

A new method of nonparametric density estimation with application in photovoltaics

Andrey Pepelyshev
Ansgar Steland, Nina Golyandina

Institut für **RWTH**
statistik **AACHEN**
& **W**irtschafts-
mathematik

Wismar
March 2, 2011

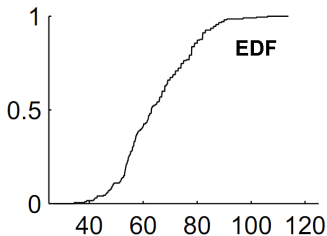
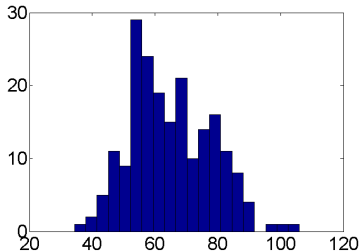
Contents

- Introduction
- A new density estimate
- A new choice of a smoothing parameter
- Application in photovoltaic

Introduction

We have a sample (x_1, \dots, x_m) from a distribution with density $p(x)$.

The problem is to estimate $p(x)$.



The classic way is to estimate $p(x)$ directly.
Another way is to smooth the EDF and then differentiate.

How to smooth the EDF

- Consider values of the EDF at equidistant points as a time series.
- Apply Singular Spectrum Analysis (SSA) for its smoothing.

Golyandina, N., Nekrutkin, V., Zhigljavsky, A. (2001) Analysis of Time Series Structure: SSA and Related Techniques. London: Chapman & Hall/CRC.

The SSA procedure for a CDF series

- Embed a series (f_1, \dots, f_N) to a matrix \mathbb{X} as
 $\mathbb{X} = (x_{i,j}) = (f_{i+j-1}), i=1, \dots, L, j=1, \dots, N-L+1.$
- Make the SVD decomposition
 $\mathbb{X} = \sum_{i=1}^L \sqrt{\lambda_i} U_i V_i^T, \lambda_1 \geq \lambda_2 \geq \dots$

$$\mathbb{X} = \begin{array}{|c|c|c|c|c|} \hline f_1 & f_2 & f_3 & \cdots & f_{N-L+1} \\ \hline f_2 & f_3 & f_4 & \cdots & f_{N-L+2} \\ \hline f_3 & f_4 & f_5 & \cdots & f_{N-L+3} \\ \hline \vdots & \vdots & \vdots & & \vdots \\ \hline f_L & f_{L+1} & f_{L+2} & \cdots & f_N \\ \hline \end{array}$$

The SSA procedure for a CDF series

- Embed a series (f_1, \dots, f_N) to a matrix \mathbb{X} as
 $\mathbb{X} = (x_{i,j}) = (f_{i+j-1}), i = 1, \dots, L, j = 1, \dots, N - L + 1.$
- Make the SVD decomposition
 $\mathbb{X} = \sum_{i=1}^L \sqrt{\lambda_i} U_i V_i^T, \lambda_1 \geq \lambda_2 \geq \dots$
- Extract the most important components
 $\mathbb{X}^{(r)} = (x_{i,j}^{(r)}) = \sum_{i=1}^r \sqrt{\lambda_i} U_i V_i^T.$
- Transform $\mathbb{X}^{(r)}$ back to the series

$$\hat{f}_j = \hat{f}_j(L, r) = \begin{cases} 0 & 1 \leq j < L, \\ \frac{1}{L} \sum_{k=1}^L x_{k,j-k+1}^{(r)} & L \leq j \leq K, \\ 1 & K < j \leq N, \end{cases}$$

Representation of the SSA procedure

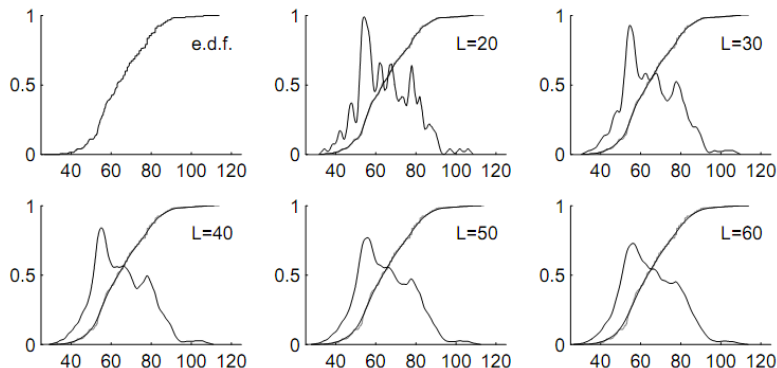
$$\hat{f}_j = \sum_{i=1}^L \sum_{l=1}^L (u_{1,i}u_{1,l} + \dots + u_{r,i}u_{r,l}) f_{j+i-l} / L$$

for $j = L, \dots, K$, where $(u_{l,1}, \dots, u_{l,L})^T = U_l$.

