

## EXAMEN FINAL, ACT6420, HIVER 2012

ARTHUR CHARPENTIER

(énoncé avec quelques éléments de correction)

Les calculatrices sont autorisées. Tous les documents sont, en revanche, interdits.

Dans les feuilles qui suivent, il y a

- 10 questions générales sur les modèles de régression sur données individuelles
- 10 questions portant sur la modélisation des matchs de basket-ball (à partir des sorties disponibles dans l'Annexe A, 9 pages)
- 10 questions générales sur les modèles de séries temporelles
- 10 questions portant sur la modélisation et la prévision (à partir des sorties disponibles dans l'Annexe B, 14 pages)

Pour chaque question, quatre réponses sont proposées, une seule est valide, et vous ne devez en retenir qu'une (au maximum), et reportez votre réponse sur la feuille jointe. Choisir l'affirmation la "*moins juste*" signifie que trois affirmations sont valides et qu'une est fausse. Il faut identifier la fausse. Choisir l'affirmation la "*plus juste*" signifie que trois affirmations sont fausses et qu'une est valide. Il faut identifier cette dernière.

Pour le décompte des points,

- vous gagnez 2.5 points par bonne réponse
- vous perdez ne perdez aucun point par mauvaise réponse

Aucune justification n'est demandée. Votre note finale est le total des points.

**Formulaire:** si  $Z$  suit une loi  $\mathcal{N}(0, 1)$ ,  $\mathbb{P}(Z > 0.25) = 40\%$ ,  $\mathbb{P}(Z > 0.52) = 30\%$ ,  $\mathbb{P}(Z > 0.84) = 20\%$ ,  $\mathbb{P}(Z > 1.28) = 10\%$ ,  $\mathbb{P}(Z > 1.64) = 5\%$  et  $\mathbb{P}(Z > 1.96) = 2.5\%$ .

Vous avez 3 heures.

## 1. RÉGRESSIONS SUR DONNÉES INDIVIDUELLES

Considérons un modèle de régression, de la forme

$$Y_i = \beta_0 + \beta_1 X_{1,i} + \varepsilon_i \quad (1.1)$$

où  $Y$  est le poids d'un individu, en kilogrammes, et  $X$  sa taille, en mètres. Ce modèle sera appelé (1.1) dans la suite.

**Question 1.** Le modèle de régression contient un terme d'erreur  $\varepsilon_i$  pour plusieurs raisons. Parmi les affirmations suivantes, une seule n'est jamais valide. Laquelle ?

- A. parce qu'une partie du comportement de  $Y$  n'a pas pu être capturé par  $X_1$ ,
- B. parce que le modèle linéaire n'est (généralement) qu'une approximation de la vraie relation liant  $Y$  à la variable  $X_1$ ,
- C. parce que l'estimateur par moindres carrés ne donne qu'une approximation numérique de l'estimateur du maximum de vraisemblance, que l'on corrige à l'aide du terme d'erreur,
- D. parce qu'on souhaite tenir compte d'une éventuelle marge d'erreur dans l'observation de la variable  $Y$ .

*l'estimateur par maximum de vraisemblance coïncide avec l'estimateur par moindres carrés, dans le cas où les résidus sont Gaussiens. Il n'y a pas d'approximation (numérique) dans le calcul de  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ . Les trois autres réponses pouvant être des manières d'expliquer les résidus dans la régression*

Pour estimer les coefficients  $\beta_0$  et  $\beta_1$ , l'approche naturelle est de considérer l'estimateur des moindres carrés ordinaires.

**Question 2.** Mais auparavant, il convient de faire des hypothèses afin d'avoir construit ainsi les meilleurs estimateurs linéaires non biaisés de  $\beta_0$  et  $\beta_1$ . Laquelle de ces hypothèses sur les résidus n'est pas nécessaire (et ne fait pas partie des *hypothèses de base* telles qu'énoncées dans le cours)

- A. les résidus doivent être centrés,
- B. les résidus sont de variance identique,

- C. les résidus doivent suivre une loi normale,
- D. les résidus doivent être non corrélés, et non corrélés avec  $X_1$ .

l'hypothèse de normalité n'est pas nécessaire pour établir ces propriétés. Elle le devient si on veut faire des tests, et que l'on a peu d'observations.

Une fois cette hypothèse faite, on calcule par la méthode des moindres carrés  $\hat{\beta}_0$  et  $\hat{\beta}_1$ , à partir de  $n$  observations. On peut alors construire la série des erreurs,  $\hat{\varepsilon}_i = Y_i - \hat{Y}_i$  où  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1,i}$ ,

**Question 3.** Laquelle des affirmations suivante n'est pas vraie

- A. la droite de régression  $y = \hat{\beta}_0 + \hat{\beta}_1 x$  passe forcément par le point  $(\bar{X}_1, \bar{Y})$  où, classiquement,  $\bar{X}_1 = \frac{1}{n} \sum_{i=1}^n X_{1,i}$  et  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ ,
- B. la somme des résidus est nulle  $\sum_{i=1}^n \hat{\varepsilon}_i = 0$ ,
- C. la somme des valeurs observées des  $Y_i$  et la somme des valeurs prédites  $\hat{Y}_i$  doivent être égales,  $\sum_{i=1}^n \hat{Y}_i = \sum_{i=1}^n Y_i$ ,
- D. il y a autant de points au dessus que de points en dessous de la droite de régression.

il faudrait faire une régression médiane pour séparer le nuage en deux, ou  $L_1$ , ou en considérant l'estimateur pour moindre valeur absolue des résidus, i.e.  $\sum |Y_i - (\beta_0 + \beta_1 X_i)|$ . On notera que B. et C. sont équivalentes: si une était fausse, l'autre serait fausse aussi.

**Question 4.** Le  $R^2$  (appelé coefficient d'ajustement) mesure une quantité parmi les quatre suivantes. Laquelle ?

- A. la corrélation (empirique) entre les  $X_{1,i}$  et  $Y_i$
- B. la covariance (empirique) entre  $X_{1,i}$  et  $Y_i$
- C. le ratio de la somme des  $\hat{\varepsilon}_i^2$  sur la somme des  $Y_i^2$
- D. le ratio de la somme des  $\hat{Y}_i^2$  sur la somme des  $Y_i^2$

C'est la définition

$$R^2 = \frac{\sum \hat{Y}_i^2}{\sum Y_i^2} = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum Y_i^2}$$

De plus, le  $R^2$  est - comme son nom le laisse entendre - positif, et donc il est difficile de l'imaginer étant une corrélation ou une covariance (ce qui excluait les premières affirmations).

