

Shape Constraints and Multiscale Methods for Density Estimation

Lutz Dümbgen, Kaspar Rufibach

(University of Berne)

Günther Walther

(Stanford University)

- I. Estimating a Log-Concave Density**
- II. Inference via Multiscale Methods**

I. Estimating a Log-Concave Density

Consider order statistics

$$X_1 < X_2 < \dots < X_n$$

of a random sample from unknown density f .

Shape constraint: f assumed log-concave,

$$f = \exp(\psi) \quad \text{with} \quad \psi : \mathbb{R} \rightarrow [-\infty, \infty) \text{ concave.}$$

Goal: Compute and analyze the NPMLE

$$\hat{f} = \exp(\hat{\psi}).$$

Why log-concavity?

- Many standard models satisfy this constraint, e.g.

$$\mathcal{N}(\mu, \sigma^2)$$

$$\text{Gamma}(a, b) \quad (a \geq 1, b > 0)$$

$$\text{Beta}(a, b) \quad (a \geq 1, b \geq 1)$$

$$\text{Weibull}(a, b) \quad (a \geq 1, b > 0)$$

Gumbel

...

- Log-concave densities are unimodal.

As opposed to NPMLE of a unimodal density,

- no trying out of many potential modes,
- no “spiking” near the estimated mode.

Estimation of $\psi = \log f$

$$\hat{F}_{\text{emp}} := n^{-1} \sum_{i=1}^n 1\{X_i \leq \cdot\} \quad (\text{empirical c.d.f.})$$

$$\hat{\psi} := \arg \max_{\psi \text{ concave}} \left(\underbrace{\int \psi d\hat{F}_{\text{emp}}}_{\text{log-likelihood}} - \underbrace{\int \exp(\psi(x)) dx}_{\text{Lagrange term}} \right)$$

Theorem 1 (Existence and uniqueness)

- $\hat{\psi}$ exists and is unique,
- $\hat{\psi}$ is piecewise linear and continuous on $[X_1, X_n]$ with knots only in $\{X_1, X_2, \dots, X_n\}$,
- $\hat{f} \equiv 0$ on $\mathbb{R} \setminus [X_1, X_n]$.

Characterisation and properties of $\hat{\psi}$, \hat{f}

By definition,

$$\frac{d}{dt} \Big|_{t=0} \left(\int (\hat{\psi} + t\Delta) d\hat{F}_{\text{emp}} - \int \exp(\hat{\psi}(x) + t\Delta(x)) dx \right) \leq 0$$

whenever $\psi + t\Delta$ is concave for some $t > 0$.

Lemma

$$\int \Delta d\hat{F}_{\text{emp}} \leq \int \Delta(x) \hat{f}(x) dx$$

whenever $\psi + t\Delta$ is concave for some $t > 0$.

In addition to \hat{F}_{emp} , $\hat{\psi}$ and \hat{f} define

$$\hat{F}(r) := \int_{-\infty}^r \hat{f}(x) dx .$$

Setting $\Delta(x) := x$ or $\Delta(x) := -x^2$ in the previous Lemma yields:

Corollary 1

$$\text{Mean}(\hat{F}) = \text{Mean}(\hat{F}_{\text{emp}}) ,$$

$$\text{Var}(\hat{F}) \leq \text{Var}(\hat{F}_{\text{emp}}) .$$

Let

$$\hat{\mathcal{S}} := \left\{ \text{knots of } \hat{\psi} \right\} \supset \{X_1, X_n\}.$$

Corollary 2

$$\int \Delta(x) \hat{F}(dx) = \int \Delta(x) \hat{F}_{\text{emp}}(dx)$$

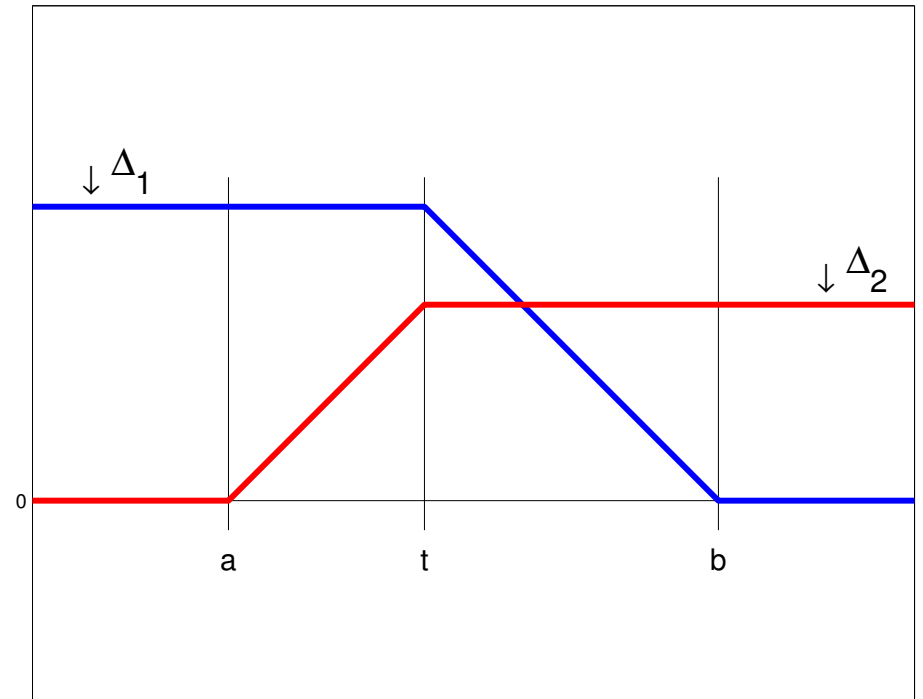
whenever $\Delta : \mathbb{R} \rightarrow \mathbb{R}$ is continuous and piecewise linear with knots only in $\hat{\mathcal{S}}$.

Corollary 3 For $a < t < b$ with $a, b \in \widehat{\mathcal{S}}$,

$$\int_a^t \widehat{F}_{\text{emp}}(x) dx \geq \int_a^t \widehat{F}(x) dx ,$$

$$\int_t^b \widehat{F}_{\text{emp}}(x) dx \leq \int_t^b \widehat{F}(x) dx ,$$

$$\int_a^b \widehat{F}_{\text{emp}}(x) dx = \int_a^b \widehat{F}(x) dx .$$



Corollary 4

$$\widehat{F}(X_1) = 0, \quad \widehat{F}(X_n) = 1$$

and

$$\widehat{F} \in \left[\widehat{F}_{\text{emp}} - n^{-1}, \widehat{F}_{\text{emp}} \right] \text{ on } \widehat{\mathcal{S}}.$$

Corollary 5

$$\left\| \widehat{F} - F \right\|_{\infty} \leq 3 \left\| \widehat{F}_{\text{emp}} - F \right\|_{\infty} + n^{-1}$$

whence

$$\left\| \widehat{F} - F \right\|_{\infty} = O_p \left(n^{-1/2} \right).$$

Conjecture

$$\sup_{\mathbb{R}} \left| \widehat{F}_{\text{emp}} - \widehat{F} \right| = o_p \left(n^{-1/2} \right).$$

Theorem 2 (Consistency of $\hat{\psi}$)

Suppose that ψ is Hölder–continuous with exponent $\beta \in [1, 2]$ on $[a, b] \subset \{f > 0\}$, i.e. for some constant L ,

$$|\psi'(x) - \psi'(y)| \leq L|x - y|^{\beta-1} \quad \text{for all } x, y \in [a, b].$$

Then

$$\sup_{[a+\delta_n, b-\delta_n]} |\hat{\psi} - \psi| = O_p \left(\left(\frac{\log n}{n} \right)^{\beta/(2\beta+1)} \right)$$

where $\delta_n \downarrow 0$.

Theorem 3 (Consistency of \hat{F})

Suppose that ψ is twice continuously differentiable on $[a, b] \subset \{f > 0\}$ with $\psi'' < 0$. Then

$$\sup_{[a+\delta_n, b-\delta_n]} \left| \hat{F} - \hat{F}_{\text{emp}} \right| = o_p \left(n^{-1/2} \right).$$

Remark 1

Rate $O_p\left((\log(n)/n)^{\beta/(2\beta+1)}\right)$ is optimal under the given conditions.

Remark 2

Integrating the log-concave density estimator \hat{f} yields an estimator \hat{F} which is essentially equivalent to \hat{F}_{emp} .

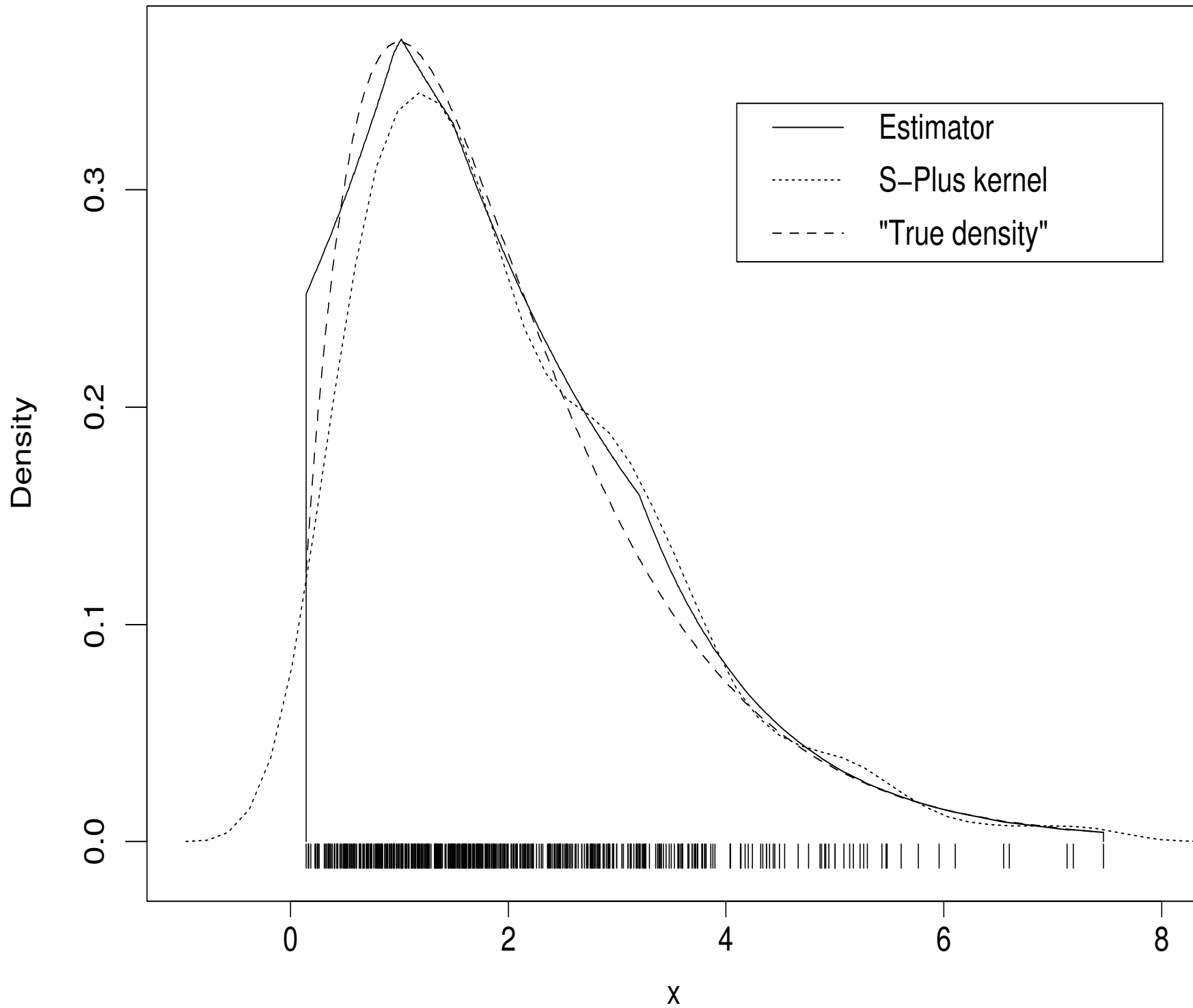
This is not true in case of a kernel density estimator \hat{f} with nonnegative kernel and optimal bandwidth of order $O(n^{-1/2})$.

