

# SINGULAR SPECTRUM ANALYSIS IN R

JAN DE LEEUW

ABSTRACT. Meet the abstract. This is the abstract.

## 1. INTRODUCTION

Singular Spectrum Analysis (SSA from now on) decomposes an observed time series into a sum of component series, in which the components hopefully capture and show the dynamics of the series more clearly. The SSA decomposition method can be thought of as one possible generalization of the singular value decomposition (SVD), or of principal component analysis (PCA), to a single time series. It can display and isolate trends and seasonal effects, as well as stationary residuals.

SSA originated and has mostly been studied and applied in geophysics and atmospheric science. A comprehensive early review paper is Ghil et al. [2002]. The technique is not very well known in statistics, despite the publication some time ago of the book by Golyandina et al. [2001]. Recent examples, perhaps showing an increasing interest, are the review paper by Hassani [2007], and the application published by Bilancia and Stea [2008].

## 2. PCA OF STATIONARY SERIES

Suppose  $\underline{x}$  is a  $T$ -dimensional vector random variable with  $E(\underline{x}\underline{x}') = \Omega$ . The cross product matrix  $\Omega$ , which we assume to be positive definite, has spectral decomposition  $\Omega = Q\Lambda^2Q'$ . Define  $\underline{z} = \Lambda^{-1}Q'\underline{x}$ . Then  $E(\underline{z}\underline{z}') =$

---

*Date:* October 21, 2008 — 19h 33min — Typeset in LUCIDA BRIGHT.

*2000 Mathematics Subject Classification.* 62M10.

*Key words and phrases.* Time Series, Singular Value Decomposition, Singular Spectrum Analysis.

$\Lambda^{-1}Q'\Omega Q\Lambda^{-1} = I$ . Also  $\underline{x} = Q\Lambda\underline{z}$ , the *Karhunen-Loève Decomposition* of  $\underline{x}$ .

In the usual time series analysis we only observe a single  $T$ -dimensional realization  $x$  of  $\underline{x}$ , and we have no idea what  $\Omega$  is. If we assume  $\underline{x}$  is stationary, however, we can estimate  $\omega_{st}$  with  $s \geq t$  by the average of the  $T - (s - t)$  products  $x_u x_v$  for which  $u - v = s - t$ . Thus

$$\hat{\omega}_{st} = \frac{1}{T - (s - t)} \sum_{v=1}^{T-(s-t)} x_v x_{v+(s-t)}.$$

Now compute  $\hat{\Omega} = \hat{Q}\hat{\Lambda}^2\hat{Q}'$  and  $\hat{z} = \hat{\Lambda}^{-1}\hat{Q}'x$ . Again  $x = \hat{Q}\hat{\Lambda}\hat{z}$ . The columns of  $\hat{Q}$  are *Empirical Orthogonal Functions* or EOF's.

More specific assumptions about the nature of the stationary process may lead to more precise estimates of  $\Omega$ , provided of course these assumptions are more or less true. For instance, we might assume that  $\underline{x}$  is **AR**(1), in which case  $\Omega$  only depends on the variance  $\sigma^2$  and the autocorrelation  $\rho$ .

```

1  dg<-function(s) {
2      n<-nrow(s); nn<-1:n; r<-rep(0,n)
3      for (k in 0:(n-1))
4          r[k+1]<-mean(s[which(outer(nn,k+nn,"=="))])
5      return(toeplitz(r))
6  }
```

### 3. THE FOUR-STEP PROGRAM

**3.1. Embedding.** Suppose  $\{x_1, x_2, \dots, x_T\}$  is our observed time series. For now, we assume there are no gaps (missing data). Choose a *window width*  $2 \leq L \leq \lfloor \frac{1}{2}T \rfloor$ . Define  $K = T - L + 1$ , so that  $K = \lceil \frac{1}{2}T \rceil + 1 \geq L + 1$ . Define the  $K \times L$  matrix  $Z$  with row  $i$  equal to  $\{x_i, x_{i+1}, \dots, x_{i+L-1}\}$ . Thus  $z_{ij} = x_{i+j-1}$ , which implies that  $Z$  is a *Hankel matrix*.  $Z$  is constant along its skew-diagonals, if  $i + j = k + \ell$  then  $z_{ij} = z_{k\ell}$ .

**3.2. Decomposition.** The singular value decomposition or SVD of  $Z$  is  $Z = U\Lambda V'$ . Here  $U$  is a  $K \times L$  orthonormal matrix,  $V$  is  $L \times L$  square orthonormal, and  $\Lambda$  is a diagonal matrix of order  $L$ . We suppose, for

identification purposes, that the diagonal elements of  $\Lambda$  are non-negative and are in non-increasing order along the diagonal. The SVD can also be written as

$$Z = \sum_{s=1}^L Z_s = \sum_{s=1}^L \lambda_s u_s v_s'.$$

Each of the  $Z_s$  is a rank-one matrix, and the  $Z_s$  are orthogonal in the sense that both  $Z_s Z_t' = 0$  and  $Z_s' Z_t = 0$  for  $s \neq t$ . Also, for the sum of squares, a.k.a. the squared *Frobenius norm*,

$$\text{SSQ}(Z) = \sum_{s=1}^L \text{SSQ}(Z_s) = \sum_{s=1}^L \lambda_s^2,$$

If we define

$$\pi_t \triangleq \frac{\lambda_t^2}{\sum_{s=1}^L \lambda_s^2}$$

then the  $\pi_t$  add up to one, and can be interpreted as the percentage of the sum of squares “explained” by  $Z_t$ , i.e. by the singular triple  $(\lambda_t, u_t, v_t)$ .

