

**Exercice 1** Lors d'un mémoire de DEA (promotion 2000-2001), Frédéric Rey et Alexandre Turpin ont étudié le taux de chômage. Ils ont recueilli un jeu de données (présenté à la fin du devoir) constitué de  $n = 34$  observations annuelles (de 1960 à 1993). Voici une description des variables fournie par les auteurs du mémoire :

- **an** : année
- **chom** : taux de chômage.
- **txpib** : taux de variation du produit intérieur brut (pib) représentant le taux de croissance de l'économie.
- **deppub** : part des dépenses publiques par rapport au pib, qui peut ainsi représenter le degré d'intervention de l'état dans l'économie.
- **pfisc** : pressions fiscales, pour voir si une imposition trop importante des entreprises nuit à leur embauche et donc au niveau du chômage.
- **salva** : la part des salaires par rapport à la valeur ajoutée permettant de connaître l'influence du coût sur l'embauche.
- **infl** : taux d'inflation afin de vérifier la relation inverse entre le chômage et l'inflation définie par la courbe de Philips.

### Partie I : modèle à un seul régresseur

On envisage un modèle linéaire expliquant la variable **chom** en fonction de la seule variable **txpib**. On tente une modélisation linéaire du type :

$$(\text{chom})_i = \beta_0 + \beta_1(\text{txpib})_i + \varepsilon_i, \quad i = 1, \dots, 34$$

**Question 1:** Pourquoi les paramètres  $\beta_0$  et  $\beta_1$  ne sont pas calculables ?

**Réponse**

Pour pouvoir évaluer les paramètres  $\beta_0$  et  $\beta_1$  il faudrait "en théorie" une infinité de données. Ces paramètres ne sont donc pas évaluables mais simplement estimables à partir d'un nombre fini de données.

Fin

**Question 2:** On s'intéresse tout naturellement à l'estimation des paramètres  $\beta_0$  et  $\beta_1$ . Déterminez les estimations obtenues par la méthode des moindres carrés.

On rappelle à titre indicatif que

- $\text{var}(\mathbf{x}) = \overline{x^2} - \bar{x}^2$
- $\text{cov}(\mathbf{x}, \mathbf{y}) = \overline{x \times y} - \bar{x} \times \bar{y}$ .

```

1 | > mean(txpib)
2 | [1] 3.488235
3 | > mean(txpib^2)
4 | [1] 16.06412
5 | > mean(chom)
6 | [1] 5.458824
7 | > mean(chom^2)
8 | [1] 42.03294
9 | > mean(txpib*chom)
10| [1] 13.85794

```

**Réponse**

En notant  $\mathbf{y} = \text{chom}$  et  $\mathbf{x} = (\mathbf{1}, \mathbf{x}^{(1)} := \text{txpib})$ ,

$$\widehat{\beta}_1(\mathbf{y}|\underline{\mathbf{x}}) = \frac{\text{cov}(\mathbf{y}, \mathbf{x}^{(1)})}{\text{var}(\mathbf{x}^{(1)})} = \frac{13.85794 - (3.488235) * (5.458824)}{16.06412 - 3.488235^2} \simeq -1.33041$$

$$\widehat{\beta}_0(\mathbf{y}|\underline{\mathbf{x}}) = \bar{\mathbf{y}} - \widehat{\beta}_1(\mathbf{y}|\underline{\mathbf{x}}) \overline{\mathbf{x}^{(1)}} = 5.458824 - (-1.33041) * 3.488235 \simeq 10.09961.$$

Vérification en R :

```

1 | > beta1Est<-(mean(txpib*chom)-mean(txpib)*mean(chom))/
2 | + (mean(txpib^2)-mean(txpib)^2)
3 | > beta1Est
4 | [1] -1.33041
5 | > mean(chom) - beta1Est*mean(txpib)
6 | [1] 10.09961

```

Fin

**Question 3:** Déterminez le coefficient de détermination linéaire ( $R^2$ ) (puis le coefficient de corrélation linéaire ( $R$ )) entre *txpib* et *chom*, et donnez-en une interprétation.

**Réponse**

Dans le cadre de la régression simple,  $R^2$  vaut le carré du coefficient de corrélation linéaire :

$$R^2 = \frac{\text{cov}(\mathbf{y}, \mathbf{x}^{(1)})^2}{\text{var}(\mathbf{y}) \times \text{var}(\mathbf{x}^{(1)})} = \frac{(13.85794 - (3.488235) * (5.458824))^2}{(16.06412 - 3.488235^2)(42.03294 - 5.458824^2)} \simeq 0.5637051.$$

D'où  $\text{corr}(\mathbf{x}, \mathbf{y}) = -\sqrt{0.5637051} \simeq -0.750803$  qui a en particulier le même signe que  $\widehat{\beta}_1(\mathbf{y}|\underline{\mathbf{x}})$ . Plus  $|\text{corr}(\mathbf{x}, \mathbf{y})|$  est proche de 1 et plus la dispersion des points autour de la droite ajustée est faible.

