

Séries chronologiques - Partie 3

Désaisonnalisation et calcul des coefficients saisonniers

BUT Science des Données, deuxième année

On dispose d'une série chronologique $(y_i)_{i=1,\dots,n}$ où les composantes présentes sont

- La **tendance**.
- La **saisonnalité** (composante périodique).
- La **composante résiduelle**.

On considèrera deux types de modèles :

- Le modèle **additif**.
- Le modèle **multiplicatif**.

Soit $(y_i)_{i=1,\dots,n}$ une série chronologique.

On considère le modèle :

$$y_i = f_i + s_i + e_i \text{ pour } i = 1, \dots, n$$

où :

- $(f_i)_{i=1,\dots,n}$ désigne la tendance.
- $(s_i)_{i=1,\dots,n}$ désigne la composante saisonnière de période p telle que :

$$\sum_{j=1}^p s_j = 0.$$

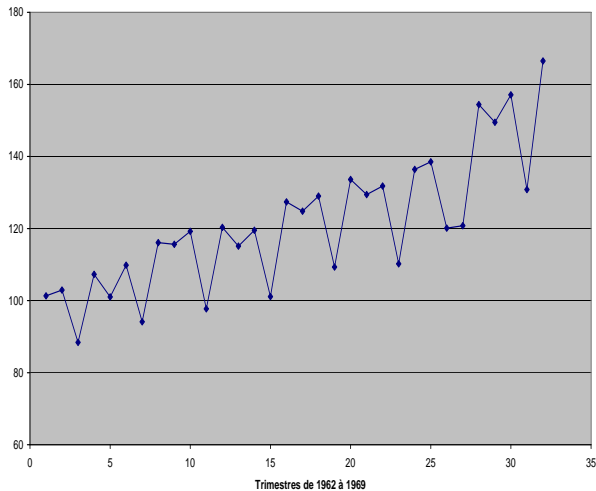
- $(e_i)_{i=1,\dots,n}$ désigne la composante résiduelle (ou bruit).

On étudie l'indice trimestriel de la production industrielle, entre 1962 et 1969, base 100 en 1962 (source : INSEE).

Année	Trimestre 1	Trimestre 2	Trimestre 3	Trimestre 4
1962	101.3	102.9	88.4	107.3
1963	101	109.8	94.1	116.1
1964	115.6	119.2	97.7	120.3
1965	115.1	119.5	101.1	127.4
1966	124.8	129	109.3	133.6
1967	129.4	131.8	110.2	136.4
1968	138.5	120.1	120.8	154.4
1969	149.5	157.1	130.8	166.5

Représentation graphique

Indice trimestriel de la production industrielle base 100 en 1962



Les différentes étapes de la désaisonnalisation

- 1 Estimation de la tendance par moyenne mobile.
- 2 Constitution de la série des différences.
- 3 Calcul des coefficients saisonniers non centrés.
- 4 Centrage des coefficients saisonniers.
- 5 Calcul de la série corrigée des variations saisonnières.

Etape 1 : Estimation de la tendance par moyenne mobile

On effectue un **lissage par la méthode des moyennes mobiles** afin d'obtenir une **première évaluation de la tendance de la série**.

On note $(\hat{f}_i)_{i=1,\dots,n}$ la série obtenue.

Soit p la période de la composante saisonnière :

- Si p est impaire alors $\hat{f}_i = MM(p)_i$.
- Si p est paire alors $\hat{f}_i = MMC(p)_i$.

