



Université Abderahmane Mira de Bejaia  
Faculté des Sciences Exactes  
Département Recherche Opérationnelle

## TP Modèles Linéaires

Djamila Ouaret / Faiza Amarouche

# Application de régression linéaire simple et multiple sur "ozone" sous R

October 14, 2023

# SOMMAIRE

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Régression linéaire simple</b>	<b>2</b>
2.1	Importer les données . . . . .	2
2.2	Tracer le nuage des points et superposer la droite de régression . . . . .	4
2.3	Comparer graphiquement la régression avec et sans intercept . . . . .	5
2.4	Les valeurs des estimateurs renvoyées par la fonction lm sont-elles exactes ? . . . . .	7
2.5	Tracer le nuage des points des valeurs aberrantes . . . . .	8
2.6	Vérifier avec QQ plot l'hypothèse gaussienne des résidus . . . . .	8
2.7	Comparer par un graphique les résidus estimés aux résidus standardisés et aux résidus studentisés . . . . .	9
2.8	Marquer les points leviers . . . . .	11
2.9	Analyser la distance de Cook des observations . . . . .	12
2.10	Intervalle de confiance et prévision . . . . .	13
<b>3</b>	<b>Régression linéaire multiple</b>	<b>16</b>
3.1	Saisir les données . . . . .	16
3.2	Estimer les coefficients de la régression multiple . . . . .	17
3.3	Intervalle de confiance . . . . .	18
3.4	Test d'hypothèse . . . . .	18
3.5	Qualité du modèle . . . . .	19
3.6	Retirer les variables non significatives . . . . .	19
3.6.1	visualisation des variables significatives . . . . .	19
3.7	Créer le nouveau modèle . . . . .	20
3.8	Définir n et k . . . . .	20
3.9	Calculer les résidus studentisés . . . . .	21
3.10	Valeurs aberrantes . . . . .	21
<b>4</b>	<b>Conclusion</b>	<b>22</b>

## 1 Introduction

**La régression linéaire** est l'une des techniques les plus fondamentales en statistiques et en apprentissage automatique. Elle est utilisée pour modéliser et prédire les relations linéaires entre une variable dépendante continue et une ou plusieurs variables indépendantes.

**La régression linéaire simple** est le cas où une seule variable indépendante est utilisée pour prédire la variable dépendante. L'objectif est de trouver une relation linéaire qui minimise l'erreur de prédiction entre les valeurs observées et les valeurs prédites par le modèle. La relation linéaire est décrite par une équation de la forme  $Y = \beta_0 + \beta_1 X + \epsilon$ , où  $Y$  est la variable dépendante,  $X$  est la variable indépendante,  $\beta_0$  et  $\beta_1$  sont les coefficients de régression, et  $\epsilon$  est l'erreur résiduelle.

**La régression linéaire multiple** est une extension de la régression linéaire simple, où plusieurs variables indépendantes sont utilisées pour prédire la variable dépendante. L'équation de régression devient

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

où  $X_1, X_2, \dots, X_n$  sont les variables indépendantes,  $\beta_1, \beta_2, \dots, \beta_n$  sont les coefficients de régression correspondants, et  $\epsilon$  est l'erreur résiduelle.

## 2 Régression linéaire simple

L'ozone troposphérique est une problématique environnementale de plus en plus répandue à l'échelle de la planète. L'augmentation démesurée de la taille des agglomérations, la croissance démographique et l'augmentation des émissions sont autant de facteurs qui contribuent à cette pollution de nature urbaine.

De ce fait, la pollution de l'air est actuellement une des préoccupations majeures de santé publique. Des associations de surveillance de la qualité de l'air existent sur tout le territoire français et mesurent la concentration des polluants ainsi que les conditions météorologiques, comme la température, la nébulosité, le vent, ...

Dans ce rapport, on dispose d'un jeu de données de concentration d'OZONE dans l'air qui contient 112 données et 13 variables, relevées durant l'été 2001 à Rennes stockées dans le fichier CSV "ozone.csv".

Nous allons commencer par la mise en œuvre de la régression linéaire simple sur le pic d'ozone (maxO3), expliqué par la température à midi (T12).

### 2.1 Importer les données

On charge nos données sous le nom "ozone" et voyons toutes les variables existantes dans ce jeu de données.

October 14, 2023

```
1 ozone=read.csv(file.choose('ozone.csv'), header = TRUE, sep= ";", encoding="UTF-8"
  )
2 ozone
3 # Nos variables
4 names(ozone)
```

	maxO3	T9	T12	T15	Ne9	Ne12	Ne15	Vx9	Vx12	Vx15	maxO3v	vent	pluie
20010601	87	15.6	18.5	18.4	4	4	8	0.6946	-1.7101	-0.6946	84	Nord	Sec
20010602	82	17	18.4	17.7	5	5	7	-4.3301	-4	-3	87	Nord	Sec
20010603	92	15.3	17.6	19.5	2	5	4	2.9544	1.8794	0.5209	82	Est	Sec
20010604	114	16.2	19.7	22.5	1	1	0	0.9848	0.3473	-0.1736	92	Nord	Sec
20010605	94	17.4	20.5	20.4	8	8	7	-0.5	-2.9544	-4.3301	114	Ouest	Sec
20010606	80	17.7	19.8	18.3	6	6	7	-5.6382	-5	-6	94	Ouest	Pluie
20010607	79	16.8	15.6	14.9	7	8	8	-4.3301	-1.8794	-3.7588	80	Ouest	Sec
20010610	79	14.9	17.5	18.9	5	5	4	0	-1.0419	-1.3892	99	Nord	Sec
20010611	101	16.1	19.6	21.4	2	4	4	-0.766	-1.0261	-2.2981	79	Nord	Sec
20010612	106	18.3	21.9	22.9	5	6	8	1.2856	-2.2981	-3.9392	101	Ouest	Sec

Table Chart

<< < 1 of 12 > >>

Rows per page 10 112 rows

'maxO3' · 'T9' · 'T12' · 'T15' · 'Ne9' · 'Ne12' · 'Ne15' · 'Vx9' · 'Vx12' · 'Vx15' · 'maxO3v' · 'vent' · 'pluie'

```
1 # Statistique descriptive des données de l'étude
2 summary(ozone)
```

```
maxO3      T9      T12      T15
Min.   : 42.00  Min.   :11.30  Min.   :14.00  Min.   :14.90
1st Qu.: 70.75  1st Qu.:16.20  1st Qu.:18.60  1st Qu.:19.27
Median : 81.50  Median :17.80  Median :20.55  Median :22.05
Mean   : 90.30  Mean   :18.36  Mean   :21.53  Mean   :22.63
3rd Qu.:106.00  3rd Qu.:19.93  3rd Qu.:23.55  3rd Qu.:25.40
Max.   :166.00  Max.   :27.00  Max.   :33.50  Max.   :35.50

Ne9      Ne12      Ne15      Vx9
Min.   :0.000  Min.   :0.000  Min.   :0.00  Min.   :-7.8785
1st Qu.:3.000  1st Qu.:4.000  1st Qu.:3.00  1st Qu.: -3.2765
Median :6.000  Median :5.000  Median :5.00  Median :-0.8660
Mean   :4.929  Mean   :5.018  Mean   :4.83  Mean   :-1.2143
3rd Qu.:7.000  3rd Qu.:7.000  3rd Qu.:7.00  3rd Qu.: 0.6946
Max.   :8.000  Max.   :8.000  Max.   :8.00  Max.   : 5.1962

Vx12      Vx15      maxO3v      vent
Min.   :-7.878  Min.   :-9.000  Min.   : 42.00  Length:112
1st Qu.: -3.565  1st Qu.: -3.939  1st Qu.: 71.00  Class :character
Median : -1.879  Median : -1.550  Median : 82.50  Mode  :character
Mean   : -1.611  Mean   : -1.691  Mean   : 90.57
3rd Qu.: 0.000  3rd Qu.: 0.000  3rd Qu.:106.00
Max.   : 6.578  Max.   : 5.000  Max.   :166.00

pluie
Length:112
Class :character
Mode  :character
```

Pour décrire la relation entre la concentration maximale en ozone maxO3 et la température à midi T12, on calcule la corrélation entre les deux variables.

```
1 # Calculer la corrélation entre 'maxO3' et 'T12'
2 cor.test(ozone$T12, ozone$maxO3)
```

Pearson's product-moment correlation

```
data: ozone$T12 and ozone$maxO3
t = 13.258, df = 110, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.7006713 0.8466150
sample estimates:
cor
0.7842623
```

Le test sur le coefficient de corrélation entre maxO3 et T12 donne une p-value  $< 2.2 \times 10^{-16}$ , inférieure au seuil  $= 0.05$ . L'hypothèse nulle d'absence de corrélation est donc rejetée au seuil de 5%.