Supposons que l'on change d'unité pour la mesure de  $Y$ , en donnant désormais le poids en livres (pour rappels, une livre vaut 0,4535 kilogramme). On notera  $\tilde{Y}$  la nouvelle variable. On a désormais un modèle  $\tilde{Y}_i = \tilde{\beta}_0 + \tilde{\beta}_1 X_{1,i} + \tilde{\varepsilon}_i$ .

**Question 5.** Tous les “coefficients” vont changer. Sauf un, lequel ?

- A. la constante de la régression, i.e.  $\tilde{\beta}_0 = \beta_0$
- B. la pente de la régression, i.e.  $\tilde{\beta}_1 = \beta_1$
- C. le  $R^2$ , coefficient d'ajustement  $\tilde{R}^2 = R^2$
- D. la variance des résidus,  $\tilde{\sigma}^2 = \sigma^2$

On devrait avoir  $\tilde{Y} = aY$ , et donc

$$aY_i = \tilde{\beta}_0 + \tilde{\beta}_1 X_{1,i} + \tilde{\varepsilon}_i$$

ou de manière équivalente

$$Y_i = \underbrace{\frac{\tilde{\beta}_0}{a}}_{\beta_0} + \underbrace{\frac{\tilde{\beta}_1}{a}}_{\beta_1} X_{1,i} + \underbrace{\frac{\tilde{\varepsilon}_i}{a}}_{\varepsilon_i} .$$

Aussi, la constante va changer, la pente aussi, et la variance sera elle aussi multipliée par  $a$ . Ces trois quantités vont donc changer. Le  $R^2$  est un ratio, dont le numérateur et le dénominateur sont subir le même chagement d'échelle, et donc le  $R^2$  sera inchangé.

À partir d'un jeu de données, on a estimé le modèle suivant  $Y = -100 + 100X_1 + \varepsilon$ , avec  $\text{Var}(\varepsilon) = 25$ .

**Question 6.** Une de ces affirmations n'est pas correcte,

- A. un individu qui pèse 70 kg pour 180 cm est léger, relativement à sa taille,
- B. un individu qui pèse 75 kg pour 170 cm est lourd, relativement à sa taille,
- C. aucune personne de moins d'un mètre ne faisait partie de l'échantillon
- D. en moyenne, un individu de 2 mètres pèse 100 kg

pour la première affirmation, le modèle prédit 80 kg, donc 70 kg peut-être vu comme léger. Pour la seconde, le modèle prédit 70 kg, donc 75 kg peut-être vu comme lourd. Enfin, pour la dernière, le modèle prédit 100 kg pour 200 cm, effectivement. C'est donc la troisième affirmation qui n'est pas correcte. En fait, on aura du mal à proposer une prédiction valide pour une personne mesurant moins d'un mètre; mais rien n'empêche de des individus de moins d'un mètre aient figuré dans la base.

Mais avant d'en tirer des conclusions, on nous suggère de faire un test de Student.

**Question 7.** Le test de Student peut être utilisé (une seule affirmation est correcte)

- A. pour tester si la constante  $\beta_0$  et la pente  $\beta_1$  sont de signe contraire
- B. pour tester si la constante  $\beta_0$  et la pente  $\beta_1$  sont opposés l'un de l'autre
- C. pour tester si la constante  $\beta_0$  est nulle ou pas,
- D. pour tester si la variance des résidus  $\sigma^2$  est nulle ou pas

Le test de Student est utilisé pour tester une hypothèse  $H_0 : \beta_k = 0$  contre l'hypothèse alternative  $H_1 : \beta_k \neq 0$ .

Au lieu de changer l'unité de  $Y$ , on va décaler les données. On va noter désormais  $\tilde{X}_1$  la taille au delà d'un mètre, i.e.  $\tilde{X}_1 = X_1 - 1$ . On a désormais un modèle

$$Y_i = \tilde{\beta}_0 + \tilde{\beta}_1 \tilde{X}_{1,i} + \tilde{\varepsilon}_i.$$

**Question 8.** Tous les "coefficients" seront identiques, sauf un. Lequel ?

- A. la constante de la régression, i.e.  $\tilde{\beta}_0 \neq \beta_0$
- B. la pente de la régression, i.e.  $\tilde{\beta}_1 \neq \beta_1$
- C. le  $R^2$ , coefficient d'ajustement  $\tilde{R}^2 \neq R^2$
- D. la variance des résidus,  $\tilde{\sigma}^2 \neq \sigma^2$

On devrait avoir  $\tilde{X} = X - a$ , et donc

$$Y_i = \tilde{\beta}_0 + \tilde{\beta}_1 [X_{1,i} - a] + \tilde{\varepsilon}_i$$

ou de manière équivalente

$$Y_i = \underbrace{\tilde{\beta}_0 - a\tilde{\beta}_1}_{\beta_0} + \underbrace{\tilde{\beta}_1}_{\beta_1} X_{1,i} + \underbrace{\tilde{\varepsilon}_i}_{\varepsilon_i}.$$

Aussi, la constante va changer. En revanche, la pente ne va pas changer, et les résidus étant inchangés, leur variance sera identique. Ces deux quantités vont donc être inchangées. Le  $R^2$  ne changera pas non plus car les  $Y$  sont inchangés, et les résidus aussi.

En fait, le modèle a été estimé sur 59 personnes. On a ainsi obtenu

$$\hat{\beta}_0 = -96.63 \text{ et } \hat{V}(\hat{\beta}_0) = 582.74 \quad \hat{\beta}_1 = 98.10 \text{ et } \hat{V}(\hat{\beta}_1) = 165.89 \quad \hat{\sigma} = 6.085$$

**Question 9.** Laquelle des affirmations suivantes est correcte (une seule l'est)

- A.  $\beta_0$  et  $\beta_1$  ne sont pas significativement non-nuls
- B.  $\beta_0$  n'est pas significativement non-nul, alors que  $\beta_1$  est significativement non-nul
- C.  $\beta_0$  est significativement non-nul, alors que  $\beta_1$  n'est pas significativement non-nul
- D.  $\beta_0$  et  $\beta_1$  sont significativement non-nuls.

