
Projet de séries temporelles

Jean-Baptiste DREVETON
Alexandre FILIOT

On s'intéresse à la modélisation et la prévision de l'indice de production industrielle observé en France. Il ne vous est pas demandé de connaître le processus de construction de l'indicateur. Vous ne travaillez que sur les données observées. A partir du répertoire des séries chronologiques de l'INSEE, www.insee.fr/fr/statistiques?debut=0&categorie=10, vous devez choisir une série agrégée brute, mensuelle, correspondant à n'importe quel secteur d'activité (à votre convenance) contenant au moins 25 années et contenant le mois de janvier 2018.

Date limite : 6 Juin 2018

Partie 1 : Les données

1. Que représente la série choisie ? (secteur, périmètre, traitements éventuels, transformation logarithmique, ...)

Dans ce projet, nous nous intéresserons au secteur du transport, et plus précisément à la fréquentation de passagers lors des vols internationaux au départ de Paris. Les données sont mensuelles et s'étendent de Janvier 1994 à Janvier 2018. Ces données sont issues de l'Insee¹. La série originale est affichée en page 3, graphique 1.

2&3. Transformer si besoin la série pour la rendre stationnaire (désaisonnalisation, différenciation, suppression de la tendance déterministe, ...). Justifier soigneusement vos choix. Représenter graphiquement la série choisie avant et après transformation.

Tout d'abord, il est souvent d'usage pour ce type de série d'opérer au préalable une transformation logarithmique. Ceci permet de réduire l'hétéroscédasticité des résidus.

Les touristes prenant leurs vacances à peu près au même moment de l'année, il en résulte que notre série présente une saisonnalité à l'ordre 12. Nous appliquons alors la transformation suivante :

$$\mathbf{Y}_t = (1 - B^{12})\log \mathbf{X}_t$$

où $(\mathbf{X}_t)_{12 \leq t \leq T}$ est la série initiale, et B l'opérateur backward.

Suite à l'arrivée de compagnies low-cost sur le secteur du transport aéronautique, et la construction de nouveaux aéroports, notamment dans les pays émergents, il y a eu une démocratisation d'accès au transport aérien. Cela se traduit sur notre série par une tendance haussière. La série a donc été différenciée à l'ordre 1 pour supprimer cette tendance déterministe. La série finale est alors :

$$\mathbf{Z}_t = (1 - B)\mathbf{Y}_t = (1 - B) \circ (1 - B^{12})\log \mathbf{X}_t$$

L'objectif de cette partie étant de stationnariser la série initiale, justifions nos transformations. Ces justifications sont également disponibles en commentaires du code. Étant donné que notre série présente une tendance croissante sur l'ensemble de la période observée, nous nous plaçons dans le cas avec constante non nulle et avec tendance du test de Dickey-Fuller augmenté (ADF). Nous utilisons en premier lieu ce test pour rejeter ou non la stationnarité de nos différentes séries. L'estimation du modèle ADF, pour une série quelconque (\mathbf{X}_t) est donnée par :

$$\Delta \mathbf{X}_t = c + at + \phi \mathbf{X}_t + \sum_{k=1}^h \alpha_k \Delta \mathbf{X}_{t-k} + \mathbf{u}_t$$

avec h le lag maximum utilisé dans l'estimation. Systématiquement nous avons relevé la valeur de la statistique de test et l'avons comparée aux valeurs critiques associées aux tests à 1%, 5% et 10%.

Tests sur la série originale

Commençons par notre série originale. Lorsque $h = 0$, nous trouvons une statistique de test de $-6,3258$ soit une p -valeur faible. En estimant le modèle simple :

$$\Delta \log \mathbf{X}_t = c + at + \phi \log \mathbf{X}_t + \mathbf{u}_t$$

nous rejetons donc aux seuils usuels l'hypothèse de racine unitaire. Quel est alors l'intérêt d'appliquer une désaisonnalisation et/ou une différenciation ? Nous allons y venir. Quand bien même l'hypothèse de non-stationnarité n'est pas rejetée à 0,1%, il faut néanmoins tester la validité du modèle en terme de persistance des résidus. Pour cela, nous utilisons un test de Ljung-Box avec un lag maximal égal à 24. Pour notre série $\log \mathbf{X}_t$, il s'avère que le test de blancheur des résidus est rejeté à 5% à tous les ordres : les résidus du modèle sont donc corrélés de l'ordre 1 à l'ordre 24, le modèle n'est donc pas valide. L'idée est donc d'augmenter le paramètre h dans l'estimation du modèle ADF afin de relever le lag minimal permettant de ne pas rejeter la corrélation des résidus de

1. <https://insee.fr/fr/statistiques/serie/010001981>

l'ordre 1 à l'ordre 24 au seuil de 5%. Cette valeur est de $h_{min} = 23$ pour la série $\log \mathbf{X}_t$, ce qui est beaucoup. Nous aimerions donc pouvoir diminuer cette valeur de h_{min} à l'aide d'une première transformation, la désaisonnalisation.

Tests sur la série désaisonnalisée

La p -valeur du test ADF pour un lag $h = 0$ appliqué à la série désaisonnalisée \mathbf{Y}_t est très faible, on rejette donc l'hypothèse de racine unité à 5% (df -statistique : -8,6601). Le test de Ljung-Box ne rejette pas à 5%, de peu, l'hypothèse de corrélation des résidus à l'ordre 1. En effet, la p -valeur est égale à 0,0542. Enfin, le lag minimal à partir duquel l'hypothèse de corrélation des résidus est rejetée à 5% jusqu'à l'ordre 24 est de 13, ce qui est nettement inférieur à 23 pour la série $\log \mathbf{X}_t$.

Tests sur la série désaisonnalisée et différenciée

La p -valeur du test ADF pour un lag $h = 0$ appliqué à la série différenciée et désaisonnalisée \mathbf{Z}_t est encore plus faible, on rejette donc logiquement l'hypothèse de racine unité à 5% (df -statistique : -24,697). Le test de Ljung-Box ne rejette pas à 5% l'hypothèse de corrélation des résidus à l'ordre 1. En effet, la p -valeur est égale à 0,3952. Enfin, le lag minimal à partir duquel l'hypothèse de corrélation des résidus est rejetée à 5% jusqu'à l'ordre 24 est lui aussi égal à 13.

Comment choisir ?

On remarque que la valeur de h_{min} est la même pour la série désaisonnalisée \mathbf{Y}_t que pour celle désaisonnalisée et différenciée \mathbf{Z}_t . Pour cette valeur de h , la p -valeur du test ADF sur \mathbf{Y}_t est de 0,021, tandis que celle du test sur la série \mathbf{Z}_t est inférieure à 0,01. Ainsi, les tests ADF rejettent (heureusement) l'hypothèse de racine unité sur les deux transformations. Il en est de même pour le test de Phillips-Perron (df -statistique respectives de -327,68 et -129,24 pour les séries \mathbf{Z}_t et \mathbf{Y}_t) où les p -valeurs sont plus petites que 0,01. Enfin, le test KPSS ne rejette pas l'hypothèse de stationnarité à 5% puisque les p -valeurs sont toutes deux égales à 0,1. Systématiquement, nous remarquons que le rejet des hypothèses de non-stationnarité (ou le non-rejet de celles de stationnarité) est plus fort, en terme de p -valeurs, pour la série \mathbf{Z}_t . De plus, nous l'avons dit, l'hypothèse de corrélation des résidus à l'ordre 1 est "plus" rejetée à 5% pour cette série que pour celle désaisonnalisée \mathbf{Y}_t . Nous aurions pu nous contenter d'utiliser la série désaisonnalisée uniquement pour la suite, mais nous estimons que notre modèle est plus robuste lorsqu'il prend en compte la série \mathbf{Z}_t en terme de stationnarité et de blancheur des résidus. Le calibrage du modèle ARMA aurait sans doute était sensiblement le même.

Enfin, nous n'aurions cependant pas pu utiliser la série $\log \mathbf{X}_t$ initiale dans la mesure où le test ADF ne rejette pas l'hypothèse de racine unité à 5% (p -valeur égale à 0,553) lorsque $h = h_{min} = 23$.

Partie 2 : Modèles ARMA

4. Choisir, en le justifiant, un modèle ARMA(p, q) (avec éventuellement une composante saisonnière) pour votre série corrigée \mathbf{Z}_t . Estimer les paramètres du modèle et vérifier sa validité.

Comme nous l'avons vu dans la partie précédente, notre série présente une période d'ordre 12 et une tendance haussière déterministe d'ordre 1. Ainsi,