Fin

**Question 4:** Rappelez brièvement à quoi correspondent chacune des quatre colonnes de la matrice "Coefficients" de la sortie ci-dessous. Retrouvez les résultats des deux questions précédentes.

```

1 | > summary(lm(chom~txpib))
2 |
3 | Call:
4 | lm(formula = chom ~ txpib)
5 |
6 | Residuals:
7 |      Min       1Q   Median       3Q      Max
8 | -6.3987 -1.5732 -0.2337  1.4000  5.6212
9 |
10 | Coefficients:
11 |             Estimate Std. Error t value Pr(>|t|)
12 | (Intercept)  10.0996     0.8293   12.18 1.48e-13 ***
13 | txpib        -1.3304     0.2069   -6.43 3.14e-07 ***
14 | ---
15 | Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16 |
17 | Residual standard error: 2.381 on 32 degrees of freedom
18 | Multiple R-squared:  0.5637,    Adjusted R-squared:  0.5501
19 | F-statistic: 41.34 on 1 and 32 DF,  p-value: 3.144e-07
20 |
21 | > sqrt(0.5637051)
22 | [1] 0.750803

```

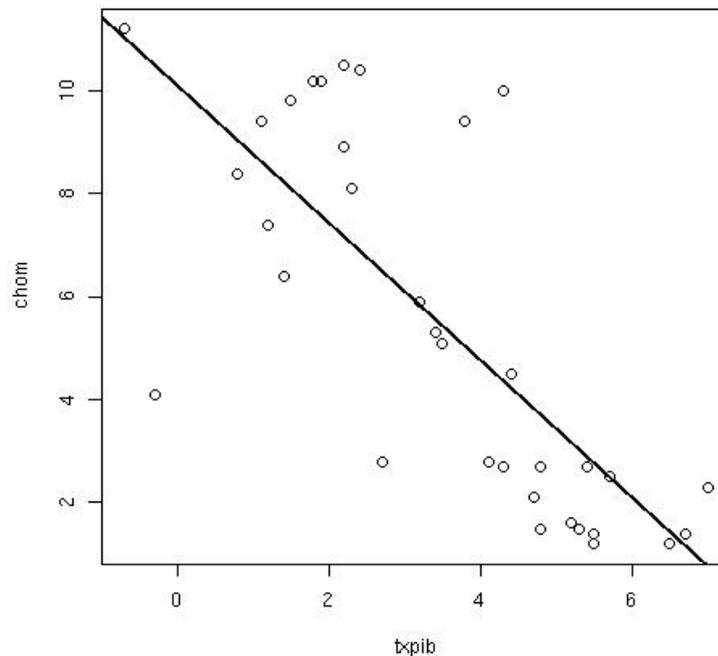
**Réponse**

On retrouve les résultats des deux questions précédentes. En particulier, la colonne **Estimate** fournit les estimations des paramètres  $\beta_0$  et  $\beta_1$  et **Multiple R-squared** le  $R^2$ . Notons que la dernière colonne du tableau **Coefficients** fournit les p-valeurs des tests de significativité locales

des paramètres  $\beta_0$  et  $\beta_1$ . La dernière laisse en particulier apparaître que la variable `txpib` semble apporter de l'information dans l'explication de la variable `chom`.

Fin

**Question 5:** Sur le graphique ci-dessous reportez la droite ajustée (même approximativement) et illustrez la notion de valeur ajustée et de résidu.



Les notions de résidu et de valeur ajustée sont illustrés dans le polycopié p.23.

**Question 6:** Peut-on penser au vu des données que la variable `txpib` apporte de l'information pour expliquer la variable `chom` (indication : fournir la  $p$ -valeur associée puis conclure.)

**Réponse**

Pour pouvoir accepter  $H_1 : \beta_1 \neq 0$ , i.e. le régresseur `txpib` apporte de l'information pour expliquer la variable `chom`, le risque à encourir est de l'ordre de  $3.144288e - 07$ .

Fin

## Partie II : modèle linéaire multiple

On envisage un modèle linéaire multiple expliquant la variable `chom` en fonction de tous les régresseurs du jeu de données (exceptée la variable `an`).

