

# TP R sur les tests et la régression

Emmanuel Rachelson and Matthieu Vignes

16 octobre 2013, SupAero - ISAE

## 1 Un petit exercice: la chaleur latente de fusion de la glace

Voici 2 séries de mesures indépendantes de chaleur latente de fusion de la glace (en cal/g) :

Méthode A	79.98	80.04	80.02	80.04	80.03	80.03	80.04	79.97	80.05
	80.03	80.02	80.00	80.02					
Méthode B	80.02	79.94	79.98	79.97	79.97	80.03	79.95	79.97	

1. Lire les données dans R.
2. Comparer les distributions à l'aide de `boxplot()`(s). Conclusions ?
3. Tester pour l'égalité des moyennes avec `t.test()`. Quelles sont les hypothèses de ce test ? Conclusions ?
4. Tester l'hypothèse d'égalité des variances à l'aide de `var.test()`. Conclusions ?
5. Appliquer un *t-test* classique qui suppose l'égalité des variances. Conclusions ?
6. Tous les tests ci-dessus supposent la normalité des deux échantillons. Appliquer un test de rang signé, `wilcox.test()`. Quelles sont les hypothèses de ce test ? Conclusions ?
7. Tester la normalité des données avec `qqplot()` et des tests appropriés (Shapiro-Wilk, Kolmogorov). Conclusions ?

## 2 Etude d'un scénario complet: étude de la concentration en ozone

### 2.1 Introduction

La pollution de l'air constitue actuellement une des préoccupations majeures de santé publique. De nombreuses études épidémiologiques ont permis de mettre en évidence l'influence sur la santé de certains composés chimiques comme le dioxyde soufre (SO<sub>2</sub>), le dioxyde d'azote (NO<sub>2</sub>), l'ozone (O<sub>3</sub>) ou des particules en suspension. Des associations de surveillance de la qualité

de l'air (Air Breizh en Bretagne depuis 1994) existent sur tout le territoire métropolitain et mesurent la concentration des polluants. Elles enregistrent également les conditions météorologiques comme la température, la nébulosité, le vent, les chutes de pluie en relation avec les services de Météo France... L'une des missions de ces associations est de construire des modèles de prévision de la concentration en ozone du lendemain à partir des données disponibles du jour : observations et prévisions de Météo France. Plus précisément, il s'agit d'anticiper l'occurrence ou non d'un dépassement légal du pic d'ozone ( $180 \mu\text{g}/\text{m}^3$ ) le lendemain afin d'aider les services préfectoraux à prendre les décisions nécessaires de prévention : confinement des personnes à risque, limitation du trafic routier. Plus modestement, l'objectif de cette étude est de mettre en évidence l'influence de certains paramètres sur la concentration d'ozone (en  $\mu\text{g}/\text{m}^3$ ) et différentes variables observées ou leur prévision. Les 112 données étudiées ont été recueillies à Rennes durant l'été 2001. Elles sont disponibles sur le site du laboratoire de mathématiques appliquées de l'Agrocampus Ouest. Les 13 variables observées sont :

- MaxO3 : Maximum de concentration d'ozone observé sur la journée en  $\mu\text{g}/\text{m}^3$
- T9, T12, T15 : Température observée à 9, 12 et 15h
- Ne9, Ne12, Ne15 : Nébulosité observée à 9, 12 et 15h
- Vx9, Vx12, Vx15 : Composante E-O du vent à 9, 12 et 15h
- MaxO3v : Teneur maximum en ozone observée la veille
- vent : orientation du vent à 12h
- pluie : occurrence ou non de précipitations

## 2.2 Exploration statistique élémentaire

- lire les données (fichier `ozone.csv` à récupérer là : [http://carlit.toulouse.inra.fr/wikiz/images/7/72/OzoneREMOVE\\_TXT\\_EXT.csv.txt](http://carlit.toulouse.inra.fr/wikiz/images/7/72/OzoneREMOVE_TXT_EXT.csv.txt); supprimer la variable inutile `obs`. Aide R: `read.csv2()`, `summary()`
- description unidimensionnelle : variables qualitatives, quantitatives. Aide R: `mean`, `sd`, `boxplot`, `hist`, `barplot`, `pie`
- description bidimensionnelle : variables quantitatives, variables qualitatives, variables qualitatives vs quantitatives. Aide R: `pairs`, `plot`, `table`, `mosaicplot`, `boxplot`

