

Statistical Estimation: Data & Non-data Information

Roger J-B Wets

University of California, Davis

& M.Casey @ Raytheon

G.Pflug @ U. Vienna,

X. Dong @ EpiRisk,

G-M You @ EpiRisk.



a little background

- Decision making under uncertainty:

$$\min f_{01}(x) + E\{Q(\xi, x)\} \text{ so that } x \in C_1 \subset \mathbb{R}^n$$

$$Q(\xi, x) = \inf \left[f_{02}(\xi, y) \mid y \in C_2(\xi, x) \right]$$

$$C_2(\xi, x) = \left\{ y \in \mathbb{R}_+^d \mid W_\xi y = d_\xi - T_\xi x \right\}$$

- = stochastic programming problems.
- Issue: realizations of $(d_\xi, T_\xi, W_\xi) = \xi \in \mathbb{R}^N, N \gg 2!$
data: 43 samples, ... never reaches the asymptotic range

Formulation

- Find F^{est} , an estimate of the distribution of a random (phenomena) variable X given **all** the information available about this random phenomena,
- i.e. such that

$$\forall x : F^{\text{est}}(x) \simeq F^{\text{true}}(x) = \text{prob.}[X \leq x].$$

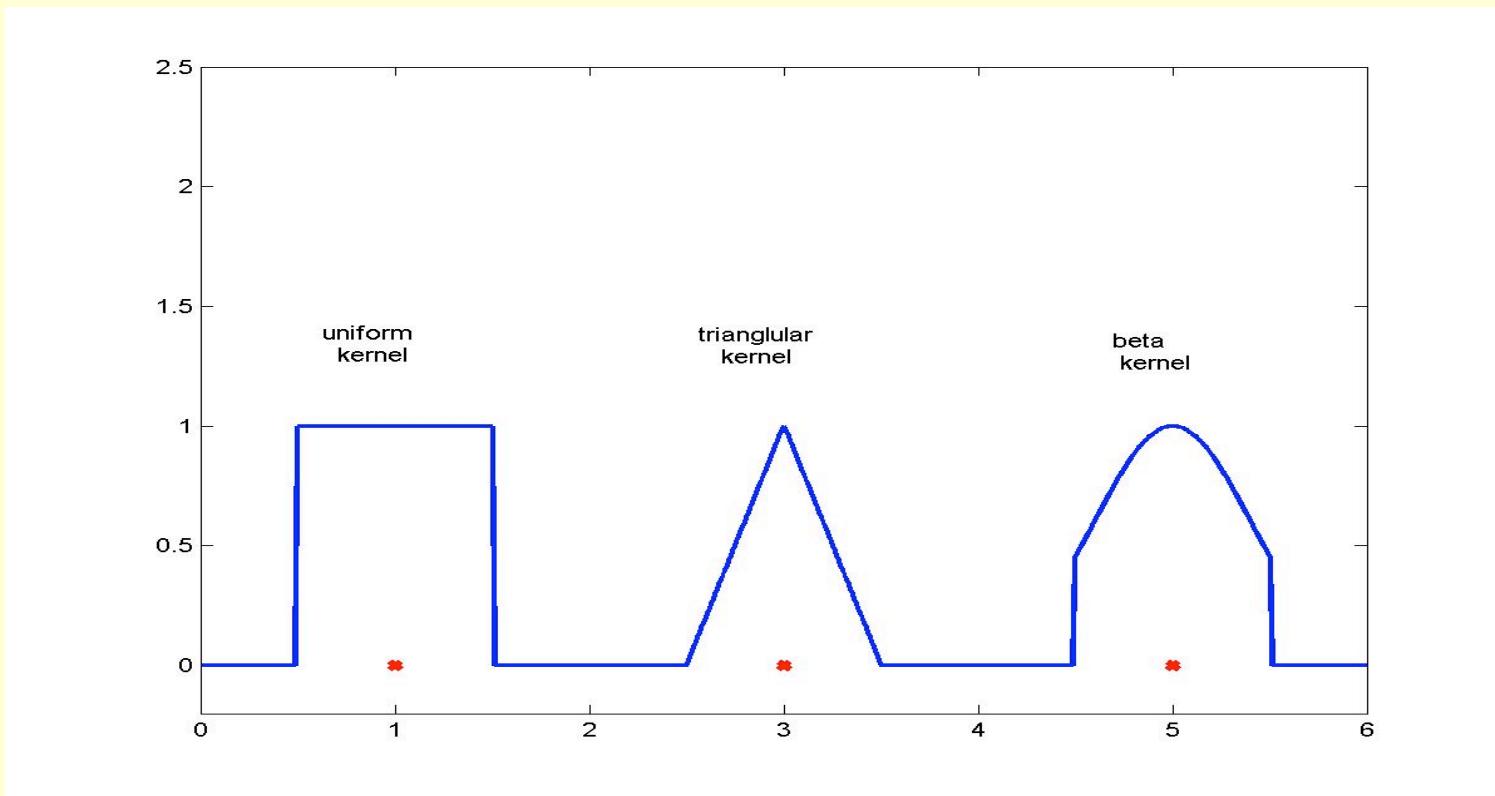
Information

- Observations (data): x_1, x_2, \dots, x_v
- Non-data facts:
 - density or discrete distribution,
 - bounds on expectation, moments,
 - shape: unimodal, decreasing, parametric class
- Non-data modeling assumptions:
 - see above + ...
 - density is smooth, (un)bounded support, ..

Applications:

- Estimating (cum.) distribution functions
- Estimating coefficient of time series
- Estimating coefficients of stochastic differential equation (SDE)
- Estimating financial curves (zero-curves)
- Dealing with lack of data: *few observations*
- Estimating density functions: h^{est}

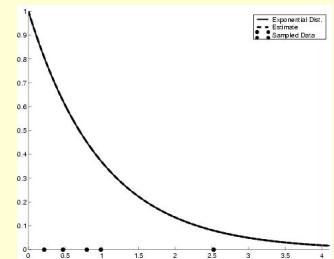
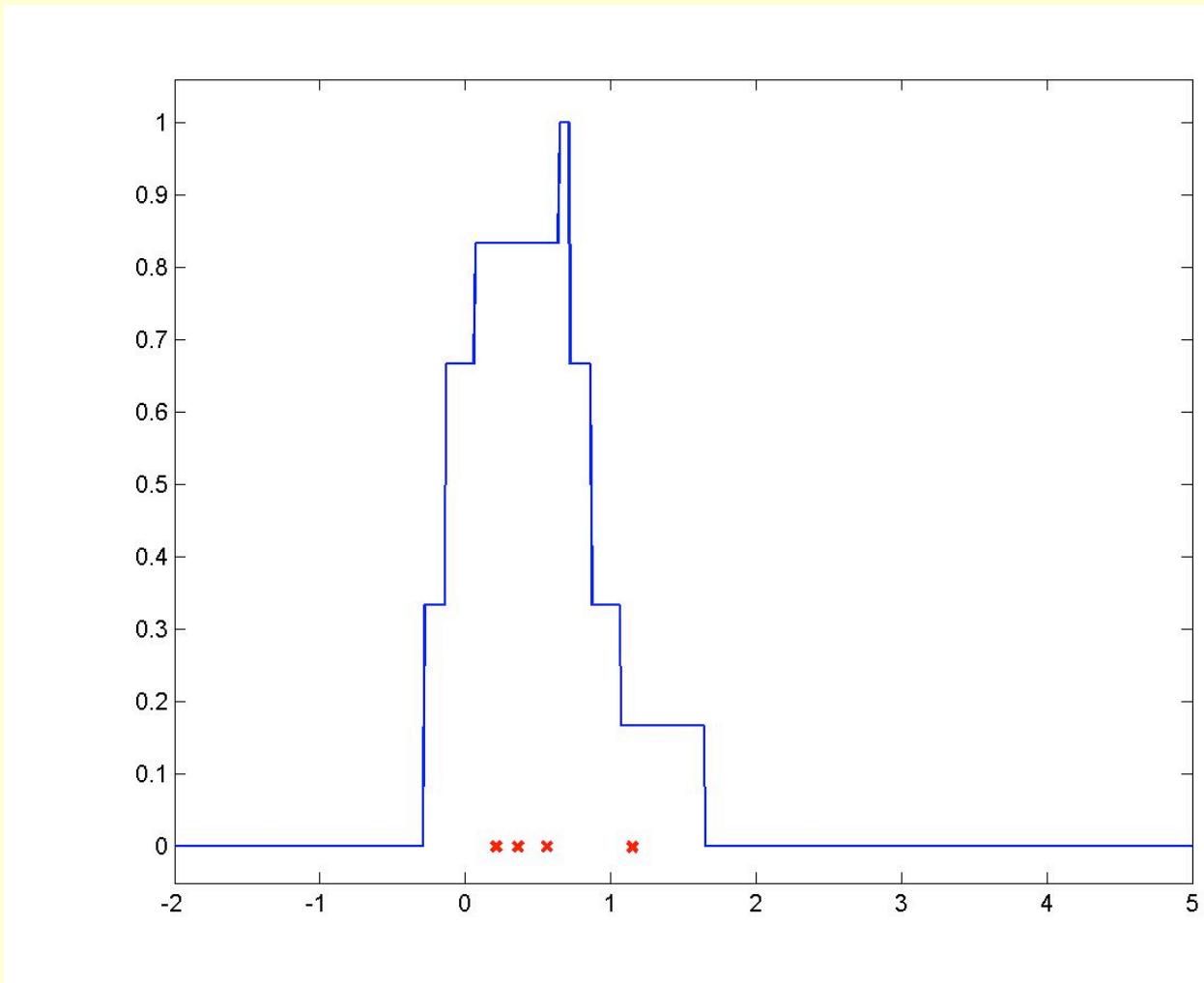
Kernel-choice



Information = Observations: x_1, x_2, \dots, x_v

Optimal bandwidth = kernel support ?