**3.3. Grouping.** Suppose  $\mathcal{I} = \{I_1, \dots, I_\Xi\}$  is a *partition* of  $\{1, 2, \dots, L\}$ , and define  $\bar{Z}_\xi = \sum_{s \in I_\xi} Z_s$ . Obviously we still have  $Z = \sum_{\xi=1}^\Xi \bar{Z}_\xi$ . Moreover  $Z_\xi Z_\mu' = 0$  and  $Z_\xi' Z_\mu = 0$  for  $\xi \neq \mu$ , and

$$\text{SSQ}(Z) = \sum_{\xi=1}^\Xi \text{SSQ}(\bar{Z}_\xi) = \sum_{\xi=1}^\Xi \left\{ \sum_{s \in I_\xi} \lambda_s^2 \right\}.$$

**3.4. Hankelization.** Suppose we have a partition  $\mathcal{I}$  into  $\Xi$  sets, and corresponding  $\bar{Z}_\xi$ . Find the Hankel matrices  $\tilde{Z}_\xi$  that minimize  $\text{SSQ}(Z - \bar{Z}_\xi)$ . They can be computed simply by *Hankelizing* or *diagonal averaging*, i.e. by replacing all elements for which  $i + j$  is constant by their average. At the same time this define a time series  $\tilde{x}_\xi$ , with element  $k$  equal to the elements of  $\tilde{Z}_\xi$  for which  $i + j = k - 1$ . Since Hankelizing is a linear operation, and since  $Z$  is Hankel already, we have

$$Z = \sum_{\xi=1}^\Xi \tilde{Z}_\xi,$$

as well as

$$x = \sum_{\xi=1}^\Xi \tilde{x}_\xi.$$

This is the SSA decomposition of  $x$ , or rather it is *a* SSA decomposition, because it depends on the choice of the window width and the grouping of the singular triples.

Note, by the way, that Hankelizing a matrix is an orthogonal projection, and thus

$$\mathbf{SSQ}(Z) \geq \sum_{\xi=1}^{\Xi} \mathbf{SSQ}(\tilde{Z}_{\xi}),$$

as well as

$$\mathbf{SSQ}(x) \geq \sum_{\xi=1}^{\Xi} \mathbf{SSQ}(\tilde{x}_{\xi}).$$

#### 4. CHOICES

**4.1. Window width.** There are many sophisticated methods to choose window width. Golyandina et al. [2001, p. 18] suggest determining the *fractal dimension* of the series, or to find an approximate order using *autoregression*. In our code we use  $L = \lfloor \frac{1}{2}T \rfloor$  as the default, and generally that seems to work rather well.

**4.2. Grouping.** We have chosen a simple grouping method based on *w-correlations* defined in Golyandina et al. [2001, p. 46–47]. These are just cosines between time series, using a weighted inner product that de-emphasizes the beginning and end of the series.

We first use the trivial grouping in which each singular triple defines a group. Compute the  $L \times L$  matrix  $R$  of *w-correlations*, and choose a cut-off quantity  $0 < \epsilon < 1$ . the adjacency matrix  $A$  defined by

$$a_{j\ell} = \begin{cases} 1 & \text{if } |r_{j\ell}| > \epsilon, \\ 0 & \text{otherwise.} \end{cases}$$

We then use *Warshall's Algorithm* [Warshall, 1962] to compute the *transitive closure* of the relation corresponding with  $A$ , and we use the *equivalence classes* of this relation to define the groups.

In principle any clustering method can be used here. In Bilancia and Stea [2008], for example, complete linkage hierarchical clustering is used. Given the almost infinite number of cluster methods that are available

there is a great deal of flexibility here [Gan et al., 2007] . Even if we limit ourselves to cluster methods available in R, we have many different choices. Also note that there are many clustering methods that do not require pairwise similarity measures first, and work directly on the data matrix.

## 5. EXAMPLES

**5.1. Nile.** Data are from the package `datasets` in base R. The series, of length 100, is the annual flow of the river Nile at Ashwan 1871–1970. A timeplot is given in Figure 1(a).

*Insert Figure 1(a) about here*

For the SSA we use window width 50 and cut-off 0.25. The first singular triple explains 97.4% and the next two 0.32% and 0.22%. The grouping gives the six groups  $\{\{1\}, \{2\}, \{\dots \text{rest} \dots\}, \{32, 33\}, \{36, 37\}, \{48, 49\}\}$ .

*Insert Figure 1(b) about here*

**5.2. Accidental Deaths.** Data are from package `MASS` and give monthly totals of accidental deaths in the USA from 1973–1978. The series has length 72, and is plotted in Figure 2(a).

*Insert Figure 2(a) about here*

For the SSA we use window width 36 and cut-off 0.25. The first singular triple explains 99.01% and the next two 0.68% and 0.12%. The grouping gives the ten groups

$$\begin{aligned} &\{\{1\}, \{2, 3\}, \{4, 5\}, \{6\}, \{7, 8\}, \{9, 10\}, \\ &\quad \{11, \dots, 18\}, \{\dots \text{rest} \dots\}, \{31\}, \{35, 36\}\}. \end{aligned}$$

*Insert Figure 2(b) about here*

**5.3. Milk.** The data from package `TSA` give the average monthly milk production per cow in the US from 1994 to 2005. The series has length 144, and is plotted in Figure 3(a).

*Insert Figure 3(a) about here*

SSA with the default settings gives 10 components. The first component explains 99.85% of the total sum of squares. The groups are

$$\{\{1\}, \{2, 3\}, \{4, 5\}, \{6\}, \{7, 8\}, \{\dots \text{rest} \dots\}, \\ \{12, 13\}, \{34, 35\}, \{44, \dots, 70\}, \{71, 72\}\}.$$

The 10 components are plotted in Figure 3(b).

*Insert Figure 3(b) about here*

**5.4. Beer Sales.** Data from package `TSA`. Monthly beer sales in millions of barrels from 1975 to 1990. The series has length 192, and is plotted in Figure 4(a).

*Insert Figure 4(a) about here*

*Insert Figure 4(b) about here*

**5.5. CO<sub>2</sub>.** Data from package `datasets`. Mauna Loa Atmospheric CO<sub>2</sub> Concentration, monthly 1959–1997. The series has length 468, and is plotted in Figure 5.

*Insert Figure 5 about here*

Window width is 234, and the algorithm with cut-off 0.25 gives 12 groups. The first eigenvalue explains 99.9955% of the total sum of squares.