## 2.3 Tests de comparaisons

Important : Lors de l'exécution de chaque test précisez vous bien :

1. la question posée,
2. l'hypothèse ( $H_0$ ) en relation avec la question et associée au test,
3. la p-valeur calculée et la décision du test,
4. la réponse à la question.

### 2.3.1 Gaussanité

Beaucoup des outils ci-dessous nécessitent de vérifier le caractère gaussien ou non de la distribution. En fait, le nombre important d'observations dans l'échantillon permet de s'affranchir de cette hypothèse mais il est utile de savoir la vérifier et éventuellement de sélectionner la transformation la plus appropriée des données notamment pour les variables de concentration d'ozone.

**Normalité d'une distribution : Shapiro-Wilks** La droite de Henri ou graphe quantile-quantile donne déjà un aperçu graphique de la normalité de la distribution avant de calculer le test.

```
# qq-plots
qqnorm(ozone$max03); qqline(ozone$max03,col=2)
qqnorm(log(ozone$max03)); qqline(log(ozone$max03),col=2)
# Test de shapiro-Wilks
shapiro.test(ozone$max03); shapiro.test(log(ozone$max03))
```

Le test de Kolmogorov-Smirnov de comparaison à une distribution théorique pourrait également être utilisé (`ks.test`).

Les variables transformées sont ajoutées dans la table.

```
ozone=data.frame(ozone,Lmax03=log(ozone$max03), Lmax03v=log(ozone$max03v))
summary(ozone)
```

**Intervalle de confiance d'une moyenne : Student** Il est important de savoir estimer l'intervalle de confiance d'une moyenne ; celui-ci permet de tester l'égalité de cette moyenne à une valeur théorique selon l'appartenance ou non de cette valeur à l'intervalle. L'effectif étant suffisamment grand, il n'est pas nécessaire de supposer la normalité des données mais la variable transformée "la plus gaussienne" est choisie. L'intervalle de confiance est calculé par défaut avec un seuil à 95% mais ce paramètre peut être précisé (`conf.level=.95`) de même que la moyenne théorique testée (`mu=0.0`, par défaut à 0).

```
t.test(log(ozone$Lmax03), conf.level=.95)
```

**Comparaison de deux variances : Fisher** On s'intéresse à l'influence de la présence de pluie sur la concentration en ozone. Tester l'égalité des deux moyennes nécessite de vérifier préalablement plusieurs points :

1. la normalité des distributions dans chaque classe (sauf si l'échantillon considéré est de taille suffisamment grande),
2. le caractère indépendant ou appariés des échantillons,
3. l'égalité ou non des variances à l'intérieur de chaque groupe.

On dispose de deux échantillons indépendants : les jours de pluie et les jours de temps sec. Testons les autres hypothèses.

```
# Normalité des distributions (facultatif)
shapiro.test(ozone[ozone$pluie=="Pluie", "LmaxO3"])
shapiro.test(ozone[ozone$pluie=="Sec", "LmaxO3"])
# égalité des variances (test de Fisher)
var.test(LmaxO3~pluie, data=ozone)
Commenter les résultats.
```

**Comparaison de deux moyennes** Le test de comparaison des moyennes à utiliser (Student *vs.* Welch) dépend du résultat précédent concernant l'égalité des variances.