Kernel estimates



Statistical Estimation from an optimization viewpoint



Criterion: Maximum Likelihood

- Find $h \in H = \text{class-fcns}(\mathbb{R})$ that maximizes the probability of observing x_1, x_2, \dots, x_v

$$\max \prod_{l=1}^v h(x_l)$$

- equivalently:

$$\max \frac{1}{v} \sum_{l=1}^v \ln h(x_l) =$$

$$\max E^v \{\ln h(x)\} = \max \int \ln h(x) P^v(dx)$$

leads to:
 ∞ -dim. optimization problem

$$\max E^\nu \left\{ \ln h(x) \right\} = \frac{1}{\nu} \sum_{l=1}^{\nu} \ln h(x_l)$$

so that $\int h(x) dx = 1$,

$$h(x) \geq 0, \quad \forall x \in \mathbb{R}$$

$$h \in \mathcal{A}^\nu \subset H$$

\mathcal{A}^ν : non-data information constraints

$$H = C^2(\mathbb{R}), L^p(S), H^l(S), \dots \quad S \subset \mathbb{R}$$

Non-data information constraints

support: $S = [\alpha, \beta]$, $S = [\alpha, \infty)$, ...

bounds on moments: $a \leq \int x h(x) dx \leq b$

unimodal: $h(x) = e^{-Q(x)}$, Q : convex

decreasing: $h'(x) \leq 0$, $\forall x \in S$

'smoothness': $h \in H_0^1$, $\int \frac{h'(x)^2}{h(x)} dx \leq \kappa$ (Fisher info.)

'Bayesian': $|h - h_a| \leq \beta$ or $\theta(|h - h_a|)$, Total variation, ...

Questions (M^+ -estimators)

- Consistency: as $v \rightarrow \infty$ does $h^{est} \rightarrow h^{true}$?
hypo-convergence: “a law of large numbers”
- Convergence rate: how large does v have to be so that $h^{est} \sim h^{true}$?
quantitative hypo-convergence: $\text{dist}(h^{est}, h^{true})$
- Solving the ∞ -dimensional optimization problem?
finite dimensional approximations
+ again convergence issues.

When is (\mathcal{P}) “near” (\mathcal{P}^a)

$$(\mathcal{P}) \quad \max f_0(x)$$

$$\text{s.t. } f_i(x) \leq 0, \ i \in I_f$$

$$x \in X_f \subset E$$

$$(\mathcal{P}^a) \quad \max g_0(x)$$

$$\text{s.t. } g_i(x) \leq 0, \ i \in I_g$$

$$x \in X_g \subset E$$

Example:

$$(\mathcal{P}^a)$$

$$(\mathcal{P}) \quad \max f_0(x)$$

$$\text{s.t. } f_i(x) \leq 0, \ i = 1, \dots, m$$

$$\max f_0(x) + \nu \sum_{i=1}^m \max[0, f_i(x)]^2$$

From: when is (\mathcal{P}) “near” (\mathcal{P}^a)

$(\mathcal{P}) \quad \max f(x) \text{ where}$

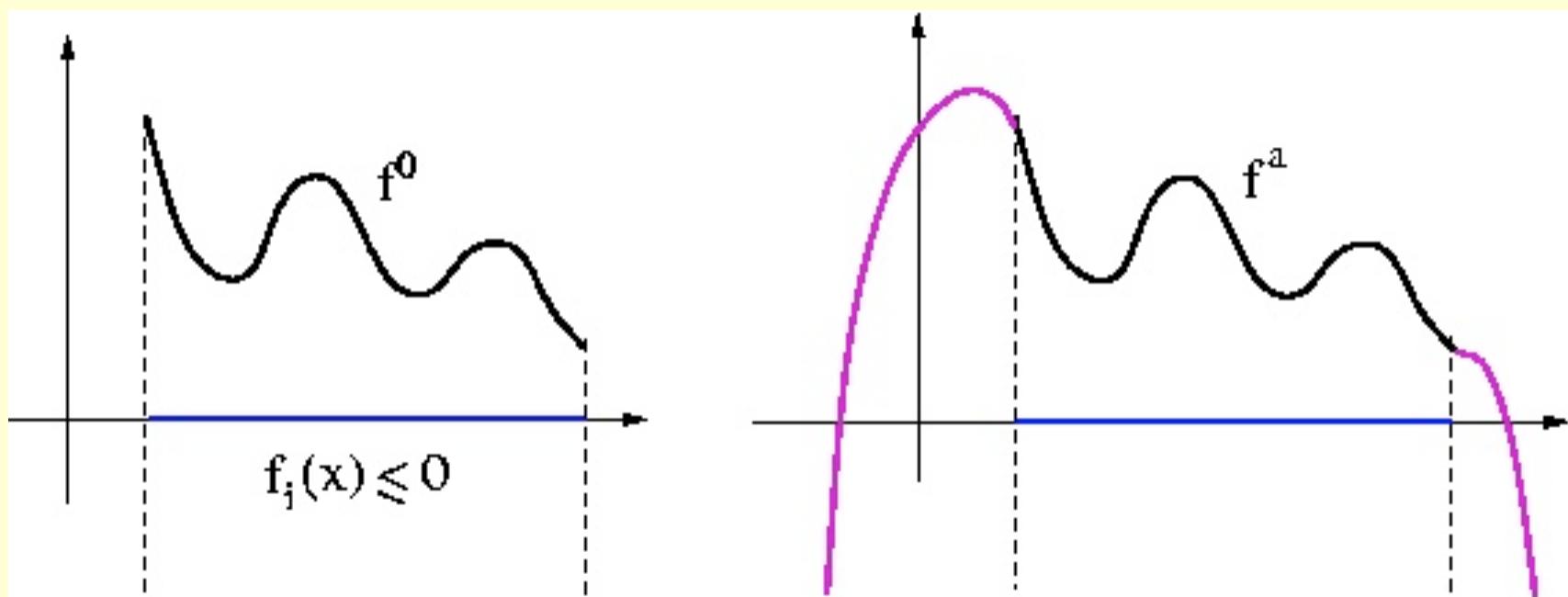
$f(x) = f_0(x) \text{ when } f_i(x) \leq 0, i \in I_f, x \in X_f \subset E$
 $- \infty \text{ otherwise}$

$(\mathcal{P}^a) \quad \max g(x) \text{ where}$

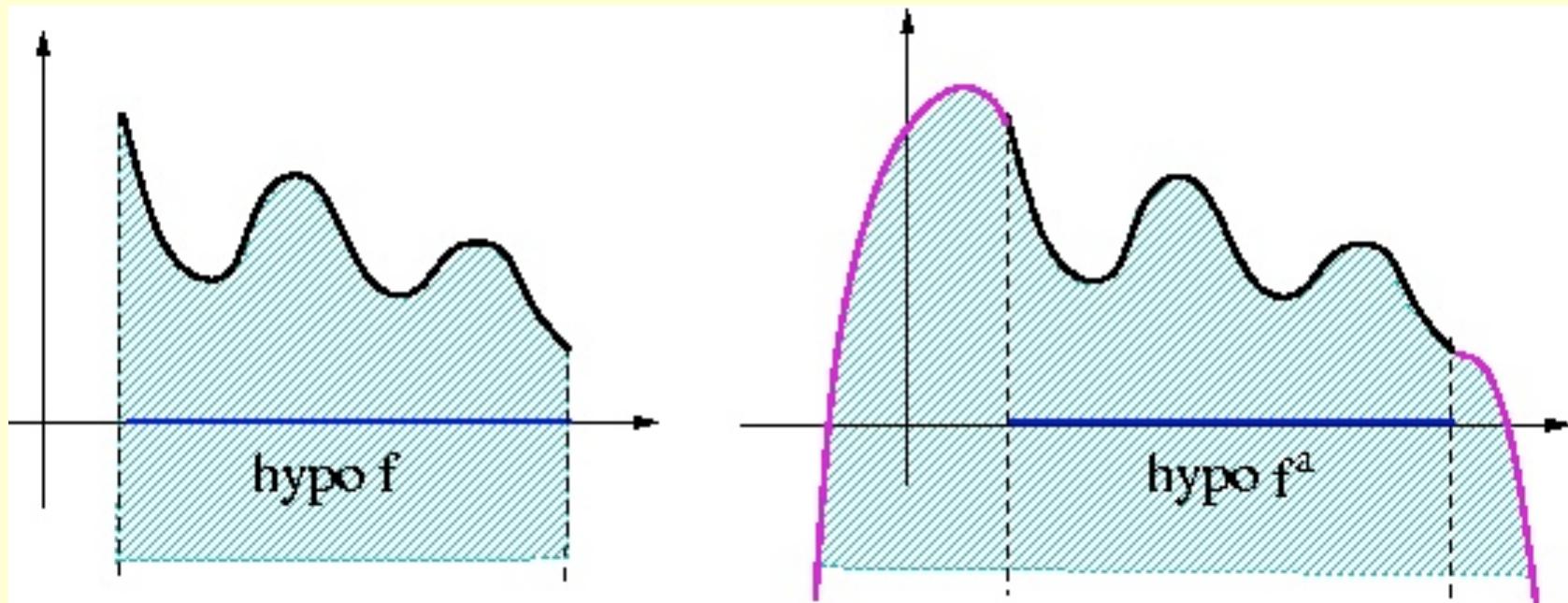
$g(x) = g_0(x) \text{ when } g_i(x) \leq 0, i \in I_g, x \in X_g \subset E$
 $- \infty \text{ otherwise}$