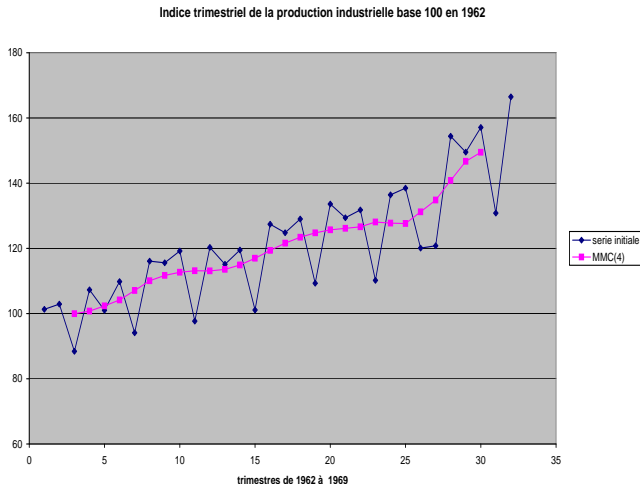
Remarque :

- Dans de nombreuses applications, la période p est paire et vaut $p = 2m$ (séries : mensuelle, trimestrielle, bi-mensuelle, semestrielle, etc.).
- Via le lissage, on perd m observations au début de la série et m autres à la fin, soit au total $p = 2m$ observations.

Etape 1 : Estimation de la tendance par moyenne mobile

date	y_i	MMC(4)
1	101,3	
2	102,9	
3	88,4	99,9375
4	107,3	100,7625
5	101	102,3375
6	109,8	104,15
7	94,1	107,075
8	116,1	110,075
9	115,6	111,7
10	119,2	112,675
11	97,7	113,1375
12	120,3	113,1125
13	115,1	113,575
14	119,5	114,8875
15	101,1	116,9875
16	127,4	119,3875
17	124,8	121,6
18	129	123,4
19	109,3	124,75
20	133,6	125,675
21	129,4	126,1375
22	131,8	126,6
23	110,2	128,0875
24	136,4	127,7625
25	138,5	127,625
26	120,1	131,2
27	120,8	134,825
28	154,4	140,825
29	149,5	146,7
30	157,1	149,4625
31	130,8	
32	166,5	

Etape 1 : Estimation de la tendance par moyenne mobile



Etape 2 : Constitution de la série des différences

On calcule la **série des différences** :

$$D_i = y_i - \hat{f}_i$$

pour $i = m + 1, \dots, n - m$.

Comme les termes \hat{f}_i correspondent à une estimation de la tendance, les différences D_i sont une **approximation des composantes saisonnière + résiduelle**.

Etape 2 : Constitution de la série des différences

date	y_i	MMC(4)	D_i
1	101,3		
2	102,9		
3	88,4	99,9375	-11,5375
4	107,3	100,7625	6,5375
5	101	102,3375	-1,3375
6	109,8	104,15	5,65
7	94,1	107,075	-12,975
8	116,1	110,075	6,025
9	115,6	111,7	3,9
10	119,2	112,675	6,525
11	97,7	113,1375	-15,4375
12	120,3	113,1125	7,1875
13	115,1	113,575	1,525
14	119,5	114,8875	4,6125
15	101,1	116,9875	-15,8875
16	127,4	119,3875	8,0125
17	124,8	121,6	3,2
18	129	123,4	5,6
19	109,3	124,75	-15,45
20	133,6	125,675	7,925
21	129,4	126,1375	3,2625
22	131,8	126,6	5,2
23	110,2	128,0875	-17,8875
24	136,4	127,7625	8,6375
25	138,5	127,625	10,875
26	120,1	131,2	-11,1
27	120,8	134,825	-14,025
28	154,4	140,825	13,575
29	149,5	146,7	2,8
30	157,1	149,4625	7,6375
31	130,8		
32	166,5		

Etape 3 : Calcul des coefficients saisonniers non centrés

On estime les p coefficients saisonniers non centrés \tilde{s}_j en moyennant les valeurs de la série $(D_i)_{i=m+1,\dots,n-m}$ sur les sous-séries qui correspondent aux différentes périodes.

Pour simplifier l'écriture, supposons que n soit multiple de la période p de la saisonnalité de telle sorte que l'on dispose de K_0 périodes complètes d'observations, $n = K_0 \times p$.