→ **Echantillons indépendants** Si les variances sont différentes, il s'agit d'un test de Welch.

```
t.test(LmaxO3~pluie, var.equal=F, data=ozone)
```

Dans le cas où elles sont considérées égales, c'est un test de Student.

```
t.test(LmaxO3~pluie, var.equal=T, data=ozone)
```

→ **échantillons appariés** On se propose d'étudier la persistance moyenne de la concentration en comparant la moyenne du jour avec celle de la veille. La mesure étant observée au même point à deux instants différents, les échantillons sont cette fois appariés.

```
t.test(ozone$maxO3, ozone$maxO3v, paired=TRUE)
```

### 2.3.2 Cas non-paramétrique

Si l'hypothèse de normalité des distributions n'est pas vérifiée et si l'échantillon est trop réduit, c'est un test non-paramétrique qu'il faut mettre en oeuvre. Les tests non-paramétriques sont basés sur les rangs des observations et donc sur les comparaisons des médianes des échantillons. Une transformation des variables par une fonction monotone (*i.e.* log) qui ne changent pas leur ordonnancement n'a donc pas d'effet sur le calcul d'un test non paramétrique.

### Comparaison de deux médianes : Wilcoxon

→ **Echantillons indépendants**

```
tapply(ozone$Lmax03, ozone$pluie, median)
wilcox.test(max03~pluie, data=ozone)
```

→ **Echantillons appariés**

```
median(ozone$Lmax03 - ozone$Lmax03v)
wilcox.test(ozone$Lmax03, ozone$Lmax03v, paired=TRUE)
```

Comparer avec les résultats des tests paramétriques.

## 2.4 Tests de liaison

**Indépendance de 2 variables qualitatives** Le test du  $\chi^2$  est adapté à ce problème.

```
chisq.test(table(ozone$pluie, ozone$vent))
```

*Rque* : un avertissement peut signaler que les effectifs théoriques (sous hypothèse d'indépendance) de certaines cellules sont trop faibles pour justifier des propriétés asymptotiques du test du  $\chi^2$ . Il est dans ce cas nécessaire de regrouper des modalités.

**Une variable qualitative, une quantitative** L'ANOVA associée à un test de Fisher adapté à cette situation est sans doute le test le plus utilisé ; il revient au test de Student lorsque la variable qualitative n'a que deux modalités. L'ANOVA nécessite de vérifier :

1. le caractère indépendant des échantillons,
2. la normalité des distributions (ou une taille suffisante d'échantillon) dans chaque classe ou plutôt la normalité des résidus au modèle,
3. l'égalité des variances internes à chaque groupe.

Même si la normalité des résidus est vérifiée *a posteriori*, c'est *a priori* qu'il faut prendre en compte ce résultat pour statuer sur la légitimité du test. Si la normalité n'est pas vérifiée pour un petit échantillon ou si l'égalité des variances n'est pas acceptable, un test non-paramétrique (Kruskal-Wallis) doit être envisagé.

### *Cas gaussien : ANOVA - Fisher*

Deux tests permettent de comparer les variances des groupes. Le test de Bartlett dans le cas gaussien, celui de Levene si l'hypothèse de normalité n'est pas admissible. Le test de Bartlett est le plus utilisé.

# test de Bartlett

```
bartlett.test(Lmax03 ~ vent, data=ozone)
```

# ANOVA à un facteur : estimation des paramètres

```
res.anova=aov(Lmax03 ~ vent, data=ozone)
```

```

# normalité des résidus au modèle d'ANOVA
qqnorm(res.anova$residuals)
qqline(res.anova$residuals)
shapiro.test(res.anova$residuals)
# Interprétation du test
summary(res.anova)
  Commenter.

```

*Cas non-paramétrique : Kruskal-Wallis*

```

kruskal.test(maxO3 ~ vent, data=ozone)
  Comparer les résultats.

```

**Deux variables quantitatives** La régression simple permet de tester l'influence éventuelle d'une variable sur une autre et plus précisément, dans le cas de cet exemple, d'expliquer et même de prévoir la concentration d'ozone en fonction de celle de la veille. La commande `lm` produit un ensemble de résultats sous la forme d'une liste de matrices et vecteurs.