L'estimation ponctuelle égale à 0,7842623 indique **une corrélation relativement forte et positive entre maxO3 et T12**.

## 2.2 Tracer le nuage des points et superposer la droite de régression

```
1 # Graphique du nuage des points
2 plot(maxO3~T12, data=ozone, col='red')
3 # Droite de la regression simple
4 reg_simple= lm(maxO3~T12, data=ozone)
5 summary(reg_simple)
6 # Superposer la droite sur le graphe de nuage des points
7 abline(reg_simple)
8 # Afficher les points de données observés pour la variable T12 en fonction des
  valeurs prédites par la régression linéaire
9 points(ozone$T12,reg_simple$fitted.values,pch=15,col='blue')
```

Call:

```
lm(formula = maxO3 ~ T12, data = ozone)
```

Residuals:

	Min	1Q	Median	3Q	Max
Residuals	-38.079	-12.735	0.257	11.003	44.671

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-27.4196	9.0335	-3.035	0.003 **
T12	5.4687	0.4125	13.258	<2e-16 ***

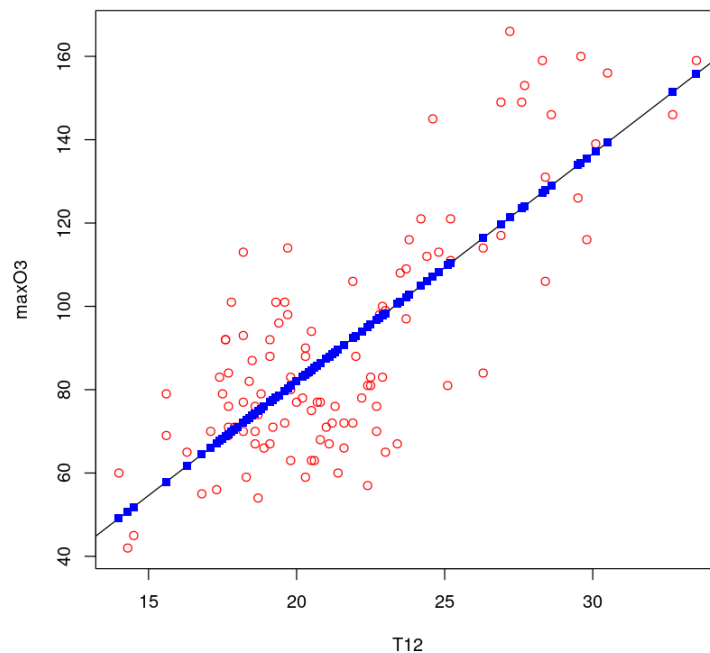
---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 17.57 on 110 degrees of freedom

Multiple R-squared: 0.6151, Adjusted R-squared: 0.6116

F-statistic: 175.8 on 1 and 110 DF, p-value: < 2.2e-16



$\beta_0 = -27.420$  et  $\beta_1 = 5.469$ , cela indique que pour chaque degré supplémentaire de température, on peut s'attendre à ce que la concentration en ozone augmente de 5.469 g/l. Par exemple, avec l'augmentation d'une unité de  $x$  (la température à midi) de 15 à 16, on voit bien que l'ozone passe de 54,6109 à 60,0796 degrés soit un 5,4687 de plus, ce qui est bel et bien la valeur de la pente  $\beta_1$ .

Les valeurs de  $p$  sont inférieures à 0.05, alors on peut rejeter l'hypothèse nulle, ce qui signifie que tout changement dans la valeur prédictive est **lié** à un changement dans la variable réponse. C'est-à-dire qu'une faible valeur de  $T12$  peut influencer significativement le modèle.

D'où,  $p = 0.003$  pour l'ordonnée à l'origine signifie que la constante (intercept  $\beta_0$ ) doit apparaître dans le modèle et  $p = 2e - 16$  pour la pente indique **une liaison significative** entre la concentration en ozone ( $\text{maxO3}$ ) et la température à midi ( $T12$ ).

### 2.3 Comparer graphiquement la régression avec et sans intercept

```

1 # Ajuster le modèle 'sans intercept'
2 reg_ajusté <- lm(maxO3 ~ T12 - 1, data = ozone)
3 summary(reg_ajusté)
4 # Tracé du nuage de points entre T12 et maxO3
5 plot(ozone$T12, ozone$maxO3, xlab = 'T12', ylab = 'maxO3')
6 # Ajout de la droite de régression linéaire du modèle reg_simple (bleue)
7 abline(reg_simple, col = 'blue')
8 # Ajout de la droite de régression linéaire du modèle reg_ajusté (rouge)
9 abline(reg_ajusté, col = 'red')

```

```

10 # Ajout d'une légende en haut à gauche
11 legend('topleft', c('régression avec intercept', 'régression sans intercept'), lty
    = 1, col = c('blue', 'red'))

```

Call:

```
lm(formula = maxO3 ~ T12 - 1, data = ozone)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-37.930	-14.290	-2.462	9.605	50.728

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
T12	4.23795	0.07855	53.95	<2e-16 ***

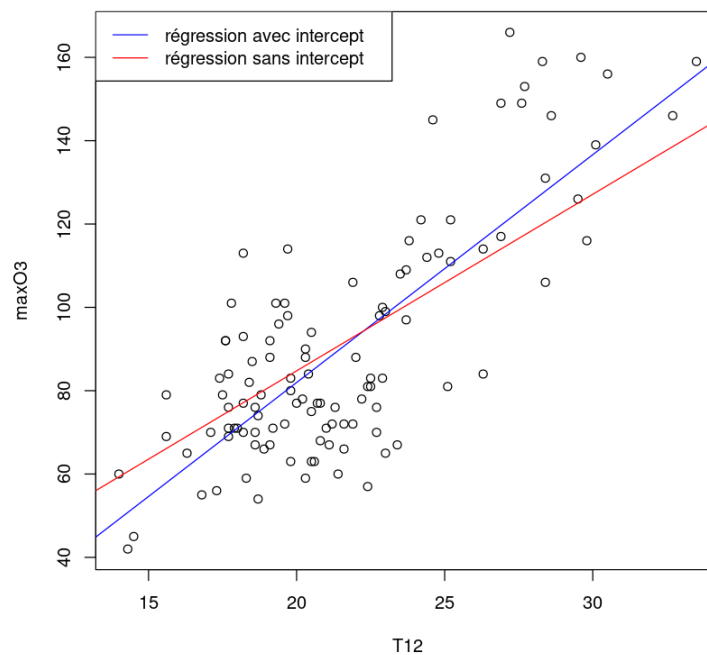
---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.21 on 111 degrees of freedom

Multiple R-squared: 0.9633, Adjusted R-squared: 0.9629

F-statistic: 2911 on 1 and 111 DF, p-value: < 2.2e-16



#### • Comparaison :

L'analyse met en évidence une relation linéaire positive entre les concentrations maximales de maxO3 et les valeurs de T12. Le nuage de points entre ces deux variables présente une tendance croissante.

L'ajustement d'un modèle de régression linéaire simple sans intercept entre maxO3 (variable dépendante) et T12 (variable explicative) montre une pente positive. Ce résultat traduit le fait qu'une augmentation de T12 est associée à une augmentation de maxO3.

La comparaison graphique avec un modèle de régression linéaire simple avec intercept indique que le modèle sans intercept (forçant une régression passant par 0) s'ajuste bien à la tendance des données.

En conclusion, cette analyse met en évidence et quantifie via un modèle statistique la relation linéaire positive existant entre les concentrations de maxO3 et les températures relevées par la variable T12 dans le jeu de données.

## 2.4 Les valeurs des estimateurs renvoyées par la fonction `lm` sont-elles exactes ?

Dans notre jeu de données, nous disposons de 112 observations et nous intéressons au couple  $(T12 ; maxO3) = (x_i; y_i)$  afin d'estimer les paramètres inconnus  $\beta_0$  et  $\beta_1$  de notre modèle.

Ce dernier se résume à notre droite :

$$y_i = \beta_0 + \beta_1 x_i$$

Et, on a:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{Cov(x, y)}{Var(x)}$$

### Calcul manuel:

Nous avons dans la sortie de la commande `summary(ozone)` le résumé statistique de toutes les variables de l'ensemble de données, alors, prenons les moyennes:  $\bar{y}$  et  $\bar{x}$ , telles que:  $\bar{y} = 90,30$ : La moyenne de la variable maxO3 et  $\bar{x} = 21,53$ : La moyenne de la variable T12, d'où :  $Cov(x, y) = 89,360$  et  $Var(x) = 16,34$ .

Ce qui nous déduit:

- $\beta_1 = \frac{89,360}{16,34}$  alors,  $\beta_1 = 5,469$
- $\beta_0 = 90,30 - 5,469 * 21,53$  donc,  $\beta_0 = -27,447$

### Sous R:

```
1 B1=cov(ozone$T12,ozone$maxO3)/var(ozone$T12)
2 B0=mean(ozone$maxO3)-mean(ozone$T12)*(cov(ozone$T12,ozone$maxO3)/var(ozone$T12))
3 cat('B0 =',B0)
4 cat('B1 =',B1)
```

B0 = -27.41964 ; B1 = 5.468685

On conclut explicitement que les valeurs des estimateurs  $\beta_0$  et  $\beta_1$  renvoyées par la fonction `lm` sont exactes en les comparant avec les résultats du calcul manuel ou du code R.



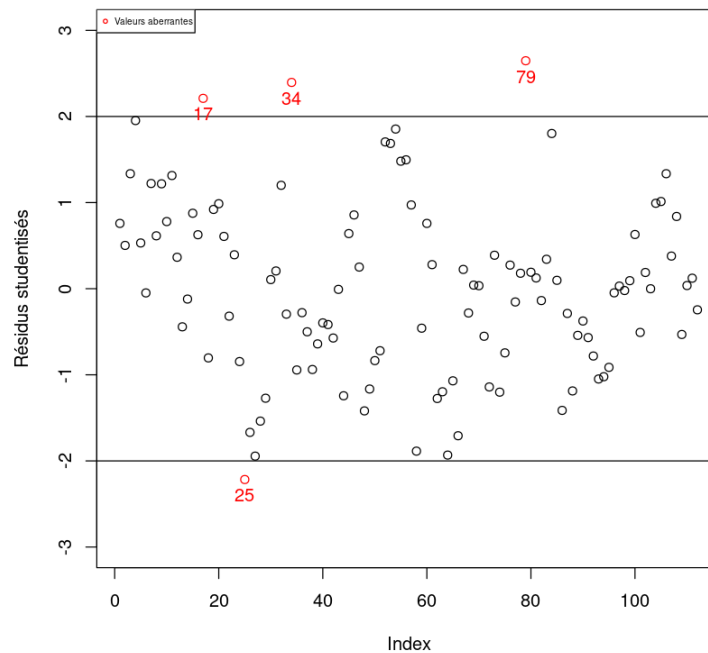
## 2.5 Tracer le nuage des points des valeurs abérrantes

Vérifions si le jeu de données contient des valeurs abérrantes.

```

1 # Résidus studentisés
2 res_stud <- rstudent(reg_simple)
3 # Identification des valeurs abérrantes
4 abr <- which(abs(res_stud) > 2)
5 # Tracé des résidus
6 plot(res_stud, ylab = "Résidus studentisés",
7       ylim = c(-3,3),
8       col = ifelse(abs(res_stud) > 2, "red", "black"))
9 # 2 Lignes horizontales (Seuils)
10 abline(-2,0)
11 abline(2,0)
12 # Légende
13 legend("topleft", c("Valeurs abérrantes"),
14        pch = 1, col = c("red"), cex = 0.5)
15 # Étiquettage des points abérrants
16 text(abr, rstudent(reg_simple)[abr], abr, col = "red", pos = 1)

```



Dans ce cas, on a **quatre** valeurs abérrantes qui se trouvent à l'extérieur de l'intervalle  $[-2;2]$  et qui sont: 17, 25, 34 et 79.

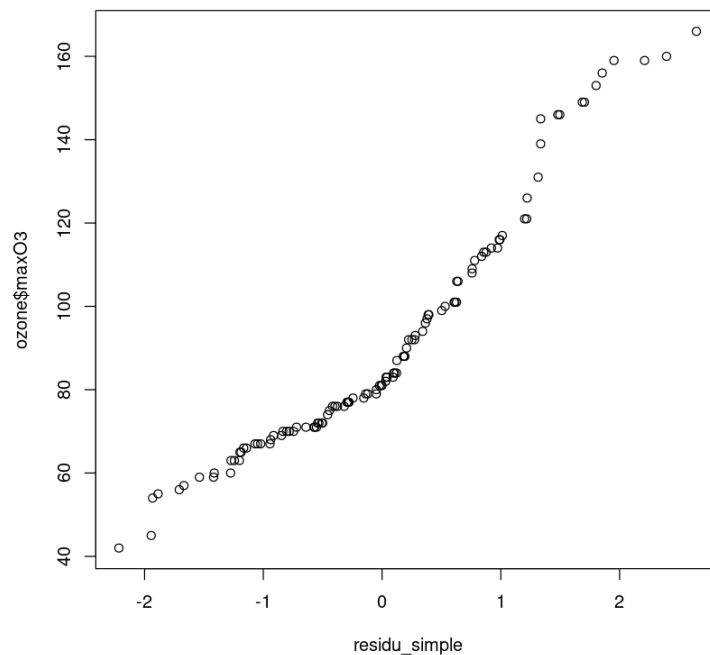
## 2.6 Vérifier avec QQ plot l'hypothèse gaussienne des résidus

On parcourt ce code afin de voir le graphique voulu.

```

1 # QQ-plot
2 qqplot(rstudent(reg_simple), ozone$maxO3)

```



Ici, la fonction `qqplot()` est utilisée pour comparer la distribution des résidus studentisés du modèle de régression (reg-simple) à la distribution des valeurs de la variable maxO3 contenues dans le dataframe ozone.

On peut remarquer que les points se distribuent approximativement le long de la diagonale, cela indique que la distribution est gaussienne et les résidus suivent une loi normale. Car, si les résidus s'écartent beaucoup de la diagonale, cela remet en cause l'hypothèse de normalité et la validité du modèle.

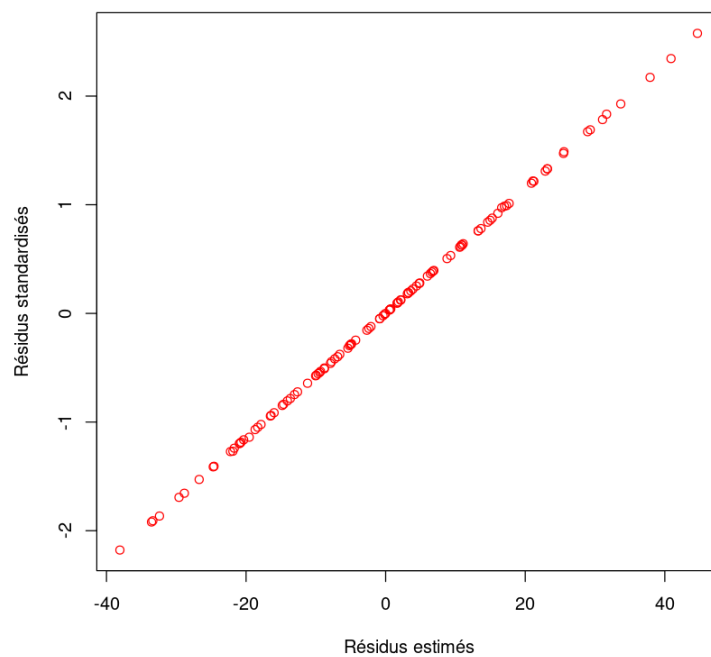
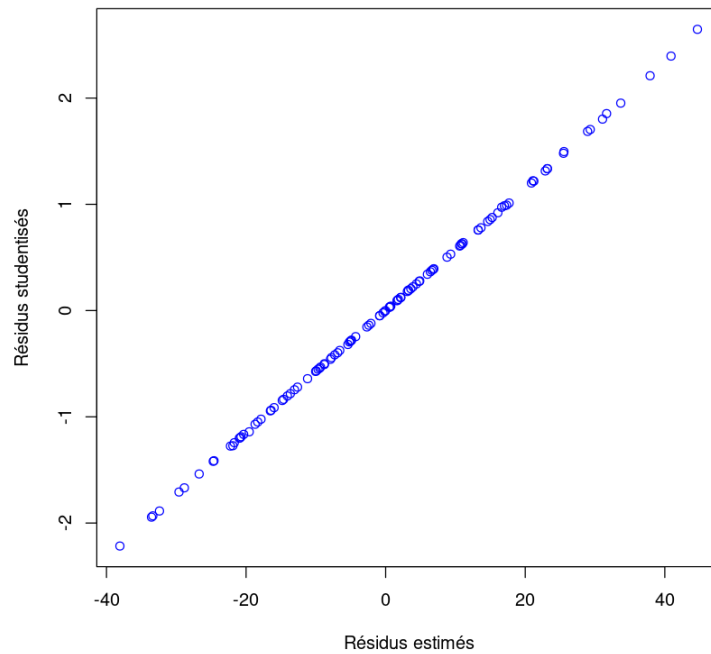
## 2.7 Comparer par un graphique les résidus estimés aux résidus standardisés et aux résidus studentisés