Il suffit de calculer les statistiques de Student  $t_k = \hat{\beta}_k / \sqrt{\hat{V}(\hat{\beta}_k)}$ , soit ici  $t_0 = -4$  et  $t_1 = 7.61$  qui sont en dehors des valeurs critiques associées à un test de Student (de l'ordre de  $\pm 2$  pour un seuil à 95%).

Une fois estimé le modèle, on s'aperçoit qu'une information supplémentaire est disponible: le sexe de l'individu. On a alors

$$Y_i = \alpha_0 + \alpha_1 X_{1,i} + \alpha_2 X_{2,i} + \eta_i,$$

où  $X_{2,i}$  est une variable qui prend la valeur 1 si l'individu  $i$  est une femme, et 0 si l'individu  $i$  est un homme.

**Question 10.** Laquelle des affirmations suivantes est correcte (une seule l'est)

- A. à taille identique, une femme pèse  $\alpha_2$  kilogrammes de plus qu'un homme
- B. à taille identique, une femme pèse  $\alpha_2$  kilogrammes de moins qu'un homme
- C. à taille identique, une femme pèse  $\alpha_0 + \alpha_2$  kilogrammes de plus qu'un homme
- D. à taille identique, une femme pèse  $\alpha_0 + \alpha_2$  kilogrammes de moins qu'un homme

Pour un homme mesurant une hauteur  $x$ , on prédit  $\alpha_0 + \alpha_1 x + 0$  alors que pour une femme de même taille, on va prédire  $\alpha_0 + \alpha_1 x + \alpha_2$ . La différence est alors  $\alpha_2$ , de plus pour une femme.

## 2. RÉGRESSIONS SUR DONNÉES INDIVIDUELLES ET MATCHS DE BASKET

Les questions suivantes portent sur les sorties données dans l'Annexe A.

**Question 11.** Dans la sortie de régression **R1**, à quoi correspond la valeur 131.3 (valeur encadrée)

- A. à la valeur de la pente dans la régression linéaire
- B. à la statistique de Student permettant de tester si la pente de la droite de régression est nulle, ou pas
- C. à la valeur  $t$  qui est (en %) la probabilité que le modèle soit linéaire
- D. au score maximum à la fin d'un match, prédit par le modèle

C'est la statistique  $t$  du test de Student, de significativité de  $\beta_1$  (pour reprendre les notations des questions précédentes).

**Question 12.** A la lecture des sorties **R2**, la quelle de ces affirmations vous semble la moins valide

- A. les résidus suivent une loi normale centrée
- B. avec ce modèle, on a environ 5% de chance de faire une erreur de prévision supérieure à 20 points
- C. il y a de l'hétéroscédasticité, car la variance des résidus n'est pas constante avec la différence de points à la mi-temps
- D. en moyenne, la différence de points à la fin du match est deux fois plus grande qu'à la mi-temps

Au contraire, si l'on regarde la Figure où sont représenté - dans le plan - les deux variables, la dispersion autour de la droite de régression semble indépendante du score à la mi-temps.

**Question 13.** À l'aide de la sortie **R3** prédire la différence de points à la fin du match en faveur d'une équipe qui joue à domicile et qui mène de 13 points à la mi-temps.

- A. environ 20 points d'avance
- B. environ 17 points d'avance
- C. environ 12 points d'avance
- D. environ 6 points d'avance

Faisons le calcul rapidement, voire grossièrement car on demande un ordre de grandeur  $-2.72 + 1.13 \times 13 + 5.44$  (car l'équipe joue à domicile) soit ici environ 17.49

**Question 14.** À l'aide de la sortie **R4** laquelle des affirmations suivante semble la **moins juste** ? (une seule l'est)

- A. les résidus sont nuls
- B. l'écart-type des résidus est nul
- C. il manque une constante dans la régression pour que le modèle soit valide
- D. le  $R^2$  de la régression vaut 1

la régression est parfaite, mais on ne fait pas *vraiment une régression* puisque l'on a une relation parfaite: la différence de points à la fin du match est rigoureusement égale à la différence de points entre les deux équipes. Les résidus sont donc toujours nuls. Par construction.

**Question 15.** À l'aide de la sortie **R5** quelle affirmation vous semble la **moins juste**.

- A. les modèles **reg4** et **reg5** sont équivalents
- B. le modèle **reg5** est meilleur que le modèle **reg6** au sens du  $R^2$  ajusté, et du critère d'Akaike
- C. le modèle **reg5** est meilleur que le modèle **reg2** au sens du  $R^2$  ajusté, et du critère d'Akaike
- D. Ces modèles ne sont pas comparables, car ils n'utilisent pas les mêmes variables explicatives.

Les modèles **reg4** et **reg5** sont équivalents, car ils sont identiques, à une transformation affine près. Et le modèle **reg5** est meilleur que les deux autres, car le critère AIC est toujours plus faible alors que le  $R^2$  ajusté est toujours plus élevé. En fait, le critère AIC



permet précisément de comparer des modèles construits sur des variables explicatives différentes. Cette dernière affirmation est donc fausse.

Dans les questions suivantes (16-18), on va se demander si ces estimations donnent les mêmes résultats pour toutes les équipes, et on va se restreindre à des sous-bases.

**Question 16.** Pour le modèle présenté dans la série des sorties **R6**, seuls les matchs de 6 équipes ont été retenus. Quelle interprétation du nombre 0.707 (encadré) retiendriez vous ?