**Question 1:** A la vue de la matrice de corrélation ci-après, quels sont les régresseurs qui vous semblent être les plus explicatifs ?

```

1 > cor(chomage[-1])
2      chom      txpib      deppub      pfisc      salva      infl
3 chom      1.0000000 -0.7508029  0.978885951  0.97923991 -0.167503533 -0.04814577
4 txpib    -0.75080295  1.0000000  -0.780822861 -0.75374536 -0.105374429 -0.30407625
5 deppub   0.97888595  -0.7808229  1.000000000  0.99212248  -0.007785484  0.06750973
6 pfisc    0.97923991  -0.7537454  0.992122482  1.00000000  -0.035551565  0.06640145
7 salva   -0.16750353  -0.1053744  -0.007785484 -0.03555157  1.000000000  0.70322004
8 infl    -0.04814577  -0.3040763  0.067509734  0.06640145  0.703220044  1.00000000

```

**Réponse**

A la vue de la matrice de corrélation, les régresseurs les plus explicatifs de la variabilité de la variable `chom` semblent être dans l'ordre `pfisc`, `deppub` et `txpib` (en ne tenant compte que des

régresseurs ayant un coefficient de corrélation avec chom en valeur absolue supérieur à 30.0%).  
Fin

**Question 2:** *Interprétez la sortie ci-dessous, en particulier les p-valeurs des tests de significativité locale, le  $R^2$ .*

```
1 > summary(lm(chom~txpib+deppub+pfisc+salva+infl))
2
3 Call:
4 lm(formula = chom ~ txpib + deppub + pfisc + salva + infl)
5
6 Residuals:
7     Min       1Q   Median       3Q      Max
8 -0.84262 -0.26630  0.05442  0.27937  1.33259
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) -5.80283     3.86147  -1.503  0.144098
13 txpib       -0.07188     0.07753  -0.927  0.361820
14 deppub      0.35855     0.12702   2.823  0.008665 **
15 pfisc       0.22751     0.14324   1.588  0.123432
16 salva      -0.19195     0.05079  -3.779  0.000757 ***
17 infl       -0.03131     0.03961  -0.791  0.435873
18 ---
19 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
20
21 Residual standard error: 0.469 on 28 degrees of freedom
22 Multiple R-squared:  0.9852,    Adjusted R-squared:  0.9826
23 F-statistic: 372.7 on 5 and 28 DF,  p-value: < 2.2e-16
```

**Réponse** \_\_\_\_\_

Les régresseurs *salva* et *deppub* sont significatifs au seuil de 5%. Enfin, le pouvoir explicatif de ce modèle, mesuré par la part de variance  $R^2$  expliquée par celui-ci, est particulièrement bon puisqu'il est de l'ordre de 98.51961%.

Fin

**Question 3:** *Rappelez les effets indésirables sur les tests de significativité locale s'il y a colinéarité entre les régresseurs.*

**Réponse** \_\_\_\_\_

La colinéarité entre régresseurs peut entraîner artificiellement une augmentation de l'erreur standard des paramètres des régresseurs et par voie de fait une diminution en valeur absolue de la statistique de test de significativité locale et donc une augmentation de la p-valeur de ce même test.

Fin

**Question 4:** *A la lumière de la matrice de corrélation associée au jeu de données, peut-on suspecter de la colinéarité entre les régresseurs ?*

**Réponse** \_\_\_\_\_

Fin

**Question 5:** *Rappelez la définition du VIF, et son interprétation générale. Ensuite, interprétez-les quant au jeu de données étudié.*

```
1 > vif(lm(chom~txpib+deppub+pfisc+salva+infl))
2     txpib  deppub   pfisc   salva   infl
3  3.620834 91.775821 81.135288  2.384586  2.712915
```

**Réponse** \_\_\_\_\_

Par définition,  $VIF_j = \frac{1}{1-R_j^2}$  où  $R_j^2$  est le coefficient de détermination multiple au carré lorsque l'on régresse le  $j$ -ème régresseur sur tous les autres. Par ailleurs, on sait d'après le cours que l'erreur standard de l'estimateur de  $\beta_j$  est proportionnelle au  $VIF_j$ . Par conséquent, plus le  $j$ -ème régresseur est corrélé avec les autres, plus  $R_j^2$  est proche de 1, plus l'erreur standard de l'estimateur de  $\beta_j$  sera grande et plus il sera difficile de montrer que ce régresseur apporte de l'information. Les VIFs de `pfisc` et `deppub` sont relativement importants, ici supérieurs à 5. Quand il existe un VIF relativement important, cela traduit une forte colinéarité entre régresseurs.