```

1 # Les résidus estimés
2 residus_estimes <- reg_simple$residuals
3 # Les résidus standardisés
4 residus_standardises <- rstandard(reg_simple)
5 # Les résidus studentisés
6 residus_studentises <- rstudent(reg_simple)
7
8 # Représentation graphique des résidus studentisés en fonction des résidus estimés
  (bleue)
9 plot(residus_estimes, residus_studentises, col='blue', xlab="Résidus estimés",
  ylab="Résidus studentisés")
10 # Représentation graphique des résidus standardisés en fonction des résidus
  estimés (rouge)

```

```
11 plot(residus_estimes, residus_standardises, col='red', xlab="Résidus estimés",  
      ylab="Résidus standardisés")
```



On voit clairement une similarité des graphiques qui est plutôt un signal positif pour la régression, indiquant que les résidus bruts (estimés) ne présentent a priori pas de problème

October 14, 2023

majeur de non-normalité.

En d'autres termes, les différentes transformations des résidus (standardisation et studentisation) suivent globalement la même tendance et n'ont dans ce cas que peu d'impact sur leur distribution.

## 2.8 Marquer les points leviers

On saisit la ligne de code ci-dessous:

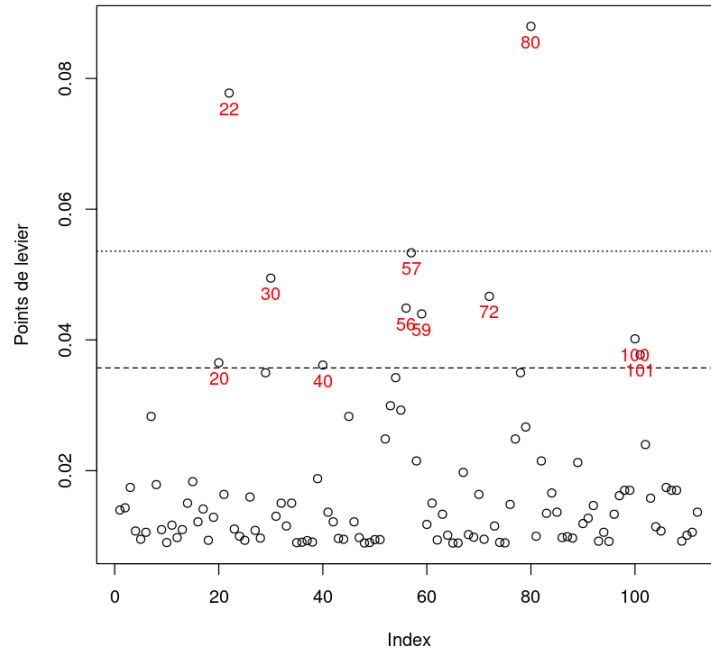
```
# Vérifier l'existence des points leviers
leviers=hatvalues(reg_simple)
leviers
```

20010601:	0.0139795884009378	20010602:	0.0143188562636468	20010603:	0.0174299602861664	20010604:	0.0107684564769346
20010605:	0.00950983768503931	20010606:	0.018572535685643	20010607:	0.0282951957268396	20010610:	0.0178684684290874
20010611:	0.0109754039660275	20010612:	0.009080536626062821	20010613:	0.0116624066201143	20010614:	0.00975833219620512
20010615:	0.0109754039660275	20010616:	0.0150304720824687	20010617:	0.0183180632698097	20010618:	0.0121755418785122
20010620:	0.0141357343180174	20010621:	0.00934896604033054	20010622:	0.0128684516921354	20010623:	0.0365120546420133
20010624:	0.0163674591988806	20010625:	0.0777576146732268	20010626:	0.0110752349621932	20010627:	0.00996823326506626
20010628:	0.00934896604033054	20010629:	0.0159679383009081	20010630:	0.0108631679348354	20010701:	0.00968744628176796
20010702:	0.0349742241093628	20010703:	0.0494516876065627	20010704:	0.013027944999619	20010705:	0.0150304720824687
20010706:	0.0115324491103129	20010707:	0.0150304720824687	20010708:	0.00898744787670104	20010709:	0.009052032820966
20010710:	0.00930544996507884	20010711:	0.00908156861864815	20010712:	0.0187785648083333	20010713:	0.0361511271943375
20010714:	0.0136513472360302	20010715:	0.0121755418785122	20010716:	0.00962857159172155	20010717:	0.00950983768503931
20010718:	0.0282951957268396	20010719:	0.0121755418785122	20010720:	0.00975833219620512	20010721:	0.00893152676642284
20010722:	0.00900536626062821	20010723:	0.00945076608967501	20010724:	0.00945076608967501	20010725:	0.024846399084154
20010726:	0.029939158225828	20010727:	0.0342218488897395	20010728:	0.029263969890814	20010729:	0.0448628091733462
20010730:	0.05332127091356	20010731:	0.0214899510355546	20010801:	0.0439781135801069	20010802:	0.0117775962310747
20010803:	0.0150304720824687	20010804:	0.0094021304761584	20010805:	0.0133341327689239	20010806:	0.0101251668034174
20010807:	0.00893152676642284	20010808:	0.00893743392595927	20010809:	0.0197327679787846	20010810:	0.0102137741964639
20010811:	0.00982232642451644	20010812:	0.0163674591988806	20010813:	0.00950983768503931	20010814:	0.0466052804532288
20010815:	0.0115324491103129	20010819:	0.00902899489877393	20010820:	0.00895692755242948	20010821:	0.0148355358177665
20010822:	0.024846399084154	20010823:	0.0349742241093628	20010824:	0.0266734835287717	20010825:	0.0879667615947204
20010826:	0.00996823326506626	20010827:	0.0214899510355546	20010828:	0.0134800396094737	20010829:	0.0165860240937285
20010830:	0.0136513472360302	20010831:	0.00975833219620512	20010901:	0.00989911949849003	20010902:	0.00968744628176796
20010903:	0.0212467729679716	20010904:	0.0119134609004126	20010905:	0.0127327839281154	20010906:	0.0146691508241571
20010907:	0.00921979615180061	20010908:	0.010572535685643	20010912:	0.0091784460350456	20010913:	0.0133341327689239
20010914:	0.0161805960422115	20010915:	0.0170024788418468	20010916:	0.0170024788418468	20010917:	0.0401630730461627
20010918:	0.0377228254416635	20010919:	0.0239906485726432	20010920:	0.0157861946924959	20010921:	0.0114223790376174
20010922:	0.0107684564769346	20010923:	0.0174299602861664	20010924:	0.0170024788418468	20010925:	0.0170024788418468
20010927:	0.00921979615180061	20010928:	0.0101251668034174	20010929:	0.010572535685643	20010930:	0.0136513472360302

Et voici le marquage des points leviers dans le nuage des points (T12,maxO3):