To: when is f “near” g

Example (with penalty)



Hypographs of f & f^a



$\text{hypo } f^a$ “near” $\text{hypo } f$

→ $\text{argmax } (\mathcal{P}^a) \sim \text{argmax } (\mathcal{P})$

$$f = h\text{-}\lim_\nu f^\nu, \quad f^\nu \xrightarrow{h} f$$

- $x^\nu \in \arg \max f^\nu, x^\nu \underset{\text{cluster}}{\rightarrow} x \Rightarrow x \in \arg \max f$

aw-topology for usc functions (Attouch-Wets, 1992)

- $\exists dl$ such that $dl(f^\nu, f) \rightarrow 0 \Leftrightarrow f^\nu \xrightarrow{h} f$

$$dl(f^\nu, f) \leq \eta \Rightarrow \text{dist}(\arg \max f^\nu, \arg \max f) \leq \kappa \psi^{-1}(\eta)$$

- $F^\nu(x) = \int f(\xi, x) P^\nu(d\xi) \quad \& \quad F(x) = \int f(\xi, x) P(d\xi)$

$$P^\nu \xrightarrow{n} P \quad \& \quad f\text{-tight} \Rightarrow F^\nu \xrightarrow{h} F$$

Quantitative hypo-convergence

- $\psi : \mathbb{R}_+ \rightarrow [0,1]$ continuous, strictly increasing
 - for example: $\psi(a) = (2/\pi) \arctan a$
- $dl_\rho(C, D) = \sup_{|x| \leq \rho} |\psi(d(x, C)) - \psi(d(x, D))|$
- $dl(C, D) = \int_0^\infty e^{-\rho} dl_\rho(C, D) d\rho$
 - metric on space of closed sets: τ_{aw} -topology
- $dl(f, g) = dl(\text{hypo } f, \text{hypo } g)$
 - metric on space of usc-fcns: τ_{aw} -topology

Quantitative hypo-convergence II

- $dl(f^\nu, f) \rightarrow 0 \Leftrightarrow f^\nu \xrightarrow{aw\text{-hypo}} f$
- Theorem:

$f, g : X \rightarrow \bar{\mathbb{R}}$, usc, concave,

$\rho_0 : \arg \max f \cap \rho_0 B \neq \emptyset$, $\sup f < \rho_0$ & for g

then $\forall \rho > \rho_0$, $\varepsilon > 0$:

$$dl_\rho(\varepsilon\text{-}\arg \max f, \varepsilon\text{-}\arg \max g) \leq (1 + \frac{4\rho}{\varepsilon}) dl_\rho(f, g)$$

“Opt”-Formulation

$$\max E^v \left\{ \ln h(x) \right\} = \frac{1}{v} \sum_{l=1}^v \ln h(x_l)$$

so that $\int h(x)dx = 1$,

$$h(x) \geq 0, \quad \forall x \in \mathbb{R}$$

$$h \in \mathcal{A}^v \subset H$$

\mathcal{A}^v : non-data information constraints

$$H = C^2(\mathbb{R}), L^p(S), H^l(S), \dots \quad S \subset \mathbb{R}$$

Convergence “rate”

$$L^v(h, x) = \ln h(x) \text{ if } h \in A^v \subset H, \int h(x) dx = 1, h \geq 0 \\ = -\infty \text{ otherwise}$$

$$L = L^v \text{ except for } A = aw\text{-}\lim A^v \quad (aw = dl)$$

$$\text{then } \int L^v(\bullet, x) dP^v = \boxed{E^v L^v \underset{aw\text{-hypo}}{\rightarrow} EL} = \int L(\bullet, x) dP^{true}$$

when H Hilbert, $H \xrightarrow[\text{embedding}]{cont.\text{-}compact} \check{H}$, \check{H} RKHS

constraint set: H -bounded, \check{H} -closed

\check{H} lin. subspace, $\forall x : x \mapsto h(x)$ cont.,

relies on the Hilbert-Schmidt Embedding Theorem

Numerical Procedures

$$1. \quad h = \sum_{k=1}^q u_k \varphi_k(\bullet)$$

Fourier coefficients, wavelets, “kernels-basis”

$$2. \quad h = \exp(s(\bullet))$$

$s(\bullet)$ constrained cubic (or quadratic) spline

Basis: Fourier coefficients

$$\max \frac{1}{\nu} \sum_{l=1}^{\nu} \ln \left[\frac{1}{\sqrt{\theta}} u_0 + \sqrt{\frac{2}{\theta}} \sum_{k=1}^q \cos\left(\frac{w_k \pi x_l}{\theta}\right) u_k \right]$$

$$\text{so that } u_0 + \sqrt{2} \sum_{k=1}^q \frac{\sin(w_k \pi)}{w_k \pi} u_k = \sqrt{\theta},$$

$$u_0 + \sqrt{2} \sum_{k=1}^q \cos\left(\frac{w_k \pi x}{\theta}\right) \geq 0, \quad \forall x \in [0, \theta],$$

$$u_k \in \mathbb{R}, k = 0, \dots, q, \quad w_k \in \mathbb{R}_+, k = 1, \dots, q$$

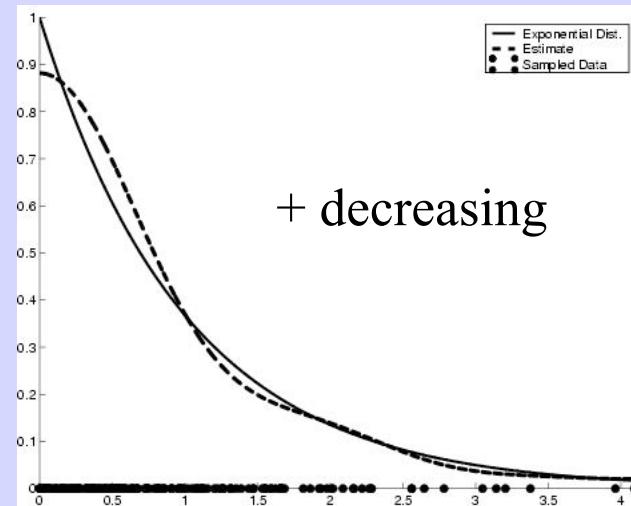
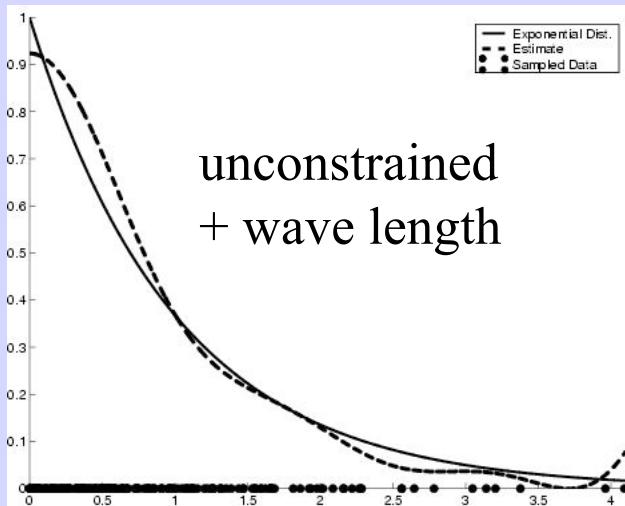
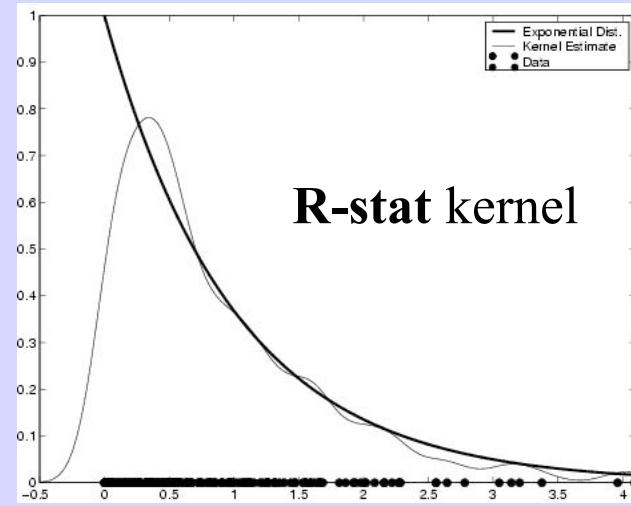
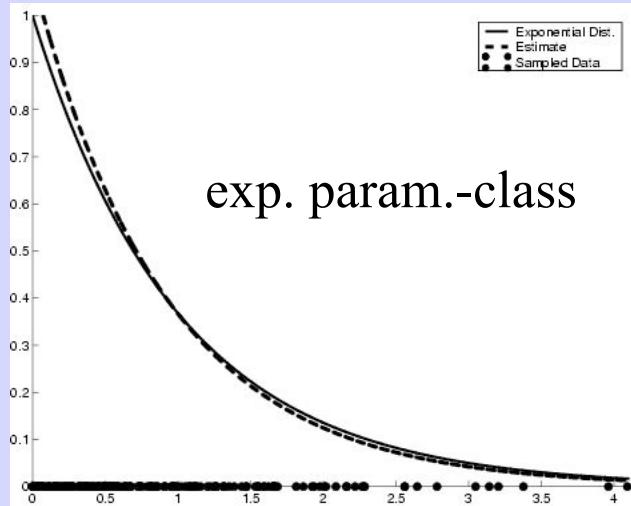
$$\sum_{k=1}^q \left[\cos\left(\frac{w_k \pi x}{\theta}\right) - \cos\left(\frac{w_k \pi x'}{\theta}\right) \right] u_k \geq 0, \quad \forall 0 \leq x \leq x' \leq \theta$$