## 6. VARIATIONS

### 6.1. To Center or Not to Center.

**6.2. From Hankel to Toeplitz.** Expanding the time series to a Hankel matrix, and decomposing this matrix into a sum of Hankel matrices, takes up a large amount of space. We can avoid all this, by working with the matrix  $C = Z'Z$  and its eigen-decomposition  $C = V\Lambda^2V'$ . Now

$Z_s = Z v_s v_s'$ , which becomes in elementwise notation

$$(Z_s)_{ik} = \sum_{j=1}^L z_{ij} v_{js} v_{ks} = \sum_{j=1}^L x_{i+j-1} v_{js} v_{ks}$$

Hankelizing means setting

$$(x_s)_v = \frac{1}{n_v} \sum \{(Z_s)_{ik} \mid i + k = v + 1\}$$

Suppose  $(i, k)$  are the  $n_v$  index pairs for which  $1 \leq i \leq K$  and  $1 \leq k \leq L$  and  $i + k = v + 1$ . This means that  $\max(1, (v + 1) - L) \leq i \leq \min(v, K)$  and  $k = (v + 1) - i$ . Thus

$$(x_s)_v = \frac{1}{n_v} \sum_{i=\max(1, (v+1)-L)}^{\min(v, K)} \sum_{j=1}^L x_{i+j-1} v_{js} v_{(v+1)-i, s}$$

The expression only involves the eigenvectors of  $C$  and the values of the original series. No *Embedding* and no *Hankelizing* is required.

### 6.3. Effect of Window Width.

## 7. EXTENSIONS

### 7.1. Gaps.

### 7.2. Multivariate Series.

## REFERENCES

- M. Bilancia and G. Stea. Timescale Effect Estimation in Time-series Studies of Air Pollution and Health: A Singular Spectrum Analysis Approach. *Electronic Journal of Statistics*, 2:432–453, 2008. URL <http://www.i-journals.org/ejs/include/getdoc.php?id=967&article=123&mode=pdf>.
- G. Gan, C. Ma, and J. Wu. *Data Clustering. Theory, Algorithms, and Applications*. ASA-SIAM Series on Statistics and Applied Probability. SIAM, Philadelphia, PA, 2007.
- M. Ghil, M.R. Allen, M.D. Dettinger, K. Ide, D. Kpndrashov, M.E. Mann, A.W. Robertson, A. Saunders, Y. Tian, F. Varadi, and P. Yiou. Advanced Spectral Methods for Climatic Time Series. *Review of Geophysics*, 40: 1–41, 2002.

- N. Golyandina, V. Nekrutkin, and A. Zhigljavsky. *Analysis of Time Series Structure: SSA and Related Techniques*. Number 90 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, 2001.
- H. Hassani. Singular Spectrum Analysis: Methodology and Comparison. *Journal of Data Science*, 5:239-257, 2007. URL <http://www.sinica.edu.tw/~jds/JDS-396.pdf>.
- S. Warshall. A Theorem on Boolean Matrices. *Journal of the Association of Computer Machinery*, 9:11-12, 1962.



## APPENDIX A. CODE

## A.1. R Code.

```

1  #
2  #  ssa package
3  #  Copyright (C) 2008 Jan de Leeuw <deleeuw@stat.ucla.edu>
4  #  UCLA Department of Statistics, Box 951554,
5  #  Los Angeles, CA 90095-1554
6  #
7  #  This program is free software; you can redistribute it
8  #  and/or modify it under the terms of the GNU General Public
9  #  License as published by the Free Software Foundation;
10 #  either version 2 of the License, or (at your option)
11 #  any later version.
12 #
13 #  This program is distributed in the hope that it will be
14 #  useful, but WITHOUT ANY WARRANTY; without even the implied
15 #  warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR
16 #  PURPOSE. See the GNU General Public License for more details.
17 #
18 #  You should have received a copy of the GNU General Public
19 #  License along with this program; if not, write to the
20 #  Free Software Foundation, Inc., 675 Mass Ave, Cambridge,
21 #  MA 02139, USA.
22 #
23 #####
24 #
25 # version 0.0.1, 2008-10-10, initial
26 # version 0.1.0, 2008-10-16, similarity and cluster parameters
27 # version 0.2.0, 2008-10-17, got rid of hankel
28 # version 0.3.0, 2008-10-17, added Toeplitz SSA
29 # version 0.3.1, 2008-10-19, C code for crossproduct
30 # version 0.3.2, 2008-10-19, preprocessing option
31 # version 0.3.3, 2008-10-19, C code for fromCross
32 # version 0.4.0, 2008-10-19, added Fourier SSA
33 #
34 # to do:
35 #
36 # -- add clustering methods
37 # -- add similarity measures
38 # -- add preprocessors
39 # -- ndim for eigen
40 # -- better algorithm for Toeplitz eigenvectors
41 #
42
43 dyn.load("ssa.so")
44
45 ssa<-function(data,L=floor(length(data)/2),b=as.list(1:L),preproc=center,
46             ndim=width,method="hankel",ssa_similarity=ssa_w_cor,ssa_cluster=ssa_trans,
47             par=.25) {
48   if (!identical(class(data),"ts")) stop("Data should be of class ts")
49   data<-preproc(data); T<-length(data); K<-T-L+1

```

```

49   if (identical(method,"toeplitz")) {
50     cross<-toeplitz(rconv(data,1:L))
51     s<-eigen(cross)$vectors
52   }
53   if (identical(method,"hanke1")) {
54     cross<-hanke1(data,L)
55     s<-eigen(cross)$vectors
56   }
57   if (identical(method,"fourier")) {
58     s<-sincos(L)
59   }
60   a<-fromCross(data,s)
61   r<-ssa_similarity(a,L,K)
62   b<-ssa_cluster(r,par)
63   a<-apply(b,function(kk) rowSums(as.matrix(a[,kk])))
64   r<-ssa_similarity(a,L,K)
65   y<-ts(a,start=start(data),end=end(data),frequency=frequency(data))
66   return(list(y=y,b=b,s=colSums(a^2),r=r))
67 }
68
69 ssa_w_cor<-function(z,l,k) {
70   ls<-min(l,k); ks<-max(l,k)
71   n<-nrow(z); w<-rep(ls,n)
72   w[1:ls]<-1:ls; w[(ks+1):n]<-n-(ks:(n-1))
73   c<-crossprod(z,w*z); d<-diag(c)
74   return(c/sqrt(outer(d,d)))
75 }
76
77 ssa_trans<-function(r,cut) {
78   s<-ifelse(abs(r)>cut,1,0)
79   v<-warshall(s)
80   h<-unique(v)
81   return(apply(h,1,function(v) which(v==1)))
82 }
83
84 warshall<-function(a) {
85   n<-nrow(a)
86   for (j in 1:n) {
87     for(i in 1:n) {
88       if (a[i,j]==1) a[i,]<-pmax(a[i,],a[j,])
89     }
90   }
91   return(a)
92 }
93
94 sincos<-function(L) {
95   x<-2*pi*(1:L)/L
96   M<-floor(L/2)
97   if (is.even(L)) N<-M-1 else N<-M
98   s<-matrix(0,L,N+M+1)
99   for (k in 1:N)
100     s[,k]<-sin(x*k)
101   for (k in 0:M)