$$\log \mathbf{X}_t \sim SARIMA(p, d = 1, q)(P, D = 1, Q)_{s=12}$$

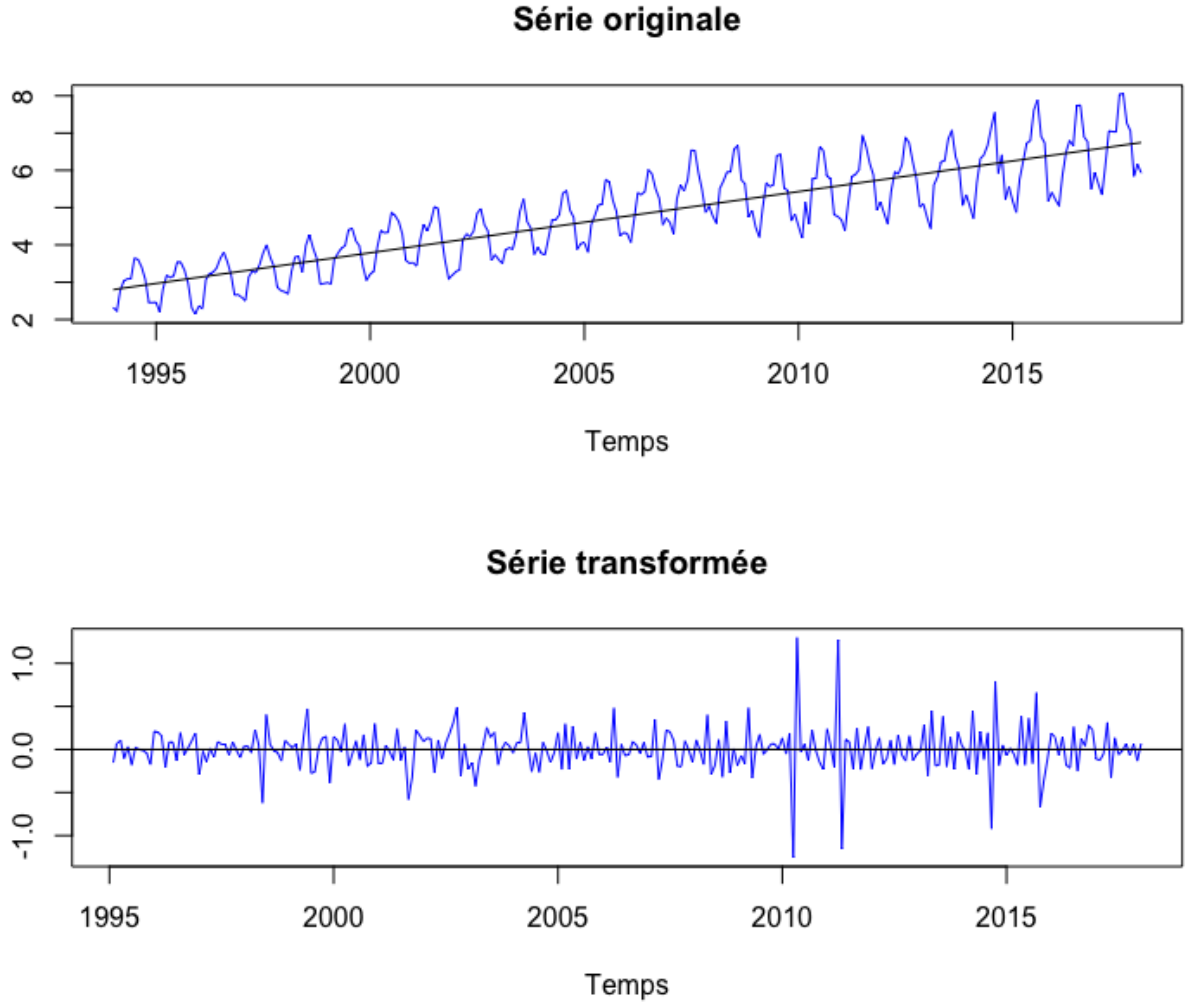
$$\mathbf{Z}_t \sim SARIMA(p, d = 0, q)(P, D = 0, Q)_{s=12}$$

Nous suspectons donc la série $\log \mathbf{X}_t$ de satisfaire l'équation :

$$(1 - B)^d(1 - B^s)^D\phi(B)\Phi(B^s)\log \mathbf{X}_t = \theta(B)\Theta(B^s)\epsilon_t$$

où p, q, P, Q sont respectivement les ordres des polynômes $\phi, \theta, \Phi, \Theta$ et sont à déterminer au même titre que la valeur des coefficients des différents polynômes. ϵ_t est lui un bruit blanc.

Graphique 1 – Différence entre série brute et série corrigée



Nous déterminons les ordres maximums des coefficients p , q , P , Q en regardant l'autocorrélogramme et l'autocorrélogramme partiel de la série différenciée et désaisonnalisée \mathbf{Z}_t (figure 2). Tout d'abord, les autocorrélations empiriques pour des lags compris entre 2 et 11 sont dans l'intervalle de confiance à 95%. Nous pouvons donc prendre $q_{max} = 3$. Pour la PACF, les autocorrélations partielles de lags 1,2,4,5,7,11,12,13,21 n'appartiennent pas aux bornes de significativité à 5%. Nous pouvons considérer que l'écart entre le lag 13 et le lag 21 permet de prendre $p_{max} = 14$.

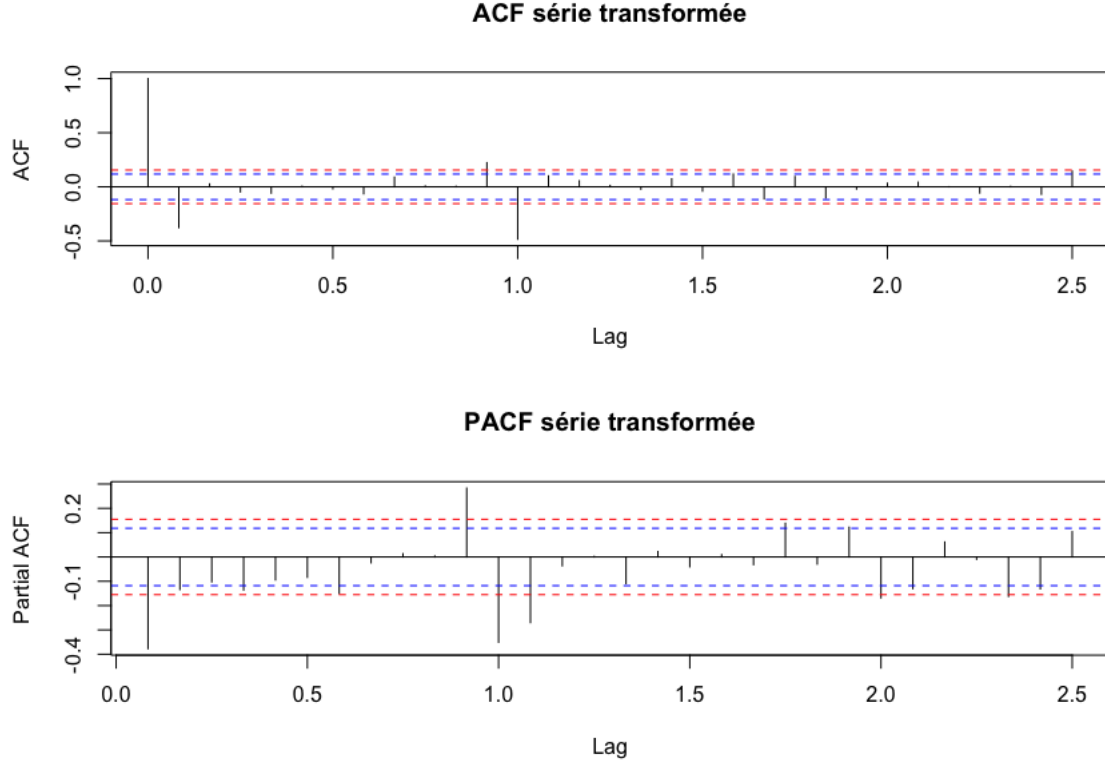
Pour ce qui est des ordres P_{max} et Q_{max} , nous nous intéressons aux lags multiples de 12. À l'aide d'une instruction simple :

```
which(abs(acf$acf) > clim95)
```

avec `clim95` la valeur de la borne de significativité supérieure de l'ACF (ou PACF), nous pouvons avoir accès aux lags pour lesquels les autocorrélations ou autocorrélations partielles empiriques sont hors des bornes de significativité à 5%. En annexe, figure 4 sont représentés les ACF et PACF jusque lag 300 avec les bornes de significativité à 99%. En contrôlant la significativité à 5% nous devrions prendre en toute rigueur $P_{max} = 4$ et $Q_{max} = 12$; ces ordres sont trop élevés pour une éventuelle modélisation *SARIMA*. Nous contrôlons donc la significativité à 1%. Dans ce cas, nous pouvons prendre $Q_{max} = 5$, ce qui est déjà élevé. En conclusion, les ordres maximums vraisemblables pour la série \mathbf{Z}_t sont :

$$(p_{max}, q_{max}, P_{max}, Q_{max}) = (12, 3, 4, 5)$$

Graphique 2 – ACF et PACF de la série stationnarisée jusque lag 30



Nous allons alors évaluer tous les modèles du type $SARIMA(p, 0, q)(P, 0, Q)_{s=12}$ avec $p \leq 5$, $q \leq 6$, $P \leq 4$, $Q \leq 5$ et tester leurs validités et significativités à l'aide de tests statistiques.

La validité du modèle se teste à partir des tests de Ljung-Box comme précédemment. La significativité se vérifie en s'assurant que les coefficients associés aux ordres les plus élevés sont significativement non nuls, i.e que nous pouvons rejeter l'hypothèse de leur nullité à un certain seuil. Nous prendrons dans tous nos tests un seuil de significativité à 5%.

En annexe, nous indiquons les modèles AR (tableau 1), MA (tableau 2), $ARMA$ (tableau 3) et $SARIMA$ (tableau 4) significatifs et valides à 5%. Au total, nous avions à l'origine retenu 266 modèles valides et significatifs (5 modèles AR ; 4 modèles MA ; 75 modèles $ARMA$; 182 modèles $SARIMA$). Ensuite, nous avons classé ces différents modèles retenus selon des indicateurs mesurant la qualité du modèle : les critères AIC et BIC. Nous avons privilégié de manière arbitraire l'AIC. Enfin, pour sélectionner le modèle final, nous avons calculé les erreurs de prédiction à horizon 4 (en tronquant la série $\log \mathbf{X}_t$) et la somme des carrés résidus sur les données existantes (RMSE). Le tableau 5 recense les 10 modèles que nous avons finalement retenus : les modèles minimiseurs de l'AIC ou BIC, et 8 autres modèles dont les performances soit en terme de prédiction soit en terme de fidélité aux données étaient meilleures que celles des 2 premiers modèles. L'utilisation du critère d'erreur de prédiction n'est pas pertinente puisque nous ne cherchons pas à faire de la prédiction (car en pratique nous ne connaissons pas les valeurs futures) mais bien de trouver un modèle valide, significatif et qui *fit* le mieux nos données observées. C'est la raison pour laquelle nous avons retenu comme modèle final un modèle de type :

$$\log \mathbf{X}_t \sim SARIMA(p = 0, d = 1, q = 3, P = 4, D = 1, Q = 2, s = 12)$$

L'estimation des différents coefficients amène à la modélisation suivante :

$$\begin{aligned} & (1 - B)(1 - B^{12})(1 + 0.0017B^{12} + 0.531B^{24} + 0.436B^{36} + 0.330B^{48}) \log \mathbf{X}_t \\ & = \\ & (1 + 0.441B + 0.048B^2 + 0.159B^3)(1 + 0.797B^{12} - 0.625B^{24}) \epsilon_t \end{aligned}$$

Partie 3 : Prévision

5&6. Écrire l'équation vérifiée par la région de confiance de niveau α sur les valeurs futures (Z_{T+1}, Z_{T+2}) . Préciser les hypothèses utilisées pour obtenir cette région (réponses en jaune).