i.e., h is decreasing

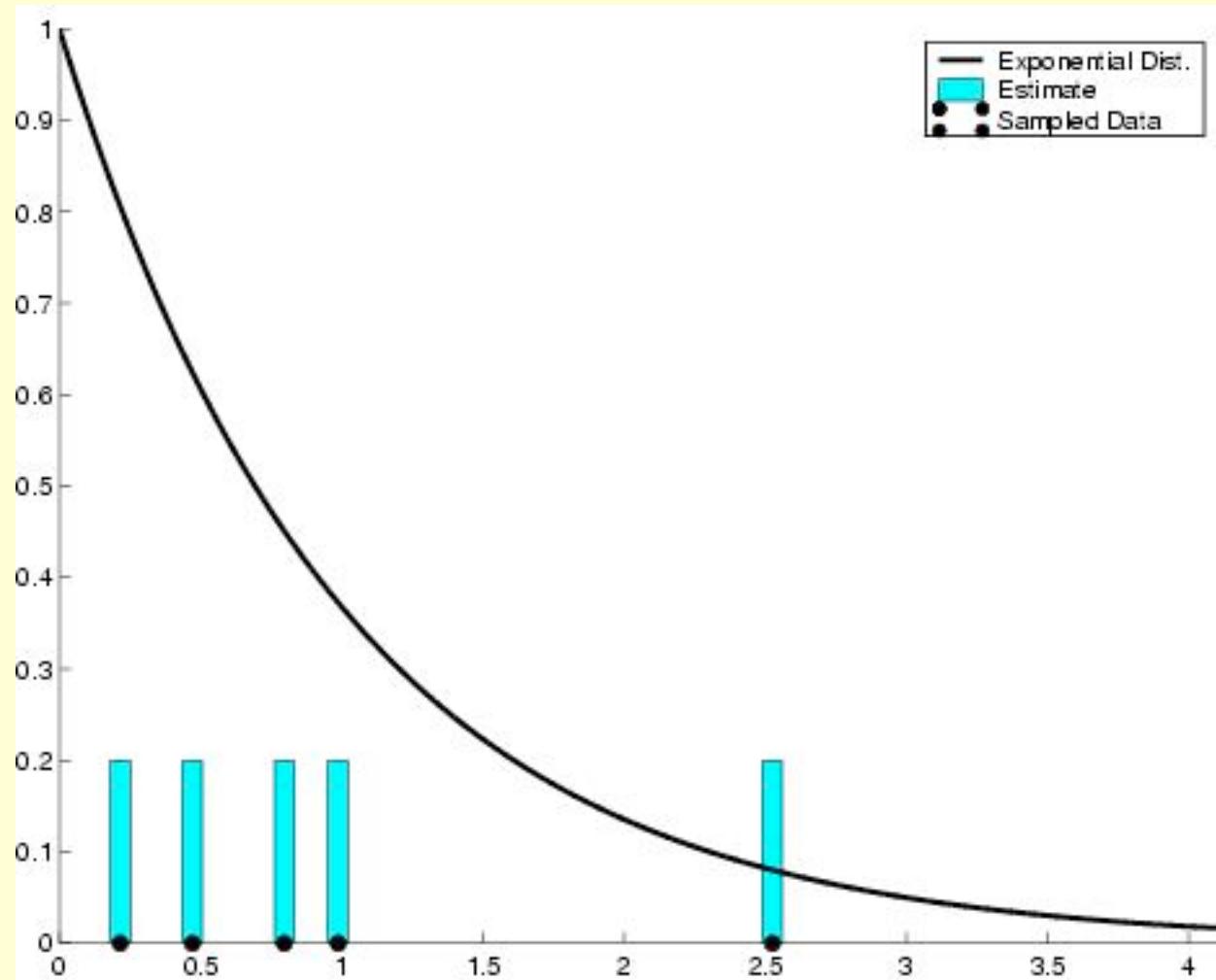
Test Case: exponential

- $h^{\text{true}}(x) = \lambda e^{-\lambda x}$ if $x \geq 0$; $= 0$ if $x < 0$ ($\lambda = 1$)
- “empirical” estimate
- kernel estimate from **R-stat**
- unconstrained with support $[0, \infty)$
- unconstrained with adaptive wave length
- constrained (h decreasing)
- parametric, i.e., $h \in \text{exp-class}$

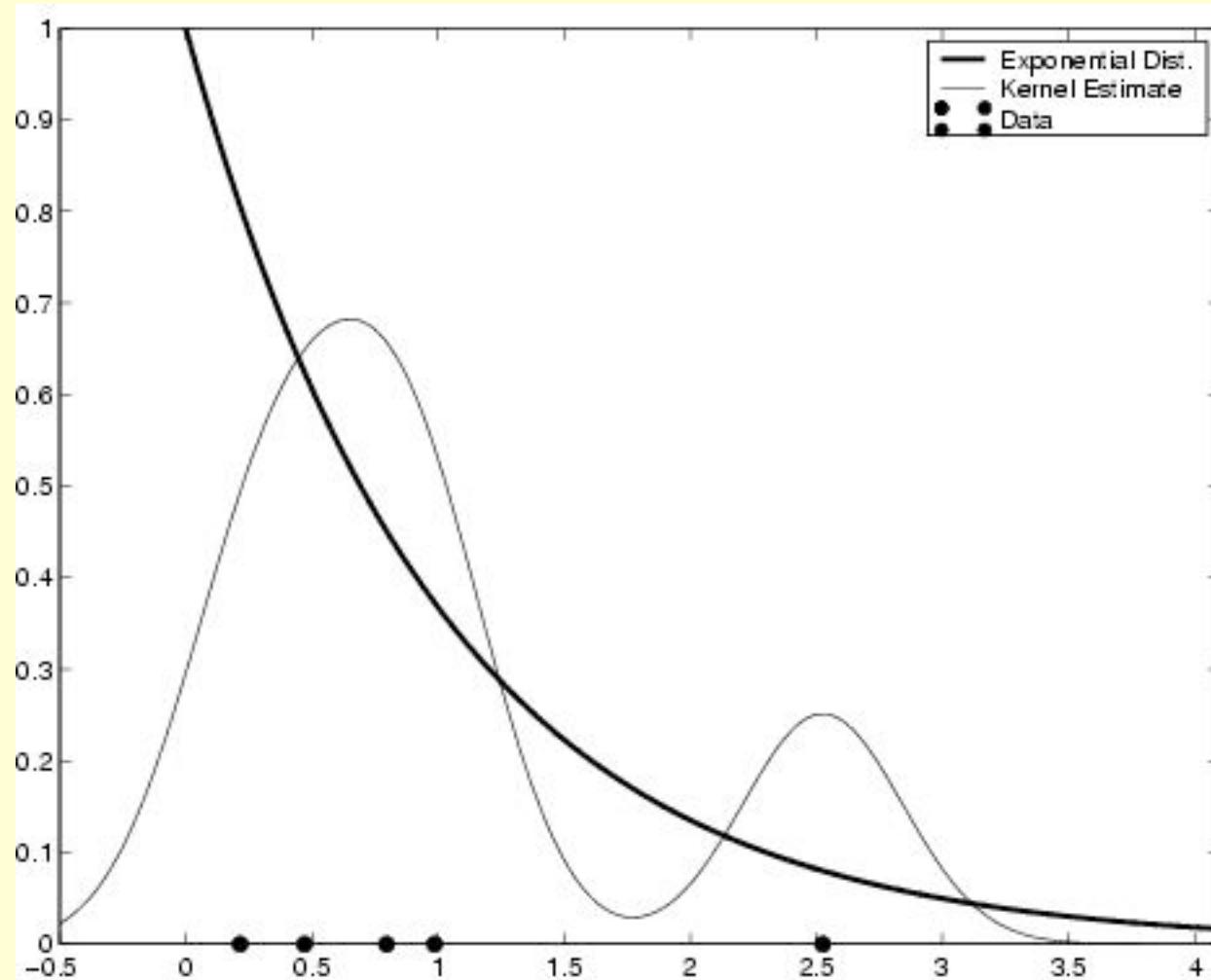
200-observations!



