) On pose  $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n 1_{X_i \leq x}$ . Calculer  $\mathbb{E}[\hat{F}_n(x)]$  et  $V[\hat{F}_n(x)]$ .
  - (b) Montrer que  $\hat{F}_n(x)$  converge presque-sûrement vers  $F(x)$ .
  - (c) Montrer que, si  $l(x, y) = (x - y)^2$  est la perte quadratique,

$$\sup_{F \in \mathcal{F}} \mathbb{E}[l(\hat{F}_n(x), F(x))] = \frac{1}{4n}.$$

- (d) En déduire que  $\hat{F}_n(x)$  converge uniformément en norme  $\mathcal{L}^2$  vers  $F(x)$ , et donc en probabilité.
4. (a) Montrer que

$$\mathbb{P}(|\hat{F}_n(x) - F(x)| > t) \leq \frac{1}{t^2} \text{Var}[\hat{F}_n(x)] \leq \frac{1}{4nt^2}$$

- (b) Soit  $\alpha \in ]0; 1[$ . Déterminer

$$t_{\alpha, n} = \inf\{t > 0 : \frac{1}{4nt^2} \leq \alpha\}$$

et en déduire un intervalle de confiance pour  $F(x)$  de niveau  $1 - \alpha$ .

- (c) Comment interpréter  $I_{n, \alpha}$  ? Quelle est sa précision ?
5. On pose  $\xi_n = \sqrt{n} \frac{\hat{F}_n(x) - F(x)}{\sqrt{\hat{F}_n(x)(1 - \hat{F}_n(x))}}$ .
  - (a) Déterminer la limite en loi de  $\xi_n$ .
  - (b) On note  $J_{n, \alpha}$  l'intervalle  $[-\phi^{-1}(1 - \frac{\alpha}{2}); \phi^{-1}(1 - \frac{\alpha}{2})]$ . Calculer la limite de  $\mathbb{P}(\xi_n \in J_{n, \alpha})$  lorsque  $n$  tend vers  $\infty$ .

- (c) Donner un intervalle de confiance asymptotique pour  $J_{n,\alpha}$ , ainsi que sa précision asymptotique.
6. Soient  $Y_j$  des variables aléatoires réelles indépendantes centrées :  $\mathbb{E}Y_j = 0$  et bornées :  $a_j \leq Y_j \leq b_j$ . On veut démontrer ce que l'on appelle l'*inégalité de Hoeffding* : pour tout  $t > 0$ ,

$$\mathbb{P}(\sum Y_j \geq t) \leq e^{-st} \prod e^{\frac{s^2(b_j - a_j)^2}{8}} \quad \forall s > 0.$$

On pose  $\Phi_Y(s) = \log \mathbb{E}[e^{s(Y - \mathbb{E}Y)}]$ .

- (a) Montrer que

$$\Phi_Y''(s) = e^{-\Phi_Y(s)} \mathbb{E}[Y^2 e^{sY}] - e^{-2\Phi_Y(s)} (\mathbb{E}[Y e^{sY}])^2.$$

- (b) On définit une nouvelle mesure de probabilité par  $\mathbb{Q}(A) = e^{-\Phi_Y(s)} \mathbb{E}[e^{sY} 1_A]$  pour tout borélien  $A$ . Comment interpréter  $\Phi_Y''(s)$  dans ce cadre ?
- (c) Montrer alors que  $\Phi_Y(s) \leq s^2 \frac{(b-a)^2}{8}$ .
- (d) En déduire l'inégalité de Hoeffding.
7. (a) Soient  $X_j$  des v.a. de Bernoulli de paramètre  $p$  et  $\bar{X}_n = \frac{1}{n} \sum_{j=1}^n X_j$ , montrer que

$$\mathbb{P}(|\bar{X}_n - p| \geq t) \leq 2e^{-2nt^2}.$$

- (b) En déduire un intervalle de confiance de niveau  $1 - \alpha$  pour  $F(x)$ .
8. Comparer les différents intervalles de confiance que vous avez obtenu.

## 2 Exercice 1

1. Que font les commandes *pnorm*, *qnorm*, *rnorm*, *dnorm*? Utilisez l'aide de *R*.
2. Tracer la fonction de densité de la loi normale. Faites varier les paramètres et afficher les différentes courbes sur le même graphique.
3. Simuler 2 vecteurs  $X$  et  $Y$  contenant chacun  $N = 100$  variables indépendantes identiquement distribuées suivant une loi normale  $\mathcal{N}(0, 1)$ .
4. Afficher les points de coordonnées  $(X[j], Y[j])$  dans le plan, pour  $j$  allant de 1 à 100.
5. Tracer la fonction de répartition empirique des  $X[j]$ .
6. Soit  $\mathcal{E}$  une v.a. de loi exponentielle de paramètre 1, et  $U$  une v.a. suivant une loi uniforme sur  $[0, 2\pi]$ . On pose

$$(X, Y) = (\sqrt{\mathcal{E}} \cos(U), \sqrt{\mathcal{E}} \sin(U)).$$

Quelle est la loi du couple  $(X, Y)$ ? (Vous pouvez le prouver, ou observer grâce à *R* ce qu'il se passe en simulant ces variables et en les traçant.)