Par exemple si la série est trimestrielle $p = 4$ et que la série est observée sur $n = 12$ trimestres, on a $K_0 = 3$ périodes complètes d'observations (3 années). Après lissage on dispose de la série D_i pour $i = 3, \dots, 10$:

- $\tilde{s}_1 = \frac{1}{2}(D_5 + D_9),$
- $\tilde{s}_2 = \frac{1}{2}(D_6 + D_{10}),$
- $\tilde{s}_3 = \frac{1}{2}(D_3 + D_7),$
- $\tilde{s}_4 = \frac{1}{2}(D_4 + D_8).$

Etape 3 : Calcul des coefficients non centrés

Pour $1 \leq j < m$:

$$\tilde{s}_j = \frac{1}{K_0 - 1} \sum_{k=1}^{K_0-1} D_{j+pk}$$

Pour $m \leq j \leq p$:

$$\tilde{s}_j = \frac{1}{K_0 - 1} \sum_{k=1}^{K_0-1} D_{j+p(k-1)}$$

Etape 3 : Calcul des coefficients non centrés

date	y_i	MMC(4)	D_i n°trimestre	
1	101,3		1	
2	102,9		2	
3	88,4	99,9375	-11,5375	3 Stilde_3 -14,7428571
4	107,3	100,7625	6,5375	4 Stilde_4 8,27142857
5	101	102,3375	-1,3375	1 Stilde_1 3,46071429
6	109,8	104,15	5,65	2 Stilde_2 3,44642857
7	94,1	107,075	-12,975	3
8	116,1	110,075	6,025	4
9	115,6	111,7	3,9	1
10	119,2	112,675	6,525	2
11	97,7	113,1375	-15,4375	3
12	120,3	113,1125	7,1875	4
13	115,1	113,575	1,525	1
14	119,5	114,8875	4,6125	2
15	101,1	116,9875	-15,8875	3
16	127,4	119,3875	8,0125	4
17	124,8	121,6	3,2	1
18	129	123,4	5,6	2
19	109,3	124,75	-15,45	3
20	133,6	125,675	7,925	4
21	129,4	126,1375	3,2625	1
22	131,8	126,6	5,2	2
23	110,2	128,0875	-17,8875	3
24	136,4	127,7625	8,6375	4
25	138,5	127,625	10,875	1
26	120,1	131,2	-11,1	2
27	120,8	134,825	-14,025	3
28	154,4	140,825	13,575	4
29	149,5	146,7	2,8	1
30	157,1	149,4625	7,6375	2
31	130,8			3
32	166,5			4

Etape 4 : Centrage des coefficients saisonniers

Rappelons que l'influence des variations saisonnières sur une année est neutre. La composante saisonnière est censée vérifier :

$$\sum_{j=1}^p s_j = 0$$

On va donc corriger les coefficients \tilde{s}_j en les centrant.

Etape 4 : Centrage des coefficients saisonniers

On calcule la moyenne des p coefficients saisonniers obtenus à l'étape 3 :

$$\bar{s} = \frac{1}{p} \sum_{j=1}^p \tilde{s}_j.$$

On centre ensuite les coefficients saisonniers $(\tilde{s}_j)_{j=1,\dots,p}$ en leur retranchant la moyenne \bar{s} . On obtient :

$$\hat{s}_j = \tilde{s}_j - \bar{s}, \quad j = 1, \dots, p.$$

On a bien :

$$\sum_{j=1}^p \hat{s}_j = \sum_{j=1}^p \tilde{s}_j - \sum_{j=1}^p \bar{s} = p \times \bar{s} - p \times \bar{s} = 0.$$