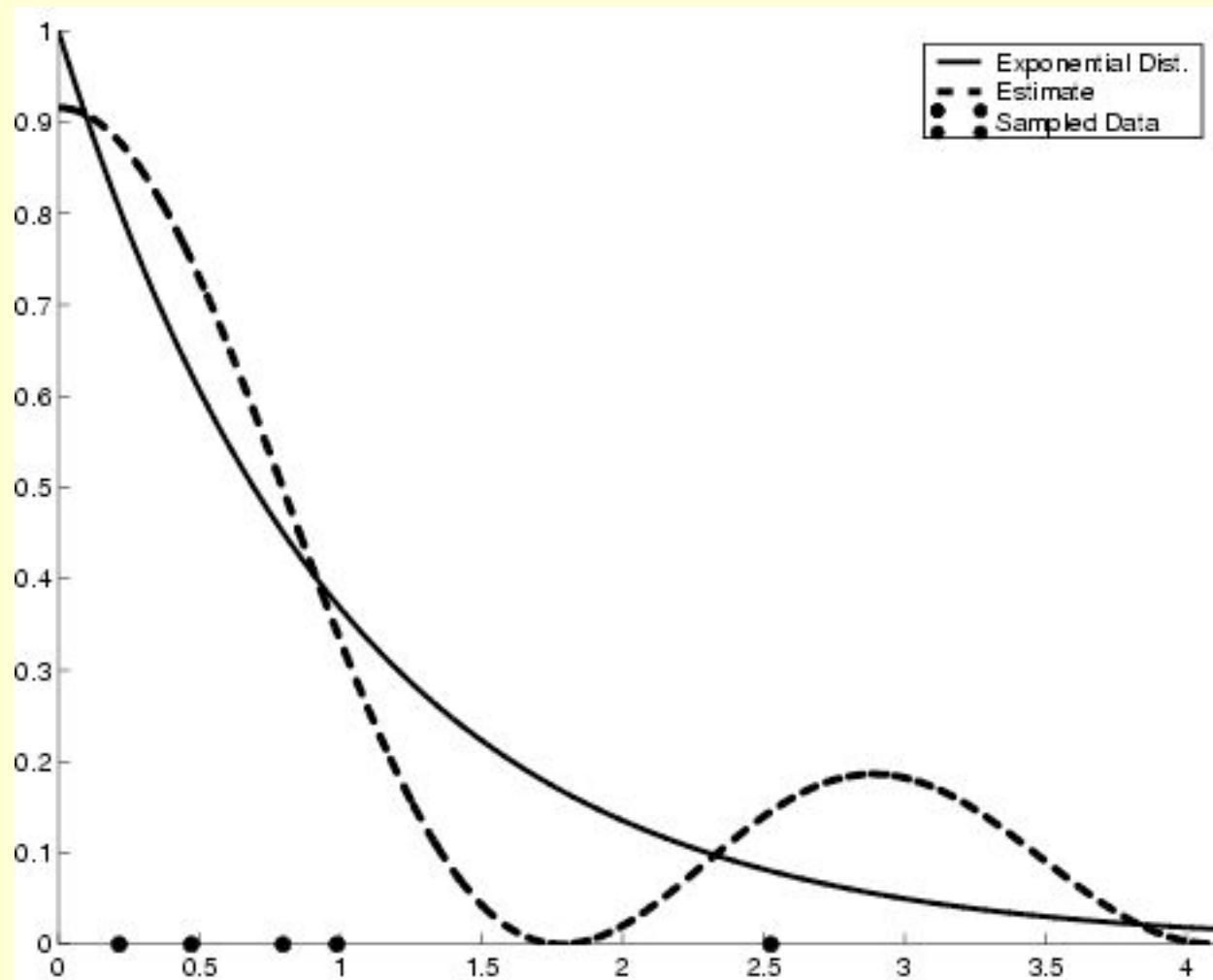
h^{est} : empirical (five observations)