```

```

102     s[,N+k+1]<-cos(x*k)
103     return(apply(s,2,function(z) z/sqrt(sum(z^2))))
104 }
105
106 center<-function(x) x-mean(x)
107
108 ident<-function(x) x
109
110 rconv<-function(x,lag) {
111   .C("cconv",
112     as.integer(lag),
113     as.integer(length(lag)),
114     as.double(x),
115     as.integer(length(x)),
116     as.double(vector("double",length(lag))))[[5]]
117 }
118
119 hankel<-function(x,L){
120   cc<-rep(0,L*L); T<-length(x)
121   cl<-.C("hankelC",
122     as.double(x),
123     as.double(cc),
124     as.integer(L),
125     as.integer(T))
126   return(matrix(cl[[2]],L,L))
127 }
128
129 fromCross<-function(x,s) {
130   L<-nrow(s); T<-length(x)
131   aa<-rep(0,T*L); s<-as.vector(s)
132   al<-.C("fromCrossC",
133     as.double(x),
134     as.double(s),
135     as.double(aa),
136     as.integer(L),
137     as.integer(T))
138   return(matrix(al[[3]],T,L))
139 }
140
141 is.even<-function(x) return((as.integer(x) %% 2) == 0)

```

## A.2. C Code.

```

1
2 void
3 cconv (int *l, int *m, double *x, int *n, double *s)
4 {
5   double *y, *z, *u;
6   int i;
7   for (i = 0; i < *m; i++)
8   {
9     y = x + (*n - 1[i]);

```

```

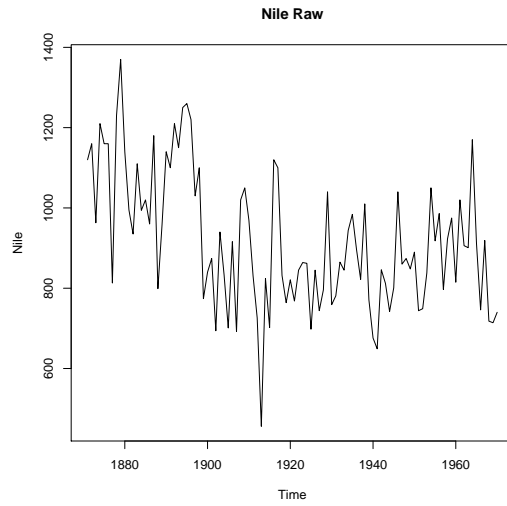
10     z = x + l[i];
11     u = x;
12     s[i] = 0.0;
13     while (u < y)
14         s[i] += *u++ * *z++;
15     }
16 }
17
18 void
19 hankelC(double *x, double *c, int *l, int *t)
20 {
21     int i, j, h, K, L, T; double s;
22     L = *l; T=*t; K = T - L + 1;
23     for (j = 1; j <= L; j++) for (h = j; h <= L; h++) {
24         s = 0.0;
25         for (i = 1; i <= K; i++) s += x[i+j-2] * x[i+h-2];
26         c[j + (h - 1) * L - 1] = c[h + (j - 1) * L - 1] = s;
27     }
28 }
29
30 void
31 fromCrossC(double *x, double *s, double *a, int *l, int *t)
32 {
33     int i, j, k, h, ilw, iup, ni, K, L, T;
34     double sv, sw;
35     L = *l; T=*t; K = T - L + 1;
36     for (i = 1; i <= T; i++) {
37         ilw = (i + 1) - L; if (ilw < 1) ilw = 1;
38         iup = i; if (iup > K) iup = K;
39         ni = iup - ilw + 1;
40         for (j = 1; j <= L; j++) {
41             sv = 0.0;
42             for (k = ilw; k <= iup; k++) {
43                 sw = 0.0;
44                 for (h = 1; h <= L; h++)
45                     sw += x[k + h - 2] * s[h + (j - 1) * L - 1];
46                 sv += sw * s[(i + 1) - k + (j - 1) * L - 1];
47             }
48             a[i + (j - 1) * T - 1] = sv / ((double) ni);
49         }
50     }
51 }

```

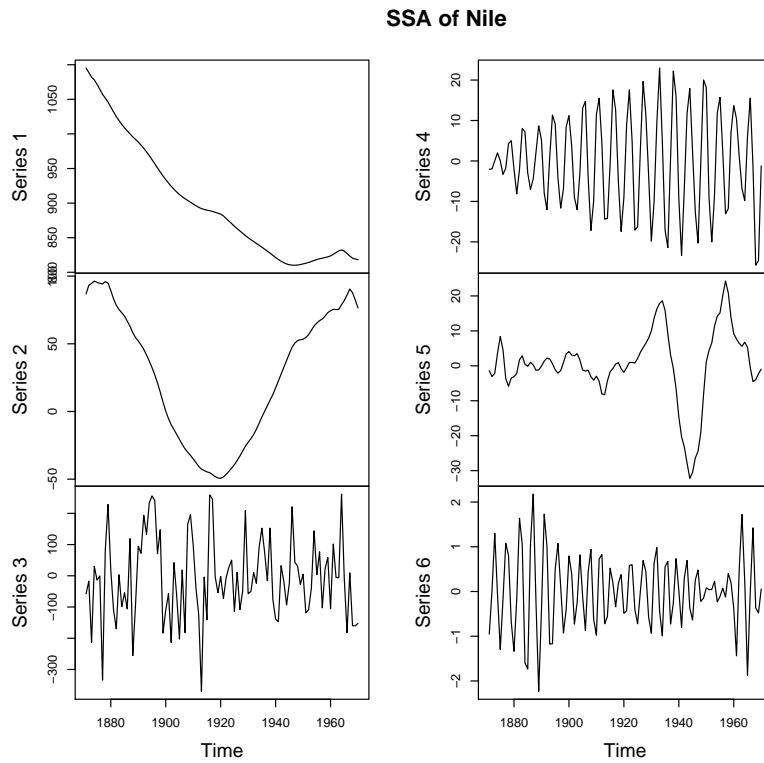
DEPARTMENT OF STATISTICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, CA 90095-1554