- A. l'équipe de basket de COLUMBIA a 70.7% de chance de gagner un match
- B. l'équipe de basket de COLUMBIA a 70.7% de chance de gagner un match, si elle gagne à la mi-temps
- C. le coefficient est non-significatif, ce qui veut dire que pour les matchs de l'équipe COLUMBIA, s'il y avait  $x$  points de différence à la mi-temps, la meilleure prédiction possible est qu'il y aura toujours  $x$  points de différence à la fin du match.
- D. le coefficient est non-significatif, ce qui veut dire que l'équipe de basket de COLUMBIA a autant de chance de gagner contre toutes les équipes

Il s'agit de la statistique de Student, pour une variable indicatrice. Les deux premières affirmations sont fausses (ou en tout cas, la statistique  $t$  n'a aucun lien avec la probabilité de gagner. La dernière affirmation n'est pas valide non plus: la non-significativité signifie que pour les matchs de l'équipe COLUMBIA, le modèle serait  $Y = \beta_1 X$ , avec ici  $\beta_1$  proche de 1. On retient en fait la troisième affirmation.

**Question 17.** À l'aide des sorties **R6** laquelle affirmation, parmi les 4 proposées vous semble la plus juste ?

- A. si l'équipe de PERDUE perd de 4 points à la mi-temps, l'équipe a plus de chance de gagner le match que de le perdre
- B. si l'équipe d'UCLA gagne de 4 points à la mi-temps, l'équipe a plus de chance de perdre le match que de le gagner

- C. si l'équipe de COLUMBIA perd de 4 points à la mi-temps, l'équipe a plus de chance de gagner le match que de le perdre
- D. si l'équipe de NORTHWESTERN gagne de 4 points à la mi-temps, l'équipe a plus de chance de perdre le match que de le gagner

Cette fois, on va faire un peu de calculs. Pour PERDUE, si l'équipe perd à la mi-temps de 4 points le score à la fin suit une loi  $\mathcal{N}(-4 \times 1 + 0, 9^2)$ . La probabilité de gagner est la probabilité que

$$\mathbb{P}(-4 + 9Z > 0) = \mathbb{P}(9Z > 4) = \mathbb{P}\left(Z > \frac{4}{9}\right) = \mathbb{P}(Z > 0.55)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . D'après le formulaire  $\mathbb{P}(Z > 0.55) = 30\%$  donc la probabilité que PERDUE gagne est de l'ordre de 30%. Donc A. est fausse.

On continue. Pour UCLA, si l'équipe gagne à la mi-temps de 4 points le score à la fin suit une loi  $\mathcal{N}(4 \times 1 - 6.5, 9^2)$ . La probabilité de gagner est la probabilité que

$$\mathbb{P}(-2.5 + 9Z > 0) = \mathbb{P}(9Z > 2.5) = \mathbb{P}\left(Z > \frac{2.5}{9}\right) = \mathbb{P}(Z > 0.27)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . D'après le formulaire  $\mathbb{P}(Z > 0.25) = 40\%$  donc la probabilité que UCLA gagne est de l'ordre de 40%. Elle a plus de chances de perdre que de gagner, donc B. est vraie (On fait tous les calculs ici, mais il va de soi qu'il suffisait de regarder les signes car on veut juste savoir si la probabilité est plus petite, ou plus grande, que 50%).

Pour COLUMBIA, on se retrouve dans la même situation que pour l'affirmation A., un effet équipe nul: si l'équipe perd, elle devrait perdre à la fin. Donc C. est fausse.

Pour la dernière affirmation, on va faire un calcul rapide, en notant qu'il convient de prendre en compte l'effet équipe. Pour NORTHWESTERN, si l'équipe gagne à la mi-temps de 4 points le score à la fin suit une loi  $\mathcal{N}(4 + 5 \times 1, 9^2)$ . La probabilité de gagner est la probabilité que

$$\mathbb{P}(9 + 9Z > 0) = \mathbb{P}(9Z > -9) = \mathbb{P}(Z > -1) = 1 - \mathbb{P}(Z > 1)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . D'après le formulaire  $\mathbb{P}(Z > 0.84) = 20\%$  et  $\mathbb{P}(Z > 1.28) = 10\%$  donc la probabilité que NORTHWESTERN gagne est de l'ordre de 75%. Donc D. est fausse.

On va ici supposer que l'équipe qui nous intéresse perd fortement à la mi-temps (de plus de 10 points d'écart).

**Question 18.** À l'aide des sorties **R7**, laquelle des affirmations suivantes vous semble la **moins juste** ? *petite coquille: il fallait lire la plus juste*

- A. si l'équipe de PERDUE perd de 10 points à la mi-temps, la probabilité de l'emporter à la fin du match est supérieure à 25%
- B. si l'équipe d'UCLA perd de 10 points à la mi-temps,, la probabilité de l'emporter à la fin du match est supérieure à 25%
- C. si l'équipe de NORTHWESTERN perd de 10 points à la mi-temps, la probabilité de l'emporter à la fin du match est supérieure à 25%
- D. si l'équipe de MICHIGAN STATE perd de 10 points à la mi-temps, la probabilité de l'emporter à la fin du match est supérieure à 25%

À nouveau, on est bon pour quelques calculs. Pour PERDUE, si l'équipe perd à la mi-temps de 10 points le score à la fin suit une loi  $\mathcal{N}(-10 \times 1.34, 9^2)$ . La probabilité de gagner est la probabilité que

$$\mathbb{P}(-13.4 + 9Z > 0) = \mathbb{P}(9Z > 13.4) = \mathbb{P}\left(Z > \frac{13.4}{9}\right) = \mathbb{P}(Z > 1.5)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . D'après le formulaire  $\mathbb{P}(Z > 1.28) = 10\%$ ,  $\mathbb{P}(Z > 1.64) = 5\%$  donc la probabilité que PERDUE, gagne est comprise entre 5% et 10%: elle est moins de 25% de chance de gagner. Donc A. est fausse.

On continue. Pour UCLA, il n'y a pas d'effet équipe, on aura les mêmes résultats numériques. Donc B. est fausse.

Pour NORTHWESTERN, on doit cette fois prendre en compte l'effet équipe, qui est significatif. Si l'équipe perd à la mi-temps de 4 points le score à la fin suit une loi  $\mathcal{N}(-13.4 + 9.6, 9^2)$ . La probabilité de gagner est la probabilité que

$$\mathbb{P}(-4 + 9Z > 0) = \mathbb{P}(9Z > 4) = \mathbb{P}\left(Z > \frac{4}{9}\right) = \mathbb{P}(Z > 0.55)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . Or selon le formulaire  $\mathbb{P}(Z > 0.52) = 30\%$  donc la probabilité que NORTHWESTERN, gagne est de l'ordre de 30%. On a trouvé la bonne réponse.