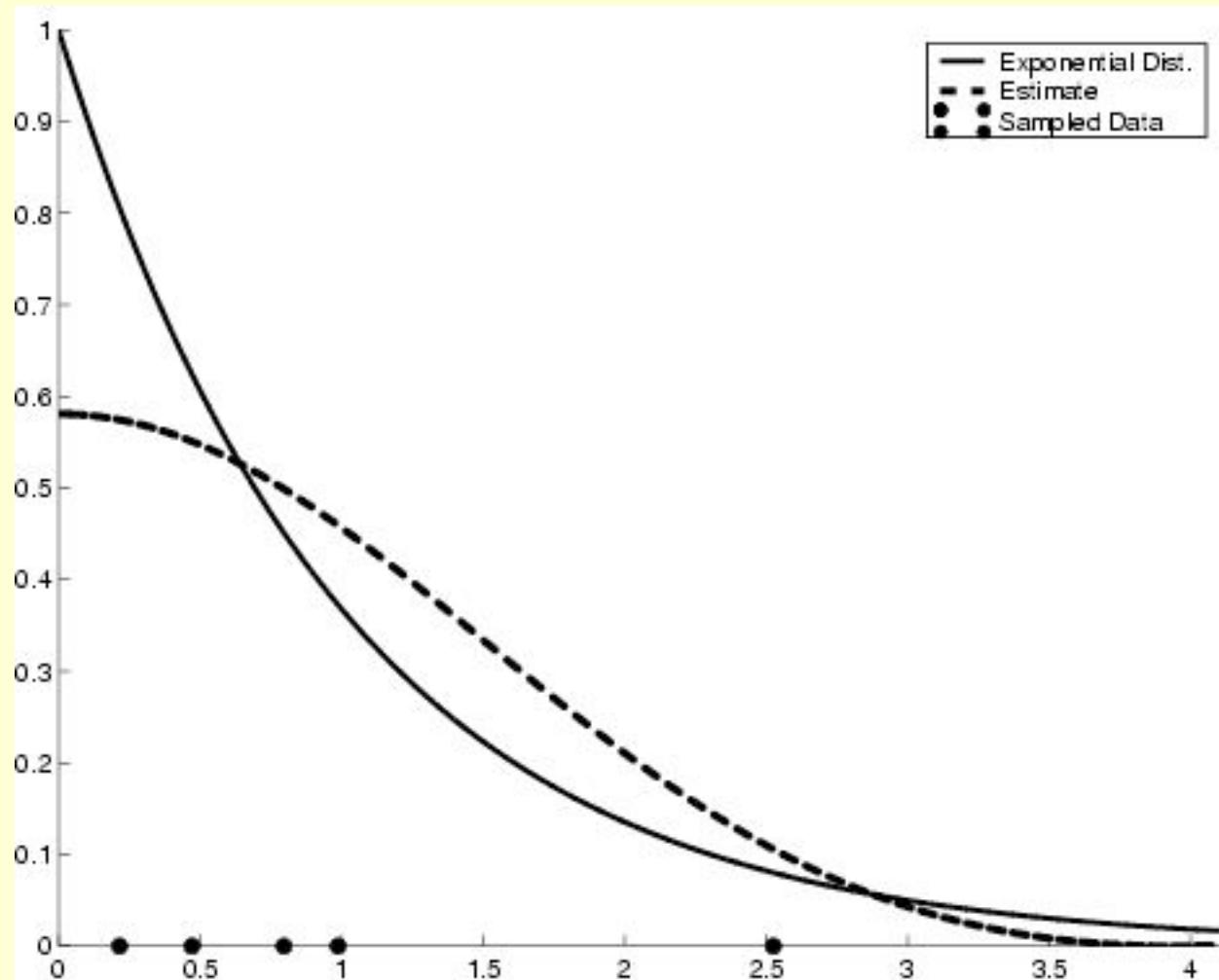
h^{est} : kernel



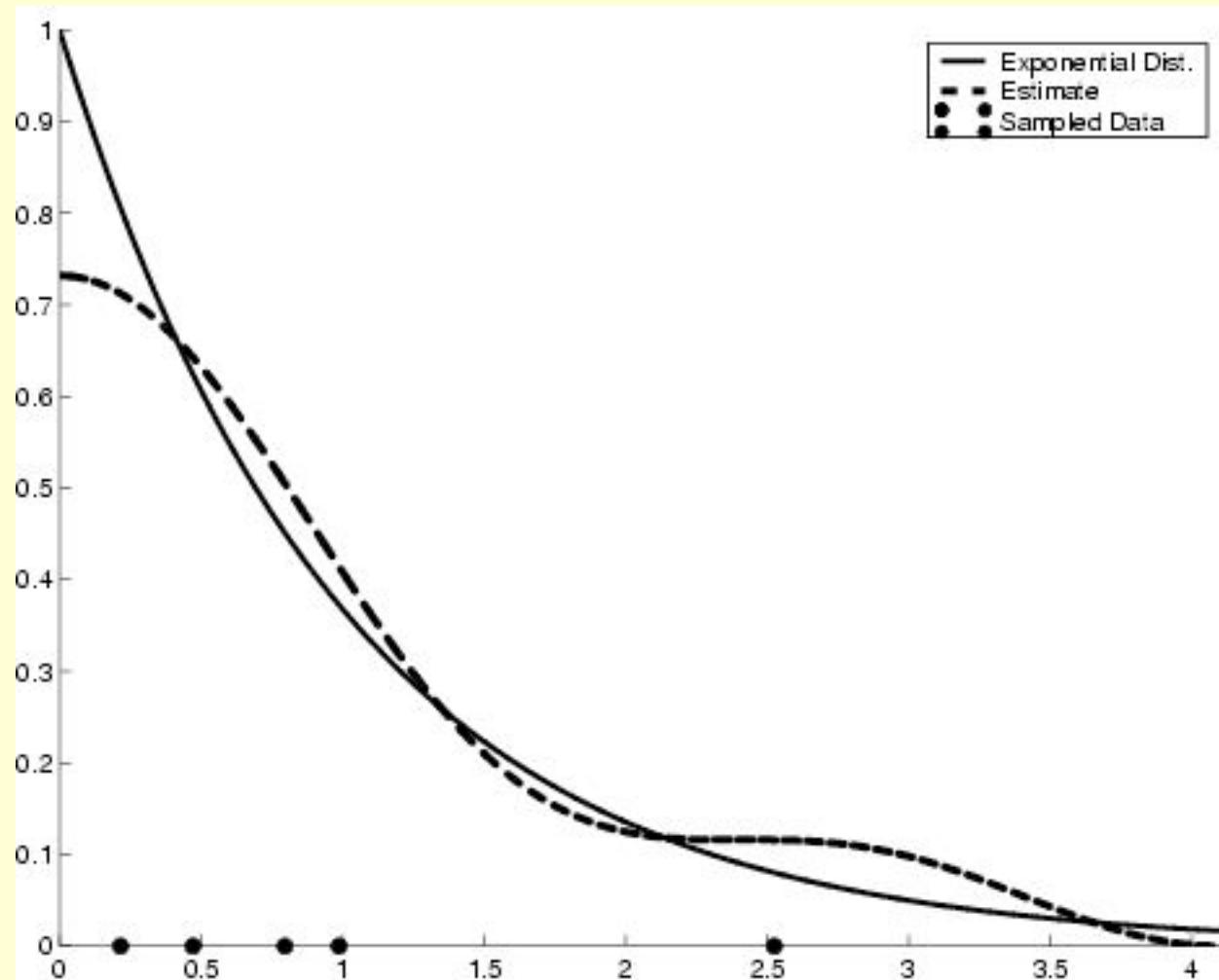
h^{est} : unconstrained



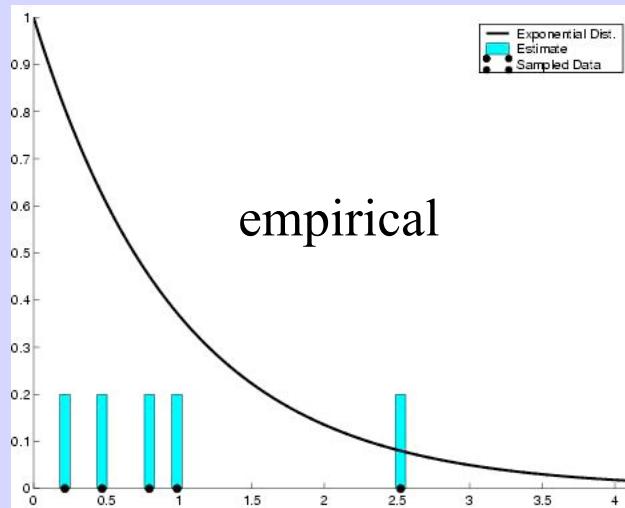
h^{est} : ..with adaptive wave #



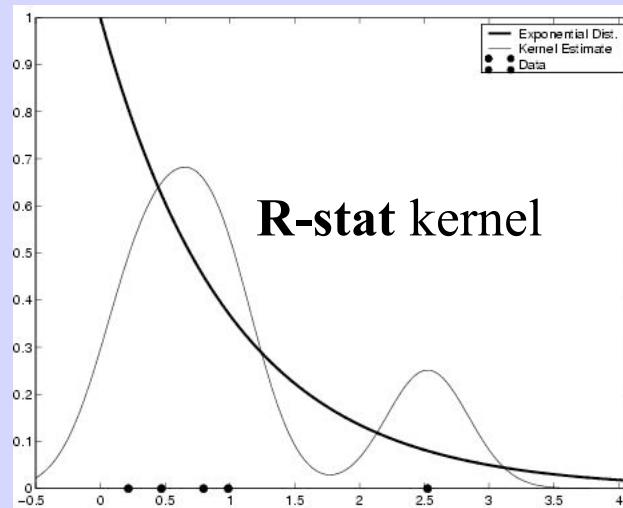
h^{est} : with decreasing constraint



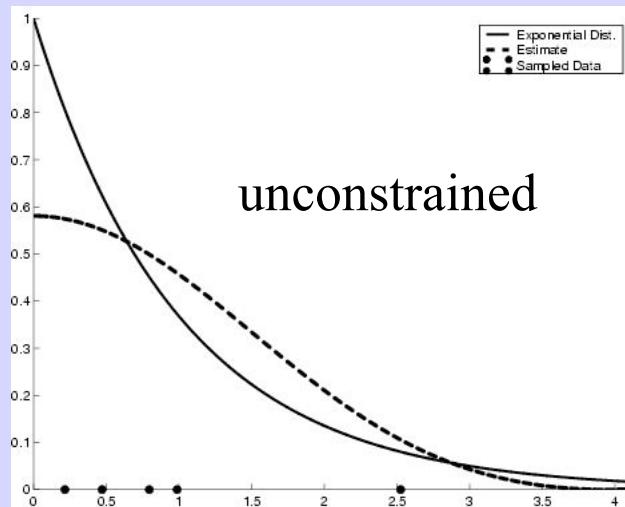
5-observations!



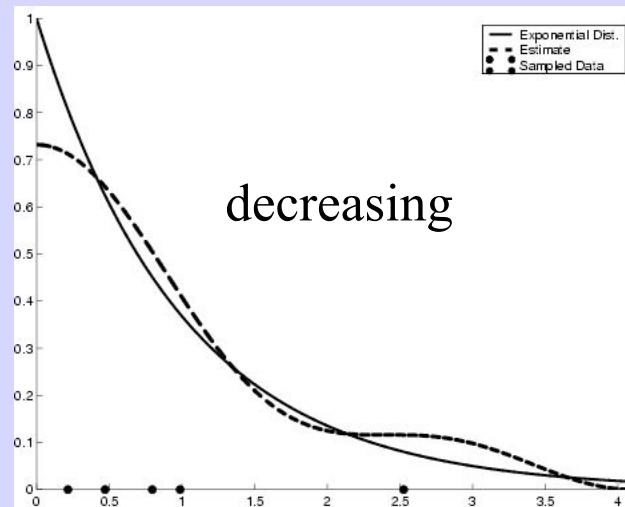
empirical



R-stat kernel

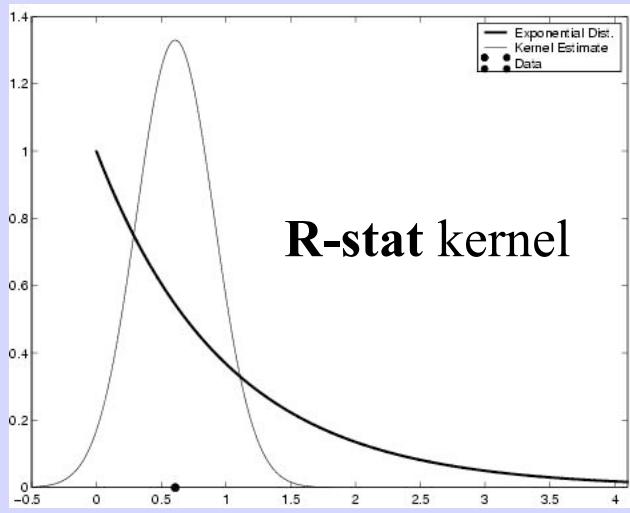


unconstrained

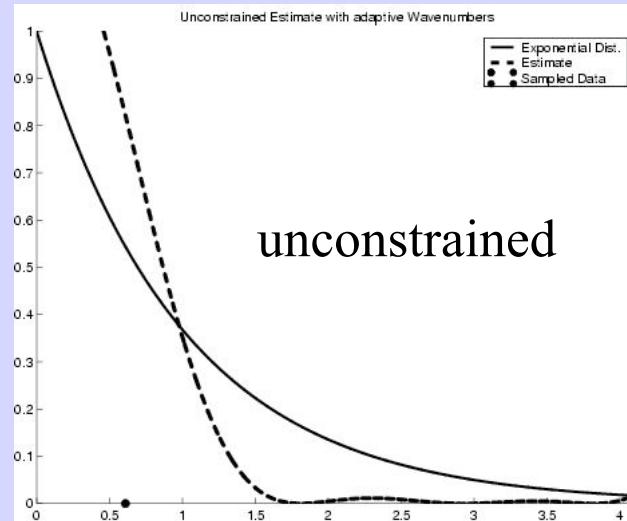


decreasing

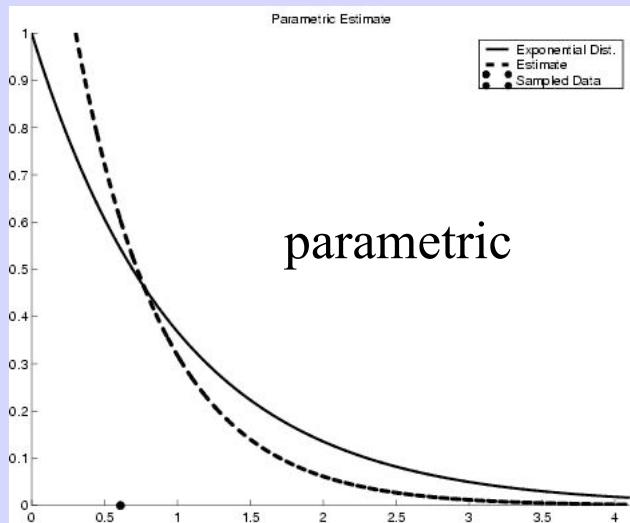
1-observation!