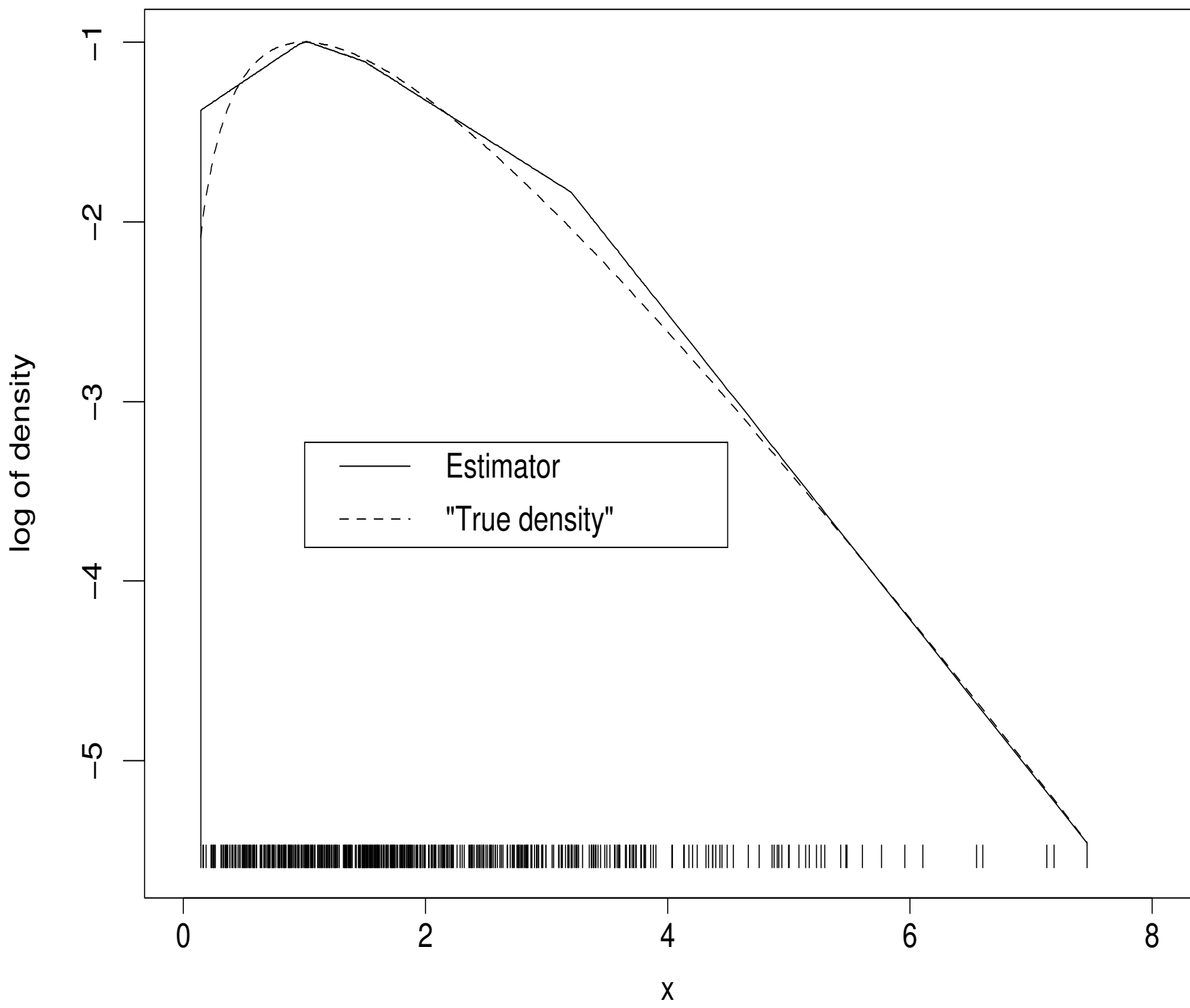
References

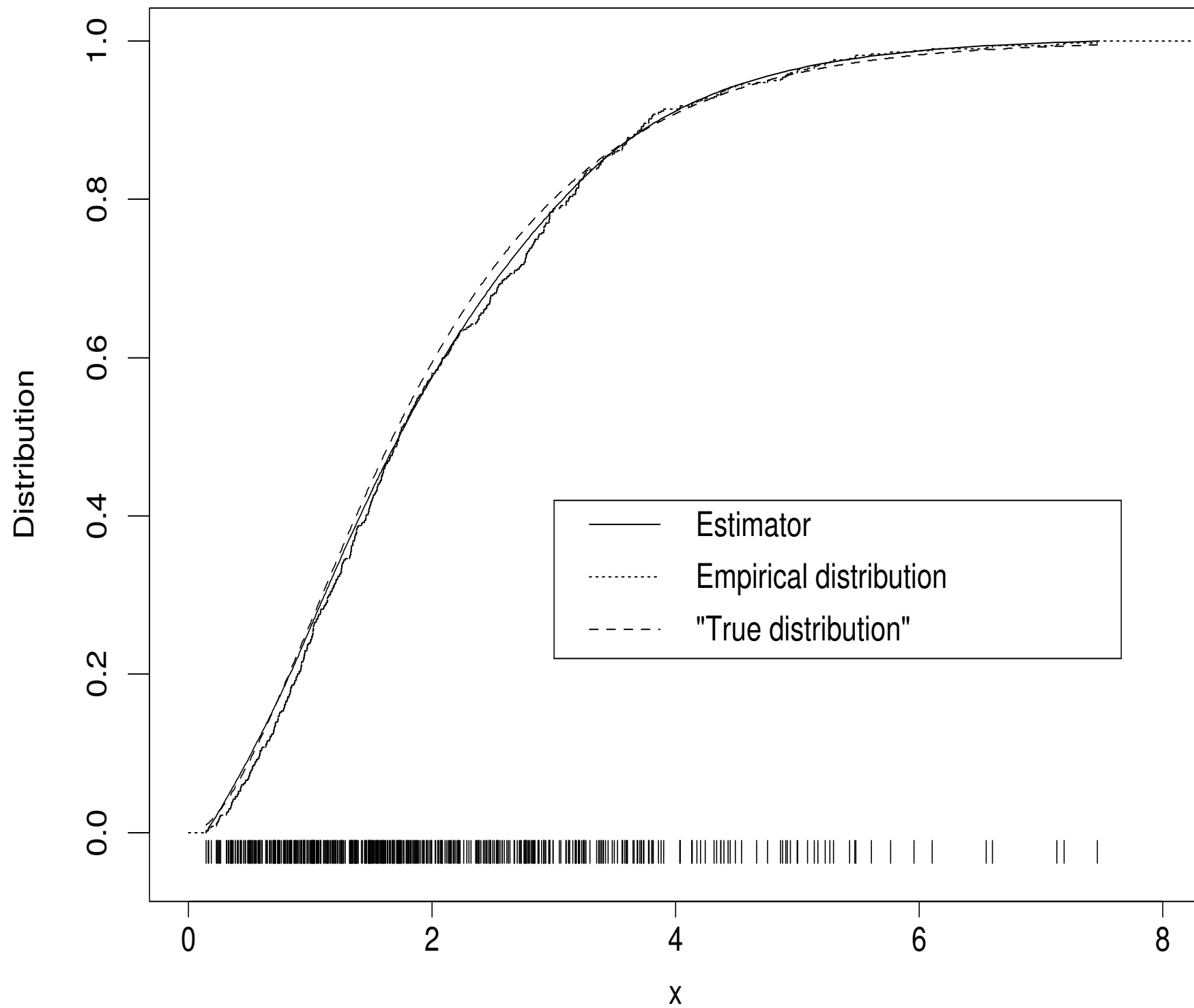
- K. Rufibach and LD (2004). *Maximum likelihood estimation of a log-concave density: basic properties and uniform consistency*. Preprint
- K. Rufibach (2004). *Computing maximum likelihood estimators of log-concave density functions*. Preprint

Numerical example

Sample of size $n = 500$ from $\text{Gamma}(2, 1)$...