La première affirmation est vraie, comme on l'a noté auparavant.

Dans les questions suivantes (19-20), on va se demander ce qui se passe si les équipes sont ex-aequo, ou presque. Pour cela, on se limite aux matchs où à la mi-temps, les équipes étaient ex-aequo (différence de points nulle) ou avec un point d'écart.

**Question 19.** À l'aide des sorties **R8** et **R3**, si une équipe perd d'un point à la mi-temps, et joue à domicile, quelle serait la différence de prédiction entre les deux modèles. Une seule des affirmations suivantes est **fausse**.

- A. avec le modèle de la sortie **R8**, peu importe que l'équipe mène ou perde d'un point, la prédiction sera significativement la même
- B. avec le modèle de la sortie **R8**, peu importe que l'équipe qui perd joue à domicile, ou pas, la prédiction sera significativement la même
- C. les deux modèles prévoient une victoire (avec un nombre de points d'écart à la fin du match positif, en moyenne)
- D. les deux modèles **reg2** et **reg10** prévoient une victoire, mais avec un point de *moins* (environ) avec le modèle **reg2** qu'avec le modèle **reg10**

Avec ce dernier modèle, la différence de points n'est pas significative, donc effectivement, la prédiction sera la même que la différence soit positive, ou négative. Donc A. est juste. En revanche, il y a un effet de jouer à domicile ou pas. Donc B. est fausse. La bonne réponse est donc ici B.

Vérifions rapidement que les deux autres affirmations sont justes. Si on joue à domicile, le dernier modèle prévoit une victoire, avec  $5.14 - 2.57$  points d'avance, soit  $2.57$  points d'avance. Avec le premier modèle (**reg2**), nous avons une prévision  $5.55 - 2.72 + 1.11 \times (-1)$  soit  $1.72$ . Bref, les deux modèles prévoient une victoire, et effectivement, le modèle **reg2** prévoit un point (environ,  $0.85$ ) de moins.

**Question 20.** À l'aide de la sortie **R8**, si une équipe perd d'un point à la mi-temps, en jouant à l'extérieur, quelle serait la probabilité qu'elle gagne le match ?

- A. 0%
- B. 20%

C. 40%

D. 80%

Pour une équipe qui joue à l'extérieur et qui perd d'un point, la différence de score à la fin suit une loi  $\mathcal{N}(-2.5 + 00, 9^2)$ . La probabilité de gagner de plus d'un point est la probabilité que

$$\mathbb{P}(-2.5 + 9Z > 0) = \mathbb{P}(9Z > 2.5) = \mathbb{P}\left(Z > \frac{1.5}{9}\right) = \mathbb{P}(Z > 0.27)$$

où  $Z$  suit une loi  $\mathcal{N}(0, 1)$ . Or selon le formulaire  $\mathbb{P}(Z > 0.27) = 40\%$  donc la bonne réponse est la réponse C.

### 3. MODÉLISATION DE SÉRIES TEMPORELLES

Dans les questions 21-26, on considère un modèle autorégressif de la forme

$$X_t = \phi X_{t-1} + \varepsilon_t, \quad (3.1)$$

où  $(\varepsilon_t)$  est un bruit blanc. On supposera que  $\phi \neq 0$ , et on supposera dans toute la section que la série  $(X_t)$  est stationnaire.

**Question 21.** Laquelle des affirmations suivant n'est pas correcte

A.  $\varepsilon_t$  et  $X_{t-1}$  sont non-corrélésB.  $X_t$  et  $\varepsilon_{t-1}$  sont non-corrélésC.  $\varepsilon_t$  et  $\varepsilon_{t-1}$  sont non-corrélésD.  $X_t$  et  $\varepsilon_t$  sont corrélés

par construction,  $X_t = \phi X_{t-1} + \varepsilon_t$  donc  $X_t$  et  $\varepsilon_{t-1}$  sont corrélés. En revanche,  $\varepsilon_t$  est indépendant du passé du processus  $(X_t)$ , c'est à dire des  $X_{t-h}$  pour  $h \geq 1$ .

**Question 22.** Laquelle des affirmations suivant est correcte

A.  $\text{cor}(X_t, X_{t-1}) = \phi$ B.  $\text{cor}(X_t, X_{t-2}) = 0$ C.  $\text{cor}(X_t, \varepsilon_t) = 0$ D.  $\text{cor}(\varepsilon_t, \varepsilon_{t-1}) = \phi$

Refaisons les calculs. Comme  $\text{cor}(\varepsilon_t, X_{t-1}) = 0$ ,

$$\text{cor}(X_t, X_{t-1}) = \text{cor}(\phi X_{t-1} + \varepsilon_t, X_{t-1}) = \phi + 0$$

En revanche,

$$\text{cor}(X_t, X_{t-2}) = \text{cor}(\phi X_{t-1} + \varepsilon_t, X_{t-2}) = \text{cor}(\phi[\phi X_{t-2} + \varepsilon_{t-1}] + \varepsilon_t, X_{t-2}) = \phi^2$$

$$\text{cor}(X_t, X_{t-2}) = \text{cor}(\phi X_{t-1} + \varepsilon_t, X_{t-2}) = \text{cor}(\phi[\phi X_{t-2} + \varepsilon_{t-1}] + \varepsilon_t, X_{t-2}) = \phi^2$$

alors que  $\text{cor}(\varepsilon_t, \varepsilon_{t-1}) = 0$  par hypothèse sur les bruits blancs.

On suppose ici que le bruit  $(\varepsilon_t)$  est de variance  $\sigma^2 = 1$ .

**Question 23.** Que vaut la variance de  $(X_t)$  ?

- A.  $\text{Var}(X_t) = 1$
- B.  $\text{Var}(X_t) = (1 - \phi^2)$
- C.  $\text{Var}(X_t) = \frac{1}{1 - \phi^2}$
- D.  $\text{Var}(X_t) = 0$

Là encore, faisons les calculs.

$$\text{Var}(X_t) = \phi^2 \text{Var}(X_{t-1}) + \text{Var}(\varepsilon_t) = \phi^2 \text{Var}(X_t) + \sigma^2, \text{ car } \text{Var}(X_{t-1}) = \text{Var}(X_t),$$

le processus étant supposé stationnaire. Donc  $(1 - \phi^2)\text{Var}(X_t) = \sigma^2$ .