- This is a data-adaptive filter.
- L is a half-length of the filter.
- r is a complexity of the filter.
- L controls the smoothness.

Influence of the parameter L

Measurements of Learn Body Mass from
the Australian Institute of Sport Data



The SSA^{1c} estimate (i.e. with $r = 1$) of the
distribution function and the density for $L=20, \dots, 60$

Validity of the SSA procedure

The estimated series $(\hat{f}_1, \dots, \hat{f}_N)$ can be transformed to the SSA estimate of $F(x)$ as follows

$$\hat{F}_m(x) = \sum_{j=2}^N \left(\hat{f}_{j-1} + (\hat{f}_j - \hat{f}_{j-1}) \frac{x - t_{j-1}}{t_j - t_{j-1}} \right) \mathbf{1}_{[t_{j-1}, t_j)}(x) + \mathbf{1}_{[t_N, \infty)}(x),$$

Lemma.

The SSA^{1c} estimate is a distribution function.

The SSA^b estimate

Let $S(F)$ be a filter (weighted moving average).

Our case: $S(F)$ is the SSA^{1c} estimate.

The problem: this filter has a bias.

The idea is to correct $S(F)$ as an estimate F by use of $S(S(F))$ and $S(S(S(F)))$.

$$S^b(\mathcal{F}) = 3S(\mathcal{F}) - 3S^2(\mathcal{F}) + S^3(\mathcal{F})$$

How to choose an estimate automatically

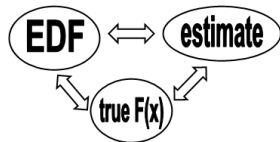
A law of the iterated logarithm.

The empirical distribution function $F_m(x)$ satisfies

$$\limsup_{m \rightarrow \infty} R_m \|F_m(x) - F(x)\|_{\infty} \leq 1$$

almost surely, where $R_m = \frac{2\sqrt{m}}{\sqrt{2 \ln \ln m}}$.

If an estimate is far from the EDF,
then this estimate is not good.



How to choose L automatically

1. Compute

$$\bar{L} = \max \left\{ L : \|\mathcal{F}_m - \hat{\mathcal{F}}_m(l)\|_\infty \leq 1/R_m \forall l \in \{1, \dots, L\} \right\}.$$

2. Compute the sequence $M_1, \dots, M_{\bar{L}}$, where M_j be the number of modes of estimated density for $L = j$. This sequence has decreasing tendency.

3. Compute $\check{M}_j = \min\{M_1, M_2, \dots, M_j\}$.

How to choose L automatically

4. Divide the set $\{1, 2, \dots, \bar{L}\}$ into groups such that the values \check{M}_j are equal to each other within each group, i.e. define $a_1, b_1, \dots, a_k, b_k$ and k such that

$$1 = a_1 \leq b_1 < \dots < a_k \leq b_k = \bar{L},$$

$a_{i+1} = b_i + 1$ and $\check{M}_j = \check{M}_j$ for all $i, j \in \{a_l, \dots, b_l\}$, $l \in \{1, \dots, k\}$.

How to choose L automatically

5. Finally, compute L_a (a 'best' value of L) as an average with weight coefficients, which are proportional to the sizes of these k groups, namely

$$L_a = \left[\sum_{i=1}^k c_i w_i \right],$$

where

$$c_i = \gamma_i a_i + (1 - \gamma_i) b_i, \quad w_i = \frac{b_i - a_i}{\sum_{j=1}^k b_j - a_j},$$

$\gamma_i = 1/2$ if $\check{M}_{a_i} = 1$ and $\gamma_i = 0.9$ otherwise.

How to choose the bandwidth

1. Compute

$$\bar{h} = \max \left\{ h : \left\| \mathcal{F}_m - \hat{\mathcal{F}}_m(\bar{h}) \right\|_{\infty} \leq 1/R_m \forall \bar{h} \in (0, h) \right\}.$$

2. Define a dense set $\{h_1, \dots, h_n\} \in (0, \bar{h})$ and compute the sequence M_1, \dots, M_n , where M_j be the number of modes of estimated density for $h = h_j$.

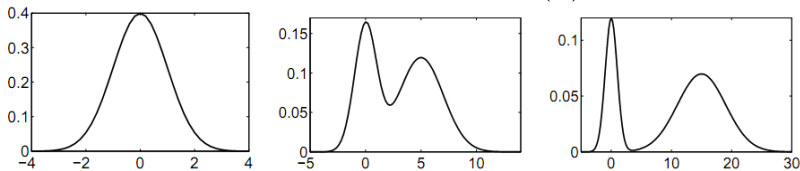
3. Compute $\check{M}_j = \min\{M_1, M_2, \dots, M_j\}$.