II. Inference via Multiscale Methods

Goal: Identify intervals $[s, t]$ on which unknown **curve** f has a certain property, for instance,

- increases $(f' \geq 0 \text{ on } [s, t])$
- decreases $(f' \leq 0 \text{ on } [s, t])$
- has a local extremum (maximum or minimum)
- bends upward $(f'' \geq 0 \text{ on } [s, t])$
- bends downward $(f'' \leq 0 \text{ on } [s, t])$

II.1 The multiscale approach in general

For (almost) any interval $[s, t]$ consider a test statistic

$$T_{s,t} = T_{s,t}(\text{data})$$

for the null hypothesis

$$H_{s,t} : f \text{ hasn't specified property on } [s, t]$$

Multiple test: For $\alpha \in (0, 1)$ let $c_{s,t}(\alpha)$ be a critical value such that

$$\mathbb{P} \left(T_{s,t} > c_{s,t}(\alpha) \text{ and } H_{s,t} \text{ true for some } [s, t] \right) \leq \alpha.$$

Claim with confidence $1 - \alpha$ that f has specified property on **any** interval $[s, t]$ such that

$$T_{s,t} > c_{s,t}(\alpha).$$

References

- P. Chaudhuri and J.S. Marron (1999). SiZer for exploration of structures in curves. *J. Amer. Statist. Assoc.* **94**, 807–823
- LD and V.G. Spokoiny (2001). Multiscale testing of qualitative hypotheses. *Ann. Statist.* **29**, 124–152
- LD (2002). Application of local rank tests to nonparametric regression. *J. Nonpar. Statist.* **14**, 511–537
- ...
- LD and G. Walther (2005). Multiscale Inference about a Density. Preprint in preparation
- LD, K. Rufibach and G. Walther (2005). Bump hunting via multiscale testing. Preprint in preparation

II.2 Conditional densities and spacings

If the support of f is known to be (contained in)

$$[a, \infty) \quad \text{or} \quad (-\infty, b] \quad \text{or} \quad [a, b]$$

with real endpoints a, b , then add

$$X_0 := a \quad \text{or} \quad X_{n+1} := b$$

(or both) to the ordered sample. After adjusting and renumbering the observations we end up with

$$\mathbf{X} = (X_i)_{i=0}^{n+1}, \quad X_0 < X_1 < \cdots < X_{n+1}.$$

Proposition

For $0 \leq j < k \leq n + 1$ with $k - j > 1$ define (random) interval

$$\mathcal{I}_{jk} := [X_j, X_k]$$

and (random) density

$$f_{jk}(x) := \frac{1\{x \in \mathcal{I}_{jk}\} f(x)}{F(X_k) - F(X_j)}$$

Then **conditional on X_j and X_k** ,

$$(X_i)_{i=j+1}^{k-1} =_{\mathcal{L}} (Y_s)_{s=1}^{k-j-1} : \text{ordered sample from } f_{jk}$$

Thus use

$$\left(X_i \right)_{i=j+1}^{k-1} \quad \text{or} \quad \left(\frac{X_i - X_j}{X_k - X_j} \right)_{i=j+1}^{k-1}$$

for inference about

shape of f on \mathcal{I}_{jk} .

II.3 Local monotonicity properties (mode hunting)

Subsequent “distribution-free” method relies on the following fact:

$$F_o := \text{c.d.f. of } f_{0,n+1},$$

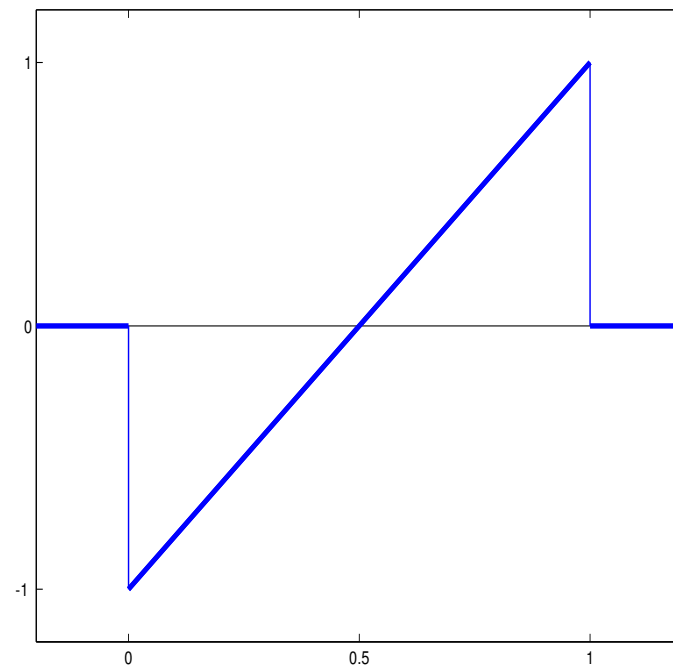
$$U = (U_i)_{i=0}^{n+1} := (F_o(X_i))_{i=0}^{n+1}.$$

Then

$$(U_i)_{i=1}^n =_{\mathcal{L}} \text{ ordered sample from } \mathcal{U}[0, 1]$$

Test statistic for an increase of f

$$\beta(x) := 1\{0 < x < 1\}(2x - 1)$$

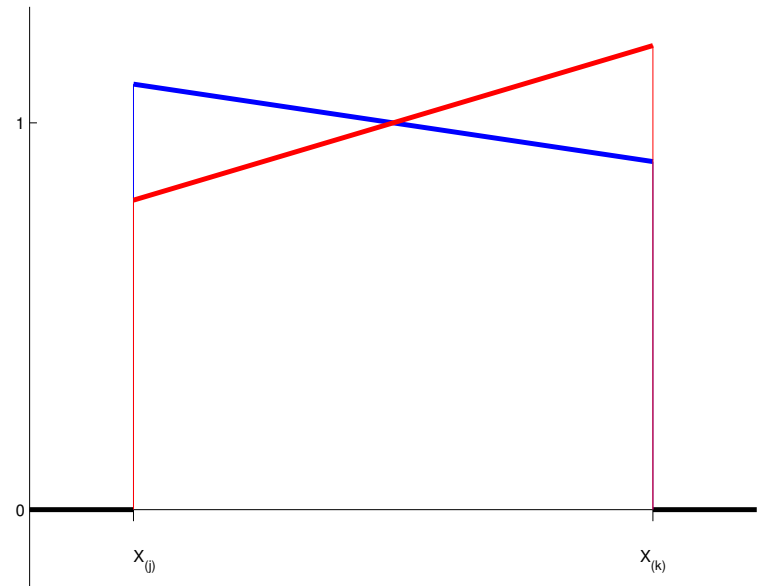


$$T_{jk}(X) := \sqrt{\frac{3}{k-j-1}} \sum_{i=j+1}^{k-1} \beta \left(\frac{X_i - X_j}{X_k - X_j} \right)$$

Possible interpretation of T_{jk}

Locally most powerful test of

“ $\lambda \leq 0$ ” versus “ $\lambda > 0$ ”



within **parametric model** where

$$f_{jk}(x) = 1\{x \in \mathcal{I}_{jk}\} \left(1 + \lambda \beta \left(\frac{x - X_j}{X_k - X_j} \right) \right)$$

Proposition

With $U = (F_o(X_i))_{i=0}^{n+1}$,

$$T_{jk}(X) \begin{cases} \geq T_{jk}(U) & \text{if } f' \geq 0 \text{ on } \mathcal{I}_{jk}, \\ \leq T_{jk}(U) & \text{if } f' \leq 0 \text{ on } \mathcal{I}_{jk}. \end{cases}$$

Application

Let $c_{jk}(\alpha)$ be critical values such that

$$\mathbb{P} (|T_{jk}(U)| > c_{jk}(\alpha) \text{ for some } (j, k)) \leq \alpha$$

Then claim with confidence $1 - \alpha$ that for arbitrary intervals \mathcal{I}_{jk} :

- $f' \not\leq 0$ on \mathcal{I}_{jk} whenever $T_{jk}(X) > c_{jk}(\alpha)$
- $f' \not\geq 0$ on \mathcal{I}_{jk} whenever $-T_{jk}(X) > c_{jk}(\alpha)$

Moreover, for arbitrary intervals $\mathcal{I}_{jk}, \mathcal{I}_{\ell m}$, the density f has a proper

- **local minimum** on \mathcal{I}_{jm} whenever

$$-T_{jk}(X) > c_{jk}(\alpha), \quad T_{\ell m}(X) > c_{\ell m}(\alpha), \quad k \leq \ell$$

- **local maximum** on \mathcal{I}_{jm} whenever

$$T_{jk}(X) > c_{jk}(\alpha), \quad -T_{\ell m}(X) > c_{\ell m}(\alpha), \quad k \leq \ell$$

- **local extremum** on $\mathcal{I}_{\min(j,\ell), \max(k,m)}$ whenever

$$\pm T_{jk}(X) > c_{jk}(\alpha), \quad \mp T_{\ell m}(X) > c_{\ell m}(\alpha)$$

Finding the critical values $c_{jk}(\alpha)$

$$T(\mathbf{U}) := \max_{k-j>1} \left(|T_{jk}(\mathbf{U})| - G \left(\frac{k-j}{n+1} \right) \right)$$

$$G(u) := \sqrt{2 \log \left(\frac{e}{u} \right)}$$

$\kappa(\alpha) := (1 - \alpha)$ – quantile of $T(\mathbf{U})$

$$c_{jk}(\alpha) := \kappa(\alpha) + G \left(\frac{k-j}{n+1} \right)$$

Theorem 1

$$T(U) \rightarrow_{\mathcal{L}} \mathbf{T} \in [0, \infty) \quad \text{as } n \rightarrow \infty$$

Theorem 2

Suppose that $\pm f' \geq c > 0$ in some neighborhood of $x \in \mathbb{R}$. Then x is localized with asymptotic probability one and precision

$$O_p \left(\left(\frac{\log n}{n} \right)^{1/3} \right).$$

Suppose that $f(x) > 0$, $f'(x) = 0$ and $\pm f'' \geq c > 0$ in some neighborhood of $x \in \mathbb{R}$. Then this local extremum is localized with asymptotic probability one and precision