## 3 Exercice 2

1. Une table est déjà en mémoire dans *R* : la table *stackloss*. Analyser la rapidement.
2. Tracer *stack.loss* en fonction de *Air.Flow*. Qu'en pensez-vous?
3. Effectuer la régression linéaire de *stack.loss* en fonction des autres variables. Quelles sont celles qui sont significatives?

## 4 Exercice 3

Le fichier *ozone.dta* contient les variables suivantes, pour une série de journées (qui sont ici nos individus) :

- l'identifiant de la journée,
- le maximum d'ozone (variable *maxO3*)
- l'heure à laquelle le maximum d'ozone a été obtenu (heure),
- les températures à 6h, 9h, 12h, 15h, 18h (resp. *T6* à *T18*)
- la nébulosité à 6h, 9h, 12h, 15h, 18h (resp. *Ne6* à *Ne18*)
- la projection du vent sur l'axe est-ouest à 12h (*Vx*),
- le maximum d'ozone de la veille (*maxO3v*).

Le but est de modéliser la valeur des pics d'ozone en fonction de grandeurs physiques facilement mesurables (température, heure, nébulosité, vent) afin d'avoir des approximations de la qualité de l'air faciles et rapides à obtenir.

1. Importer la table, et afficher un résumé de ce qu'elle contient.
2. Tracer *maxO3* en fonction de *T12*, puis effectuer une régression linéaire. Ajouter la droite de régression sur le graphique. Soignez la présentation.
3. Afficher les résultats de la régression.

4. Extraire les résidus et tracer leur densité estimée.
5. Effectuer la régression de *maxO3* sur toutes les variables, et supprimer récursivement celles qui ne sont pas significatives, jusqu'à ce qu'elles le soient toutes.

## 5 Exercice 4

1. Sur le site *data.gouv.fr*, vous pourrez trouver des tables de données publiques en libre accès. Choisissez un thème qui vous intéresse, puis une table en conséquence. Télécharger là.
2. Les tables sont souvent au format *.xls* : vous aurez besoin d'installer un package pour pouvoir les lire. La commande pour ce faire est *install.packages("nom du package")*. Installer le package *gdata*.
3. Analyser votre table.

## 6 Exercice 5

1. Télécharger et installer le package *ISwR* (Introductory Statistics with R).
2. Utiliser la commande *summary* pour analyser rapidement la table *bp.obese*. L'échantillon provient d'un échantillon de population mexicaine en Californie, et la table décrit 3 variables : le sexe (femme = 1, homme = 0), le ratio d'obésité (*obese*) et la pression sanguine systolique en *mm* de mercure (*bp*).
3. Représenter les données dans un graphe, en utilisant des symboles différents pour les hommes et les femmes.
4. Expliquer la pression sanguine en fonction du ratio d'obésité, puis du ratio d'obésité et du sexe.
5. Tracer sur un même graphe les courbes correspondant aux régressions dans les 2 modèles. Soignez la présentation (couleurs différentes, légende,...)

## 7 Exercice 6

1. Télécharger la librairie *MASS*.
2. Analyser rapidement la table *cats* et afficher les variables les unes en fonctions des autres par paires.
3. Effectuer une régression linéaire selon le modèle  $Hwt \sim Bwt * Sex$ . Cela apporte-t-il quelque-chose par rapport au modèle  $Hwt \sim Bwt + Sex$  ?
4. Visualiser les composantes de votre régression.
5. En extraire les prédictions, les coefficients, les résidus, les résidus studentisés, et la formule du modèle.
6. Tracer le *qqplot* des résidus studentisés ainsi que la première bissectrice.
7. Tracer le graphe des résidus contre les prédictions.
8. Tracer le graphe des distance de Cook.
9. Observer les attributs que vous donne *summary*

10. Afficher le  $R^2$  ajusté de la régression, le nombre de degrés de liberté résiduels, la matrice de variance-covariance des paramètres estimés.

## 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  | $Pr(>  t )$ |
|-----------|-----------------|------------------|--|-------------|
| $\beta_j$ | $\hat{\beta}_j$ | $\hat{\sigma}_j$ | $\hat{t}_j = \frac{\hat{\beta}_j}{\hat{\sigma}_j}$ | $p$ -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  $\alpha$  qui fait basculer le test. (Rappelez vous que, mécaniquement, si  $\alpha$  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  $\sigma$  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  $k$  :

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

## 8 Théorie des tests

### 8.1 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 ?

### 8.2 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       |

## 9 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.

### 9.1 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  $vent$ .

### 9.2 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.