Etape 4 : Centrage des coefficients saisonniers

y_i	MMC(4)	D_i n°trimestre		
101,3			1	
102,9			2	
88,4	99,9375	-11,5375	3	Stilde_3 -14,7428571
107,3	100,7625	6,5375	4	Stilde_4 8,27142857
101	102,3375	-1,3375	1	Stilde_1 3,46071429
109,8	104,15	5,65	2	Stilde_2 3,44642857
94,1	107,075	-12,975	3	
116,1	110,075	6,025	4	
115,6	111,7	3,9	1	moyenne sbz 0,10892857
119,2	112,675	6,525	2	
97,7	113,1375	-15,4375	3	^S3 -14,8517857
120,3	113,1125	7,1875	4	^S4 8,1625
115,1	113,575	1,525	1	^S1 3,35178571
119,5	114,8875	4,6125	2	^S2 3,3375
101,1	116,9875	-15,8875	3	
127,4	119,3875	8,0125	4	
124,8	121,6	3,2	1	
129	123,4	5,6	2	
109,3	124,75	-15,45	3	
133,6	125,675	7,925	4	
129,4	126,1375	3,2625	1	
131,8	126,6	5,2	2	
110,2	128,0875	-17,8875	3	
136,4	127,7625	8,6375	4	
138,5	127,625	10,875	1	
120,1	131,2	-11,1	2	
120,8	134,825	-14,025	3	
154,4	140,825	13,575	4	
149,5	146,7	2,8	1	
157,1	149,4625	7,6375	2	
130,8			3	
166,5			4	

Etape 5 : Calcul de la série corrigée des variations saisonnières

La série corrigée des variations saisonnières, ou désaisonnalisée, notée $(CVS_i)_{i=1,\dots,n}$, s'obtient en retranchant à la série initiale la suite des coefficients saisonniers centrés obtenus à l'étape 4 :

$$CVS_i = y_i - \hat{s}_i, \quad i = 1, \dots, n.$$

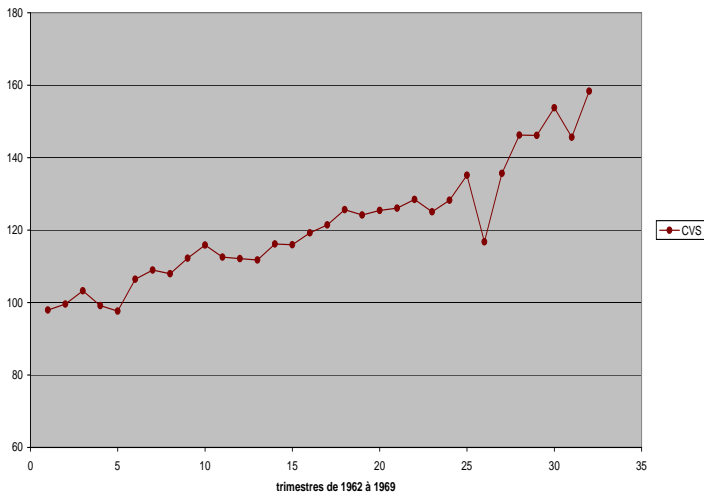
Les données CVS sont directement comparables : on a retiré l'effet saisonnier.

Etape 5 : Calcul de la série corrigée des variations saisonnières

date	y_i	n°trimestre		^S	CVSt
1	101,3	1		3,35178571	97,9482143
2	102,9	2		3,3375	99,5625
3	88,4	3	Stilde_3	-14,7428571	103,251786
4	107,3	4	Stilde_4	8,27142857	99,1375
5	101	1	Stilde_1	3,46071429	97,6482143
6	109,8	2	Stilde_2	3,44642857	106,4625
7	94,1	3		-14,8517857	108,951786
8	116,1	4		8,1625	107,9375
9	115,6	1	moyenne sb:	0,10892857	112,248214
10	119,2	2		3,3375	115,8625
11	97,7	3	^S3	-14,8517857	112,551786
12	120,3	4	^S4	8,1625	112,1375
13	115,1	1	^S1	3,35178571	111,748214
14	119,5	2	^S2	3,3375	116,1625
15	101,1	3		-14,8517857	115,951786
16	127,4	4		8,1625	119,2375
17	124,8	1		3,35178571	121,448214
18	129	2		3,3375	125,6625
19	109,3	3		-14,8517857	124,151786
20	133,6	4		8,1625	125,4375
21	129,4	1		3,35178571	126,048214
22	131,8	2		3,3375	128,4625
23	110,2	3		-14,8517857	125,051786
24	136,4	4		8,1625	128,2375
25	138,5	1		3,35178571	135,148214
26	120,1	2		3,3375	116,7625
27	120,8	3		-14,8517857	135,651786
28	154,4	4		8,1625	146,2375
29	149,5	1		3,35178571	146,148214
30	157,1	2		3,3375	153,7625
31	130,8	3		-14,8517857	145,651786
32	166,5	4		8,1625	158,3375