Comme notre série $(Z_t)_{t=1,\dots,T}$ est un **SARIMA(0, 0, 3, 4, 0, 2, 12) stationnaire**, celle-ci admet une représentation ARMA canonique obtenue en supprimant les racines communes et en inversant les racines de module < 1 . De fait, Z_t est solution de

$$\Phi(B)Z_t = \Psi(B)\epsilon_t$$

avec B l'opérateur backward, Φ et Ψ polynômes de degré respectifs $p \leq P \times s = 48$ et $q \leq q' + Q \times s = 27$ n'ayant pas de racine commune ni de racines à l'intérieur du disque unité, et ϵ_t **l'innovation linéaire** de Z_t telle que **$\epsilon_t \sim \mathcal{N}(0, \sigma^2)$** et **$(\epsilon_t)_{t=1,\dots,T}$ iid**. Nous pouvons réécrire :

$$Z_t = \epsilon_t + \sum_{j=1}^{+\infty} c_j \epsilon_{t-j}, \quad \sum_{j=1}^{+\infty} c_j^2 < \infty$$

De fait, comme nous ne disposons que des valeurs de la série entre $t = 1$ et T , nous avons plutôt :

$$Z_t = \epsilon_t + \sum_{j=1}^{T-1} c_j \epsilon_{t-j}, \quad \sum_{j=1}^{T-1} c_j^2 < \infty$$

En projetant sur le passé fini de Z_T , on a donc pour tout $h \in \mathbb{N}$:

$$Z_{T+h|Z_T, \dots, Z_1} := \hat{Z}_{T+h|T} = \sum_{j=h}^{T+h-1} c_j \epsilon_{T+h-j} := EL(Z_{T+h}|Z_T, \dots, Z_1) = EL(Z_{T+h}|\epsilon_T, \dots, \epsilon_1)$$

où $\hat{Z}_{T+h|T}$ est la prévision linéaire optimale *tronquée* d'horizon $h > 0$ (au sens des moindres carrés). On en déduit au passage l'erreur de précision :

$$e(h) = Z_{T+h} - \hat{Z}_{T+h|T} = \sum_{j=0}^{h-1} c_j \epsilon_{T+h-j} \quad (1)$$

et l'erreur quadratique moyenne :

$$MSE(h) = \mathbb{E} \left(Z_{T+h} - \hat{Z}_{T+h|T} \right)^2 = \sigma^2 \sum_{j=0}^{h-1} c_j^2 \quad (2)$$

Étant donné que les résidus sont gaussiens, les équations (1) et (2) permettent de conclure que, conditionnellement à $Z_u, u \leq T$:

$$Z_{T+h} \sim \mathcal{N} \left(\hat{Z}_{T+h|T}, \sigma^2 \sum_{j=0}^{h-1} c_j^2 \right)$$

Pour finalement obtenir un intervalle de confiance de la valeur prédite de Z_{T+h} à 95%, $\forall h \in \mathbb{N}$ du type :

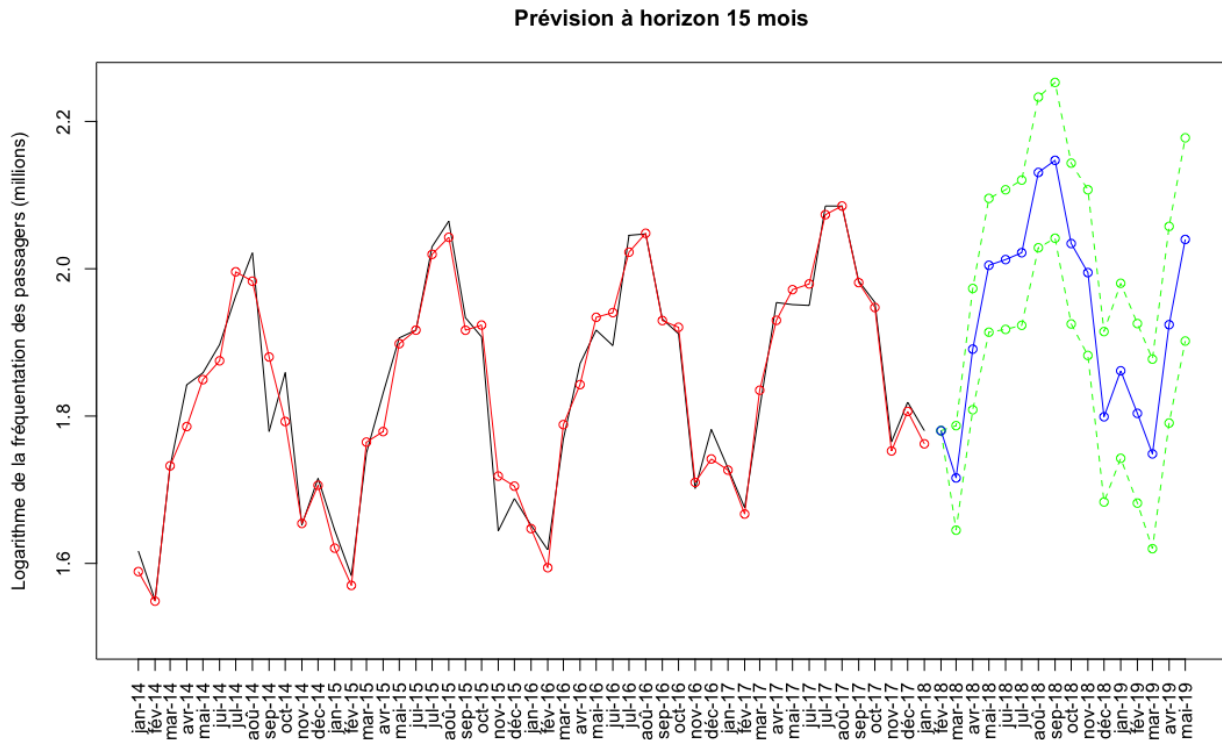
$$Z_{T+h|Z_T, \dots, Z_1} \in \left[\hat{Z}_{T+h|T} \pm \sigma \sqrt{\sum_{j=0}^{h-1} c_j^2} \times q_{1-\alpha/2}^N \right], \quad q_{1-\alpha/2}^N = 1,96$$

On peut également remplacer σ et les c_i par leurs estimateurs respectifs $\hat{\sigma}$ et \hat{c}_i .

7. Déterminer graphiquement cette région pour $\alpha = 95\%$. Commenter.

Nous avons affiché ci-dessous (figure 3) les prévisions à horizon 15 mois de notre modèle *SARIMA*. Est tracé en vert l'intervalle de confiance à 95% des valeurs prédites. Nous remarquons que l'intervalle de confiance s'élargit au fur et à mesure que la prédiction est lointaine. Ceci normal étant donné la dépendance en h au travers de l'expression $\sum_{j=0}^{h-1} c_j^2$ dans l'intervalle de confiance théorique. Vous trouverez en annexe (figure 5) l'évolution de cet intervalle de confiance en fonction des 15 mois à prédire.

Graphique 3 – Prévisions



8. Question ouverte : soit Y_t une série stationnaire disponible de $t = 1$ à T . On suppose que Y_{T+1} est disponible plus rapidement que Z_{T+1} . À quelles conditions cette information permet-elle d'améliorer la prévision de Z_{T+1} ?

Se demander à quelles conditions la prévision de Y_{T+1} permet d'améliorer la prévision de Z_{T+1} est équivalent à savoir si (Y_t) cause instantanément (Z_t) au sens de Granger, c'est à dire :

$$\mathbb{E}(Z_{T+1} | \{Z_u, Y_u, u \leq T\} \cup \{Y_{T+1}\}) \neq \mathbb{E}(Z_{T+1} | \{Z_u, Y_u, u \leq T\})$$

Cette inégalité est vérifiée si et seulement si les innovations de Z_{T+1} et Y_{T+1} sont non corrélées (en effet, d'après les hypothèses, la série $(X_t) := (Z_t, Y_t)$ est stationnaire). Or, nous avons supposé que les résidus de (Z_t) sont gaussiens. Étant donné que (ϵ_t) est un bruit blanc fort, les innovations sont donc indépendantes et identiquement distribuées de loi $\mathcal{N}(0, \sigma_\epsilon^2)$. Il faut donc distinguer les cas selon la loi des innovations de la série (Y_t) .

Dans le cas où les innovations de la série Y_t suivraient une loi gaussienne $\mathcal{N}(0, \sigma_y^2)$, un simple test de Student permet de tester la corrélation entre les innovations de la série Z_t et celles de la série Y_t .

Si nous n'avons pas d'information sur la loi des innovations de la série Y_t , mais que nous savons que les résidus sont i.i.d. (s'il s'agit un bruit blanc fort, par exemple), alors nous pouvons effectuer des tests de corrélations ne s'appuyant pas sur la distribution des variables (les tests de corrélation de Spearman ou de Kendall, par exemple).

En revanche, si nous n'avons pas d'information sur la loi des innovations de la série Y_t , et qu'en plus elles ne sont pas i.i.d., il est difficile de conclure sur l'amélioration de la prévision, car nous ne pouvons tester la corrélation entre les innovations de la série Z_t , et celles de la série Y_t .