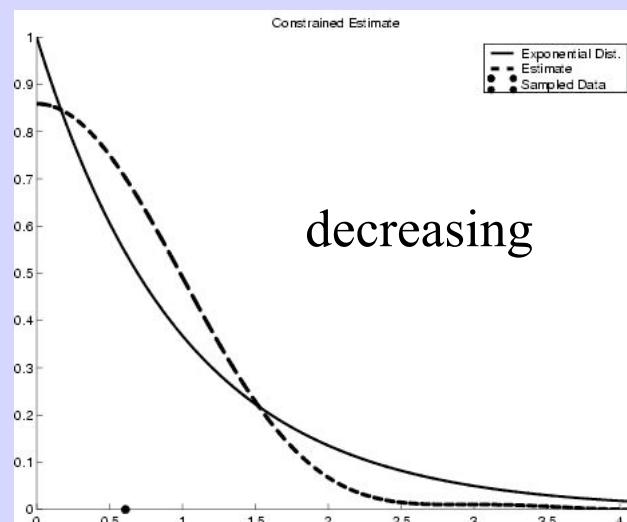
R-stat kernel



unconstrained



parametric



decreasing

Exponential (quadratic) spline



$$h(x) = e^{-s(x)}$$

$$\begin{aligned} s(x) &= s_0 + v_0 x + \int_0^x dr \int_0^r dt z(t), \quad z(t) \equiv z_k \text{ on } (x_k, x_{k+1}] \\ &= s_0 + v_0 x + \sum_{j=1}^k a_{kj} z_k \text{ when } x \in (x_k, x_{k+1}] \end{aligned}$$

$$\max E^\nu \{ \ln h(x) \} \sim \min \frac{1}{\nu} \sum_{l=1}^\nu s(x_l)$$

$$\text{so that } \int e^{-s(x)} dx \leq 1, \quad (h \geq 0)$$

$z_k \in [-\kappa_l, \kappa_u]$ 'constrained'-spline

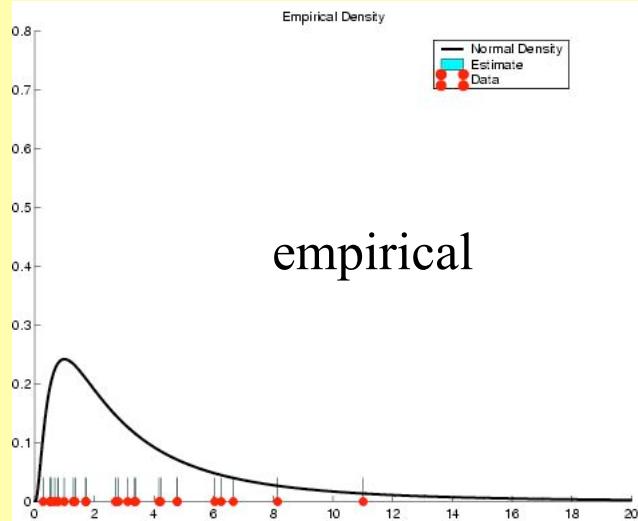
unimodal: $\kappa_l = 0$

Test Case: log-normal, $V=25$

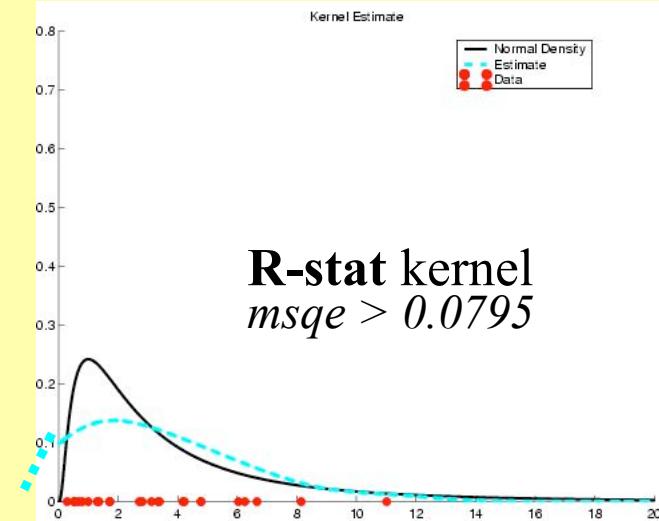
- $h^{\text{true}}(x) = \left(x\tau\sqrt{2\pi}\right)^{-1} e^{-(\ln x - \theta)^2 / 2\tau^2}$ for $x \geq 0$.
- “empirical” estimate
- kernel estimate from **R-stat**
- unconstrained with support $[0, \infty)$
- constrained (h unimodal)

25 observations!