*E-mail address*, Jan de Leeuw: `deleeuw@stat.ucla.edu`

*URL*, Jan de Leeuw: `http://gifi.stat.ucla.edu`

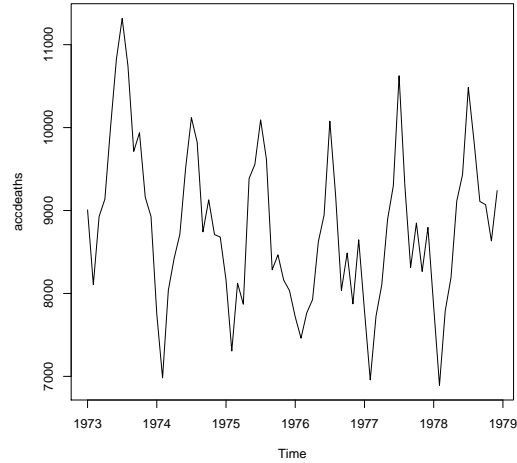


(a) Nile Data

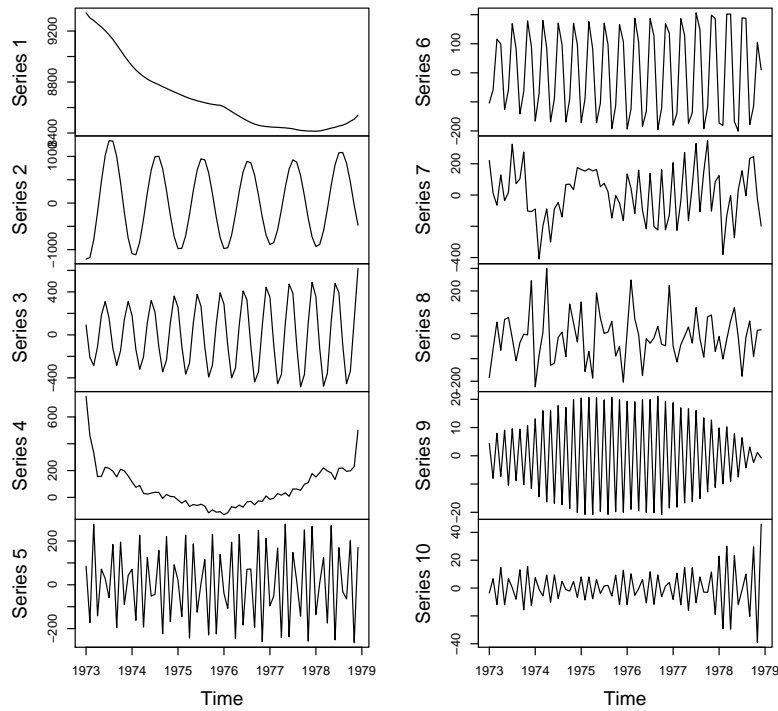


(b) Nile Data, window width 50, Cut-off 0.25

FIGURE 1. Nile Data

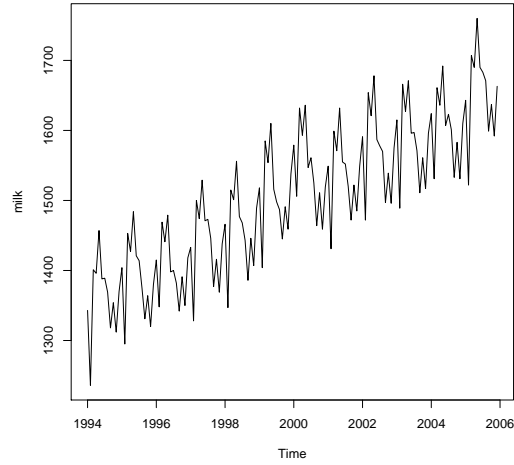


(a) Accidental Deaths Data

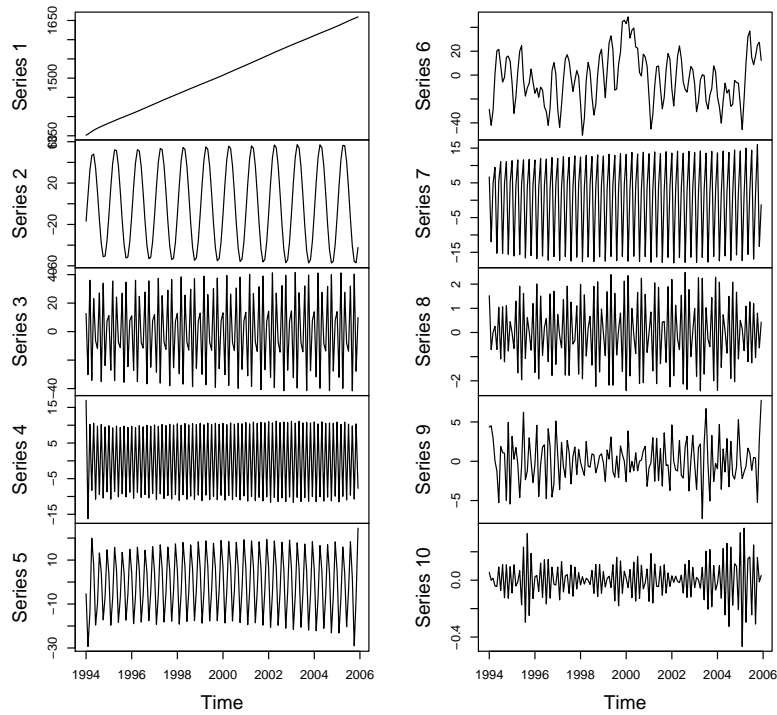
**SSA of accdeaths**

(b) Accidental Deaths Data, window width 36, Cut-off 0.25

FIGURE 2. Accidental Deaths



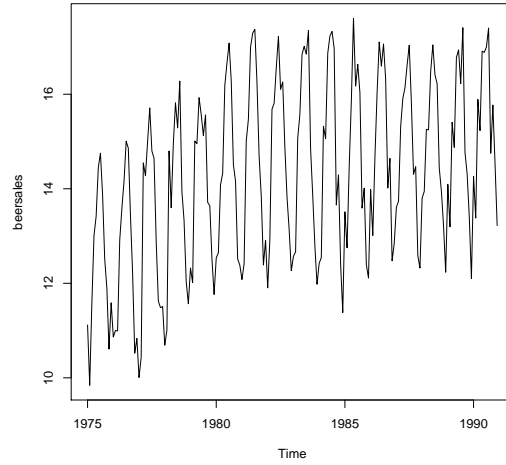
(a) Milk Yield Data

**SSA Milk Production**

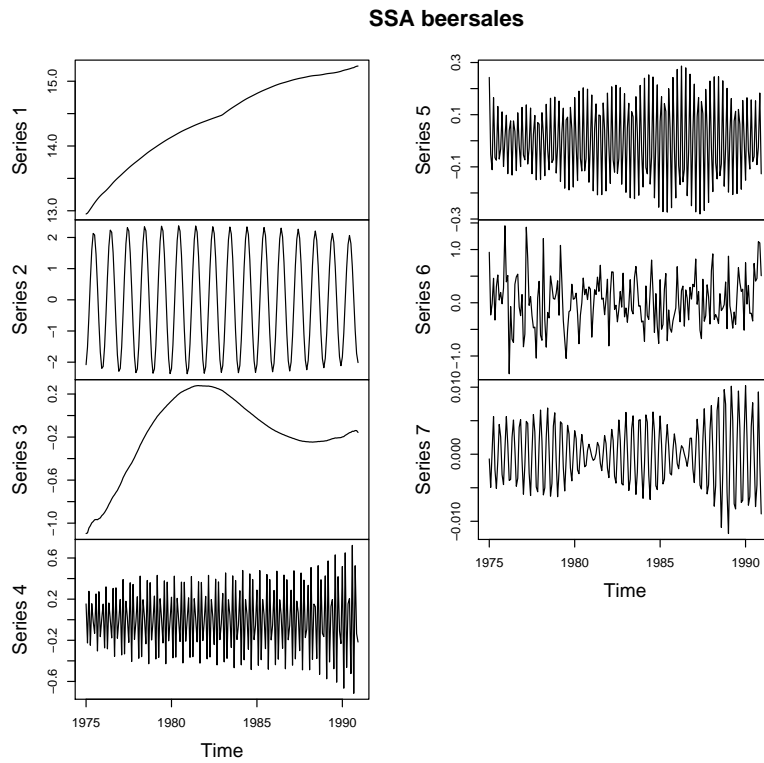
(b) Milk Yield Data, window width 72, Cut-off 0.25