Annexes

Graphique 4 – ACF et PACF de la série stationnarisée jusque lag 300

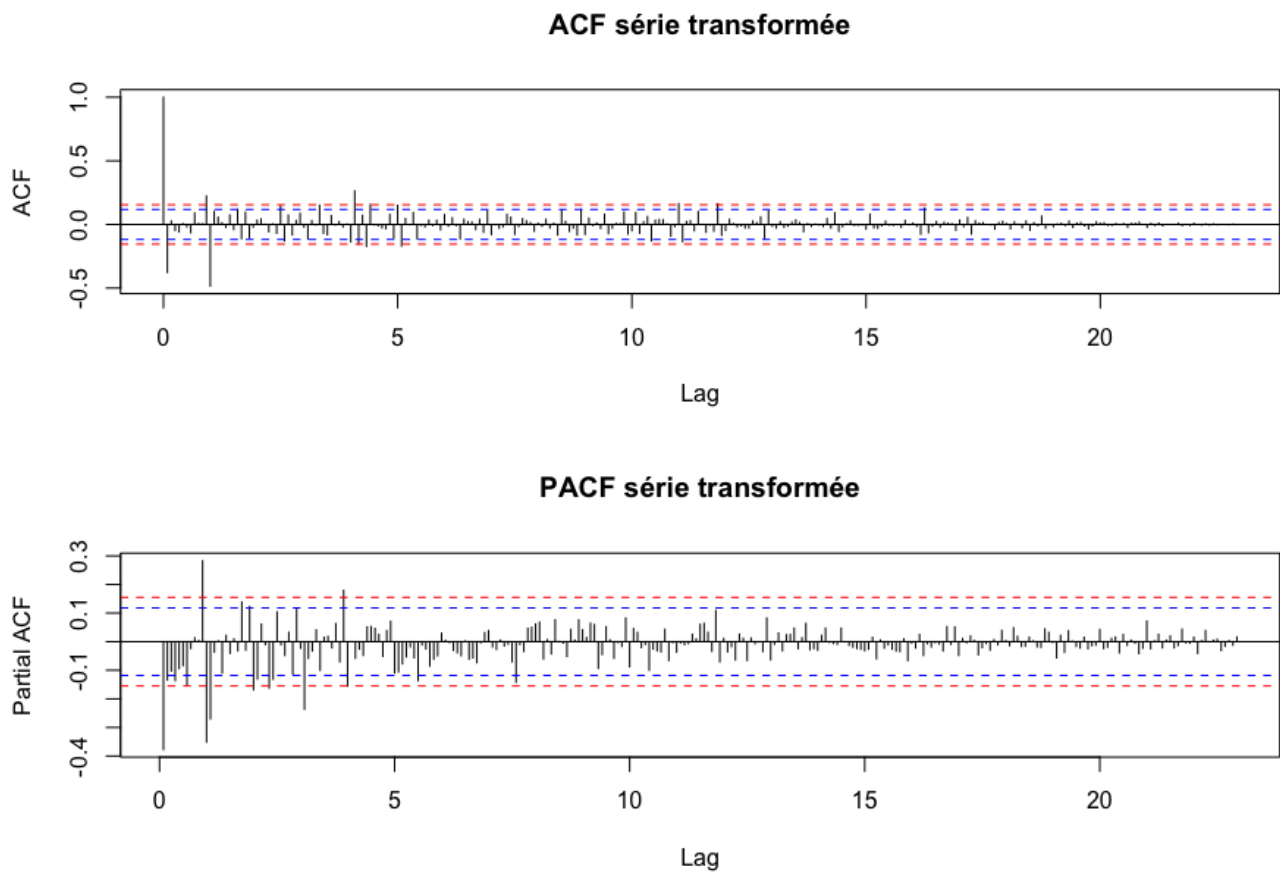


Tableau 1 – Modèles AR retenus

p	q	AIC	BIC
29	0	-106.70	5.53
25	0	-104.60	-6.85
28	0	-104.57	4.05
24	0	-100.23	-6.10
21	0	-98.99	-15.72

¹ $p_{max} = 30$.

Tableau 2 – Modèles MA retenus

p	q	AIC	BIC
0	13	-145.07	-90.76
0	20	-139.97	-60.32
0	27	-134.39	-29.40
0	26	-134.37	-33.00

¹ $q_{max} = 30$.

Tableau 3 – 20 premiers modèles ARMA retenus (sur 75)

p	q	AIC	BIC
1	13	-146.66	-88.73
2	15	-146.48	-77.69
6	19	-146.41	-48.66
5	16	-145.62	-62.35
8	20	-145.37	-36.76
0	13	-145.07	-90.76
18	14	-145.03	-21.93
9	20	-144.87	-32.64
3	15	-144.66	-72.26
12	20	-144.44	-21.35
7	12	-141.18	-65.15
3	14	-141.01	-72.23
19	18	-140.90	0.30
5	18	-140.78	-50.27
9	20	-139.97	-60.32
16	15	-139.91	-20.44
2	12	-139.43	-81.50
4	17	-139.21	-55.94
17	20	-138.38	2.82
4	16	-137.42	-57.77
⋮	⋮	⋮	⋮
3	10	-85.78	-31.48

¹ $p_{max} = 20, q_{max} = 20$.

Tableau 4 – 20 premiers modèles SAR(I)MA retenus (sur 182)

p	d	q	P	D	Q	s	AIC	BIC
0	0	3	4	0	2	12	-159.39	-119.56
4	0	5	4	0	1	12	-159.28	-101.35
4	0	5	4	0	2	12	-158.85	-97.30
4	0	5	0	0	5	12	-157.57	-99.65
0	0	3	2	0	4	12	-157.40	-117.58
3	0	3	4	0	2	12	-157.07	-106.38
5	0	4	2	0	4	12	-156.61	-95.06
0	0	3	3	0	5	12	-156.33	-109.26
5	0	2	2	0	4	12	-156.17	-101.86
5	0	4	3	0	3	12	-155.88	-94.33
5	0	5	0	0	5	12	-155.81	-94.27
3	0	2	0	0	5	12	-155.68	-112.23
3	0	5	2	0	4	12	-155.50	-97.58
4	0	3	3	0	4	12	-155.04	-97.12
0	0	1	4	0	2	12	-155.01	-122.43
0	0	1	5	0	3	12	-154.62	-114.80
0	0	3	0	0	5	12	-154.53	-118.32
0	0	3	3	0	3	12	-154.42	-114.60
4	0	2	4	0	2	12	-154.42	-103.73
4	0	1	4	0	3	12	-154.35	-103.67
0	0	3	5	0	0	12	-154.28	-118.08
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
5	0	4	1	0	0	12	-111.55	-68.11

¹ $p_{max} = 5, q_{max} = 5, P_{max} = 5, Q_{max} = 5$.

Tableau 5 – Erreurs de prédiction à horizon 4 dates et RMSE des 10 modèles retenus

Modèles	Critère minimisé	$\mathcal{E}_{prediction}$	RMSE
$SARIMA(0, 0, 3, 4, 0, 2, 12)$	AIC	0.02326554	0.6010692
$SARIMA(2, 0, 1, 0, 0, 1, 12)$	BIC	0.01947327	0.6296436
$SARIMA(0, 0, 1, 4, 0, 2, 12)$	-	0.02467597	0.6135300
$SARIMA(0, 0, 3, 2, 0, 4, 12)$	-	0.02277452	0.6032207
$SARIMA(0, 0, 1, 2, 0, 4, 12)$	-	0.01623702	0.6582885
$SARIMA(0, 0, 3, 0, 0, 5, 12)$	-	0.02662866	0.6212872
$SARIMA(0, 0, 3, 5, 0, 0, 12)$	-	0.02478279	0.6079626
$SARIMA(2, 0, 1, 0, 0, 5, 12)$	-	0.02810212	0.6227271
$SARIMA(0, 0, 1, 0, 0, 5, 12)$	-	0.01918966	0.6952696
$SARIMA(0, 0, 3, 4, 0, 0, 12)$	-	0.02534345	0.6243116

¹ $p_{max} = 5, q_{max} = 5, P_{max} = 5, Q_{max} = 5$.

² Erreurs minimales en jaune.

Graphique 5 – Largeur de l'intervalle de confiance à 95% dans le cadre de la prévision