$$O_p \left(\left(\frac{\log n}{n} \right)^{1/5} \right).$$

Numerical example

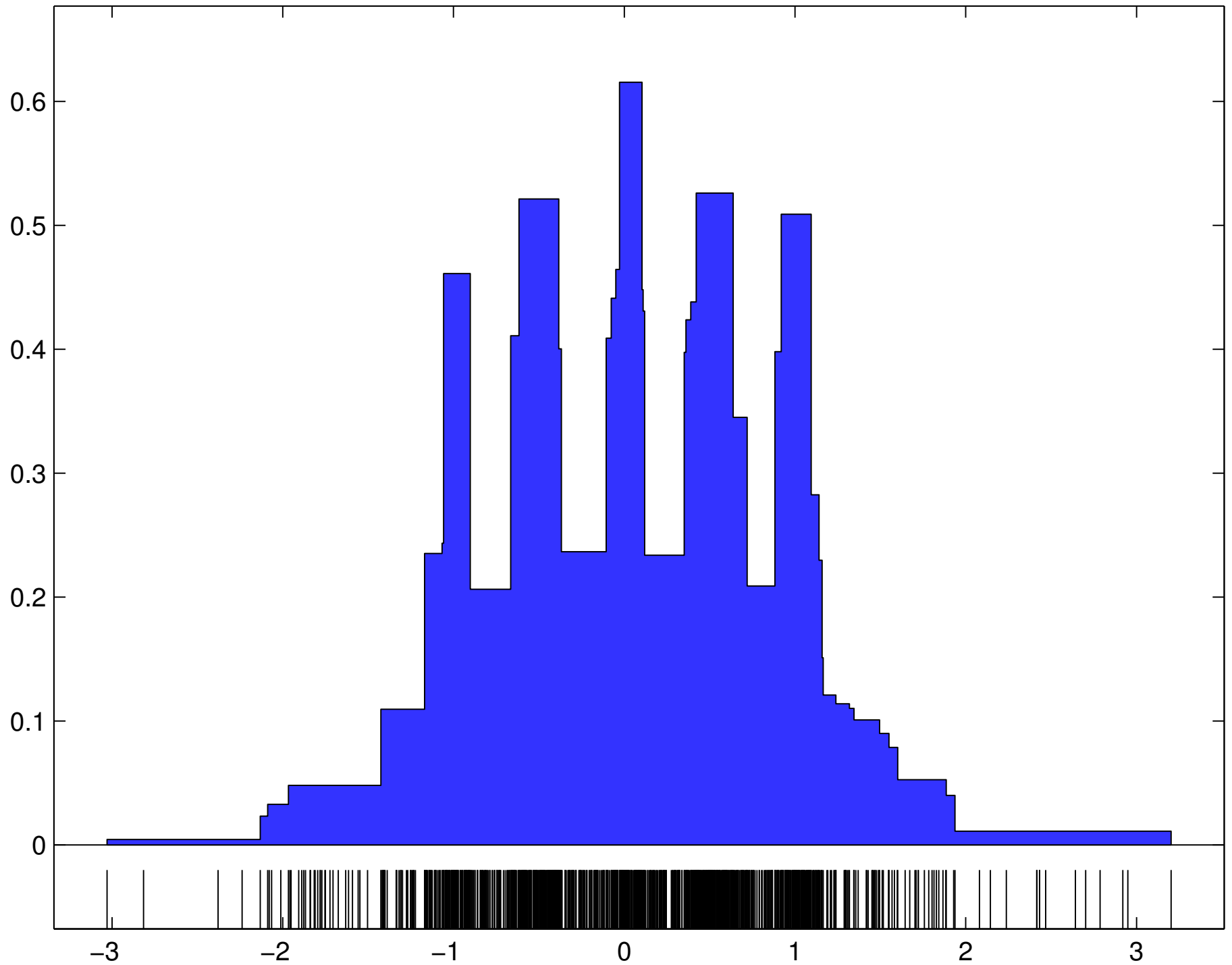
For a simulated sample of size $n = 1000$ from the “claw density” show

- a histogram estimator of f (Davies and Kovac)
- the function

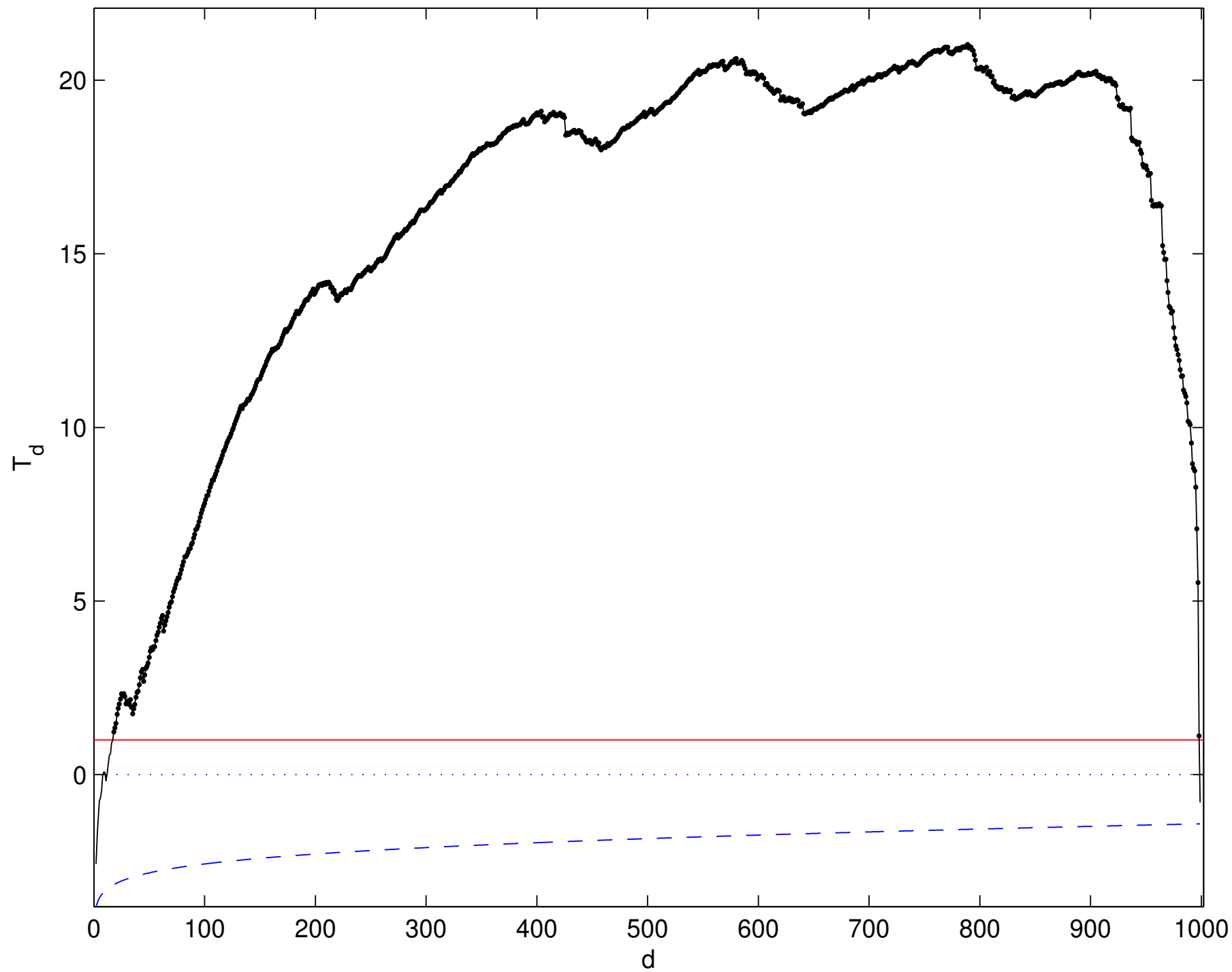
$$d \mapsto \max_{k-j=d} \left(\left| T_{jk}(X) \right| - G \left(\frac{d}{n+1} \right) \right)$$

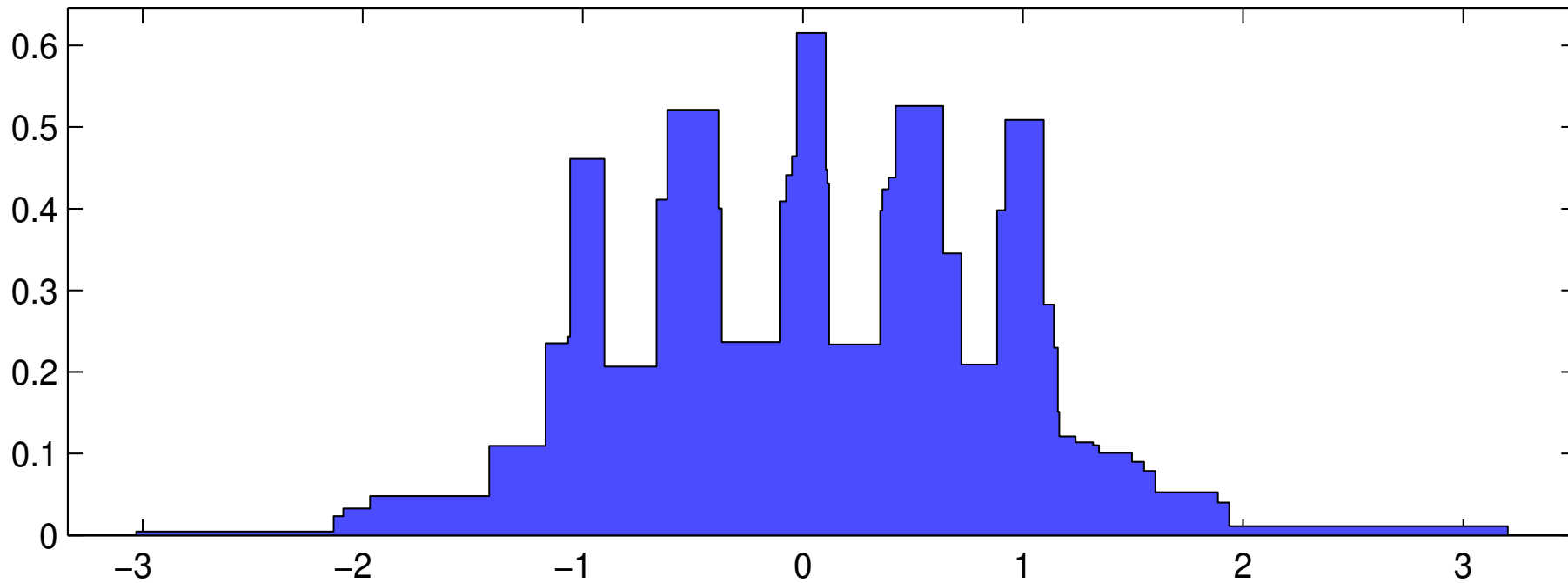
- minimal intervals with essential increases/decreases or local extrema