Etape 5 : Calcul de la série corrigée des variations saisonnières

indice trimestriel de la production industrielle base 100 en 1962



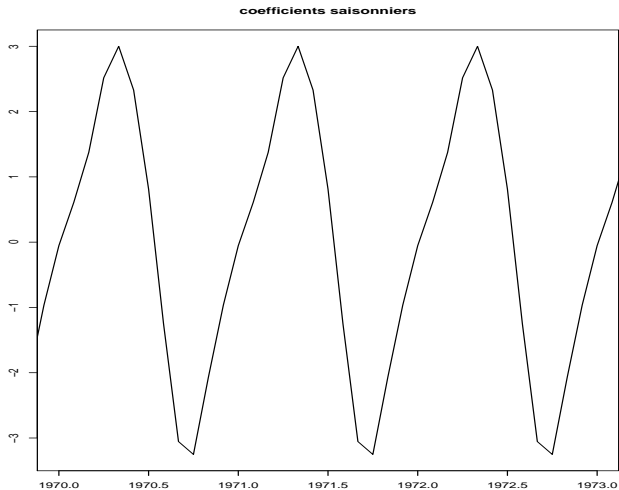
La **série ajustée** \hat{y}_i , ou série lissée des prédictions, est obtenue en additionnant la tendance \hat{f}_i et la composante saisonnière \hat{s}_i :

$$\hat{y}_i = \hat{f}_i + \hat{s}_i .$$

La **série ajustée** représente **l'évolution** qu'aurait subi la grandeur observée si les **variations saisonnières avaient été parfaitement périodiques** (i.e. s'étaient répétées à l'identique d'une période à une autre) et s'il n'y avait **pas eu de composante résiduelle**.

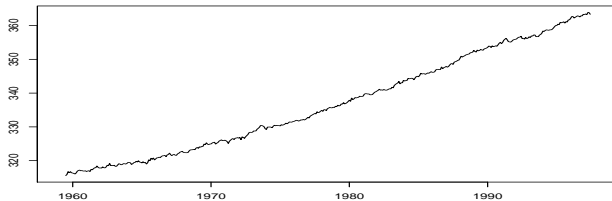
Remarque: Lorsque la tendance \hat{f}_i obtenue par moyenne mobile est encore trop irrégulière et représente mal la tendance, on peut en donner une seconde estimation (par exemple en approchant la série CVS_i par une courbe paramétrique à l'aide des moindres carrés).

Retour à l'exemple du taux de CO2

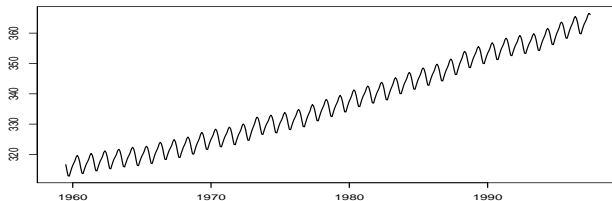


Retour à l'exemple du taux de CO₂

série corrigée des variations saisonnières



série lissée des prédictions



La différence entre la série initiale y_i et la série ajustée \hat{y}_i représente l'erreur d'ajustement :

$$\hat{e}_i = y_i - \hat{y}_i.$$

La moyenne des valeurs ajustées $\hat{y}_i = \hat{f}_i + \hat{s}_i$ est en général un peu différente de la moyenne des observations y_i : $\bar{\hat{y}} \neq \bar{y}$.