```
# Tracé des points de levier en fonction de leur index
plot(1:n, leviers, xlab='Index', ylab='Points de levier')
# Calcul du seuil théorique 1
p = reg_simple$rank # Nombre de paramètres du modèle
seuil1 = 2*p/n # Seuil à 2p/n
# Calcul du seuil théorique 2
seuil2 = 3*p/n # Seuil à 3p/n
# Ajout de lignes horizontales aux seuils
abline(seuil1,0, lty=2)
abline(seuil2,0, lty=3)
# Identification des points au dessus du seuil 1
lev = (1:n)[leviers > seuil1]
# Etiquetage des points de levier au dessus du seuil
text(lev, leviers[lev], lev, col='red', pos=1)
```

October 14, 2023



Le graphique des points de levier en fonction de leur index montre que 9 points dépassent le seuil théorique de  $2p/n$ , et 2 points dépassent le seuil de  $3p/n$ . Ces points avec un effet de levier élevé sont potentiellement influents sur le modèle.

## 2.9 Analyser la distance de Cook des observations

Sous R, on obtient la distance de Cook par la fonction `cooks.distance`.

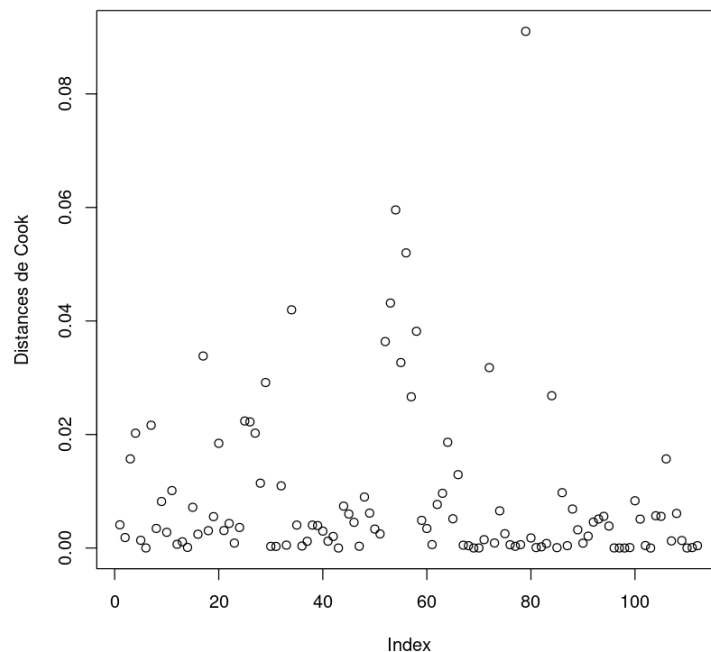
```
1 cook=cooks.distance(reg_simple)
2 cook
```

20010601:	0.00408918653573329 20010602:	0.00184731050572267 20010603:	0.015703649320165 20010604:	0.0202311629927259
20010605:	0.0013616685801166 20010606:	1.29504999381566e-05 20010607:	0.0216319873789547 20010610:	0.00344744705770462
20010611:	0.00819589358876562 20010612:	0.00277825908126665 20010613:	0.0101207536312791 20010614:	0.008661495708667857
20010615:	0.00109652163254591 20010616:	0.000111794079539862 20010617:	0.00717551581074413 20010618:	0.00243172584515188
20010620:	0.033828553416731 20010621:	0.00305921743470014 20010622:	0.00553046458257112 20010623:	0.0184489091134327
20010624:	0.00308458354724597 20010625:	0.0043292530868707 20010626:	0.000874927984835927 20010627:	0.00361552099006556
20010628:	0.0223798416404685 20010629:	0.0222281020252553 20010630:	0.020244953805004 20010701:	0.0114254187964979
20010702:	0.0291576436544318 20010703:	0.000291209539752326 20010704:	0.00028212163763763 20010705:	0.0109530854188215
20010706:	0.000514730341494657 20010707:	0.0419664791729556 20010708:	0.0040444887438058 20010709:	0.000357327723050524
20010710:	0.00118470172816609 20010711:	0.00404135514383491 20010712:	0.00395577141544445 20010713:	0.00298100916293773
20010714:	0.00121076853852777 20010715:	0.00203458298459526 20010716:	3.18600513714913e-07 20010717:	0.00738717432049722
20010718:	0.0059987060972919 20010719:	0.00452889720768916 20010720:	0.000312896775636571 20010721:	0.00899887225670057
20010722:	0.00614901971632422 20010723:	0.00333812786539842 20010724:	0.00248760413897949 20010725:	0.0363713210602906
20010726:	0.0431619431029163 20010727:	0.0595671989030089 20010728:	0.0326748993946178 20010729:	0.0519960646071453
20010730:	0.0266413209293377 20010731:	0.0381879982878962 20010801:	0.004873366965255 20010802:	0.00343801667417981
20010803:	0.000600093581653987 20010804:	0.00767479516952771 20010805:	0.00964204645899498 20010806:	0.0186314698297323
20010807:	0.00515369748387096 20010808:	0.012925411141489 20010809:	0.000507363502932897 20010810:	0.000414549711710327
20010811:	8.7358103460528e-06 20010812:	1.015740169414835e-05 20010813:	0.00147410234059285 20010814:	0.0317849382702789
20010815:	0.000887301983009746 20010819:	0.00655011326543949 20010820:	0.00252135291366663 20010821:	0.000569691816433924
20010822:	0.0003058959619106364 20010823:	0.000588109340489449 20010824:	0.0918274316747239 20010825:	0.00177503416572847
20010826:	7.87897170954975e-05 20010827:	0.000210633109069629 20010828:	0.000803476296419307 20010829:	0.0260348621633292
20010830:	6.5861825861833e-05 20010831:	0.009752713595617414 20010901:	0.000416861780049919 20010902:	0.00687027491849585
20010903:	0.00321183293977904 20010904:	0.000855722869793804 20010905:	0.00209048340037949 20010906:	0.00456584038956178
20010907:	0.00511204955323856 20010908:	0.00558136363696103 20010912:	0.00387046441298211 20010913:	1.58362782347353e-05
20010914:	7.61299339395888e-06 20010915:	4.0321141523236e-06 20010916:	7.51767727214779e-05 20010917:	0.00832694905515087
20010918:	0.00509091989743456 20010919:	0.000438986647553955 20010920:	7.356196464981593e-09 20010921:	0.00568512016819358
20010922:	0.00557691326733047 20010923:	0.015703649320165 20010924:	0.00125079782801956 20010925:	0.0009466233143995
20010927:	0.0013243053331595 20010928:	6.85479352666282e-06 20010929:	0.01045979861869e-05 20010930:	0.000419931203851446

October 14, 2023

**Note:** Le seuil critique pour la distance de Cook à partir duquel on considère que l'observation est **trop influente** est le quantile  $f(n, n-p)_{0.5}$ . Une distance de Cook en dessous de  $f(n, n-p)_{0.1}$  est considérée comme **souhaitable**. Traçons son graphique:

```
1 # Taille des données
2 n <- length(cook)
3 # Nombre de paramètres
4 p <- length(coef(reg_simple))
5 # Graphique des distances de Cook
6 plot(1:n, cook, xlab='Index', ylab='Distances de Cook')
7 # Seuil à 50%
8 s1 = qf(0.5, p, n-p)
9 # Seuil à 10%
10 s2 = qf(0.1, p, n-p)
11 # Lignes horizontales aux seuils
12 abline(h = s1, lty=3)
13 abline(h = s2, lty=2)
```



D'après le graphique ci-haut, les observations de l'analyse de régression n'ont pas une influence excessive sur les estimations des paramètres du modèle, car les distances de Cook sont toutes en dessous du seuil critique  $f(n, n-p)_{0.1}$ . Cela n'est qu'une indication que notre modèle de régression est relativement **robuste aux observations atypiques ou influentes**.

## 2.10 Intervalle de confiance et prévision

- Ajouter à la figure du nuage de points et de la droite de régression les intervalles de prédiction et de confiance en tout points  $x_i$  observé.

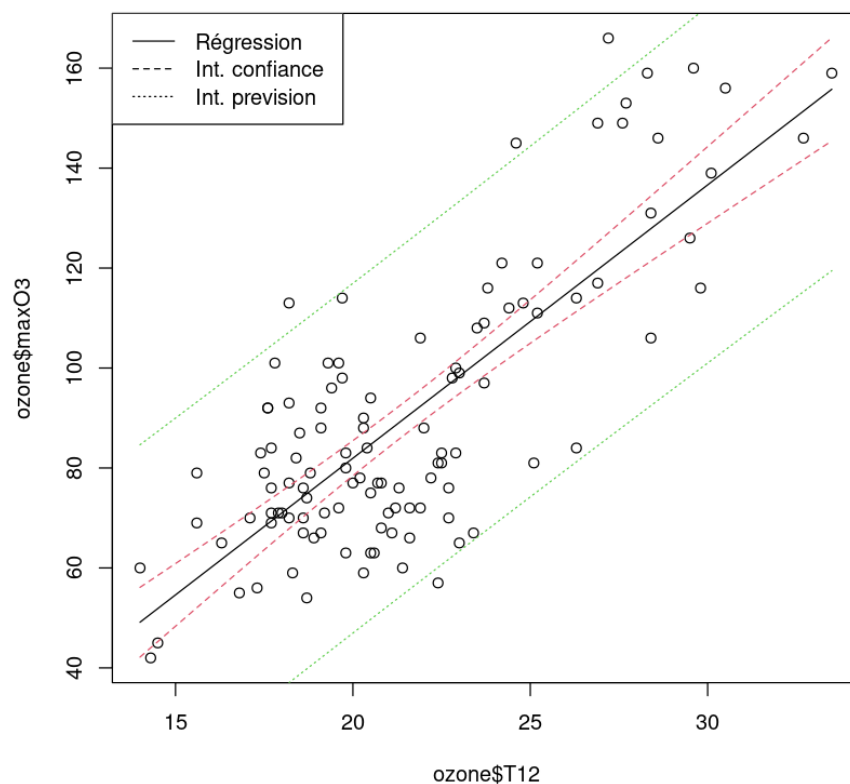
On parcourt le code suivant:

```

1 plot(ozone$maxO3 ~ ozone$T12, ylim = c(min(ozone$maxO3), max(ozone$maxO3)))
2
3 seq.x <- seq(min(ozone$T12), max(ozone$T12), length=3*n)
4 grid.x <- data.frame(seq.x)
5 dimnames(grid.x)[[2]] <- "T12"
6
7 ICconf <- predict(reg_simple, new=grid.x, interval="confidence", level=0.95)
8 ICprev <- predict(reg_simple, new=grid.x, interval="prediction", level=0.95)
9
10 matlines(grid.x, cbind(ICconf, ICprev[, -1]), lty=c(1,2,2,3,3), col= c(1, 2, 2, 3,
11 3))
12 legend("topleft", lty = 1:3, c("Régression", "Int. confiance", "Int. prévision") )

```

On obtient ce graphique:



### • Interprétation du graphique

L'écart des points par rapport à la ligne de régression est en accord avec la discussion précédente sur les résidus et la variance.

Les intervalles des prévision sont toujours plus amples que ceux de confiance.

- Ceci est interprétable comme suit : Les intervalles de confiance concernent l'espérance de

la concentration O3, pour une température T12 donnée, alors que les intervalles de prévision essaient de prendre en compte toute concentration possible dans le seuil fixé.

- **Comparaison entre les deux intervalles**

**Remarque:** Les observations atypiques (ou valeurs aberrantes) peuvent avoir un impact important sur les intervalles de confiance et de prédiction (seront plus larges).

Dans notre dessin, on observe que l'intervalle de prédiction est plus large que l'intervalle de confiance et cela signifie que l'estimation de la dispersion des valeurs futures est plus large que l'estimation de la dispersion des valeurs moyennes.

De plus, l'intervalle de confiance est d'une forme hyperbolique, plus une observation est éloignée du centre de gravité, moins nous avons d'information sur elle. Lorsque  $x_i$  est proche de  $\bar{x}$ , le terme dominant de la variance est  $1/n$ , mais dès que  $x_i$  s'éloigne de la moyenne, le terme dominant est le terme au carré.

Egalement, l'intervalle de confiance a 95% de chance de contenir la vraie valeur et l'intervalle de prédiction a 95% de contenir une observation future.



### 3 Régression linéaire multiple

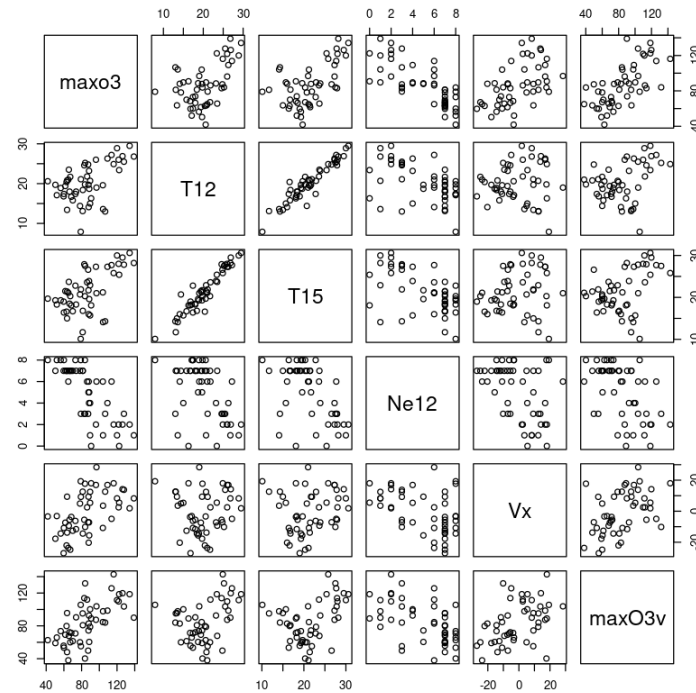
Dans cette partie de la régression linéaire multiple, on utilisera des données saisies à la main, composées de 50 observations et 6 variables (MaxO3, maxO3v, Vx, Ne12, T15, T12)

#### 3.1 Saisir les données