Data and Taut-String Histogram

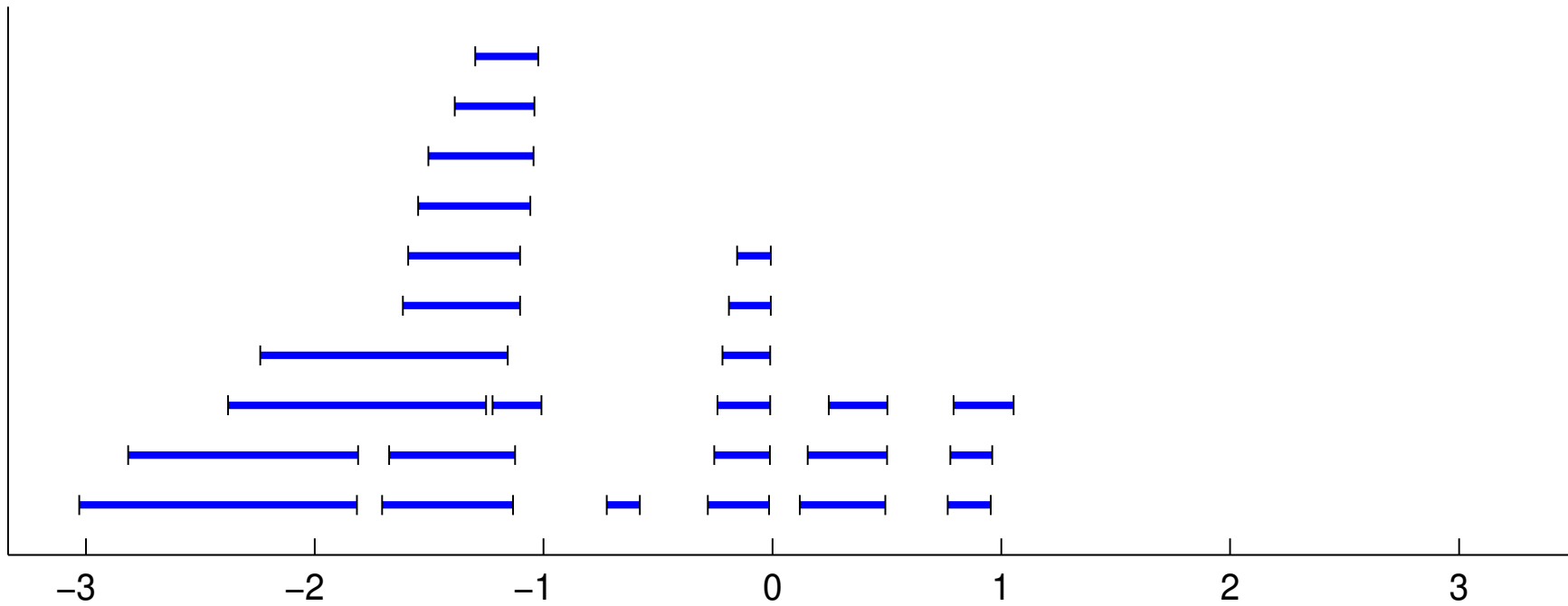


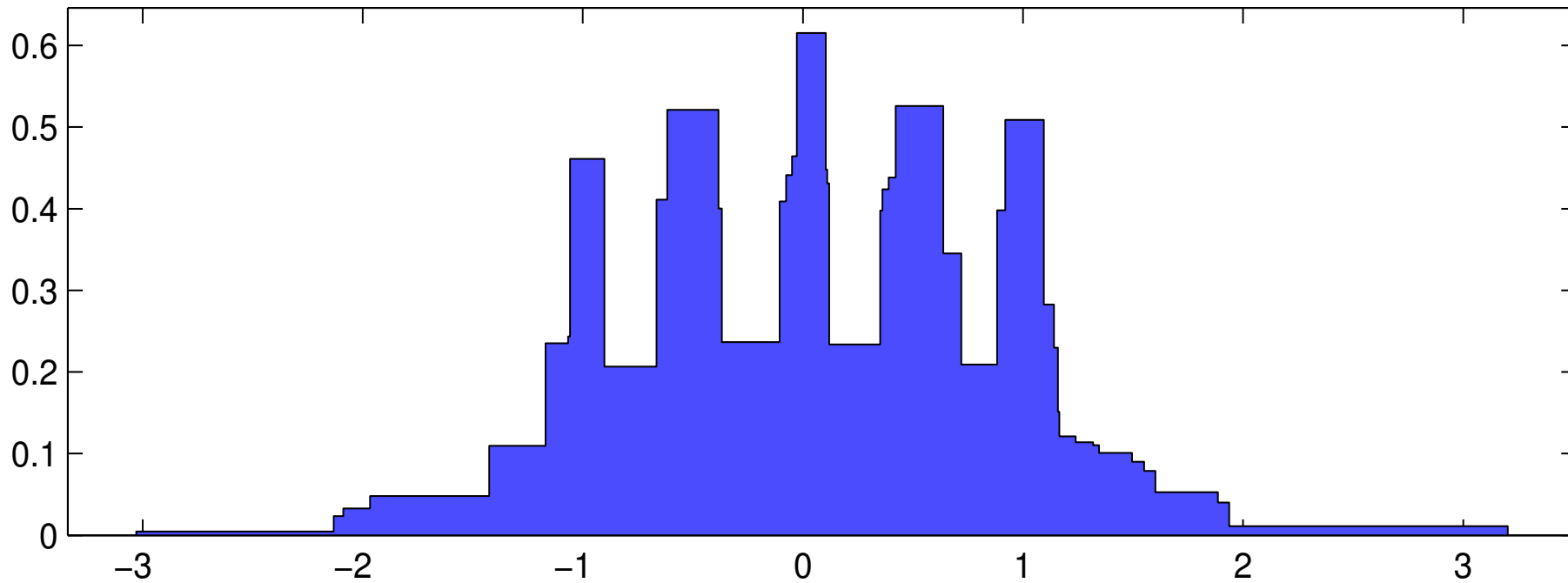
Single scale test statistics and global critical value



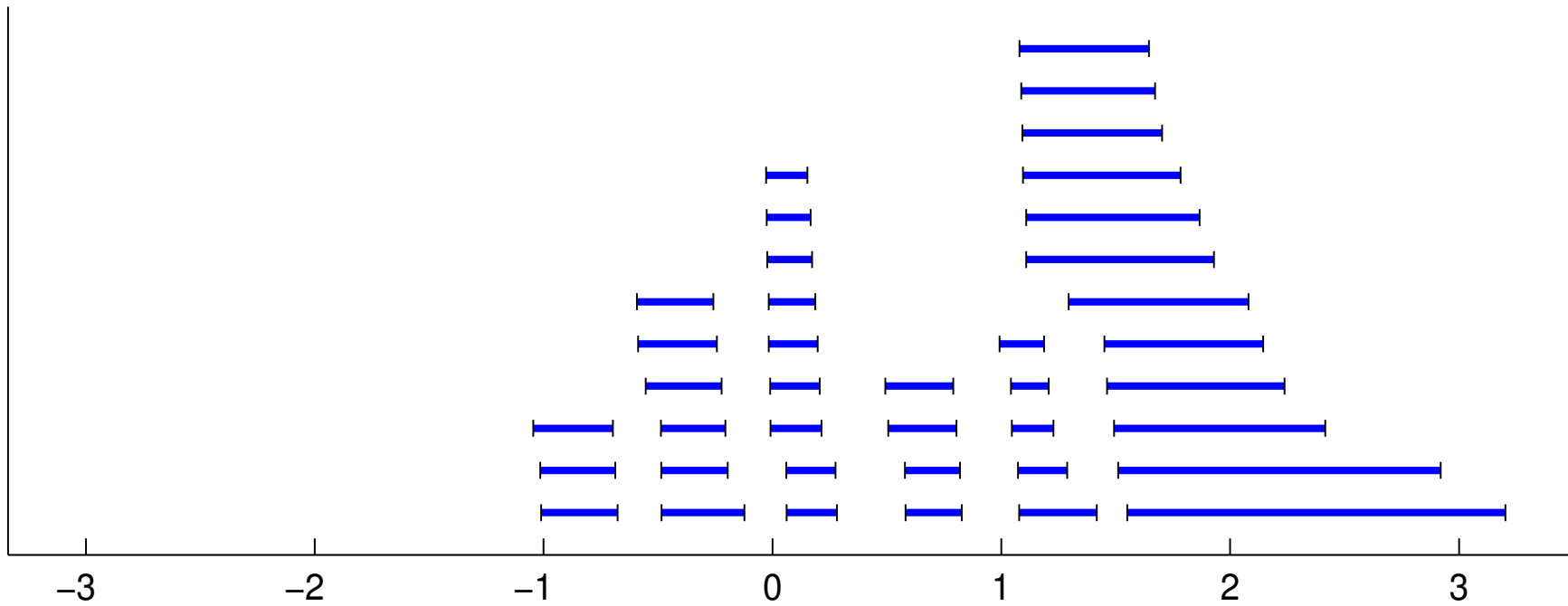


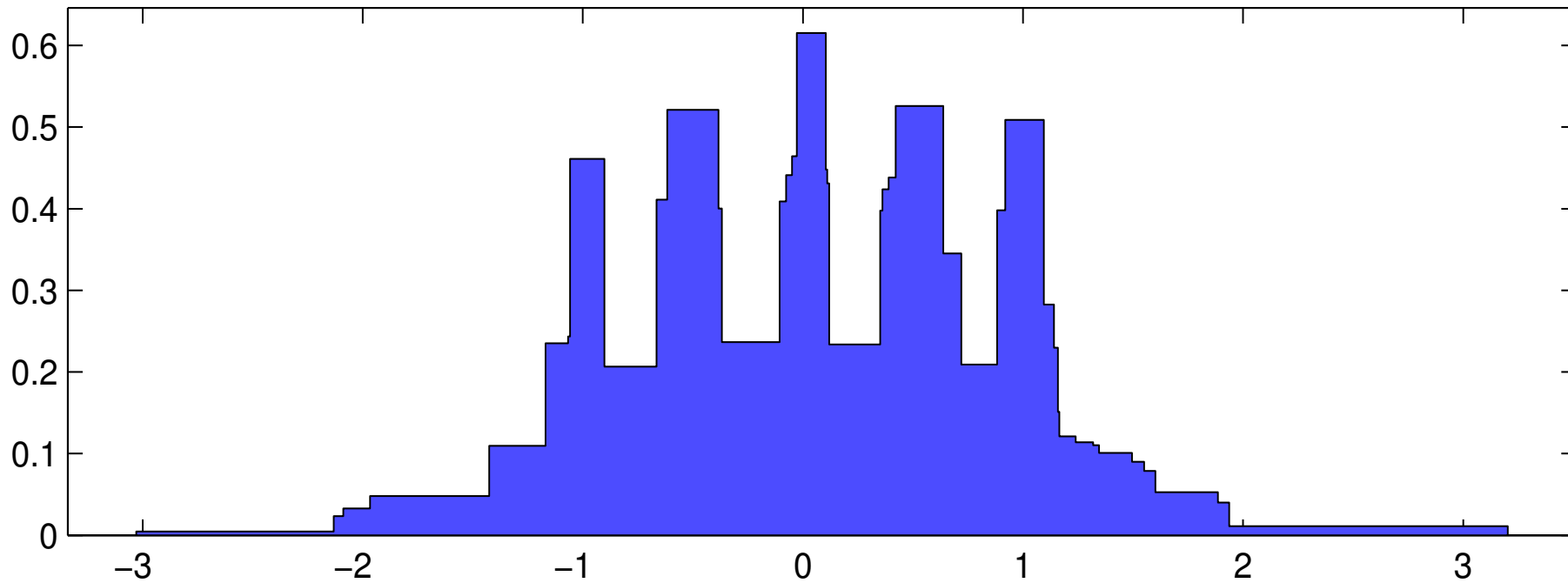
Minimal intervals with significant increase of density