On suppose que l'on travaille sur des données journalières. On a observé  $-1$  hier, et  $+1$  aujourd'hui

**Question 24.** Quelle prévision feriez vous pour demain, à l'aide du modèle (3.1) ?

- A.  $\hat{X} = \phi$
- B.  $\hat{X} = -\phi$
- C.  $\hat{X} = 0$
- D.  $\hat{X} = \phi + \varepsilon$

Comme on a un processus  $AR(1)$ , la meilleure prévision que l'on puisse faire à la date  $T$  pour un horizon 1 est  ${}_T\hat{X}_{T+1} = \phi X_T$ , et comme  $X_{T-1} = -1$  (mais ça, on s'en moque) et  $X_T = +1$ , la prévision sera  $\phi$ .

**Question 25.** Quel serait la variance de l'erreur de prévision à deux jours (pour après demain) ?

- A.  $\phi$
- B. 1
- C.  $1 + \phi^2$
- D.  $\phi^2$

Cette fois, il faut aller un peu plus loin dans l'écriture

$$X_{T+2} = \phi X_{T+1} + \varepsilon_{T+2} = \phi[\phi X_T + \varepsilon_{T+1}] + \varepsilon_{T+2} = {}_T\hat{X}_{T+2} + \phi\varepsilon_{T+1} + \varepsilon_{T+2}$$

L'erreur de prévision est ici la variable  $\phi\varepsilon_{T+1} + \varepsilon_{T+2}$  dont la variance est  $(1 + \phi^2)\sigma^2$ .

**Question 26.** Une seule des affirmations suivantes est correcte. Laquelle ?

- A. La série  $Y_t = t + X_t$  est stationnaire
- B. La série  $Y_t = X_t - X_{t-1}$  est stationnaire
- C. La série  $Y_t$  vérifiant la relation  $Y_t = Y_{t-1} + X_t$  est stationnaire
- D. La série  $Y_t$  vérifiant la relation  $Y_t = \phi^{-1}Y_{t-1} + \varepsilon_t$  est stationnaire

La série  $X_t$  est ici stationnaire.  $Y_t^{(A)} = t + X_t$  n'est pas stationnaire, car  $\mathbb{E}(Y_t^{(A)})$  sera alors affine en  $t$  (et donc non constante).  $Y_t^C = Y_{t-1}^C + X_t$  peut se récrire

$$Y_t^C = Y_0^C + \sum_{k=1}^t X_k = Y_0^C + \sum_{k=1}^t \sum_{h=1}^k \varepsilon_{k-h+1} = 1^k \phi^h \varepsilon_{k-h},$$

donc

$$\text{Var}(Y_t^C) = \sum_{k=1}^t \sum_{h=1}^k \varepsilon_{k-h+1}^2 = 1^k \phi^h \sum_{k=1}^t (t-k) \phi^k$$

qui n'est pas du tout constante. Donc  $Y_t^{(C)}$  n'est pas stationnaire. Enfin, pour la dernière affirmation, comme  $|\phi| < 1$  alors  $|\phi^{-1}| > 1$ . Donc  $Y_t^{(D)}$  n'est pas stationnaire.

Même s'il ne reste plus que B., essayons de montrer que cette série est effectivement stationnaire. D'un côté,  $Y_t^B = X_t - X_{t-1}$  peut se récrire  $Y_t^B = (\phi - 1)X_{t-1} + \varepsilon_t$ . De l'autre, on peut écrire

$$Y_{t-1}^B = X_{t-1} - X_{t-2} = X_{t-1} - \phi^{-1}[X_{t-1} - \varepsilon_{t-1}]$$

donc on peut écrire

$$X_{t-1} = [1 - \phi^{-1}]^{-1} Y_{t-1}^B - \phi^{-1} \varepsilon_{t-1}$$

En substituant dans la première expression

$$Y_{t-1}^B = (\phi - 1)[1 - \phi^{-1}]^{-1} Y_{t-1}^B - \phi^{-1} \varepsilon_{t-1} + \varepsilon_t$$

on obtient que  $Y_t^B$  peut s'écrire sous forme  $ARMA(1, 1)$ , avec comme coefficient autorégressif  $\phi$ . Donc la série  $Y_t^B$  est stationnaire.

On suppose maintenant que la série  $(X_t)$  suit un processus  $ARMA(1, 1)$

$$X_t = \phi X_{t-1} + \eta_t + \theta \eta_{t-1} \quad (3.2)$$

où  $(\eta_t)$  est un bruit blanc de variance  $\sigma^2$ .

**Question 27.** Quelle est la condition nécessaire et suffisante, sur  $\phi$ ,  $\theta$  et  $\sigma$ , pour que  $(X_t)$  soit stationnaire

- A.  $|\phi| < 1$ ,  $|\theta| < 1$  et  $\sigma^2 < 1$ ,
- B.  $|\phi| < 1$  et  $|\theta| < 1$
- C.  $|\theta| < 1$
- D.  $|\phi| < 1$ .

C'est du cours.

**Question 28.** Parmi les affirmations suivantes, une seule est toujours vraie. Laquelle ? (un seul choix est possible)

- A. Toutes les prédictions faites avec le modèle (3.2) seront identiques avec celles faites à l'aide du modèle (3.1)
- B. Toutes les prédictions faites avec le modèle (3.2) auront une variance (conditionnelle à l'information disponible aujourd'hui) identiques à celles du modèle (3.1)
- C. L'autocorrélation à l'ordre 1 est toujours non-nulle  $\rho_X(1) \neq 0$
- D. À partir, d'un moment, on retrouve la propriété vérifiée par les  $AR(1)$  sur les autocorrélations, à savoir  $\rho_X(h) = \phi \rho_X(h-1)$  (pour  $h \geq 1$ ).