```

1 maxO3<-c(63.6,89.6, 79, 81.2, 88, 68.4, 139, 78.2, 113.8, 41.8, 65, 73, 126.2,
2         127.8, 61.6, 63.6, 134.2, 67.2, 87.8, 96.8, 89.6, 66.4, 60, 90.8, 104.2, 70,
3         96.2, 65.6, 109.2, 86.2, 87.4, 84, 83, 59.7, 52, 73.8, 120, 122.4, 106.6,
4         121.8, 116.2, 81.4, 88.6, 63, 104, 88.4, 83.8, 56.4, 50.4, 79.2)
5
6
7 T12 <- c(13.4, 15, 7.9, 13.1, 14.1, 16.7, 26.8, 14.4, 27.2, 20.6, 21, 17.4, 26.9,
8         25.5, 19.4, 20.8, 29.5, 21.7, 19.7, 19, 20.7, 18, 17.4, 16.3, 13.6, 15.8, 26,
9         23.5, 26.3, 21.8, 24.8, 25.2, 24.6, 16.8, 17.1, 18, 28.9, 23.4, 13, 26, 24.9,
10        18.4, 18.7, 20.4, 19.6, 23.2, 19.8, 18.9, 19.7, 21.1)
11
12
13 T15 <- c(15, 15.7, 10.1, 11.7, 16, 18.1, 28.2, 20.7, 27.7, 19.7, 21.1, 22.8, 29.5,
14         27.8, 21.5, 21.4, 30.6, 20.3, 21.7, 21, 22.9, 18.5, 16.4, 18.1, 14.1, 16.7,
15         27.3, 23.7, 27.3, 23.6, 26.6, 27.5, 27.9, 19, 18.3, 18.3, 30, 25.4, 14.3, 28,
16         25.8, 16.8, 19.6, 16.6, 21.2, 23.9, 20.3, 19.3, 19.3, 21.9)
17
18
19 Ne12 <- c( 7, 4, 8, 7, 6, 7, 1, 7, 6, 8, 6, 8, 2, 3, 7, 7, 2, 7, 5, 6, 1, 7, 8,
20         0, 1, 7, 2, 7, 4, 6, 3, 3, 3, 7, 8, 7, 1, 0, 3, 2, 2, 7, 5, 7, 6, 4, 8, 8, 7,
21         3)
22
23
24 Vx <- c(9.35, 5.4, 19.3, 12.6, -20.3, -3.69, 8.27, 4.93, -4.93, -3.38, -23.68,
25         -6.24, 14.18, 13.79, -7.39, -13.79, 1.88, -24.82, 9.35, 28.36, 12.47,
26         -5.52, -10.8, 18, 3.55, -12.6, 16.91, -9.35, 16.91, 2.5, -7.09, -10.15,
27         -5.52, -27.06, -3.13, -11.57, 8.27, 5.52, 12.6, 2.5, 18, -14.4, -15.59,
28         -22.06, -10.8, -7.2, 17.73, -14.4, -17.73, 9.26)
29
30
31 maxO3v <- c(95.6, 100.2, 105.6, 95.2, 82.8, 71.4, 90, 60, 125.8, 62.6, 38, 70.8,
32         119.8, 103.6, 69.2, 48, 118.6, 60, 74.4, 103.8, 78.8, 72.2, 53.4, 89, 97.8,
33         61.4, 87.4, 67.8, 98.6, 112, 49.8, 131.8, 113.8, 55.8, 65.8, 90.4, 111.4,
34         118.6, 84, 109.8, 142.8, 80.8, 60.4, 79.8, 84.6, 92.6, 40.2, 73.6, 59, 55.2)
35
36
37 # Créer un dataframe
38 df_ozone=data.frame(maxO3, T12, T15, Ne12, Vx, maxO3v)
39
40 # Visualiser les données
41 plot(df_ozone)

```



L'objectif de cette application est de modéliser le pic d'ozone journalier en fonction de toutes les autres variables météorologiques; il s'agira là d'une **régression linéaire multiple** avec 5 variables explicatives.

La régression linéaire assume que la relation entre les variables explicatives et la variable à expliquer (variable numérique continue) va être linéaire, du type :  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_5 X_5 + \epsilon$

où :

- $y$ : est la variable à expliquer (ici :maxO3).
- $X_1, \dots, X_5$  sont les variables explicatives (ici : maxO3v, Vx, Ne12, T15, T12).
- $\epsilon$ : est le terme de erreur, assumé être distribué selon une loi Normale  $N(0, \sigma^2)$ .

On souhaite donc estimer les coefficients du modèle:

$$\text{maxO3} = \beta_0 + \beta_1 * T12 + \beta_2 * T15 + \beta_3 * Ne12 + \beta_4 * Vx + \beta_5 * \text{maxO3v} + \epsilon$$

### 3.2 Estimer les coefficients de la régression multiple

```
1 reg_multiple = lm(maxo3 ~ T12 + T15 + Ne12 + Vx + maxO3v)
2 coef(reg_multiple)
```

Call:  
lm(formula = maxo3 ~ T12 + T15 + Ne12 + Vx + maxO3v)

Coefficients:  
(Intercept)      T12      T15      Ne12      Vx      maxO3v  
62.51384      0.08446      0.93121      -3.94166      0.31653      0.26146

L'équation de la droite de régression est donc:

$$\text{maxO3} = 62.51 + 0.08 * T12 + 0.93 * T15 + -3.94 * Ne12 + 0.31 * Vx + 0.26 * \text{maxO3v}$$

### 3.3 Intervalle de confiance

```
1 # Intervalle de confiance (à 95%) des coefficients
2 confint(reg_multiple)
```

A matrix: 6 × 2 of type dbl

	2.5 %	97.5 %
(Intercept)	32.27516866	92.7525159
T12	-2.48979050	2.6587052
T15	-1.55276369	3.4151849
Ne12	-5.98499300	-1.8983294
Vx	-0.02411048	0.6571717
maxO3v	0.07410465	0.4488183

### 3.4 Test d'hypothèse

```
1 summary(reg_multiple)
```

Call:

```
lm(formula = maxo3 ~ T12 + T15 + Ne12 + Vx + maxO3v)
```

Residuals:

Min	1Q	Median	3Q	Max
-25.673	-8.006	0.501	5.896	25.755

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	62.51384	15.00405	4.166	0.000142	***
T12	0.08446	1.27731	0.066	0.947581	
T15	0.93121	1.23252	0.756	0.453953	
Ne12	-3.94166	1.01388	-3.888	0.000338	***
Vx	0.31653	0.16902	1.873	0.067761	.
maxO3v	0.26146	0.09296	2.812	0.007318	**

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 13 on 44 degrees of freedom

Multiple R-squared: 0.7275, Adjusted R-squared: 0.6966

F-statistic: 23.5 on 5 and 44 DF, p-value: 2.026e-11

En utilisant une valeur seuil de 0,05, on rejette l'hypothèse nulle  $H_0: B_j = 0$  pour les variables Ne12 et maxO3v, car leurs P-values sont inférieures à 0,05.

Cela signifie que ces variables sont statistiquement significatives à un niveau de confiance de 95%.

Cependant, on ne peut pas rejeter l'hypothèse nulle pour les variables T12, T15 et Vx, car leurs P-values sont supérieures à 5%.

Cela signifie que ces variables ne sont pas statistiquement significatives à un niveau de confiance de 95 %.

### 3.5 Qualité du modèle

La qualité du modèle de la régression linéaire multiple est mesurée par le coefficient de détermination (R-Squared ou  $R^2$ ), qui se définit comme la part de variation dans la variable y qui est expliquée par des variations dans les variables explicatives (souvent exprimé en %).

Formellement:  $R^2 = 1 - \frac{SSR}{SST}$

où:

**SSR** représente la somme des carrés des résidus et **SST** la somme des écarts à la moyenne des valeurs observées.

Plus sa valeur est proche de 1, et plus l'adéquation entre le modèle et les données observées va être forte. Cependant, cette valeur est fortement influencée, entre autres, par le nombre de variables explicatives incluses dans la régression.

Le  $R^2_{ajust}$  (Adjusted R-Squared) va alors tenir compte de ce nombre et sera donc plus correct.

Formellement:  $R^2_{adjusted} = 1 - \frac{SSR/n-k}{SST/n-1}$ .

Alors,

Dans les résultats obtenus avec `summary(reg_mmultiple)` le `Rcarrest` : 0.7275 et le `Rcarrajusted` : 0.6966. Il est donc de 70%, ce qui est plus que correct.

Si on compare avec les valeurs trouvées en régression simple qui sont respectivement de 0.6151 et 0.6116, donc inférieur à 70%, on conclut que le modèle de régression multiple est plus **intéressant**.

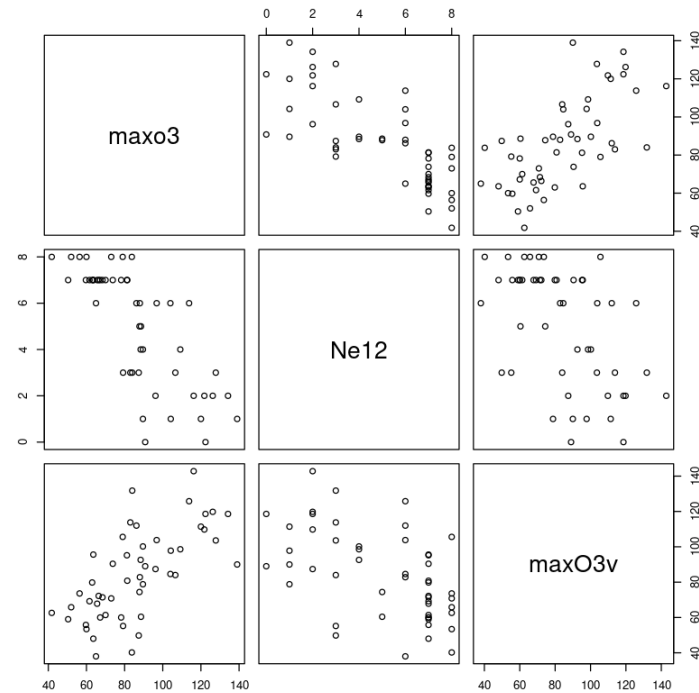
De même, cela indique que la variable T12 seule ne peut pas bien expliquer la variation de maxO3 et que les variables qu'on a rajoutées dans le cas multiple expliquent mieux la variation de maxO3.

### 3.6 Retirer les variables non significatives

D'après les résultats obtenus dans la question 4 par rapport à la P-value de chaque variable, on a constaté que T12, T15 et Vx sont des variables **non significatives**, on peut donc les **retirer**.

#### 3.6.1 visualisation des variables significatives