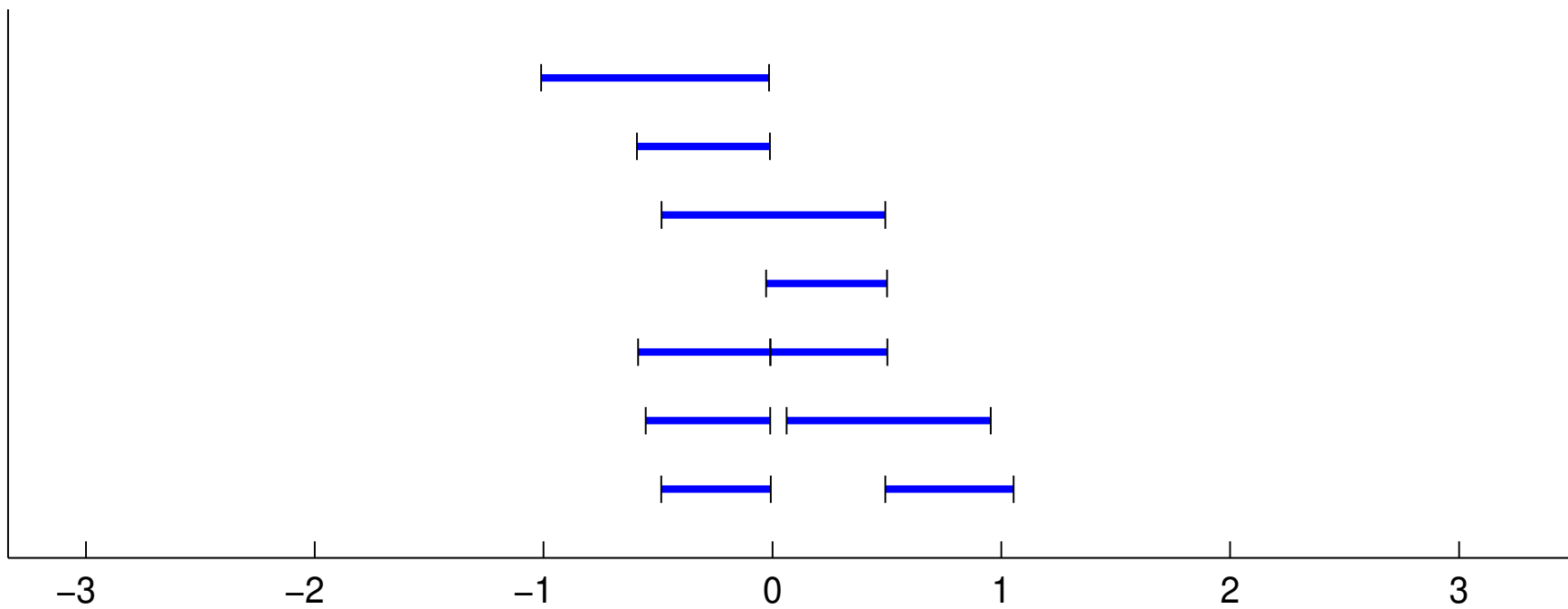


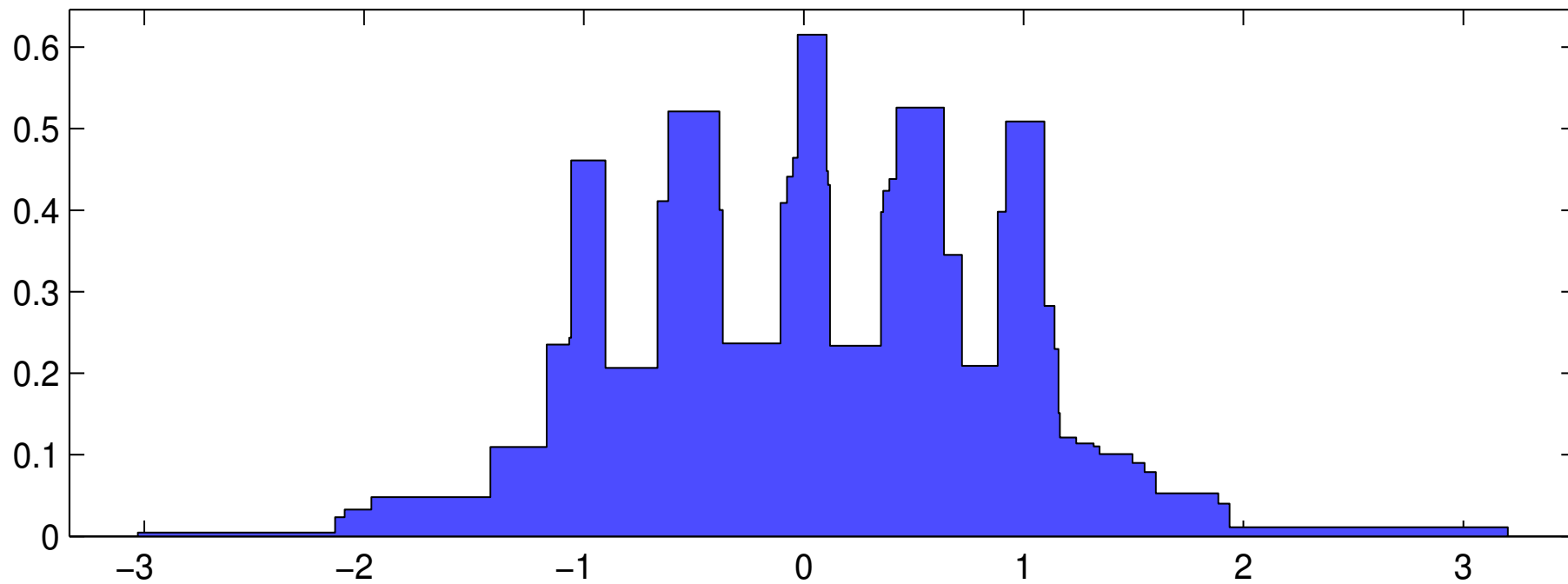
Minimal intervals with significant decrease of density



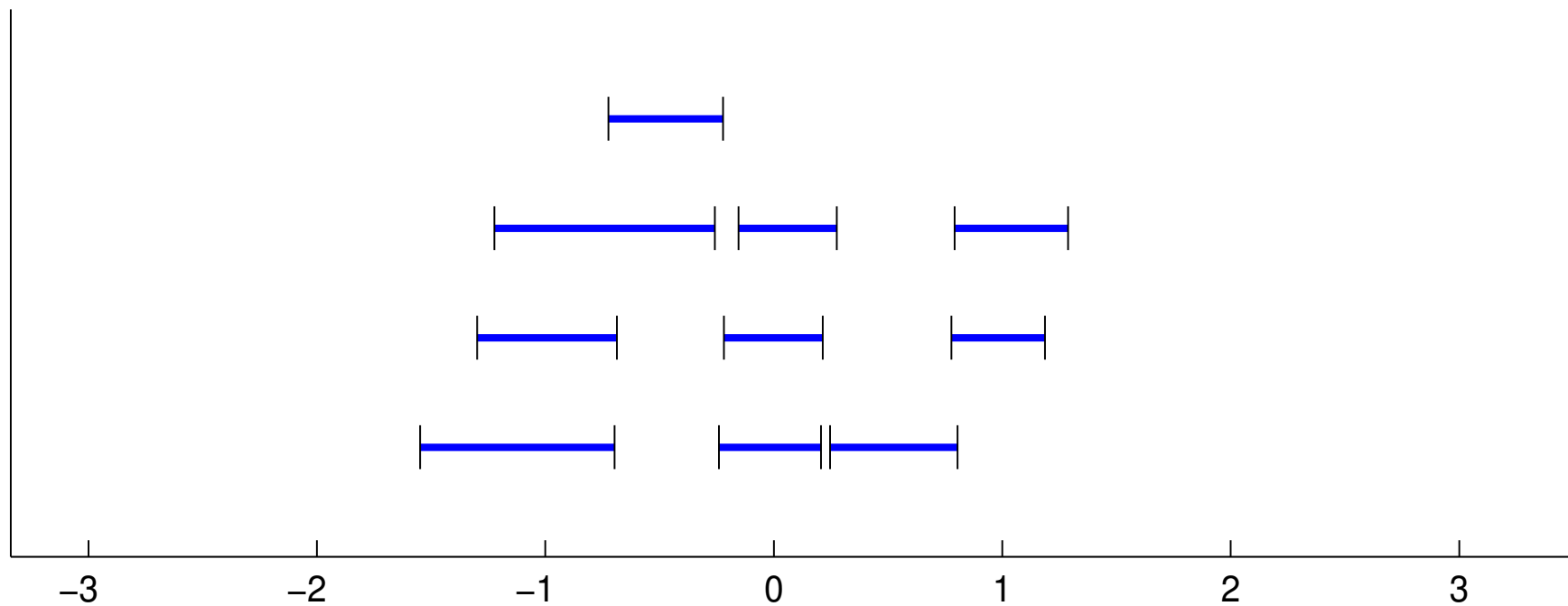


Minimal intervals with significant proper local minimum





Minimal intervals with significant proper local maximum



II.3 Auxiliary results

Theorem 1 follows from general results about stochastic processes.

Here a simple version:

Z_n : stochastic process on finite set $\Pi_n \subset \Pi$

$$\Pi := \{(s, t) : 0 \leq s < t \leq 1\}$$

Define

$$T_n := \max_{(s,t) \in \Pi_n} \left(\frac{|Z_n(s, t)|}{\sqrt{t - s}} - G(t - s) \right)$$

Theorem A

Assumption 1 (subgaussian variables):

$$\mathbb{P} \left\{ \frac{|Z_n(s, t)|}{\sqrt{t-s}} \geq \eta \right\} \leq 2 \exp \left(-\frac{\eta^2}{2} \right)$$

Assumption 2 (subexponential increments):

$$\mathbb{P} \left\{ \frac{|Z_n(s, t) - Z_n(u, v)|}{\sqrt{|s-u| + |t-v|}} \geq \eta \right\} \leq K \exp \left(-\frac{\eta}{K} \right)$$

Then

$$\sup_n \mathbb{P} \{T_n \geq \eta\} \rightarrow 0 \text{ as } \eta \rightarrow \infty.$$

Theorem B

Suppose that the assumptions of Theorem A hold. In addition suppose that

- “ $\Pi_n \rightarrow \Pi$ ” .
- the finite-dimensional distributions of Z_n (suitably extended) converge to those of a centered gaussian process Z_∞ on Π .
- $\text{Cov}(Z_\infty(s, t), Z_\infty(u, v)) = 0$ whenever $t \leq u$,
 $\text{Var}(Z_\infty(s, t)) = t - s$.