Fin

**Question 6:** (*Relation avec la matrice de corrélation*) Justifier l'ordre de grandeur des VIFs des covariables `deppub` et `pfisc` en utilisant l'instruction suivante.

```
1 > 1/(1-(0.992122482)^2)
2 [1] 63.72276
```

**Réponse**

Il est connu que lorsque l'on ajoute des régresseurs le coefficient  $R^2$  augmente nécessairement. Ainsi par exemple, le  $R^2$  obtenu en régressant `deppub` sur tous les autres ou `pfisc` sur tous les autres est nécessairement supérieurs au  $R^2$  obtenu en régressant simplement `deppub` sur `pfisc` (de l'ordre de 0.992122482<sup>2</sup>). Par conséquent les VIF associées à ces deux régresseurs sont nécessairement supérieurs à 63.72276.

Fin

**Question 7:** *Quelle est la stratégie qui a été adoptée dans la série d'instructions ci-dessous ? A la dernière étape, précisez l'équation du modèle sélectionné et analysez brièvement les sorties.*

```
1 > ## Etape 1
2 > summary(lm(chom~txpib+deppub+pfisc+salva))
3
4 Call:
5 lm(formula = chom ~ txpib + deppub + pfisc + salva)
6
7 Residuals:
8     Min       1Q   Median       3Q      Max
9 -0.80961 -0.31312 -0.00256  0.26502  1.41147
10
11 Coefficients:
12             Estimate Std. Error t value Pr(>|t|)
13 (Intercept) -3.93872     3.03798  -1.296  0.20503
14 txpib        -0.04277     0.06779  -0.631  0.53302
15 deppub         0.39774     0.11619   3.423  0.00186 **
16 pfisc          0.18756     0.13315   1.409  0.16959
17 salva        -0.22177     0.03379  -6.563 3.44e-07 ***
18 ---
19 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
20
21 Residual standard error: 0.4659 on 29 degrees of freedom
22 Multiple R-squared:  0.9849,    Adjusted R-squared:  0.9828
23 F-statistic: 471.8 on 4 and 29 DF,  p-value: < 2.2e-16
24
25 > vif(lm(chom~txpib+deppub+pfisc+salva))
26     txpib  deppub  pfisc  salva
27  2.804297 77.798129 71.033914  1.069126
28 > ## Etape 2
29 > summary(lm(chom~deppub+pfisc+salva))
30
31 Call:
32 lm(formula = chom ~ deppub + pfisc + salva)
```

```

33
34 Residuals:
35      Min       1Q   Median       3Q      Max
36 -0.79521 -0.27194 -0.02836  0.26418  1.42664
37
38 Coefficients:
39             Estimate Std. Error t value Pr(>|t|)
40 (Intercept) -4.63681    2.80081  -1.656 0.108244
41 deppub      0.42511    0.10670   3.984 0.000399 ***
42 pfisc       0.16761    0.12804   1.309 0.200456
43 salva      -0.21907    0.03318  -6.603 2.62e-07 ***
44 ---
45 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
46
47 Residual standard error: 0.4612 on 30 degrees of freedom
48 Multiple R-squared:  0.9847,    Adjusted R-squared:  0.9831
49 F-statistic: 641.8 on 3 and 30 DF,  p-value: < 2.2e-16
50
51 > vif(lm(chom~deppub+pfisc+salva))
52      deppub      pfisc      salva
53 66.949894 67.030557  1.051973
54 > ## Etape 3
55 > summary(lm(chom~deppub+salva))
56
57 Call:
58 lm(formula = chom ~ deppub + salva)
59
60 Residuals:
61      Min       1Q   Median       3Q      Max
62 -0.73827 -0.36129 -0.03607  0.29649  1.37843
63
64 Coefficients:
65             Estimate Std. Error t value Pr(>|t|)
66 (Intercept) -2.83284    2.46623  -1.149   0.259
67 deppub      0.56373    0.01319  42.741 < 2e-16 ***
68 salva      -0.22872    0.03272  -6.990 7.6e-08 ***
69 ---
70 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
71
72 Residual standard error: 0.4665 on 31 degrees of freedom
73 Multiple R-squared:  0.9838,    Adjusted R-squared:  0.9827
74 F-statistic: 940.2 on 2 and 31 DF,  p-value: < 2.2e-16
75
76 > vif(lm(chom~deppub+salva))
77      deppub      salva
78 1.000061 1.000061
79

```

Réponse \_\_\_\_\_

Fin

**Question 8:** *A partir de cette question, nous ne considérerons que le modèle final. Peut-on montrer au vu des données que le paramètre  $\beta_1 < 1$  au seuil de 5% ?*

```

1 > (0.56373-(1))/0.01319
2 [1] -33.07582

```

Réponse \_\_\_\_\_

D'après ce qui précède, la statistique du test  $\mathbf{H}_1 : \beta_1 < 1$  évaluée sur les données correspond justement au calcul fourni et vaut donc approximativement  $-33.07582$ . Sous  $\mathbf{H}_0$  on sait également que la statistique de test sur les futures données suit approximativement une loi  $\mathcal{N}(0, 1)$ .