C'est la dernière, cf. cours. Si on n'est pas convaincu, les deux premières sont fausses, sinon je l'aurais dit en cours (un résultat pareil, ça se remarque). Pour s'en convaincre, il faut se souvenir que quand on a un  $ARMA(p, q)$ , la prévision se fait en passant par la forme  $AR(\infty)$  (si le polynôme moyenne mobile est inversible). Par exemple pour le processus  $ARMA(1, 1)$ ,

$$(1 - \phi L)X_t = (1 + \theta L)\varepsilon_t \text{ soit } (1 - \phi L)(1 + \theta L)^{-1}X_t = \varepsilon_t$$

qui va s'écrire ici

$$(1 - \phi L)(1 + \theta L)^{-1}X_t = (1 - \phi L) \left( \sum_{k=0}^{\infty} (-\theta)^k L^k \right) \varepsilon_t$$

On voit que la prévision fera intervenir non seulement la valeur en  $T$  mais aussi en  $T - 1$ , en  $T - 2$ , etc.

Quant à la troisième affirmation, un peu de calcul permet de montrer que, de manière générale,

$$\gamma(0) = \sigma^2 \frac{\theta^2 + 2\phi\theta + 1}{1 - \phi^2} \text{ et } \gamma(1) = \phi\gamma(0) + \theta\sigma^2.$$

Bref, l'expression est moche, mais on peut l'écrire

$$\rho(1) = \frac{\gamma(1)}{\gamma(0)} = \phi + \frac{\theta(1 - \phi^2)}{\theta^2 + 2\phi\theta + 1}$$

Rien ne permet de garantir que cette expression soit toujours non-nulle...

**Question 29.** Un voisin suggère que dans le modèle (3.2), on pourrait avoir  $\phi + \theta = 0$ .

Dans ce cas, quelle affirmation parmi les suivantes serait correcte

- A.  $(X_t)$  suit une marche aléatoire
- B.  $(X_t)$  suit un bruit blanc
- C.  $(X_t)$  est, en fait, un processus  $AR(1)$
- D.  $(X_t)$  est, en fait, un processus  $MA(1)$

Comme  $(1 - \phi L)X_t = (1 + \theta L)\varepsilon_t$ , si  $\theta = -\phi$ , on peut simplifier, et donc  $X_t = \varepsilon_t$ . Bref,  $X_t$  est un bruit blanc.

On suppose désormais que la série  $(X_t)$  suit un processus  $AR(p)$  stationnaire,

$$X_t = \phi_1 X_{t-1} + \cdots + \phi_p X_{t-p} + u_t, \quad (3.3)$$

où  $(u_t)$  est un bruit blanc. On supposera que  $\phi_p \neq 0$ .

**Question 30.** Parmi les affirmations suivantes, une seule est (de manière générale) vraie. Laquelle ?

- A. la série des autocorrélations  $\rho(h) = \text{cor}(X_t, X_{t-h})$  est nulle pour  $h > p$
- B. la série des autocorrélations partielles  $\psi(h) = \text{cor}(X_t - \text{EL}(X_t | X_{t-1}, \dots, X_{t-h+1}), X_{t-h} - \text{EL}(X_{t-h} | X_{t-1}, \dots, X_{t-h+1}))$  est nulle pour  $h > p$
- C. l'autocorrélation d'ordre  $p$  est non nulle,  $\rho(p) \neq 0$
- D. l'autocorrélation d'ordre  $p$  vaut 1,  $\rho(p) = 1$

C'est du cours. C'est en regardant les autocorrélations partielles que l'on identifie l'ordre d'une  $AR(p)$ .

#### 4. MODÉLISATION DE SÉRIES TEMPORELLES

On cherche ici à modéliser, puis à proposer une prévision pour une série de ventes, dont les sorties sont proposées dans l'Annexe B.

On notera  $\rho_X$  la fonction d'autocorrélations de la série  $(X_t)$  et  $\psi_X$  sa fonction d'autocorrélations partielles. Pour un processus  $ARMA(p, q)$ , on utilisera la notation suivante,

$$X_t = \alpha + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \cdots + \phi_p X_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \cdots + \theta_q \varepsilon_{t-q},$$

où  $(\varepsilon_t)$  est un bruit blanc de variance  $\sigma^2$ , et où  $\phi_p \neq 0$  et  $\theta_q \neq 0$ .

**Question 31.** À l'aide de la sortie **T1**, parmi les affirmations suivantes, laquelle n'est pas correcte

- A.  $\rho(12)$  est significativement non-nulle
- B.  $\rho(6)$  est significativement non-nulle
- C.  $\rho(3)$  est significativement non-nulle
- D.  $\psi(12)$  est significativement non-nulle

“significativement non-nulle” veut dire que les points sont en dehors de la bande critique du test de significativité, sur les autocorrélogrammes. Ici, seule  $\rho(6)$  est dans cette bande, donc  $\rho(6)$  est significativement nulle.

**Question 32.** À l’aide des sorties **T1** et **T2**, parmi les affirmations suivantes, laquelle est correcte (avec un seuil d’acceptation à 95%)

- A. la série  $(X_t)$  est stationnaire
- B. la série  $(X_t)$  n’est pas stationnaire
- C. la série  $(X_t)$  est un bruit blanc
- D. la série  $(X_t)$  est une marche aléatoire

On va accepter l’hypothèse de série stationnaire, i.e. on refuse automatiquement A. et D. Pour le bruit blanc, il n’y a pas eu de test formel. Mais on a quand même beaucoup d’autocorrélations largement significatives, donc il serait surprenant que l’on considère avoir ici un bruit blanc.

**Question 33.** À l’aide des sorties **T3**, sur la modélisation par un processus  $AR(18)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble la plus juste,

- A.  $\phi_{18}$  est significativement non-nulle
- B. la plus grande valeur  $p$  pour laquelle  $\phi_p$  est significativement non-nulle est 12
- C. la constante est significative
- D. le bruit associé  $(\varepsilon_t)$  n’est sûrement pas un bruit blanc

A. étant vraie, B. sera fausse.

**Question 34.** À l’aide des sorties **T4**, sur la modélisation par un processus  $AR(12)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble **la moins juste**,

- A.  $\phi_{12}$  est significativement non-nulle
- B. le bruit associé  $(\varepsilon_t)$  peut être considéré comme étant un bruit blanc
- C. le bruit associé  $(\varepsilon_t)$  n’est pas un bruit blanc,
- D. même si on n’a pas vérifié le bruit associé  $(\varepsilon_t)$  est centré, il l’est forcément car on a mis une constante dans le modèle

Là encore, comme on accepte  $B$ , on va rejeter  $A$  (car oui, c'est un bruit blanc).