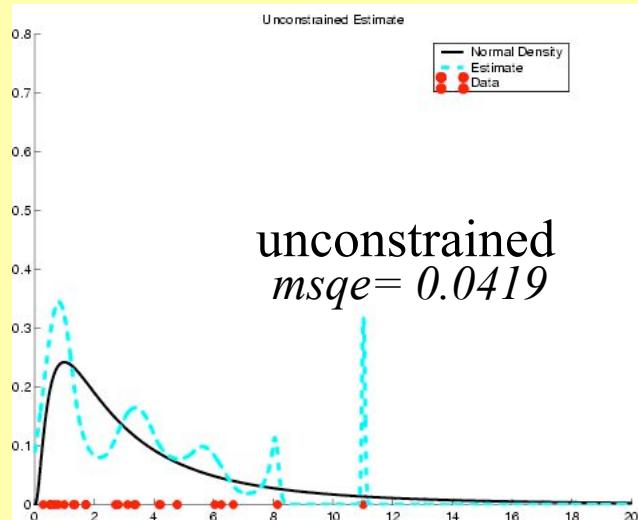
msqe = mean square error



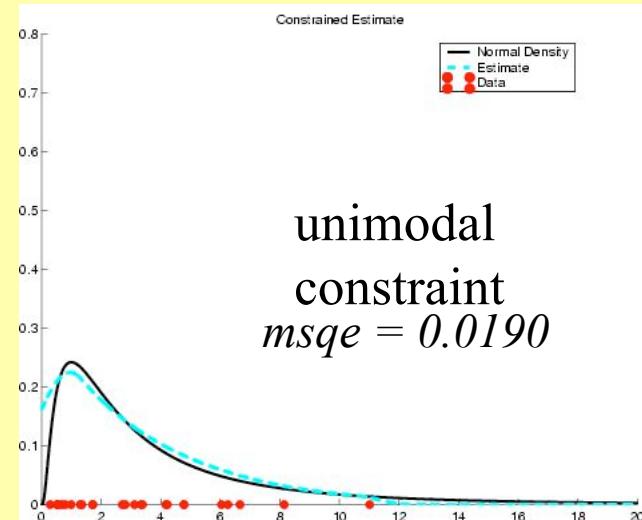
empirical



R-stat kernel
 $msqe > 0.0795$



unconstrained
 $msqe = 0.0419$

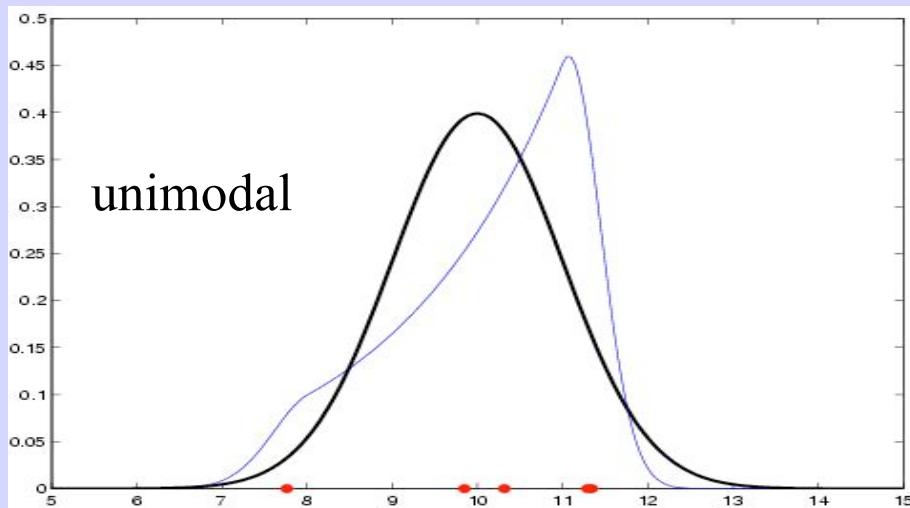
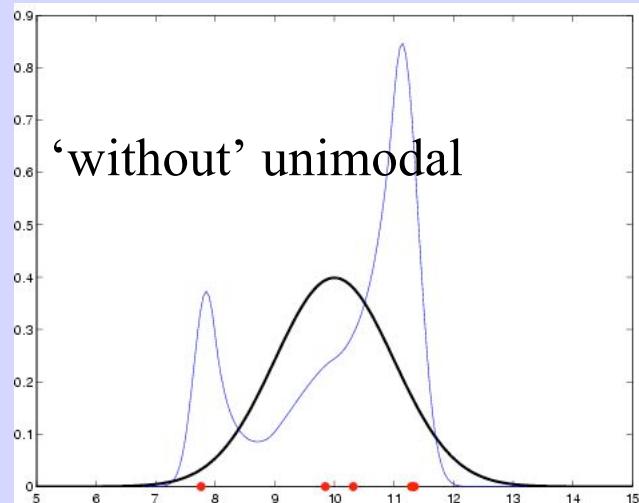
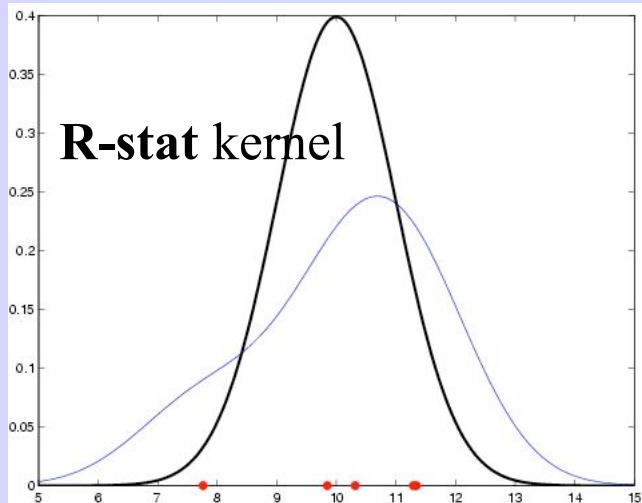


unimodal
constraint
 $msqe = 0.0190$

Test Case: $\mathcal{N}(10,1)$, V= 5

- $h^{\text{true}}(x) = (1 / \sqrt{2\pi}) e^{-(x-10)^2/2}$, $x \in \mathbb{R}$
- kernel estimate from **R-stat**
- unconstrained with support $(-\infty, \infty)$
- light-unimodality constraint on h
- stricter-unimodality constraint on h
-
- h : unimodal

5-observations: $\mathcal{N}(10,1)$

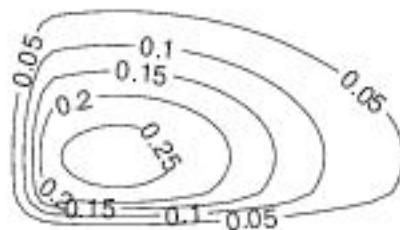


Higher Dimensions

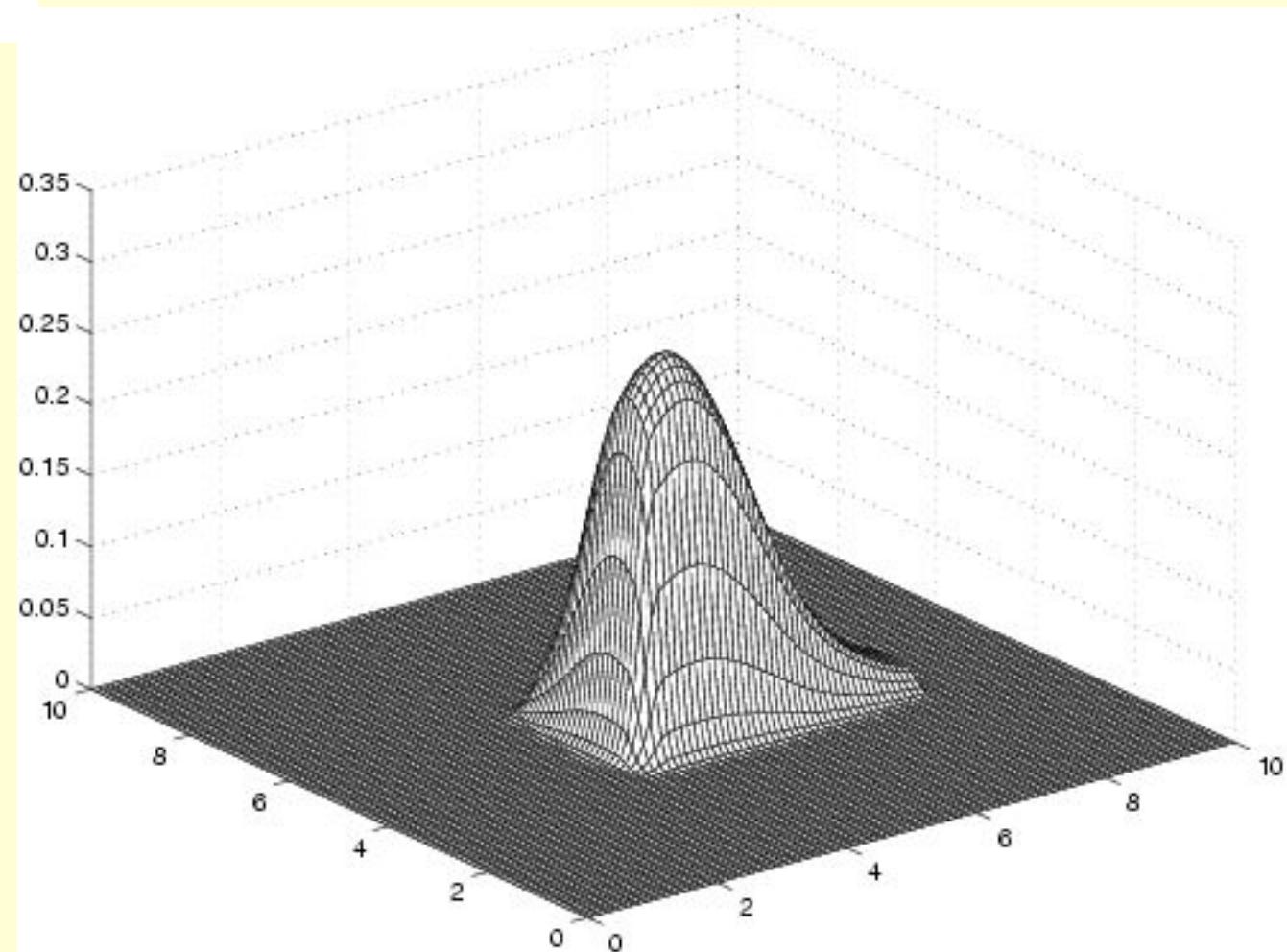
- 2-dim. (& dim. >2) $h(x, y) = \exp(z(x, y))$

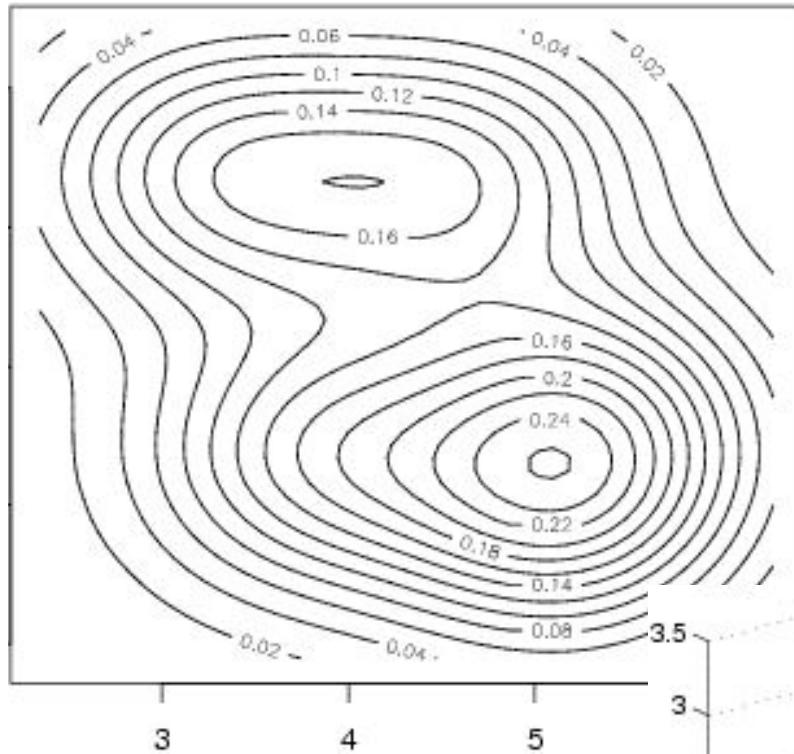
$$\begin{aligned} z(x, y) = z_0 + v_0^x x + \int_0^x d\tau \int_0^\tau ds a^x(s) \\ + \int_0^y d\tau \int_0^\tau dr a^y(r) \\ + \int_0^x ds \int_0^y dr a^{x,y}(s, r). \end{aligned}$$

- Set $a^x(t) = a_k$ on (s_{k-1}, s_k) , $a^y(t) = a_l$ on (r_{l-1}, r_l)
 $a^{x,y}(t) = a_{k,l}$ on $(s_{k-1}, s_k) \times (r_{l-1}, r_l)$



20 samples & unimodal





15 samples