Par conséquent, puisque  $\widehat{\delta}_{\beta_1,1}(\mathbf{y}|\mathbf{x}) \simeq -33.07582 < \delta_{lim,5\%}^- \stackrel{R}{=} \text{qnorm}(.05) \simeq -1.644854$ , on peut plutôt penser avec un risque de 5% que l'assertion d'intérêt  $\mathbf{H}_1 : \beta_1 < 1$  est vraie.

Fin

**Question 9:** *A partir de l'instruction R suivante, que peut-on avancer au vu des données comme assertion(s) d'intérêt au seuil 5% ?*

```
1 > pnorm((-0.22872-(-.3))/0.03272)
2 [1] 0.985315
```

Réponse

L'instruction fournie correspond à la p-valeur (gauche) du test  $\mathbf{H}_1 : \beta_2 < -.3$ . Par conséquent puisque la p-valeur droite associée au test  $\mathbf{H}_1 : \beta_2 > -.3$  qui vaut approximativement 1.468376% est inférieure à 5%, on peut plutôt penser que cette assertion est vraie. Notons, qu'il est possible de penser que  $\mathbf{H}_1 : \beta_2 \neq -.3$  est vraie car la p-valeur de ce test vaut approximativement  $2 \times 1.468376\% = 2.936751\%$ .

Fin

**Question 10:** *En vous aidant de l'instruction R ci-dessous, proposez un intervalle de confiance à 95% pour le paramètre  $\beta_2$  et interprétez-le (via l'approche expérimentale). Quelle relation y-a-t-il entre cet intervalle et le test de significativité locale du paramètre  $\beta_2$  ?*

```
1 > -0.22872+c(-1,1)*qnorm(0.975)*0.03272
2 [1] -0.2928500 -0.1645900
```

Réponse

Fin

**Question 11:** *Supposons que l'on ne connaisse pas la valeur de chom en 1993. Pourriez-vous prévoir sa valeur, calculer un intervalle de prévision au niveau 95% ? Quelle était la valeur observée de chom en 1993 et est-ce surprenant ?*

```
1 > xTau <- data.frame(chom=11.2,deppub=52.2,salva=68.6)
2 > predict(lm(chom~deppub+salva),xTau,interval="prediction")
3         fit         lwr         upr
4 1 10.9039  9.872652 11.93515
```

Réponse

Fin

Jeu de données :

```
1 > chomage
2       an chom txpib deppub pfisc salva infl
3 1 1960  1.4   5.5  34.6  34.9  72.8  3.4
4 2 1961  1.2   5.5  35.7  36.2  72.8  3.4
5 3 1962  1.4   6.7  37.0  36.3  72.8  4.7
6 4 1963  1.5   5.3  37.8  37.1  72.8  6.4
7 5 1964  1.2   6.5  38.0  38.0  72.8  4.1
8 6 1965  1.5   4.8  38.4  38.4  73.3  2.7
9 7 1966  1.6   5.2  38.5  38.4  73.3  2.9
10 8 1967  2.1   4.7  39.0  38.2  73.3  3.2
11 9 1968  2.7   4.3  40.3  38.8  73.3  4.2
12 ...
13 25 1984  9.8   1.5  52.5  49.8  75.6  7.3
14 26 1985 10.2   1.8  52.7  49.9  74.6  5.8
15 27 1986 10.4   2.4  52.2  49.4  72.1  5.3
16 28 1987 10.5   2.2  51.7  49.8  71.3  3.0
17 29 1988 10.0   4.3  50.8  49.2  70.2  3.1
18 30 1989  9.4   3.8  49.8  48.7  69.1  3.5
```

19	31	1990	8.9	2.2	50.3	48.9	69.6	3.1
20	32	1991	9.4	1.1	50.6	48.7	69.6	3.1
21	33	1992	10.2	1.9	51.3	48.5	68.8	2.9
22	34	1993	11.2	-0.7	52.2	49.0	68.6	2.9