4. Divide the set $\{h_1, \dots, h_n\}$ into groups similarly.

5. Compute $h_a = \sum_{i=1}^k c_i w_i$, where c_i and w_i are defined in the same manner.

Simulation study

We take several models of the density $p(x)$



and simulate 10000 samples of size 100.

$$D_{\text{ISE}}(\hat{p}) = \int (\hat{p}(x) - p(x))^2 dx$$

$$D_{\text{KS}}(\hat{F}) = \|\hat{F}(x) - F(x)\|_{\infty} = \max_x |\hat{F}(x) - F(x)|$$

$$D_{\text{H}}(\hat{p}) = \int \left(\sqrt{\hat{p}(x)} - \sqrt{p(x)} \right)^2 dx$$

Results on bandwidth selection

	ED_{ISE}	ED_{KS}	ED_H
model $N(0, 1)$			
Kernel est. with h_{LSCV}	0.0071	0.0551	0.0143
Kernel est. with h_{SJPI}	0.0066	0.0536	0.0131
Kernel est. with h_{ICV}	0.0075	0.0546	0.0146
Kernel est. with h_a	0.0063	0.0546	0.0131
model $0.4N(0, 1) + 0.6N(5, 2^2)$			
Kernel est. with h_{LSCV}	0.0058	0.0617	0.0231
Kernel est. with h_{SJPI}	0.0052	0.0609	0.0210
Kernel est. with h_{ICV}	0.0055	0.0614	0.0229
Kernel est. with h_a	0.0053	0.0611	0.0236
model $0.3N(0, 1) + 0.7N(15, 4^2)$			
Kernel est. with h_{LSCV}	0.0048	0.0672	0.0394
Kernel est. with h_{SJPI}	0.0069	0.0733	0.0602
Kernel est. with h_{ICV}	0.0049	0.0670	0.0396
Kernel est. with h_a	0.0050	0.0674	0.0440

Results on the SSA estimates

	ED_{ISE}	ED_{KS}	ED_H	EL_a
model $N(0, 1)$				
Kernel est. with h_a	0.0063	0.0546	0.0131	
SSA ^{1c} est.	0.0061	0.0537	0.0128	107.6
SSA ^{2c} est.	0.0060	0.0503	0.0142	145.7
SSA ^b est.	0.0052	0.0488	0.0141	141.7
model $0.4N(0, 1) + 0.6N(5, 2^2)$				
Kernel est. with h_a	0.0053	0.0611	0.0236	
SSA ^{1c} est.	0.0051	0.0607	0.0215	67.6
SSA ^{2c} est.	0.0047	0.0610	0.0195	86.4
SSA ^b est.	0.0052	0.0617	0.0206	89.7
model $0.3N(0, 1) + 0.7N(15, 4^2)$				
Kernel est. with h_a	0.0050	0.0674	0.0440	
SSA ^{1c} est.	0.0053	0.0679	0.0451	40.9
SSA ^{2c} est.	0.0046	0.0660	0.0380	60.0
SSA ^{3c} est.	0.0043	0.0643	0.0349	73.8
SSA ^b est.	0.0047	0.0670	0.0391	48.9

Sampling plans for quality control

Let measurements be distributed according to a distribution $G(x)$ with mean a and variance σ^2 . The asymptotically optimal sampling plan (n, c) is

$$n = \left\lceil \frac{(\Phi^{-1}(\alpha) - \Phi^{-1}(1 - \beta))^2}{(F^{-1}(\text{AQL}) - F^{-1}(\text{RQL}))^2} \right\rceil,$$

$$c = -\frac{\sqrt{n}}{2} (F^{-1}(\text{AQL}) + F^{-1}(\text{RQL})),$$

where $F(x) = G((x - a)/\sigma)$, AQL is the acceptable quality level, RQL is the rejectable quality level, α is the producer risk and β is the consumer risk.

Comparison of sampling plans

Means and standard deviations of sampling plan size using the empirical distribution function, the kernel density estimate and the SSA^{1c} , SSA^b and SSA^{2c} estimates for samples of size m from $0.4N(0, 1) + 0.6N(5, 2^2)$.

The true sampling plan size is 392.

	$m = 250$	$m = 500$	$m = 1000$
EDF	469.7(468.9)	462.7(278.9)	422.0(160.1)
Kernel h_{LSCV}	322.3(102.8)	324.0(79.7)	337.8(59.2)
Kernel h_a	302.1(77.2)	321.2(65.2)	339.1(52.8)
SSA^{1c} est.	301.4(71.8)	316.3(61.3)	329.1(46.9)
SSA^b est.	409.6(112.9)	408.9(88.8)	403.6(67.7)
SSA^{2c} est.	390.4(107.8)	393.1(81.4)	392.8(62.1)

Thank you for your attention!