```
1 df_ozone = data.frame(maxO3, Ne12, maxO3v)
2 plot(df_ozone)
```



### 3.7 Créer le nouveau modèle

On nomme notre nouveau modèle: *nv-reg-mult*

```
1 nv_reg_mult = lm(maxo3 ~ Ne12 + maxO3v)
2 nv_reg_mult
```

Call:

```
lm(formula = maxo3 ~ Ne12 + maxO3v)
```

Coefficients:

(Intercept)	Ne12	maxO3v
84.551	-5.378	0.339

On redéfinit notre modèle comme suit:

$$\text{maxO3} = 84.55 - 5.38 * \text{Ne12} + 0.34 * \text{maxO3v}$$

### 3.8 Définir n et k

Dans notre cas:

- $n = 50$  observations ( $n$  reste comme il est),
- $k = 2$  variables explicatives; on le trouve avec ces deux lignes de code:

```
1 k = nv_reg_mult$rank-1
2 cat('Nombre de variables explicatives : k=', k)
```

### 3.9 Calculer les résidus studentisés

Les résidus standardisés sont une mesure importante dans la vérification des hypothèses de base des modèles linéaires. Si les résidus standardisés sont très éloignés de zéro pour certaines observations, cela peut indiquer une violation de ces hypothèses ou la présence de points aberrants dans les données.

La formule pour calculer les résidus studentisés est la suivante :

$$\text{résidus}_{\text{studentisés}} = \frac{\text{residu}}{(\text{racine}(\text{variance}_{\text{residu}} * (1 - \text{influence}_{\text{pointlevier}})))}$$

```
1 # Méthode 1
2 residuals = residuals(nv_reg_mult)
3 variance_residuals = var(residuals) * (length(residuals) - 1) / df.residual(nv_reg_mult)
4 residus_studentises1 = residuals / sqrt(variance_residuals*(1 - hatvalues(nv_reg_mult)))
5 residus_studentises1
```

1:	-1.18734398365277 2:	-0.552252457210504 3:	0.130757957167447 4:	0.152990520072306 5:
0.569200152473702 6:	-0.202356983303802 7:	2.24945889476824 8:	0.822729616954312 9:	
1.48086733944581 10:	-1.57981154041567 11:	-0.0126532526377764 12:	0.563300676503301 13:	
0.897637731835963 14:	1.81733802498235 15:	-0.655167003147869 16:	0.0322891866858626 17:	
1.53568664071732 18:	-0.00320424487622826 19:	0.366194500421134 20:	0.704106656017736 21:	
-1.26414502638652 22:	-0.37210044562978 23:	0.0282527396616604 24:	-1.875628505113 25:	
-0.620703562659155 26:	0.17124904650488 27:	-0.546503413634592 28:	-0.320675889088662 29:	
0.949363422701758 30:	-0.308941974131871 31:	0.165231622811397 32:	-2.24874078185593 33:	
-1.80911126085669 34:	-0.460920027690909 35:	-0.891682917795432 36:	-0.281902312838475 37:	
0.233996678922768 38:	-0.182762851978861 39:	0.727837035468222 40:	0.814903575495996 41:	
-0.473651503514237 42:	0.531721209933437 43:	0.787791753831578 44:	-0.819393890979457 45:	
1.71488700900626 46:	-0.440518055180621 47:	2.20210498565613 48:	-0.759604433797985 49:	
-1.24003021987137 50:	-0.615896497591558			

La commande `rstandard()` en R calcule les résidus standardisés, également appelés résidus studentisés. Ces résidus mesurent la différence entre les valeurs observées et les valeurs prédites par un modèle statistique, divisée par l'estimation de l'écart-type de l'erreur.

```
1 # Méthode 2
2 residus_studentises <- rstandard(nv_reg_mult)
3 residus_studentises
```

1:	-1.18734398365277 2:	-0.552252457210504 3:	0.130757957167447 4:	0.152990520072306 5:
0.569200152473702 6:	-0.202356983303802 7:	2.24945889476824 8:	0.822729616954312 9:	
1.48086733944581 10:	-1.57981154041567 11:	-0.0126532526377764 12:	0.563300676503301 13:	
0.897637731835963 14:	1.81733802498235 15:	-0.655167003147869 16:	0.0322891866858626 17:	
1.53568664071732 18:	-0.00320424487622826 19:	0.366194500421134 20:	0.704106656017736 21:	
-1.26414502638652 22:	-0.37210044562978 23:	0.0282527396616604 24:	-1.875628505113 25:	
-0.620703562659155 26:	0.17124904650488 27:	-0.546503413634592 28:	-0.320675889088662 29:	
0.949363422701758 30:	-0.308941974131871 31:	0.165231622811397 32:	-2.24874078185593 33:	
-1.80911126085669 34:	-0.460920027690909 35:	-0.891682917795432 36:	-0.281902312838475 37:	
0.233996678922768 38:	-0.182762851978861 39:	0.727837035468222 40:	0.814903575495996 41:	
-0.473651503514237 42:	0.531721209933437 43:	0.787791753831578 44:	-0.819393890979457 45:	
1.71488700900626 46:	-0.440518055180621 47:	2.20210498565613 48:	-0.759604433797985 49:	
-1.24003021987137 50:	-0.615896497591558			

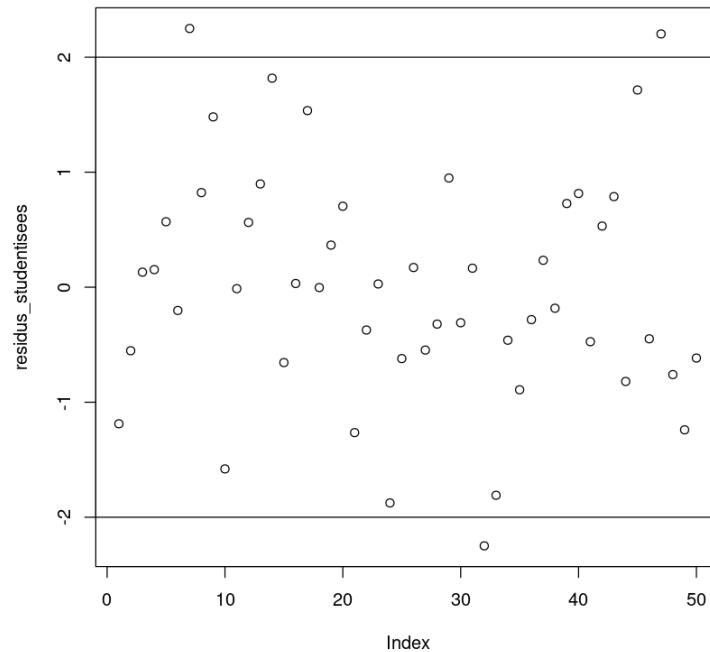
### 3.10 Valeurs aberrantes

Afin d'identifier d'éventuelles valeurs aberrantes dans les données, on utilise la représentation graphique suivante:

```

1 plot(residus_studentises)
2 abline(-2,0)
3 abline(2,0)

```



## 4 Conclusion

Cette étude avait pour objectif d'étudier les facteurs influençant les concentrations d'ozone troposphérique à partir des données météorologiques et de pollution atmosphérique du jeu ozone.

Différents modèles de régression linéaire simple et multiple ont été ajustés, en prenant la concentration d'ozone maxO3 comme variable dépendante à expliquer par différents prédictors environnementaux.

La régression linéaire simple a permis de quantifier les effets d'un facteur clé qui est la température à midi (variable T12), qui ressort comme un déterminant important de niveau d'ozone.

La régression multiple incluant plusieurs prédictors météorologiques (vent, humidité, pression, etc) se révèle plus précise, avec un R2 ajusté de 70%. Les validations sur les résidus confirment la qualité de l'ajustement.

Ces résultats illustrent l'intérêt des méthodes de régression linéaire pour modéliser les phénomènes complexes de pollution de l'air à partir de données environnementales. Les modèles développés permettent d'expliquer et de prédire les taux d'ozone de manière fiable dans ce contexte d'étude.

Cette analyse ouvre des perspectives pour affiner ces modèles statistiques et les appliquer à d'autres contextes géographiques et climatiques.