Then

$$T_n \rightarrow_{\mathcal{L}} T_\infty \quad \text{and} \quad 0 \leq T_\infty < \infty .$$

Theorems A–B apply to

$$\Pi_n := \left\{ \left(\frac{j}{n+1}, \frac{k}{n+1} \right) : 0 \leq j < k \leq n+1 \right\}$$

$$Z_n \left(\frac{j}{n+1}, \frac{k}{n+1} \right) := \sqrt{\frac{k-j-1}{n+1}} T_{jk}(U)$$

II.4 Detecting convexity/concavity of $\log f$

Local parametric models

Local coordinates:

$$\langle x \rangle_{jk} := \beta \left(\frac{x - X_j}{X_k - X_j} \right), \quad [X_j, X_k] \rightarrow [-1, 1].$$

Parametric models:

$$g_{\theta}(y) := \mathbf{1}\{-1 < y < 1\} G_{\theta} \exp(\theta y),$$

$$g_{\theta, \eta}(y) := \mathbf{1}\{-1 < y < 1\} G_{jk, \theta, \eta} \exp(\theta y + \eta y^2).$$

Local score test statistics

Norming constants $A_\theta, B_\theta, C_\theta$ s.t.

$$\int (A_\theta y^2 - B_\theta y - C_\theta) y^j g_\theta(y) dy = 0 \quad \text{for } j = 0, 1,$$
$$\int (A_\theta y^2 - B_\theta y - C_\theta)^2 g_\theta(y) dy = 1.$$

For $k - j > 2$:

$$T_{jk} := \sqrt{\frac{1}{k - j - 1}} \sum_{i=j+1}^{k-1} \left(A_{\hat{\theta}} \langle X_i \rangle_{jk}^2 - B_{\hat{\theta}} \langle X_i \rangle_{jk} - C_{\hat{\theta}} \right)$$

with

$$\hat{\theta} = \hat{\theta}_{jk} := \text{local MLE } \dots$$

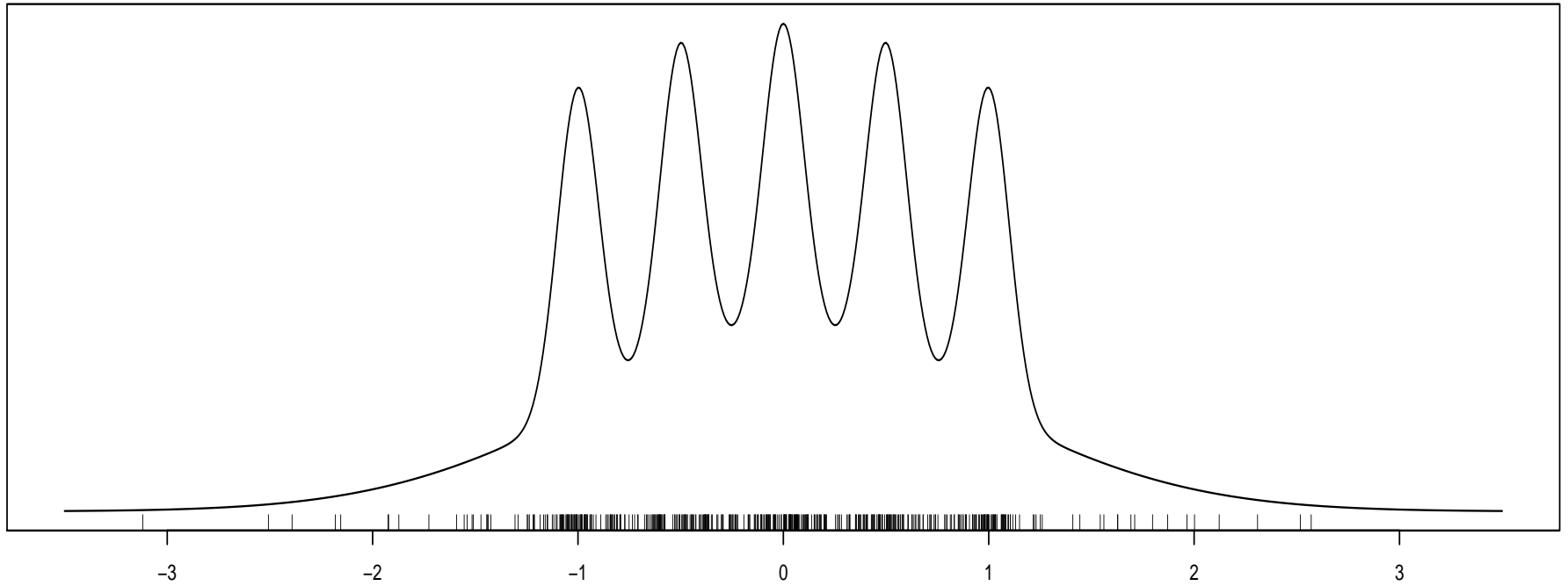
So far some theory about the local power of the single tests . . .

Conjecture that simulating critical values from uniform distribution yields **asymptotically valid** multiple test . . .

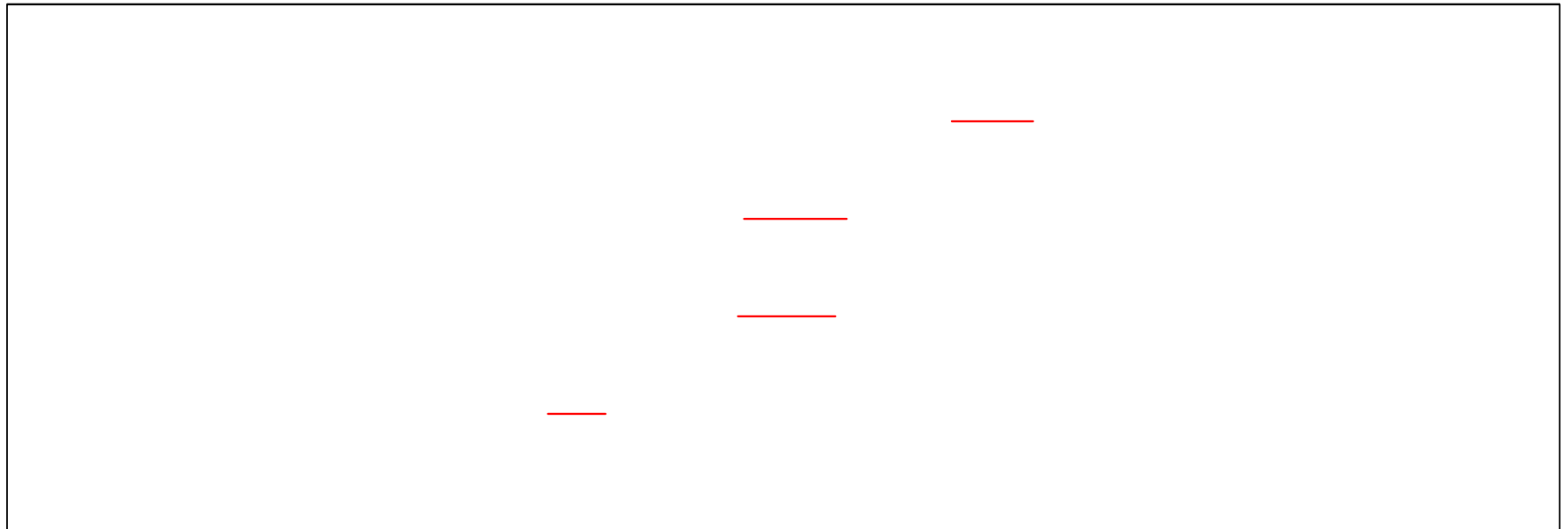
Simulations and numerical examples indicate that the method works . . .

Main difficulty: Apparently no useful transformation group connected to the parametric model $(g_\theta)_{\theta \in \mathbb{R}}$. . .

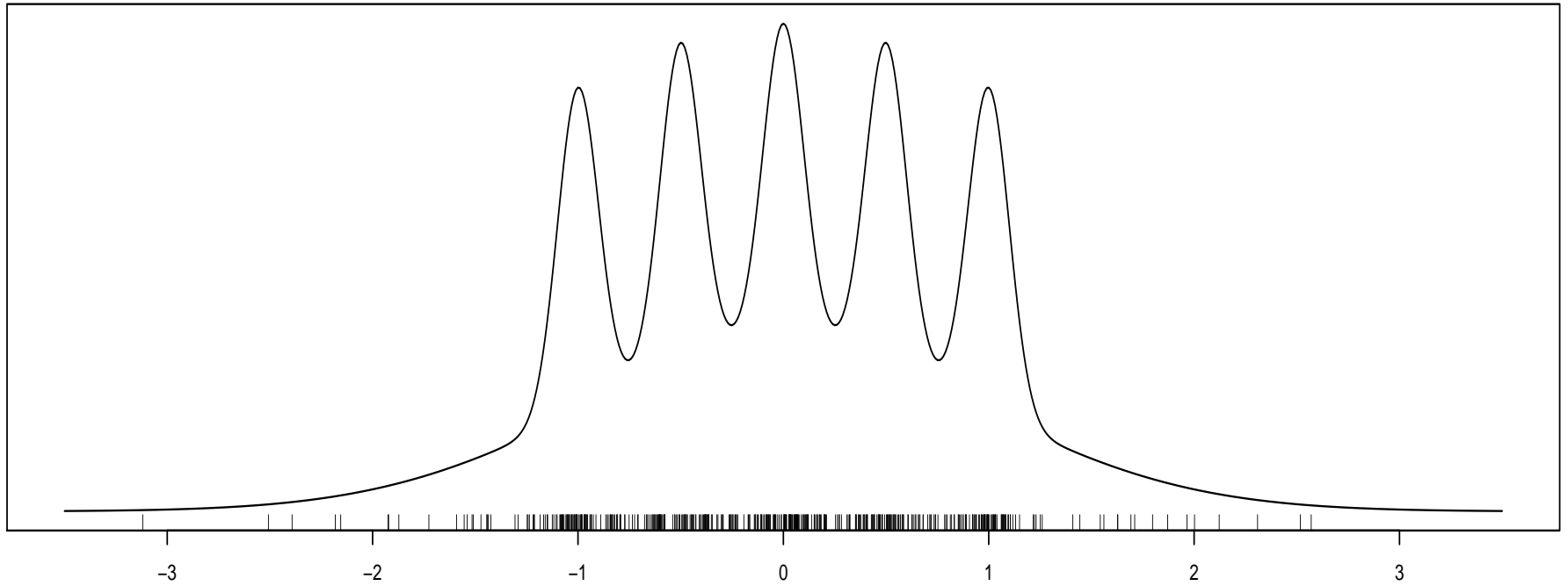
Bump hunting for claw density
n=500 / minimal lag=3 / maximal lag=250