*Estimation du modèle*

```

# retracer le nuage de point
plot(LmaxO3 ~ LmaxO3v, data=ozone)
# estimation du modèle
res1.reg=lm(LmaxO3 ~ LmaxO3v, data = ozone)
# liste des résultats
names(res1.reg)

```

*Diagnostic des résidus*

Des graphiques précédents permettent de s'assurer de la validité du modèle ; statuer sur l'homoscédasticité des résidus, leur normalité, la bonne linéarité du modèle.

```

# nuage de point, normalité des résidus
qqnorm(res1.reg$residuals)
qqline(res1.reg$residuals)
shapiro.test(res1.reg$residuals)
# Repérage d'une structure particulière du nuage ou de la présence
de "grands" résidus
res.student=rstudent(res1.reg)
ychap=res1.reg$fitted.values
plot(res.student~ychap, ylab="Résidus")
# ajouter des lignes
abline(h=c(-2,0,2), lty=c(2,1,2))
# repérage des points influents

```

```

cook=cooks.distance(res1.reg)
plot(cook~ychap,ylab="Distance de Cook")
abline(h=c(0,1),lty=c(1,2))

```

Les résidus sont "grands" si, une fois normalisés ou plutôt "studentisés", ils sont de valeur absolue plus grande que 2. Une observation est influente si elle a un grand résidu est associée à une grande valeur sur la diagonale de la hat matrix. Cela correspond à une valeur élevée (plus grande que 1) de la distance de Cook.

### *Significativité du modèle*

```
summary(res1.reg)
```

Que dire de l'influence de seuil d'ozone de la veille ? Que dire également de la présence d'observations à effet levier potentiel ? Que dire de la qualité d'ajustement de ce modèle et donc de la qualité attendue de la prévision ?

## 2.5 ACP et régression multiple

**ACP** Cette description élémentaire permet de se familiariser avec la structure de corrélation particulière des variables. Il faut sélectionner les seules variables quantitatives et l'ACP est réduite.

```

res.pca=prcomp(ozone[,c(2:10,14,15)],scale=T)
# décroissance des valeurs propres
plot(res.pca)
# parts de variance expliquée
summary(res.pca)
# biplot du premier plan principal
biplot(res.pca)

```

Comment s'interprètent les axes 1 et 2 ?

## Régression multiple

### *Modèle linéaire complet*

La régression linéaire simple ne conduit pas à un modèle bien ajusté. Le modèle linéaire multiple ci-dessous, plus complexe, recherche un meilleur ajustement des données.

```

# estimation
res2.reg=lm(Lmax03~Lmax03v+T9+T12+T15+Ne9+Ne12+Ne15+Vx9+Vx12+Vx15,
data = ozone)
# diagnostics
plot(res2.reg)
# résultats
summary(res2.reg)

```

Commenter les résultats obtenus sur la validité du modèle et la qualité de l'ajustement par rapport au modèle précédent. Que dire à propos de la significativité des tests de Student sur la nullité des paramètres ? Que penser alors de la présence de variables présentant de fortes colinéarités ?

### *Sous-modèle*

Une procédure de sélection de modèle non détaillée (stepwise) conduit à considérer le modèle ci-dessous :

```
res3.reg=lm(LmaxO3~LmaxO3v+T12+Ne9+Vx9,data=ozone)
# diagnostics plot(res3.reg)
# résultats
summary(res3.reg)
```

Commenter à nouveau les résultats.

### *Meilleure prévision*

L'objectif est de rechercher le meilleur modèle de prévision de la concentration en ozone. Ceux-ci sont comparés en considérant le PRESS (predicted residual sums of squares) ou leave one out cross validation. Une fonction élémentaire est définie pour calculer le PRESS dans le cas élémentaire de la régression linéaire.

```
# définition de la fonction PRESS
press=function(model) {
h=influence(model)$hat
e=influence(model)$wt.res
n=length(e)
sum((e/(1-h))^2)/n
}
# application aux différents modèles
press(res1.reg)
press(res2.reg)
press(res3.reg)
```

Le meilleur modèle de prévision est-il celui qui ajuste le mieux les données ? Attention, cette analyse se limite volontairement aux outils les plus élémentaires. D'autres modèles seraient à tester, notamment une analyse de covariance associant les variables qualitatives au modèle, la présence ou non d'interaction... pour tenter d'améliorer la qualité de prévision.