Dans le cas où $\hat{f}_i = MM(p)_i$ ou $MMC(p)_i$, la différence provient des m premières valeurs de la série (la différence est donc d'autant plus petite que n est grand).

Par conséquent $\bar{\hat{e}} \neq 0$, ce qui veut dire que les erreurs ne vérifient pas les contraintes de centrage.

Question: comment rectifier ce problème de façon simple?

Il suffit d'ajouter $\bar{\hat{e}}$ à la tendance: $\tilde{f}_i = \hat{f}_i + \bar{\hat{e}}$.

On obtient donc les nouvelles valeurs ajustées : $\tilde{y}_i = \tilde{f}_i + \hat{s}_i$, et les erreurs d'ajustement $\tilde{e}_i = y_i - \tilde{y}_i = \hat{e}_i - \bar{\hat{e}}$.

On vérifie bien que ces erreurs sont centrées $\bar{\tilde{e}} = 0$.

En pratique, les logiciels n'effectuent pas cette correction, qui est inutile lorsque n est assez grand: pour le taux de CO2, on a $\bar{\hat{e}} = 0.0017$.

On étudie les ventes trimestrielles d'un grand magasin parisien.

La série $(y_i)_{i=1,\dots,10}$ des ventes (en milliers d'euros) est donnée du premier trimestre 1995 au deuxième trimestre 1997 dans le tableau ci-dessous :

t_i	1	2	3	4	5	6	7	8	9	10
y_i	662	742	683	842	717	792	742	875	767	805

- 1 Calculer la série des moyennes mobiles centrées d'ordre 4.
- 2 On choisit un modèle additif. Calculer les coefficients saisonniers. Calculer leur moyenne. Calculer les coefficients saisonniers centrés.
- 3 Calculer la série corrigée des variations saisonnières.
- 4 Calculer la série des valeurs ajustées.
- 5 Calculer $\bar{\hat{e}}$.

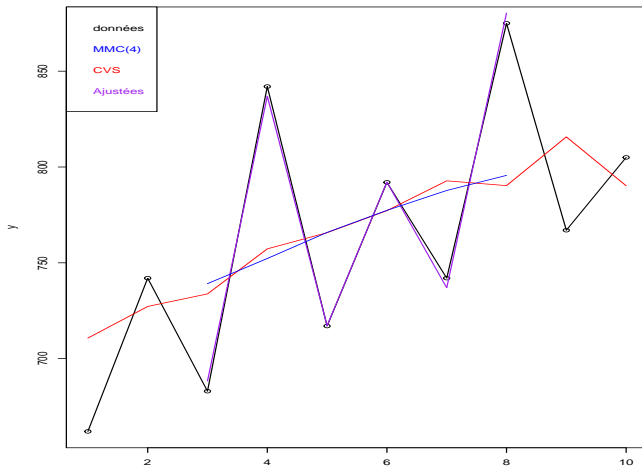
t_i	1	2	3	4	5	6	7	8	9	10
y_i	662	742	683	842	717	792	742	875	767	805

t_i	3	4	5	6	7	8
$MMC(4)_i$	739.125	752.250	765.875	777.375	787.750	795.625
D_i	-56.125	89.75	-48.875	14.625	-45.75	79.375

$$\tilde{s}_1 = -48.875 \quad \tilde{s}_2 = 14.625 \quad \tilde{s}_3 = -50.9375 \quad \tilde{s}_4 = 84.5625$$

$$\hat{s}_1 = -48.71875 \quad \hat{s}_2 = 14.78125 \quad \hat{s}_3 = -50.78125 \quad \hat{s}_4 = 84.71875$$

$$\bar{\hat{e}} = -0.156$$



Soit $(y_i)_{i=1,\dots,n}$ une série chronologique dont les termes sont positifs.