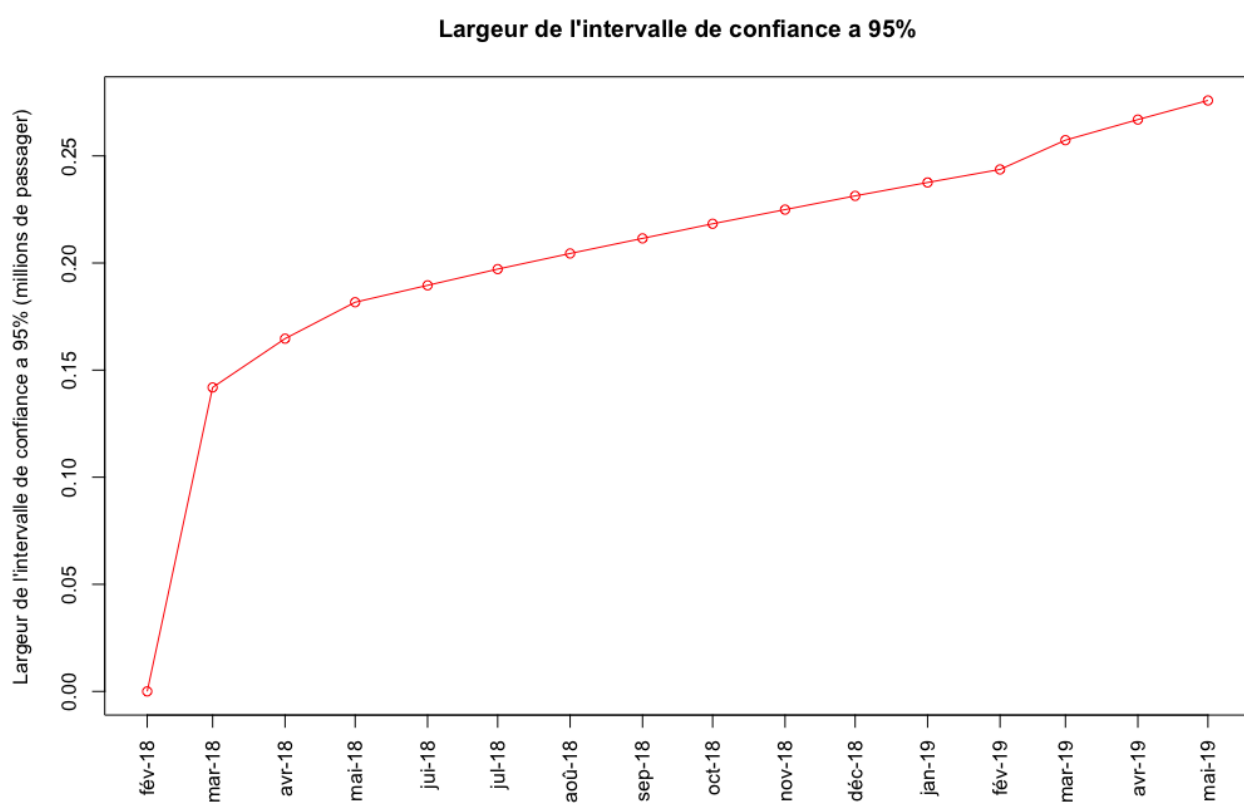


Table des figures

1	Différence entre série brute et série corrigée	3
2	ACF et PACF de la série stationnarisée jusque lag 30	4
3	Prévisions	6
4	ACF et PACF de la série stationnarisée jusque lag 300	8
5	Largeur de l'intervalle de confiance à 95% dans le cadre de la prévision	11
6	Logarithme de la fréquentation des passagers (millions) - Comparaison avant après	16

Liste des tableaux

1	Modèles AR retenus	8
2	Modèles MA retenus	9
3	20 premiers modèles ARMA retenus (sur 75)	9
4	20 premiers modèles SAR(I)MA retenus (sur 182)	10
5	Erreurs de prédiction à horizon 4 dates et RMSE des 10 modèles retenus	10

Code

Partie 1

```
1 #####
2 ##### Importation des packages #####
3 #####
4 library(zoo)
5 library(tseries)
6 library(ggplot2)
7 library(fUnitRoots)
8 library(scales)
9 library(lm)
10 library(tseries)
11 library(forecast)
12 library(xtable)
13
14 #####
15 ##### Importation et traitement des donnees #####
16 #####
17
18
19 setwd(dir="/Users/bfiliot/Desktop/ENSAE/S2/stprojet")
20 table = read.csv("valeurs_mensuelles.csv", head=T, sep=";")
21
22 names(table)[2] = "Frequentation de passagers vols internationaux - Paris"
23
24 table = table[-2,] # suppression des deux premieres lignes
25 table = table[-1,] # elles ne correspondent pas a des valeurs
26
27 table$year = substring(table[,1],1,4) # date au format annee / mois
28 table$month = substring(table[,1],6,7)
29
30 data = table[order(table[4],table[5]),] # ordonnancement du plus vieux au plus recent
31
32 xm = as.zoo(ts(data[[2]],
33               start = c(1994,1),
34               end = c(2018,1),
35               frequency = 12)) # transformation en objet zoo
36
37 dates = seq(as.POSIXct("1994-01-01"), by = "month", length.out = 289)
38
39 #####
40 ##### Affichage graphique de la serie #####
41 #####
42
43 # Utilisation de ggplot (uniquement ici)
44
45 p = ggplot(data = xm, aes(x = dates, y = xm)) +
46   geom_line(colour = 'blue') +
47   scale_x_datetime(labels = date_format("%Y"), breaks = date_breaks("years")) +
48   theme(axis.text.x = element_text(angle = 90)) +
49   labs(x = "Temps") +
50   labs(y = "Fréquentation des passagers (millions)") +
51   theme(text = element_text(size=20))
52   #labs(title = "Serie temporelle d'origine") +
53
54 p # graphe de la serie temporelle
55
56 p2 = ggplot(data = xm[1:36], aes(x = seq(as.POSIXct("1994-01-01"), by = "month", length.out = 36),
57   y = xm[1:36])) +
58   geom_line(colour = 'blue') +
59   scale_x_datetime(labels = date_format("%Y-%m"), breaks = date_breaks("months")) +
60   theme(axis.text.x = element_text(angle = 90)) +
61   labs(x = "Temps") +
62   labs(y = "Fréquentation des passagers (millions)")
63
64 p2 # graphe de la serie temporelle les 3 premieres annees
```

```

65
66 #####
67 ##### Stationnarisation de la serie #####
68 #####
69
70 # Pour ce genre de serie temporelle il est utile
71 # de transformer en log pour reduire l'heteroscedasticite
72 # des residus et/ou problemes de non-linearite.
73
74 xmlog = log(xm)
75
76 # Comme nous l'avons vu, la serie presente une tendance haussiere
77 # ainsi qu'une saisonnalite annuelle. La tendance haussiere laisse
78 # presager d'une serie stationnaire.
79
80 # Tout d'abord verifions la nature de la tendance deterministe.
81
82 T = length(xmlog)
83 time = 1:T
84 fit1 = lm(xmlog ~ poly(time,1))
85 fit2 = lm(xmlog ~ poly(time,2))
86 pred1 = predict(fit1)
87 pred2 = predict(fit2)
88
89 plot(dates,
90      xmlog,
91      type = 'l',
92      xlab = 'Temps',
93      ylab = 'Frequentation des passagers (millions)')
94
95 lines(dates, pred1, type = "l", col = "red", lty = 1, lwd = 2)
96 lines(dates, pred2, type = "l", col = "green", lty = 1, lwd = 2)
97
98 legend("bottomright",
99      legend = c("Série", "Approx. linéaire", "Approx. quadratique"),
100      col = c("black", "red", "green"),
101      text.col = c("black", "red", "green"),
102      lty = c(1,1,1),
103      ncol = 1,
104      cex = 0.5)
105
106 # Conclusion : tendance lineaire en a + bt. On se place donc dans le cas "ct".
107 # Notre serie log est-elle stationnaire ?
108
109 adf = adfTest(xmlog, lags=0, type="ct")
110 adf@test$p.value # < 0.01
111 # On rejette fortement la racine unite a 5%, la serie est donc stationnaire
112 # a 5%. Ce qui finalement n'est pas si contre-intuitif etant donne la tendance
113 # haussiere deterministe. Le modele est-il valide ?
114
115 # Regardons la validite du modele avec la persistance des residus.
116
117
118 Qtest <- function(x, kmax=24){
119   t(apply(matrix(1:kmax), 1, FUN=function(k){
120     pv = Box.test(x, lag=k, type="Ljung-Box")$p.value
121     return(c("lag"=k, "pval"=pv, "Non autocorrees a 5%" = (pv>0.05)))
122   }))
123 }
124
125 resid = adf@test$lm$residuals
126 plot(resid, type="l")
127 Qtest(adf@test$lm$residuals)
128 # Ici on voit que le test de blancheur des residus est rejete
129 # a 5% a tous les ordres. Les residus sont correles jusqu'a l'ordre 24.
130 # Le modele n'est donc pas valide si l'on considere un lag = 0 dans
131 # la formule de l'ADF.
132 # Mais nous pourrions augmenter ce lag de telle sorte a voir a partir duquel
133 # les residus ne sont pas correles (a 5%) et ceci pour tous les 24 ordres.
134 # Ceci est verifie pour le lag 23. C'est trop, il faut pouvoir diminuer

```

```

135 # ce lag. Ce que nous decidons de faire est de garder un test adf avec un lag = 0
136 # et tester la blancheur des residus a l'ordre 1. Des qu'une transformation
137 # (elle existe par anticipation...) de la serie xmlog sera telle que
138 # le residu est non correle a l'ordre 1 a 5%, on choisira cette transformation
139 # et on regardera le lag minimum a partir duquel l'ADF/le modele est valide.
140 # On prend tout d'abord en compte la saisonnalite.
141
142 xdesaison = diff(xmlog,12)
143 adf = adfTest(xdesaison, lags=0, type="ct")
144 adf@test$p.value
145 # On rejette aussi la racine unite a 5%, la serie differenciee est stationnaire
146 # ce qui est logique.
147 # Ici on remarque que la p-valeur est superieure a 5% mais tres proche (0.05417)
148 resid = adf@test$lm$residuals
149 plot(resid, type="l")
150 Box.test(resid, lag=1, type="Ljung-Box")$p.value
151 # On decide donc de differencier egalement la serie pour avoir un rejet
152 # de la racine unitaire plus net.
153
154 # Appliquons donc une desaisonnalisation et une differenciation.
155
156 xdd = diff(desaison,1)
157 adf = adfTest(xdd, lags=0, type="c") #lag 11 pour 12 premiers ; lag 13 pour 24 premiers
158 adf@test$p.value
159 # On rejette la racine unite a 5%, la serie differenciee et desaisonnalisee
160 # est stationnaire ce qui est logique.
161 resid = adf@test$lm$residuals
162 plot(resid, type="l")
163 Box.test(resid, lag=1, type="Ljung-Box")$p.value
164 # Et cette fois l'absence d'autocorrelation n'est pas rejetee a l'ordre 1 a 5%
165 # car elle est egale a 0.39516. Nous gardons donc cette transformation.
166 # Verifions a partir de quel lag (dans la formule de l'ADF) le modele est valide,
167 # i.e les residus sont non-autocorreles jusqu'a l'ordre 24 a 5%.
168 # Au passage on peut verifier qu'a lag 0 le modele n'est pas valide meme
169 # si pour certains ordres les residus ne sont pas autocorreles a 5%.
170
171 # On s'intéresse au lag minimum tel que les residus ne sont pas autocorreles jusqu'a l'ordre 24.
172
173 lagmin <- function(serie){
174   i = 0
175   if (serie == xdd){
176     type = "nc" # la serie desaisonnalisee et differenciee est centree
177   }
178   if (serie == xdesaison){
179     type = "c" # la serie desaisonnalisee est non centree
180   }
181   if (serie == xmlog){
182     type = "ct" # la serie xmlog presente une tendance lineaire deterministe
183   }
184   adf = adfTest(serie, lags=i, type=type)
185   bbxtest = Qtest(adf@test$lm$residuals)
186   while (sum(bbxtest[,3]) < 24) {
187     i = i+1
188     adf = adfTest(serie, lags=i, type=type)
189     bbxtest = Qtest(adf@test$lm$residuals)
190   }
191   return("lagminADF"=i)
192 }
193
194 lagmin(xdd)
195 # On trouve un lag min de 13 ce qui est bien mieux que le 23 pour
196 # la serie non transformee.
197 # On a donc un modele ADF valide lorsque celui-ci prend en compte
198 # 13 retards.
199 adfTest(xdd, lags=13, type="nc")
200 pp.test(xdd) # L'hypothese de racine unitaire est tres fortement rejetee.
201 kpss.test(xdd) # L'hypothese de stationnarite n'est pas rejetee a 5%.
202
203 # On se contentera donc de notre transformation dans la mesure ou
204 # elle garantit la non-autocorrelation des residus jusqu'a l'ordre 24
205 # pour un lagADF egal a 13, tout ceci pour des seuils de 5%.
206 # En soit nous aurions juste desaisonnaliser la serie...
207 # Car on remarque que le lag minimal de l'ADF pour lequel les residus