**FIGURE 3. Milk Yields**





(a) Beer Sales Data



(b) Beer Sales SSA, window width 96, Cut-off 0.25

FIGURE 4. Beer Sales

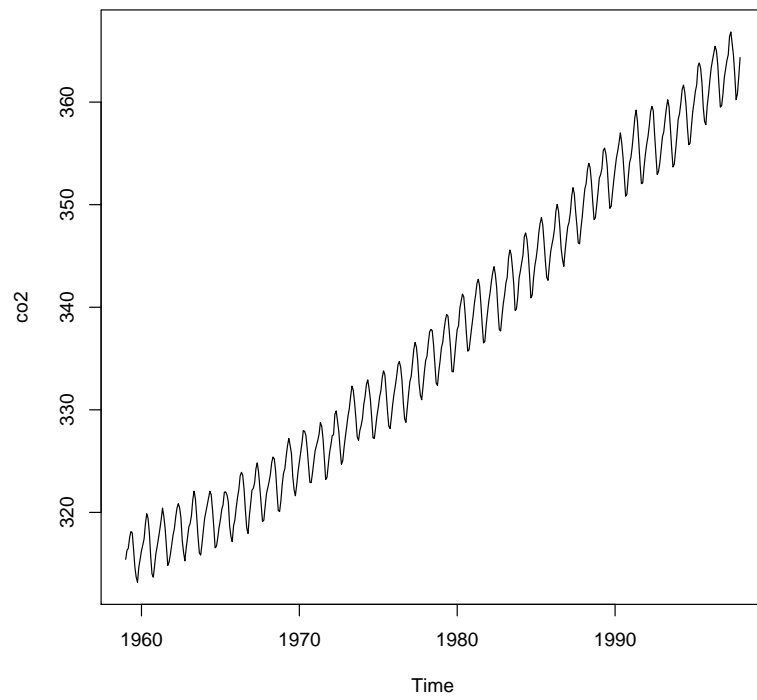
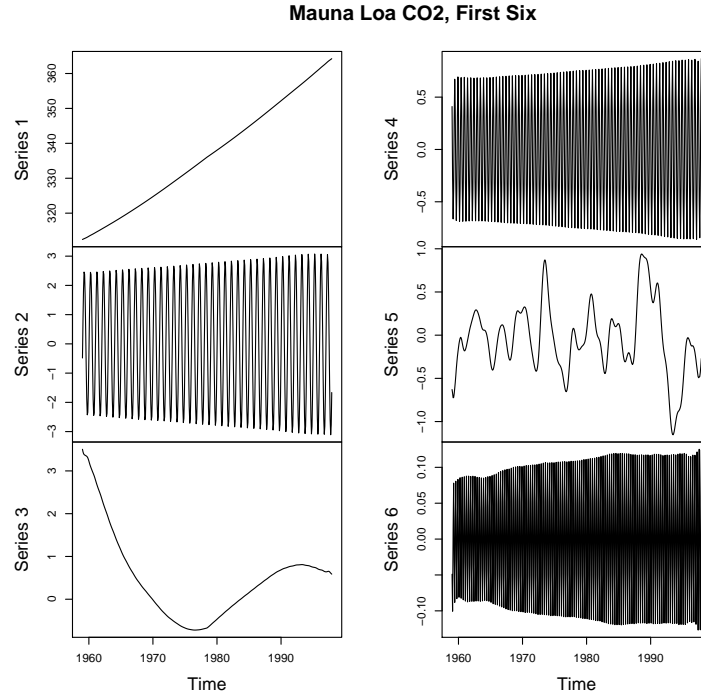
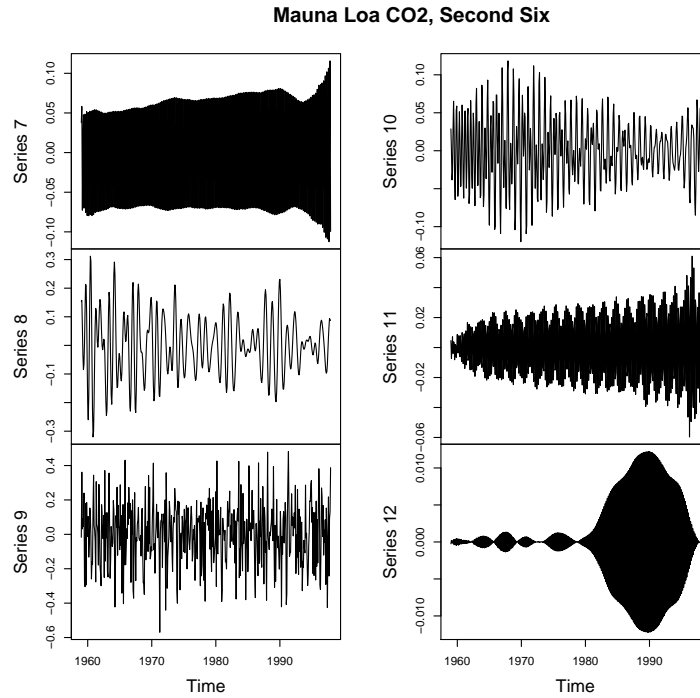
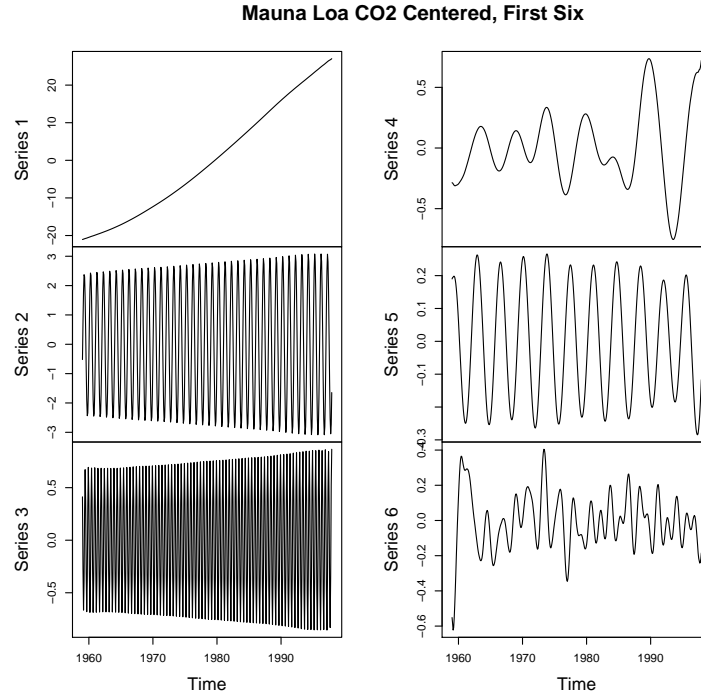
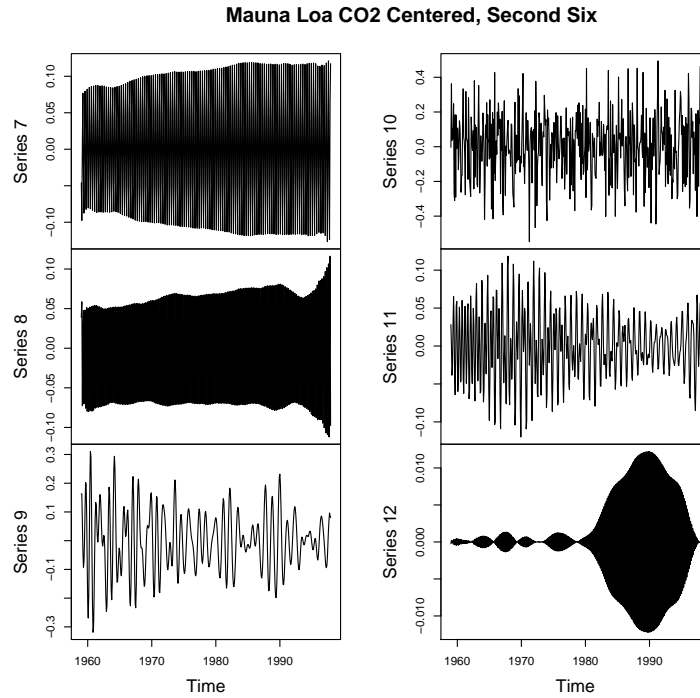