On considère le modèle :

$$y_i = F_i \times S_i \times E_i \quad \text{pour } i = 1, \dots, n$$

où :

- $(F_i)_{i=1,\dots,n}$ désigne la tendance.
- $(S_i)_{i=1,\dots,n}$ désigne la composante saisonnière de période p telle que :

$$\prod_{j=1}^p S_j = 1.$$

- $(E_i)_{i=1,\dots,n}$ désigne la composante résiduelle (ou bruit).

En prenant le logarithme des observations on obtient

$$z_i = \log(y_i) = \log(F_i) + \log(S_i) + \log(E_i),$$

c'est à dire un modèle additif avec $f_i = \log(F_i)$, $s_i = \log(S_i)$ et $e_i = \log(E_i)$.

On calcule ensuite les coefficients saisonniers de ce modèle additif en suivant les étapes décrites précédemment:

- On approche la tendance par $\hat{f}_i = MM(p)_i$ ou $\hat{f}_i = MMC(p)_i$.
- On forme la série des différences $D_i = z_i - \hat{f}_i$.
- On calcule les coefficients \tilde{s}_k en prenant les moyennes des sous-séries de D_i .
- On recentre les coefficients \tilde{s}_k pour obtenir les coefficients centrés \hat{s}_k .

Les coefficients saisonniers du modèle multiplicatif s'écrivent alors

$$\hat{S}_i = \exp(\hat{s}_i) .$$

Par construction,

$$\sum_{i=1}^p \hat{s}_i = 0 ,$$

et donc

$$\prod_{i=1}^p \hat{S}_i = \exp\left(\sum_{i=1}^p \hat{s}_i\right) = 1 .$$

La contrainte sur les coefficients saisonniers est donc bien satisfaite.

La série corrigée des variations saisonnières, ou désaisonnalisée, notée $(CVS_i)_{i=1,\dots,n}$, s'obtient en divisant la série initiale par la suite des coefficients saisonniers centrés :

$$CVS_i = \frac{y_i}{\hat{S}_i}, \quad i = 1, \dots, n.$$

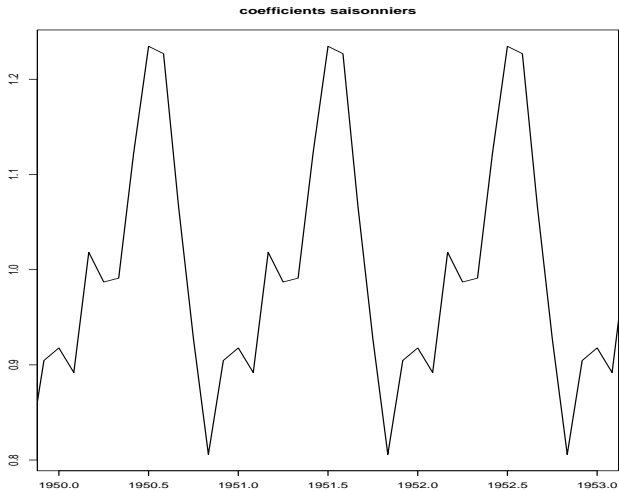
Les données CVS sont directement comparables : on a retiré l'effet saisonnier.

La **série ajustée** \hat{y}_i , ou série lissée des prédictions, est obtenue en multipliant la tendance $\hat{F}_i = \exp(\hat{f}_i)$ par la composante saisonnière \hat{S}_i :

$$\hat{y}_i = \hat{F}_i \times \hat{S}_i.$$

Remarque: Lorsque la tendance \hat{F}_i obtenue par moyenne mobile est encore trop irrégulière et représente mal la tendance, on peut en donner une seconde estimation (par exemple en approchant la série CVS_i par une courbe paramétrique à l'aide des moindres carrés).

Retour à l'exemple du trafic aérien



Retour à l'exemple du trafic aérien