**Question 35.** À l'aide des sorties **T5**, sur la modélisation par un processus  $AR(3)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble la **plus juste**,

- A. le bruit associé  $(\varepsilon_t)$  peut être considéré comme étant un bruit blanc
- B. le bruit associé  $(\varepsilon_t)$  n'est probablement pas un bruit blanc,
- C. on ne peut pas savoir si le bruit associé  $(\varepsilon_t)$  est un bruit blanc, car on n'a pas fait de test de racine unité sur ce dernier.
- D. on ne peut pas savoir si le bruit associé  $(\varepsilon_t)$  est un bruit blanc, car on n'a pas vérifié que ce dernier était centré

Un peu de logique permettait de voir qu'il s'agissait de la seule réponse possible (sinon, deux réponses seraient possibles). Maintenant, en regardant les tests de Box-Pierce, on rejette clairement l'hypothèse de bruit blanc.

**Question 36.** À l'aide des sorties **T6**, sur la modélisation par un processus  $ARMA(3, 12)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble la **moins juste**,

- A.  $\phi_3$  et  $\theta_{18}$  sont significativement non-nulles (petite coquille: il fallait lire  $\theta_{12}$ )
- B. le bruit associé  $(\varepsilon_t)$  peut être considéré comme étant un bruit blanc
- C. le bruit associé  $(\varepsilon_t)$  n'est probablement pas un bruit blanc,
- D. en utilisant le critère AIC, on ne retiendrait pas ce modèle car on en a trouvé des plus simples, et meilleurs auparavant.

On a ici un bruit blanc, donc on va rejeter la troisième affirmation.

**Question 37.** À l'aide des sorties **T7**, sur la modélisation par un processus  $AR(9)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble la **moins juste**,

- A.  $\phi_9$  est significativement non-nulle
- B. comme  $|\phi_6|$  est plus grand que  $|\phi_9|$  on rejette un  $AR(9)$  au profit d'un modèle  $AR(6)$
- C. le bruit associé  $(\varepsilon_t)$  n'est probablement pas un bruit blanc,

- D. bien que le modèle est suggéré (en premier choix) par la fonction `armaselect`, sur la base d'un critère SBC minimal parmi tous les  $ARMA(p, q)$ , on ne devrait pas le retenir.

Les affirmations A., C. et D. étant juste, on rejettera B. Effectivement, cette affirmation ne repose sur rien de tangible.

**Question 38.** À l'aide des sorties **T8**, sur la modélisation par un processus  $AR(6)$  pour  $(X_t)$ , laquelle de ces affirmations suivantes vous semble la **plus juste**,

- A. le modèle ne peut pas être correct car le bruit associé  $(\varepsilon_t)$  n'est probablement pas un bruit blanc,
- B. le modèle ne peut pas être correct car comme  $\rho(12)$  était significativement non-nulle, il faut forcément un modèle  $ARMA(p, q)$  avec  $p + q \geq 12$ ,
- C. le modèle ne peut pas être correct car la variance du bruit associé  $(\varepsilon_t)$  est plus grande que 1, et donc le modèle n'est pas stationnaire,
- D. le modèle ne peut pas être correct car  $|\phi_1 + \phi_2 + \dots + \phi_6| \approx 1$ , et donc on est en présence d'une racine unité,

la première affirmation est correcte. On va donc la retenir. Pour la dernière, on aurait une racine unité si la somme était nulle. Pour l'avant-dernière, la variance peut prendre n'importe quelle valeur positive. Enfin, pour la seconde, là encore cette affirmation ne repose sur rien.

**Question 39.** À l'aide des sorties **T9** et **T10**, sur les modélisations par des processus  $AR(10)$  et  $ARMA(10, 1)$  de la série  $(X_t)$ , laquelle de ces affirmations suivante vous semble la **plus juste**,

- A. le processus  $ARMA(10, 1)$  est préférable au processus  $AR(10)$ , car le bruit associé  $(\varepsilon_t)$  est de variance plus faible
- B. le processus  $ARMA(10)$  est préférable au processus  $AR(10, 1)$  car le coefficient  $\theta_1$  n'est pas significatif (petite coquille: il fallait lire  $AR(10)$  puis  $ARMA(10, 1)$ )

- C. le modèle  $ARMA(10, 1)$  ne peut être retenu car comme la somme des coefficients vaut  $-\pi$ , il faut un modèle avec une composante cyclique,
- D. le modèle  $ARMA(10)$  ne peut être retenu car  $\rho_\varepsilon(12)$  était significativement non-nulle, il faudrait tester un modèle  $MA(12)$  sur les résidus (petite coquille: il fallait lire  $AR(10)$ )

Pour avoir un processus  $ARMA(10, 1)$  valide il faut que  $\phi_{10} \neq 0$  et  $\theta_1 \neq 0$ . Or  $\theta_1$  n'est pas significatif. Donc pour cette raison, on va rejeter le modèle  $ARMA(10, 1)$ .

**Question 40.** À l'aide des sorties **T11**, parmi les modèles suivants, lequel retriendriez-vous pour faire une prévision sur un an (pour les 12 prochaines observations) ? On retiendra le meilleur modèle sur les observations récentes,

- A. un processus  $AR(9)$
- B. un processus  $AR(10)$
- C. un processus  $AR(12)$
- D. un processus  $AR(12, 3)$  (petite coquille: il fallait lire  $ARMA(12, 3)$ )

Je retiendrais le premier modèle car c'est celui qui a donné le moins d'erreurs sur la prévision pour les 12 derniers mois. Maintenant, c'est aussi le modèle avec le plus mauvais critère d'Akaike. Je préfère un modèle bon en prévision, donc mon choix pencherait malgré tout pour ce modèle.