```

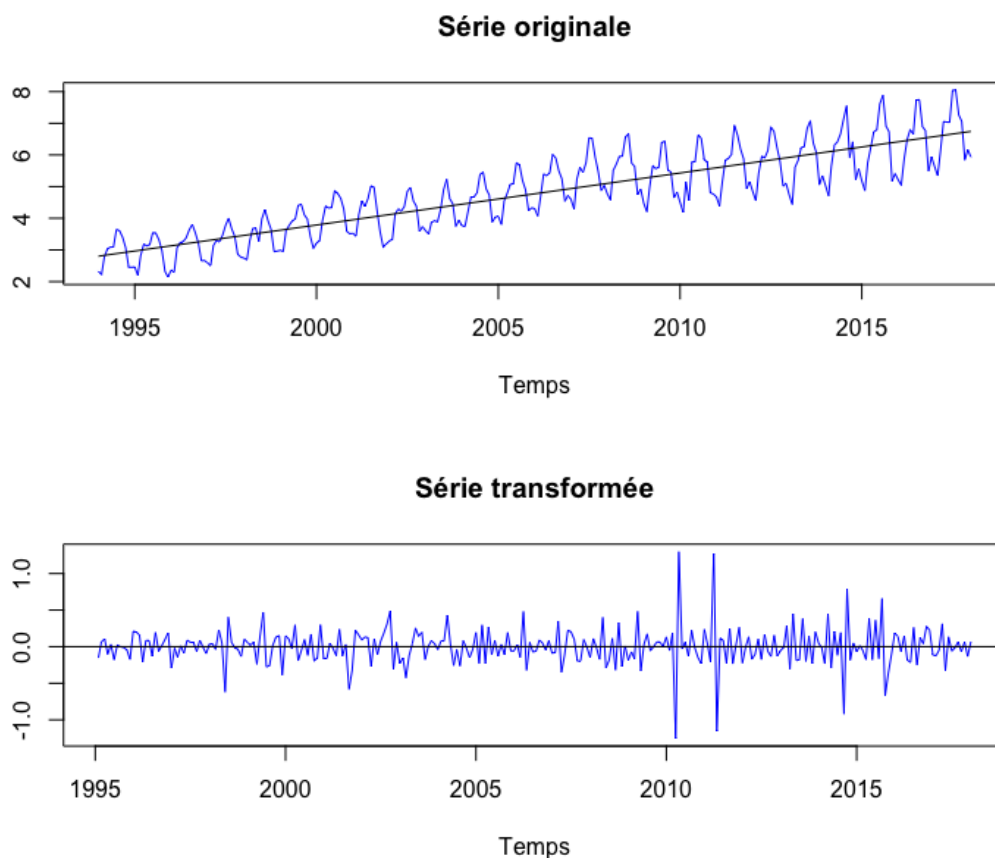


```

205 # sont non-autocorrelés à tous les ordres est... ->
206 lagmin(xdesaison,50)
207 adfTest(xdesaison, lags=13, type="c")
208 pp.test(xdesaison) # L'hypothèse de racine unitaire est très fortement rejetée.
209 kpss.test(xdesaison)
210 #... -> 13, comme pour la série désaisonnalisée. On retiendra donc qu'a
211 # lagADF = 0, les résidus de la série désaisonnée sont "presque"
212 # auto-correlés à l'ordre 1 au niveau 5% (le terme "presque" est
213 # peu rigoureux) tandis que non pour la série xdd désaisonnée et différenciée.
214
215 # Différence avant / après
216
217 par(mfrow = c(1,2)) #ou (2,1)
218 plot(dates,
219      xm,
220      type = 'l',
221      xlab = "Temps",
222      ylab = 'Logarithme de la fréquentation des passagers (millions)',
223      col = 'blue',
224      main = 'Série originale')
225
226 plot(dates[14:289],
227      xdd,
228      type = 'l',
229      xlab = "Temps",
230      ylab = 'Logarithme de la fréquentation des passagers (millions)',
231      col = 'blue',
232      main = 'Série transformée')
233 abline(a=0, b=0)

```

Graphique 6 – Logarithme de la fréquentation des passagers (millions) - Comparaison avant après



Partie 2

```
1 #####
2 ##### Calibrage d'un modèle ARMA #####
3 #####
4
5 ##### Fonctions prealables #####
6
7 # Fonction de test de l'absence d'autocorrelation des residus
8 Qtest <- function(x, kmax=24){
9   t(apply(matrix(1:kmax), 1, FUN=function(k){
10     pv = Box.test(x, lag=k, type="Ljung-Box")$p.value
11     return(c("lag"=k, "pval"=pv, "Non autocorrees a 5%" = (pv>0.05)))
12   })))
13 }
14
15 # Fonction de test des significativites individuelles des coefficients
16 signif <- function(estim){
17   coef = estim$coef
18   se = sqrt(diag(estim$var.coef))
19   t = coef/se
20   pval = (1-pnorm(abs(t)))*2
21   return(rbind(coef, se, pval))
22 }
23
24 # Fonction annexe
25 tf_binaire <- function(signif){
26   if (is.na(signif)) return(1)
27   else return(signif)
28 }
29
30 # Selection des valeurs pmax, qmax, Pmax, Qmax #
31 # Pour cela nous sommes allés voir comment sont calculés les
32 # limites de l'intervalle de confiance a 95% pour l'acf (et donc pacf)
33 # via la fonction
34 getS3method("plot", "acf")
35 # et nous avons retrouvé la variable clim égale à la borne
36 par(mfrow=c(2,1))
37 clim95 = qnorm((1 + 0.95)/2)/sqrt(length(xdd))
38 acf = acf(xdd, lag.max = 300, main = "ACF série transformée")
39 pacf = pacf(xdd, lag.max = 300, main = "PACF série transformée")
40 # Les 2 valeurs coïncident.
41
42 # On veut déterminer explicitement quels sont les lags pour lesquels les
43 # valeurs observées sont hors de la zone de significativité à 5%.
44
45 which(abs(acf$acf) > clim95)
46 which(abs(pacf$acf) > clim95)
47
48 # On remarque que les valeurs des autocorrélations pour des lag
49 # entre q=2 et q=11 sont dans l'intervalle de confiance à 95% donc on peut
50 # supposer légitimement que qmax = 3.
51 # Pour la PACF, c'est un peu plus délicat. En effet les autocorrélations
52 # partielles de lag 1,2,4,5,7,11,12,13,21 ne sont pas dans l'intervalle de
53 # confiance à 95%. On peut considérer que l'écart entre lag 13 et lag 21
54 # permet de prendre pmax = 14.
55 # Au niveau des ordres Pmax et Qmax, qu'en est-il ?
56 # Etant donné que les multiples de 12 sont : 12,24,36,48,60,72,84,96,108,120,132,144
57 # on peut prendre Pmax = 4 et on pourrait aller jusqu'à Qmax = 12...
58
59
60
61
62
63
64
65
66
67
68
```

```

69 # Si l'on regarde avec un intervalle de confiance a 99%... ->
70
71 clim99 = qnorm((1 + 0.99)/2)/sqrt(length(xdd))
72 par(mfrow = c(2,1))
73 acf = acf(xdd, lag.max = 300, main = "ACF série transformée")
74 abline(clim99, 0, col = "red", lty = 2)
75 abline(-clim99, 0, col = "red", lty = 2)
76 pacf = pacf(xdd, lag.max = 300, main = "PACF série transformée")
77 abline(clim99, 0, col = "red", lty = 2)
78 abline(-clim99, 0, col = "red", lty = 2)
79
80 which(abs(acf$acf) > clim99)
81 which(abs(pacf$acf) > clim99)
82
83 #... -> on peut prendre Qmax = 5.
84
85 # On a donc les ordres maximums vraisemblables
86 # (p*,d*,q*,P*,D*,Q*,s) = (14,0,3,4,0,5)
87
88 # La fonction suivante automatise le processus
89 # de significativité et validation puis calcul
90 # des AIC, BIC.
91
92 modelchoice <- fonction(p,q,P,Q,serie ,k=24){
93
94   print(paste0("SARIMA (p=",p," ,d=",0," ,q=",q," ,P=",P," ,D=",0," ,Q=",Q," ,s=",12,")"))
95
96   estim = try(arima(serie , order=c(p,0,q) , seasonal=list(order = c(P,0,Q) , period=12) ,
97                 method="ML" ,
98                 optim.control=list(maxit=20000)))
99
100  if (class(estim)=="try-error") return(c("p"=p,"d"=0,"q"=q,"P"=P,"D"=0,"Q"=Q,"s"=12,
101                                         "ARsignif" = NA,
102                                         "MAsignif" = NA,
103                                         "ARSAISONsignif" = NA,
104                                         "MASAISONsignif" = NA,
105                                         "ModeleSignif" = NA,
106                                         "ModeleValide" = NA,
107                                         "Selectionne" = NA,
108                                         "AIC" = NA,
109                                         "BIC" = NA))
110
111  arsignif = if (p==0) NA else signif(estim)[3,p]<=0.05
112  masignif = if (q==0) NA else signif(estim)[3,p+q]<=0.05
113  ARsaisonsignif = if (P==0) NA else signif(estim)[3,p+q+P]<=0.05
114  MASaisonsignif = if (Q==0) NA else signif(estim)[3,p+q+P+Q]<=0.05
115
116  ModeleSignif = prod(apply(matrix(c(arsignif , masignif , ARsaisonsignif , MASaisonsignif) , 1,
117                                     tf_binaire)))
118
119  resnocorr = sum(Qtest(estim$residuals , k)[,3])=k
120
121  aic = AIC(estim)
122  bic = BIC(estim)
123
124  return(c("p"=p,"d"=0,"q"=q,"P"=P,"D"=0,"Q"=Q,"s"=12,
125          "ARsignif" = arsignif ,
126          "MAsignif" = masignif ,
127          "ARSAISONsignif" = ARsaisonsignif ,
128          "MASAISONsignif" = MASaisonsignif ,
129          "ModeleSignif" = ModeleSignif ,
130          "ModeleValide" = resnocorr ,
131          "Selectionne" = ModeleSignif*resnocorr ,
132          "AIC" = aic ,
133          "BIC" = bic))
134
135
136
137