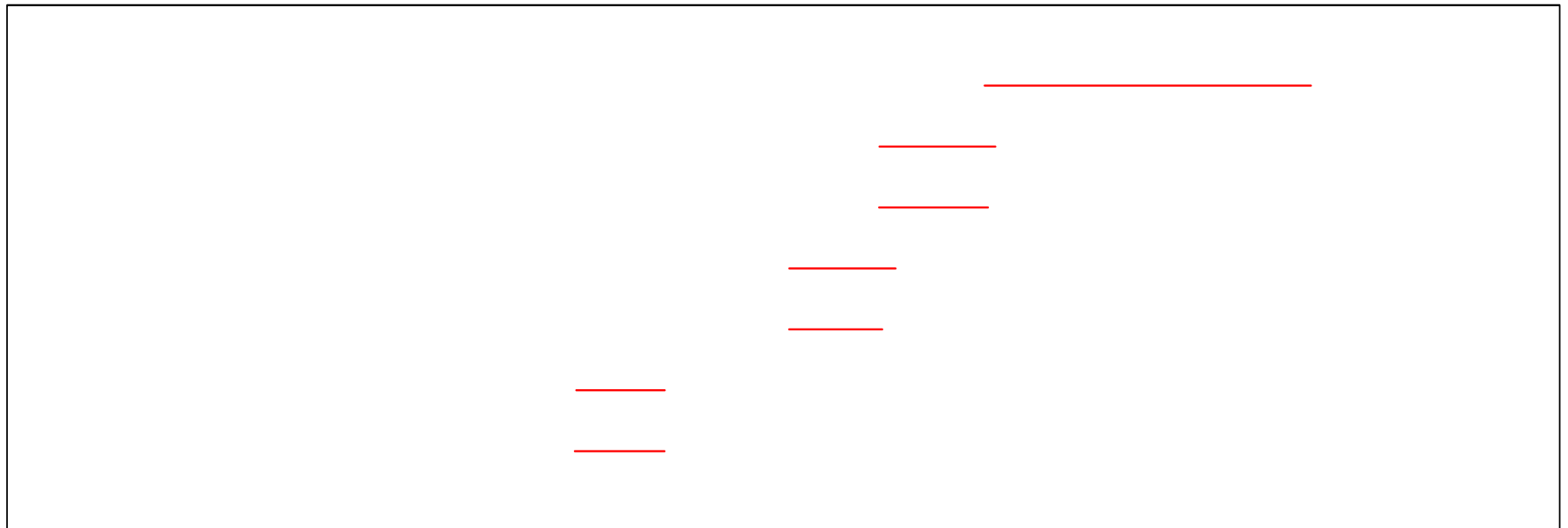
intervals on which $\log(f)$ is non-convex



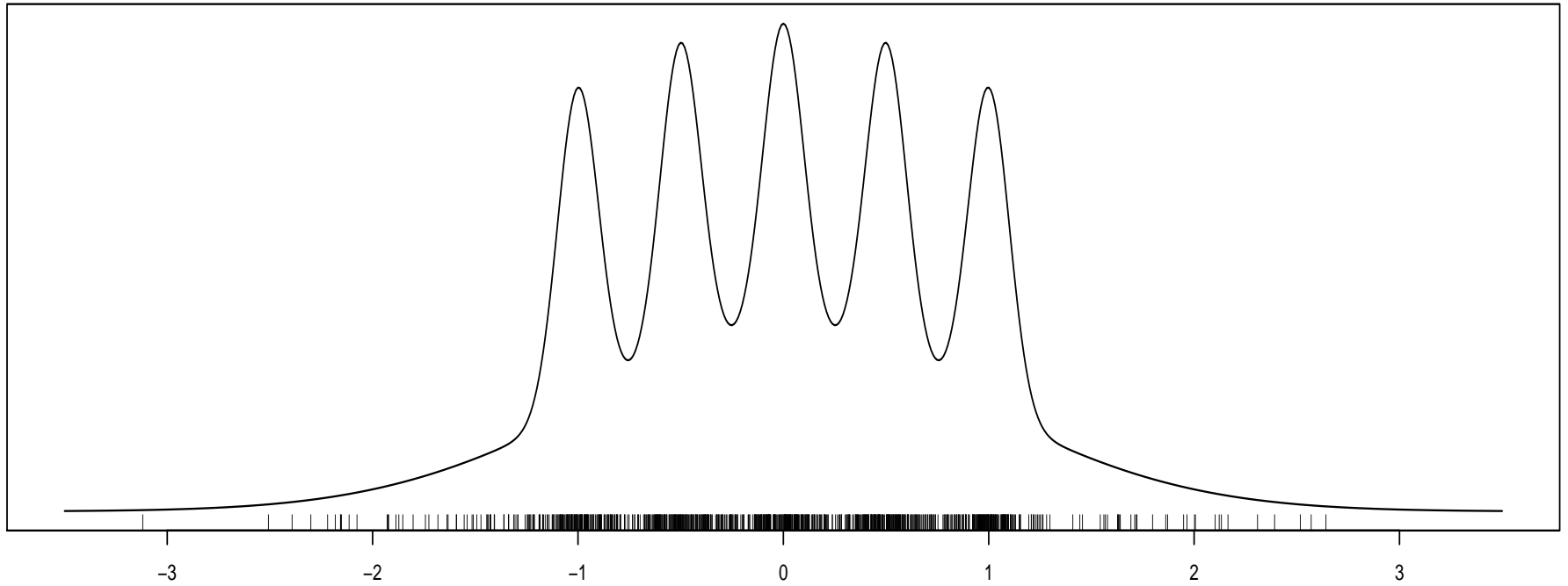
Bump hunting for claw density
n=500 / minimal lag=3 / maximal lag=250



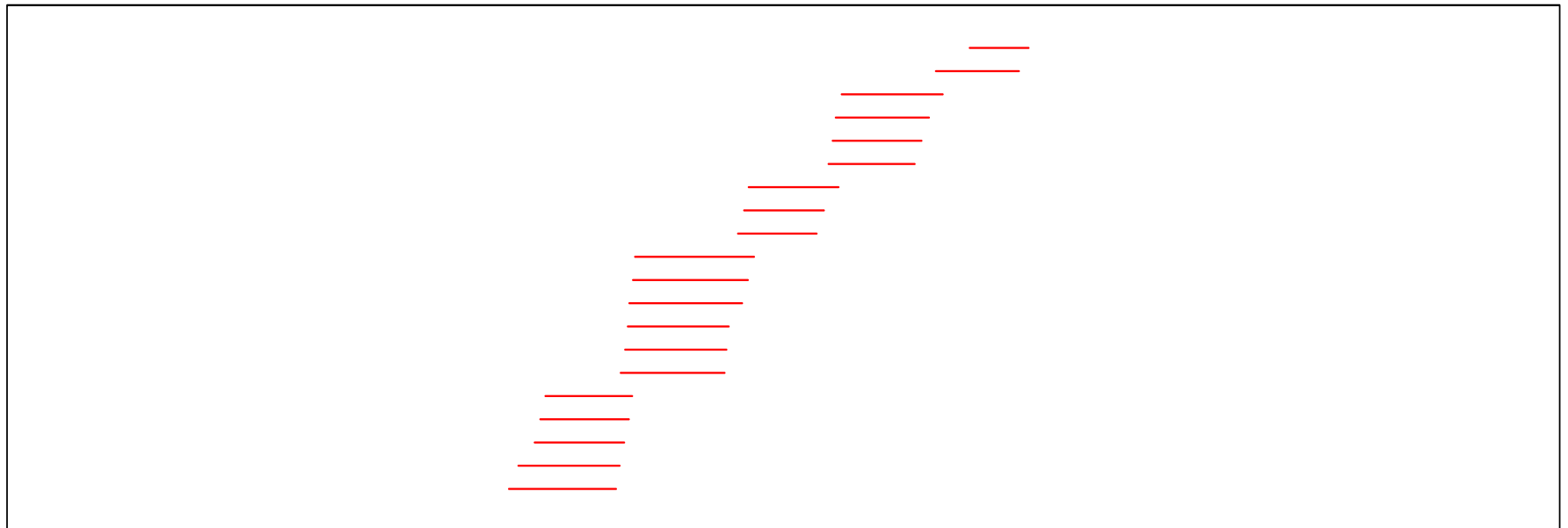
intervals on which $\log(f)$ is non-concave



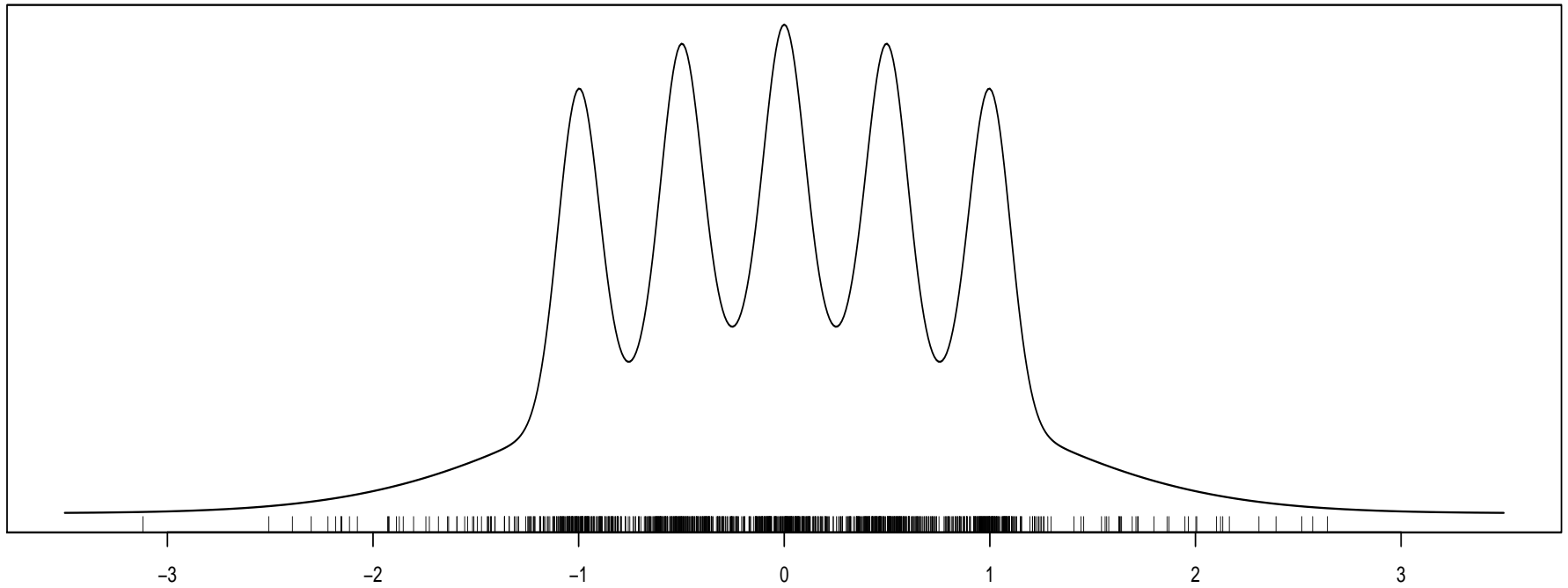
Bump hunting for claw density
n=1000 / minimal lag=3 / maximal lag=500



intervals on which $\log(f)$ is non-convex



Bump hunting for claw density
n=1000 / minimal lag=3 / maximal lag=500



intervals on which log(f) is non-concave