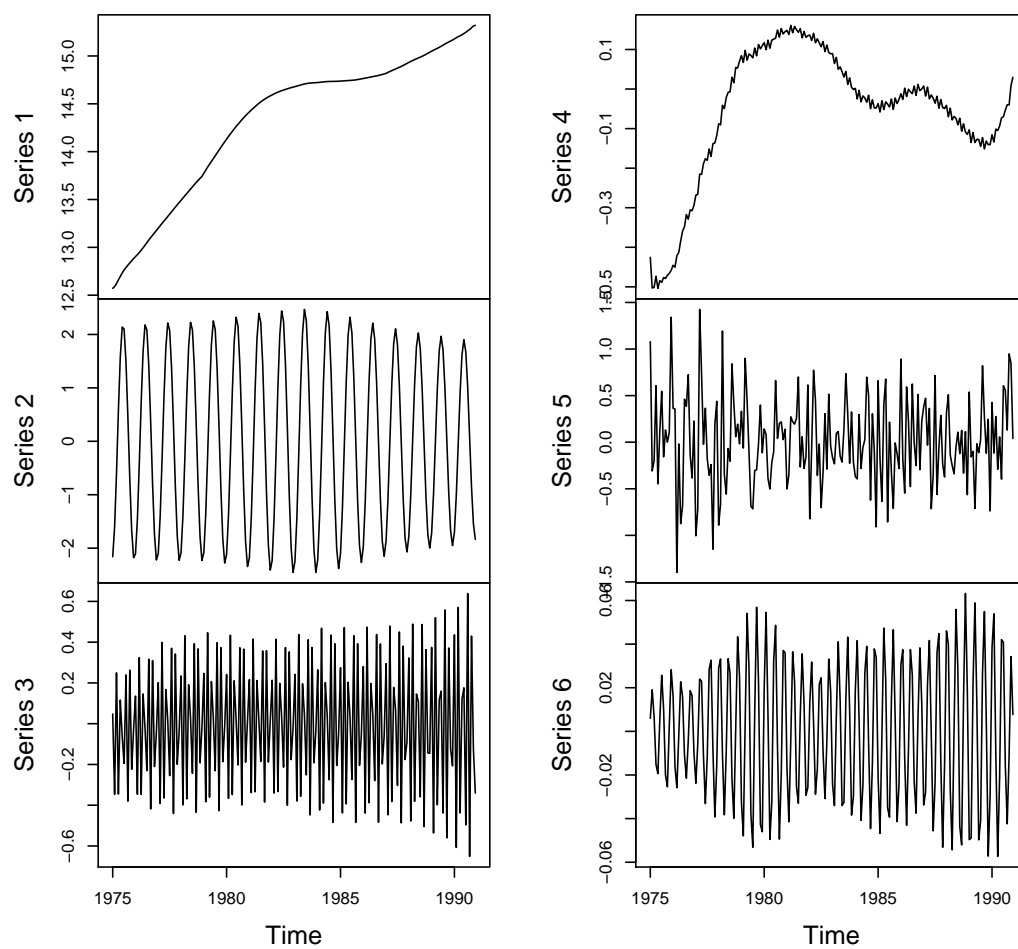
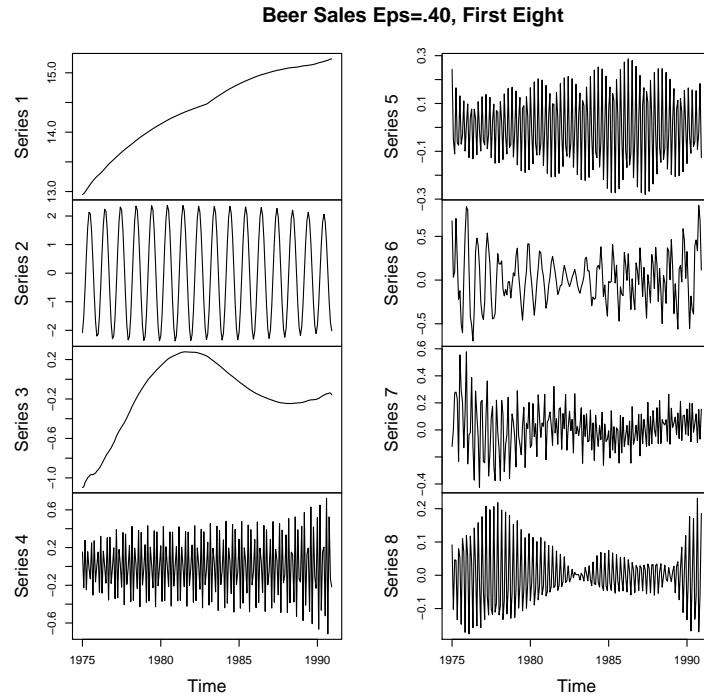
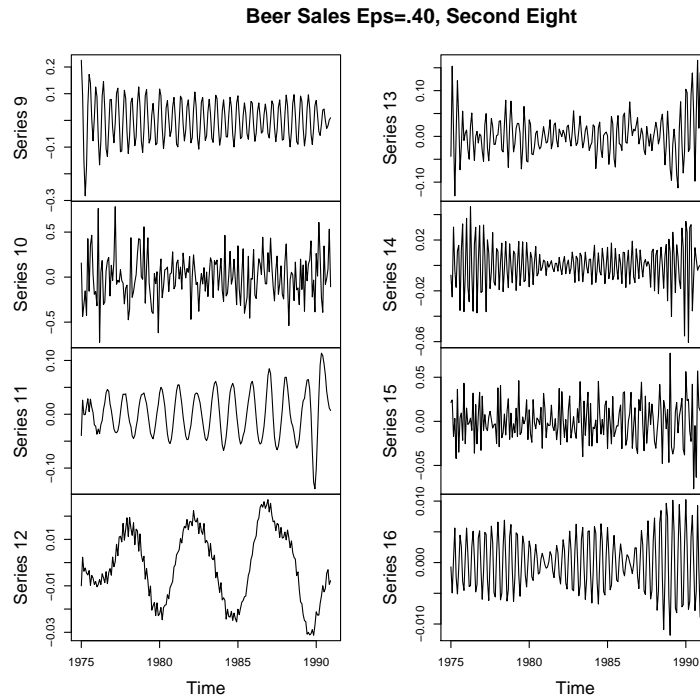


FIGURE 5. Mauna Loa  $CO_2$  Data

(a) Mauna Loa  $CO_2$ ,  $L = 234$ ,  $\epsilon = 0.25$ (b) Mauna Loa  $CO_2$ ,  $L = 234$ ,  $\epsilon = 0.25$ FIGURE 6. Mauna Loa  $CO_2$

(a) Mauna Loa  $CO_2$  Centered,  $L = 234$ ,  $\epsilon = 0.25$ (b) Mauna Loa  $CO_2$  Centered,  $L = 234$ ,  $\epsilon = 0.25$ FIGURE 7. Mauna Loa  $CO_2$  Centered

**SSA beersales, L=48**FIGURE 8. Beer Sales,  $L = 48$

(a) Beer Sales  $\epsilon = 0.40$ , First Eight(b) Beer Sales  $\epsilon = 0.40$ , Second EightFIGURE 9. Beer Sales  $\epsilon = 0.40$