```

```

138 # On verifie si des modeles AR sont selectionnes...
139
140 pmax = 30 #par extreme securite
141 resultAR = data.frame(t(mapply(function(p) modelchoice(p,0,0,0,xdd), c(0:pmax))))
142 resultAR = resultAR[with(resultAR, order(resultAR$Selectionne, decreasing = TRUE)),]
143 View(resultAR)
144 ARselec = resultAR[resultAR$Selectionne == 1,]
145 ARselec = ARselec[with(ARselec, order(ARselec$AIC, decreasing = FALSE)),]
146 View(ARselec)
147
148 # On selectionne AR(29), AR(25), AR(28), AR(24), AR(21)... dans l'ordre des
149 # AIC les plus faibles. Les ordres p sont beaucoup trop eleves pour
150 # ce type de modelisation. On ecarte donc une modelisation de la serie
151 # xdd par un AR. D'autant plus que les BIC sont tres grands.
152
153 # Des modeles MA ?
154
155 qmax = 30 #par extreme securite
156 resultMA = data.frame(t(mapply(function(q) modelchoice(0,q,0,0,xdd), c(0:qmax))))
157 resultMA = resultMA[with(resultMA, order(resultMA$Selectionne, decreasing = TRUE)),]
158 View(resultMA)
159 Maselec = resultMA[resultMA$Selectionne == 1,]
160 Maselec = Maselec[with(Maselec, order(Maselec$AIC, decreasing = FALSE)),]
161 View(Maselec)
162
163 # On selectionne MA(13), MA(20), MA(27), MA(26)... dans l'ordre des
164 # AIC les plus faibles. On retient donc le MA(13) puisqu'il minimise
165 # a la fois l'AIC et le BIC.
166
167 # Des modeles ARMA ?
168
169 pmax = 20
170 qmax = 20
171 resultARMA = data.frame(t(mapply(function(p,q) modelchoice(p,q,0,0,xdd),
172                                c(apply(matrix(c(0:pmax)),1, function(k) rep(k,(qmax+1)))),
173                                c(0:qmax))))
174
175 resultARMA = resultARMA[with(resultARMA, order(resultARMA$Selectionne, decreasing = TRUE)),]
176 View(resultARMA)
177 ARMAselec = resultARMA[resultARMA$Selectionne == 1,]
178 ARMAselec = ARMAselec[with(ARMAselec, order(ARMAselec$AIC, decreasing = FALSE)),]
179 View(ARMAselec)
180 length(which(is.na(ARMAselec[,16])==0))
181
182 # On garde 75 modeles ARMA au total.
183 # Les deux modeles ARMA optimaux sont ARMA(0,13) et ARMA(1,13). Ils minimisent
184 # chacun un critère.
185
186 # Des modeles SARIMA ?
187
188 pmax = 5
189 qmax = 5
190 Pmax = 5
191 Qmax = 5
192 resultSARIMA = data.frame(t(mapply(function(p,q,P,Q) modelchoice(p,q,P,Q,xdd),
193                                c(apply(matrix(c(0:pmax)),1, function(k) rep(k,(qmax+1)*(Pmax+1)*(Qmax+1)))),
194                                c(apply(matrix(c(0:qmax)),1, function(k) rep(k,(Pmax+1)*(Qmax+1)))),
195                                c(apply(matrix(c(0:Pmax)),1, function(k) rep(k,(Qmax+1)))),
196                                c(0:Qmax))))
197 resultSARIMA = resultSARIMA[with(resultSARIMA, order(resultSARIMA$Selectionne, decreasing = TRUE)
198                                ),]
199 View(resultSARIMA)
200 SARIMAselec = resultSARIMA[resultSARIMA$Selectionne == 1,]
201 SARIMAselec = SARIMAselec[with(SARIMAselec, order(SARIMAselec$AIC, decreasing = FALSE)),]
202 View(SARIMAselec)
203 length(which(is.na(SARIMAselec[,16])==0))
204
205 # On selectionne 182 modeles. Les 5 premiers sont par ordre croissant d'AIC:
206 # SARIMA(0,0,3,4,0,2,12), SARIMA(4,0,5,4,0,1,12), SARIMA(4,0,5,4,0,2,12)
207 # SARIMA(4,0,5,0,0,5,12), SARIMA(0,0,3,2,0,4,12)

```

```

207 # On concatene tous les modeles retenus.
208
209 resultats_retenus      = rbind(ARselec,MAselec,ARMAselec,SARIMAsselec)
210 resultats_retenus      = resultats_retenus[with(resultats_retenus,
211                                     order(resultats_retenus$AIC, decreasing = FALSE)),]
212 resultats_retenus["AIC+BIC"] = resultats_retenus$AIC + resultats_retenus$BIC
213 resultats_retenus["AIC+BIC"] = resultats_retenus$AIC + resultats_retenus$BIC
214 View(resultats_retenus)
215 length(which(is.na(resultats_retenus[,16])==0))
216
217 # Au total, ce sont 266 modeles qui sont retenus. Bien sur,
218 # il reste maintenant a selectionner le ou les meilleurs du point
219 # de la minimisation des criteres.
220
221 # Pour le latex :
222 print(xtable(ARselec,      type = "latex", tabular.environment = "longtable"),
223       file = "ARselec.tex")
224 print(xtable(MAselec,      type = "latex", tabular.environment = "longtable"),
225       file = "MAselec.tex")
226 print(xtable(ARMAselec,    type = "latex", tabular.environment = "longtable"),
227       file = "ARMAselec.tex")
228 print(xtable(SARIMAsselec, type = "latex", tabular.environment = "longtable"),
229       file = "SARIMAsselec.tex")
230
231 #####
232 ##### Selection d'un modèle ARMA #####
233 #####
234
235 # Conclusion : il nous reste a comparer les dix premiers modeles suivants :
236 # celui maximisant l'AIC : SARIMA(0,0,3,4,0,2,12)
237 # celui maximisant le BIC : SARIMA(2,0,1,0,0,1,12)
238 # 8 autres modèles ayant une somme de criteres AIC + BIC parmi les plus petites :
239 # SARIMA(0,0,1,4,0,2,12) ; SARIMA(0,0,3,2,0,4,12)
240 # SARIMA(0,0,1,2,0,4,12) ; SARIMA(0,0,3,0,0,5,12)
241 # SARIMA(0,0,3,5,0,0,12) ; SARIMA(2,0,1,0,0,5,12)
242 # SARIMA(0,0,1,0,0,5,12) ; SARIMA(0,0,3,4,0,0,12)
243 # Ces modeles sont bien sur valides et bien ajustes.
244 # Les modèles ne sont pas imbriques, on ne peut donc pas faire un LR test
245 # (rapport de vraisemblance). Il nous reste un critere de prevision pour choisir
246 # entre les 10 modeles. On choisit de comparer les erreurs de prediction
247 # pour selectionner un unique modele.
248
249 T      = length(xmlog)
250 trend  = 1:(T-4)
251 xmlog_tronq = xmlog[trend]
252 obs     = xmlog[(T-3):T]
253 lt      = lm(xmlog[trend] ~ poly(trend,2))
254
255 # Ici on rebascule sur xmlog, on doit donc le prendre en compte : d = 1, D = 1
256 # faisant reference a notre differenciation et desaisonnalisation.
257 sarima1 = arima(xmlog_tronq, order=c(0,1,3), seasonal=list(order = c(4,1,2), period=12), method=
258 "ML", optim.control=list(maxit=20000))
259 sarima2 = arima(xmlog_tronq, order=c(2,1,1), seasonal=list(order = c(0,1,1), period=12), method=
260 "ML", optim.control=list(maxit=20000))
261 sarima3 = arima(xmlog_tronq, order=c(0,1,1), seasonal=list(order = c(4,1,2), period=12), method=
262 "ML", optim.control=list(maxit=20000))
263 sarima4 = arima(xmlog_tronq, order=c(0,1,3), seasonal=list(order = c(2,1,4), period=12), method=
264 "ML", optim.control=list(maxit=20000))
265 sarima5 = arima(xmlog_tronq, order=c(0,1,0), seasonal=list(order = c(2,1,4), period=12), method=
266 "ML", optim.control=list(maxit=20000))
267 sarima6 = arima(xmlog_tronq, order=c(0,1,3), seasonal=list(order = c(0,1,5), period=12), method=
268 "ML", optim.control=list(maxit=20000))
269 sarima7 = arima(xmlog_tronq, order=c(0,1,3), seasonal=list(order = c(5,1,0), period=12), method=
270 "ML", optim.control=list(maxit=20000))
271 sarima8 = arima(xmlog_tronq, order=c(2,1,1), seasonal=list(order = c(0,1,5), period=12), method=
272 "ML", optim.control=list(maxit=20000))
273 sarima9 = arima(xmlog_tronq, order=c(0,1,0), seasonal=list(order = c(0,1,5), period=12), method=
274 "ML", optim.control=list(maxit=20000))
275 sarima10 = arima(xmlog_tronq, order=c(0,1,3), seasonal=list(order = c(4,1,0), period=12), method=
276 "ML", optim.control=list(maxit=20000))

```

```

267
268 modelComp = list(sarima1, sarima2, sarima3, sarima4, sarima5, sarima6, sarima7, sarima8, sarima9,
269                 sarima10)
270 # On introduit la fonction rmserror qui calcule la somme des carres des
271 # residus sur l'ensemble de la periode tronquee.
272
273 rmserror <- function(model) sqrt(sum(model$residuals^2))
274
275 # Et on les calcule.
276
277 i = 0
278 predmodel = rep(list(list()),10)
279 erreur = c("SARIMA(0,0,3,4,0,2,12)minAIC"= 0,
280           "SARIMA(2,0,1,0,0,1,12)minBIC"= 0,
281           "SARIMA(0,0,1,4,0,2,12)" = 0, "SARIMA(0,0,3,2,0,4,12)" = 0,
282           "SARIMA(0,0,1,2,0,4,12)" = 0, "SARIMA(0,0,3,0,0,5,12)" = 0,
283           "SARIMA(0,0,3,5,0,0,12)" = 0, "SARIMA(2,0,1,0,0,5,12)" = 0,
284           "SARIMA(0,0,1,0,0,5,12)" = 0, "SARIMA(0,0,3,4,0,0,12)" = 0)
285
286 RMSE = erreur
287
288 for (model in modelComp){
289   i = i+1
290   print(model)
291
292   predmodel[[i]] = as.zoo(ts(predict(model,4)$pred,
293                               start = c(2017,10),
294                               end = c(2018,1),
295                               frequency = 12))
296
297   erreur[i] = sqrt(sum((predmodel[[i]]-obs)^2)/4)
298   RMSE[i] = rmserror(model)
299 }
300
301 # Erreur correspond a la somme des carres des erreurs de prediction
302 # sur les 4 derniers mois de la serie xmlog pour chaque modele.
303 # RMSE regroupe les RMSE des differents modeles.
304
305 erreur = data.frame(erreur)
306 colnames(erreur) = c("Erreur prediction")
307
308 RMSE = data.frame(RMSE)
309 colnames(RMSE) = c("RMSE")
310
311 resultsPRED = cbind(erreur,RMSE) # resultsPRED regroupe l'ensemble des erreurs
312 View(resultsPRED)
313
314 print(xtable(resultsPRED, type = "latex", tabular.environment="longtable"), file = "resultsPRED.
315         tex")
316
317 # On remarque que l'erreur de prediction est minimisee par
318 # le SARIMA(0,0,1,2,0,4,12). Le RMSE est minimise
319 # par le SARIMA(0,0,3,4,0,2,12), i.e le SARIMA minimisant
320 # l'AIC. On remarque que l'erreur de prediction a horizon 4 periodes
321 # est plus petite pour le SARIMA minimisant le BIC. Cependant
322 # le RMSE de ce meme modele est plus faible.
323 # On a donc d'un cote un modele minimisant l'AIC et
324 # qui fit mieux la serie sur les donnees observees (en fait,
325 # parmi les modeles retenus, il fit LE mieux la serie xmlog) mais predit moins bien.
326 # De l'autre, un autre modele minimisant le BIC et
327 # predit mieux mais qui fit moins bien la serie sur les donnees observees.
328 # Enfin, il existe 2 autres modeles ne minimisant aucun des 2 criteres
329 # mais qui pourtant predisent mieux a horizon 4 periodes.
330
331
332
333
334

```

```

335 # On affiche les erreurs de prediction.
336
337 pred = plot(dates[(T-3):T],
338             obs,
339             type = 'o',
340             xlab = "Temps",
341             ylim = c(1.7,2),
342             ylab = 'Logarithme de la fréquentation des passagers (millions)',
343             col = 'black',
344             main = 'Prévision',
345             xaxt = "n")
346
347 axis(side=1, at=dates[(T-3):T], labels=format(dates[(T-3):T], '%b-%y'))
348
349 lines(predmodel[[1]], type = "o", col = "red", lty = 1, lwd = 2)
350 lines(dates[(T-3):T], predmodel[[2]], type = "o", col = "blue", lty = 1, lwd = 2)
351 lines(dates[(T-3):T], predmodel[[3]], type = "o", col = "green", lty = 1, lwd = 1)
352 lines(dates[(T-3):T], predmodel[[4]], type = "o", col = "gray", lty = 1, lwd = 1)
353 lines(dates[(T-3):T], predmodel[[5]], type = "o", col = "orange", lty = 1, lwd = 1)
354 lines(dates[(T-3):T], predmodel[[6]], type = "o", col = "darkblue", lty = 1, lwd = 1)
355 lines(dates[(T-3):T], predmodel[[7]], type = "o", col = "cyan", lty = 1, lwd = 1)
356 lines(dates[(T-3):T], predmodel[[8]], type = "o", col = "brown", lty = 1, lwd = 1)
357 lines(dates[(T-3):T], predmodel[[9]], type = "o", col = "yellow", lty = 1, lwd = 1)
358 lines(dates[(T-3):T], predmodel[[10]], type = "o", col = "pink", lty = 1, lwd = 1)
359
360 legend(y = 2.0,
361        x = dates[[T-1]],
362        bty = "n",
363        legend = c("Valeurs observées", "SARIMA(0,0,3,4,0,2,12) minimiseur AIC",
364                  "SARIMA(2,0,1,0,0,1,12) minimiseur BIC",
365                  "SARIMA(0,0,1,4,0,2,12)", "SARIMA(0,0,3,2,0,4,12)",
366                  "SARIMA(0,0,1,2,0,4,12)", "SARIMA(0,0,3,0,0,5,12)",
367                  "SARIMA(0,0,3,5,0,0,12)", "SARIMA(2,0,1,0,0,5,12)",
368                  "SARIMA(0,0,1,0,0,5,12)", "SARIMA(0,0,3,4,0,0,12)"),
369        col = c("black", "red", "blue", "green", "gray", "orange", "darkblue", "cyan",
370               "brown", "yellow", "pink"),
371        text.col = c("black", "red", "blue", "green", "gray", "orange", "darkblue", "cyan",
372                    "brown", "yellow", "pink"),
373        lty = c(1,1,1,1,1,1,1,1,1,1),
374        pch = c(NA,1,1,1,1,1,1,1,1,1),
375        y.intersp = 0.5,
376        ncol = 1)
378

```

Partie 3

```
1 #####
2 ##### Predictions #####
3 #####
4
5 # Modele retenu
6
7 model = arima(xmlog, order=c(0,1,3), seasonal=list(order = c(4,1,2), period=12), method="ML",
8             optim.control=list(maxit=20000))
9
10 # On realise une prevision sur 15 mois
11 lag = 15
12 lastobs = xmlog[T]
13 forecasts = forecast(model, h=lag, level=95)
14 prediction = as.zoo(ts(c(lastobs, forecasts$mean), frequency = 12, start=c(2018,1)))
15 up = as.zoo(ts(c(lastobs, forecasts$upper), frequency = 12, start=c(2018,1)))
16 down = as.zoo(ts(c(lastobs, forecasts$lower), frequency = 12, start=c(2018,1)))
17 datesPRED = seq(as.POSIXct("2018-02-01"), by = "month", length.out = lag+1)
18
19 plot(c(dates[(T-48):T], datesPRED),
20      c(xmlog[(T-48):T], as.zoo(ts(rep(NA, length(datesPRED)), start = c(2018,2), frequency = 12))),
21      type = 'l',
22      xlab = '',
23      ylab = 'Logarithme de la fréquentation des passagers (millions)',
24      col = "black",
25      ylim = c(1.5, 2.25),
26      xaxt = "n",
27      main = paste0("Prévision à horizon ", lag, " mois"))
28
29 axis(side=1, at=c(dates[(T-48):T], datesPRED), labels=format(c(dates[(T-48):T], datesPRED),
30                        '%b-%y'), las=2)
31
32 lines(dates[(T-48):T], fitted(model)[(T-48):T], type = "o", col = "red", lty = 1, lwd = 1)
33 lines(datesPRED, up, type = "o", col = "green", lty = 2, lwd = 1)
34 lines(datesPRED, down, type = "o", col = "green", lty = 2, lwd = 1)
35 lines(datesPRED, prediction, type = "o", col = "blue", lty = 1, lwd = 1)
36
37 legend("topleft",
38       bty = "n",
39       legend = c("Série observée", "Modèle SARIMA", "Prédiction",
40                 "Borne supérieures et inférieures à 95%"),
41       col = c("black", "red", "blue", "green"),
42       text.col = c("black", "red", "blue", "green"),
43       lty = c(1, 1, 1, 2),
44       pch = c(NA, 1, 1, 1, 1),
45       ncol = 1,
46       y.intersp = 0.5)
47
48 # On verifie bien que l'intervalle de confiance grandit selon que la
49 # prevision se fait a horizon plus large.
50
51 plot(datesPRED,
52      up-down,
53      type = 'o',
54      xlab = '',
55      ylab = "Largeur de l'intervalle de confiance a 95% (millions de passager)",
56      col = "red",
57      xaxt = "n",
58      main = "Largeur de l'intervalle de confiance a 95%")
59
60 axis(side=1, at=datesPRED, labels=format(datesPRED, '%b-%y'), las=2)
```