

# Inference on a Distribution Function from Ranked Set Samples

Lutz Dümbgen (Univ. of Bern)  
Ehsan Zamanzade (Univ. of Isfahan)

October 17, 2013  
Swiss Statistics Meeting, Basel

# I. The Setting

In some situations, **ranking observations** is much easier than obtaining their **precise values** ...

Consider  $n \cdot k$  independent random variables

$$X_{ij}, \quad 1 \leq i \leq n, \quad 1 \leq j \leq k,$$

with unknown continuous c.d.f.  $F$ .

Instead of complete  $i$ -th sample  $(X_{ij})_{j=1}^k$  observe only

- one of its elements,  $X_i$ , and
- its rank  $R_i \in \{1, 2, \dots, k\}$  within  $(X_{ij})_{j=1}^k$ .

## Example 1: Ranked set sampling (RSS)

(McIntyre 1952)

For  $1 \leq i \leq n$  let

$$X_{i:1} < X_{i:2} < \cdots < X_{i:k}$$

be the order statistics of  $(X_{ij})_{j=1}^k$ .

- Pick a rank  $R_i \in \{1, 2, \dots, k\}$
- Obtain precise value  $X_i := X_{i:R_i}$ .

Special case:  $k = n$  and  $R_i = i$ .

**Example 2: Judgement post-stratification (JPS)**  
(MacEachern et al. 2004)

- $X_i := X_{i1}$  ,
- $R_i := \sum_{j=1}^k 1_{[X_{ij} \leq X_{i1}]}$ .

## General setting:

Independent observations  $(X_1, R_1), (X_2, R_2), \dots, (X_n, R_n)$ .

Conditional on  $R_i = r$ ,

$$X_i \sim F_r(x) := B_r(F(x))$$

with

$$B_r(p) := \sum_{i=r}^k \binom{k}{i} p^i (1-p)^{k-i} = \int_0^p \beta_r(u) du$$

$$\beta_r(u) := k \binom{k-1}{r-1} u^{r-1} (1-u)^{k-r}$$

(c.d.f. and density of  $\text{Beta}(r, k+1-r)$ ).

## Basic fact

$$\frac{1}{k} \sum_{r=1}^k \text{Beta}(r, k+1-r) = \text{Unif}[0, 1],$$

whence

$$\frac{1}{k} \sum_{r=1}^k F_r = F.$$

## II. Estimation paradigms

### **Paradigm 0: Naive empirical distribution function**

$$\hat{F}(x) := \frac{1}{n} \sum_{i=1}^n 1_{[X_i \leq x]}.$$

In JPS setting: Okay but inefficient.

In unbalanced RSS settings: Potentially severely biased.

## Paradigm 1: Stratification (Stokes and Sager 1988)

Groupwise empirical c.d.f.s:

$$\hat{F}_r(x) := \frac{1}{N_r} \sum_{i: R_i=r} 1_{[X_i \leq x]} \quad \text{with} \quad N_r := \#\{i : R_i = r\}.$$

Stratified estimator

$$\hat{F}^S := \frac{1}{k} \sum_{r=1}^k \hat{F}_r.$$

Potential problems with small/zero numbers  $N_r$ .



## Paradigm 2: Likelihood (Kvam and Samaniego 1994)

$$\hat{F}^L(x) := \arg \max_{p \in [0,1]} L(x, p)$$

with (conditional) log-likelihood function

$$\begin{aligned} L(x, p) &:= \sum_{i=1}^n [1_{[X_i \leq x]} \log B_{R_i}(p) + 1_{[X_i > x]} \log(1 - B_{R_i}(p))] \\ &= \sum_{r=1}^k N_r [\hat{F}_r(x) \log B_r(p) + (1 - \hat{F}_r(x)) \log(1 - B_r(p))] \end{aligned}$$

### Paradigm 3: Moments

$$\mathbb{E}(\hat{F}(x) \mid \mathbf{R}) = \sum_{r=1}^k \frac{N_r}{n} B_r(F(x)).$$

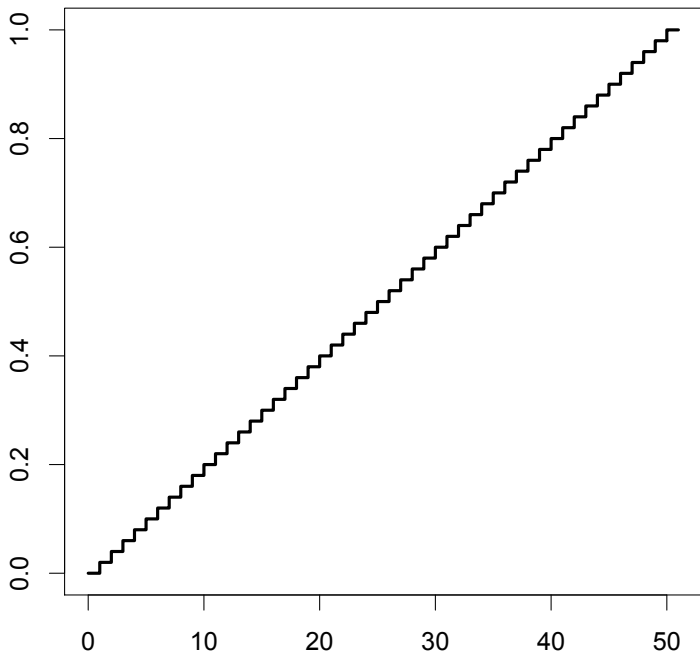
$$\hat{F}(x) \stackrel{!}{=} \sum_{r=1}^k \frac{N_r}{n} B_r(\hat{F}^M(x)).$$

Paradigms 2-3 deal with small/zero numbers  $N_r$  properly.

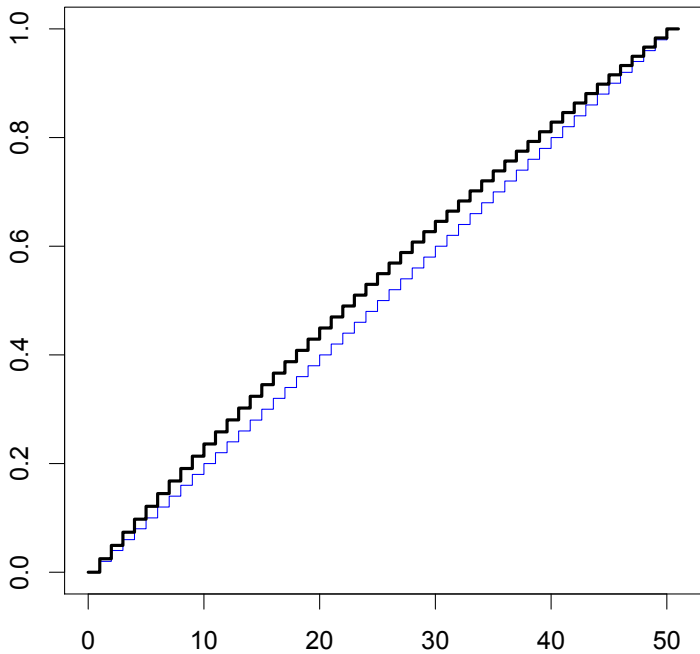
Estimator  $\hat{F}^M$  when

$$k = 2 \quad \text{and} \quad \mathbf{N} = (m, 50 - m)^\top$$

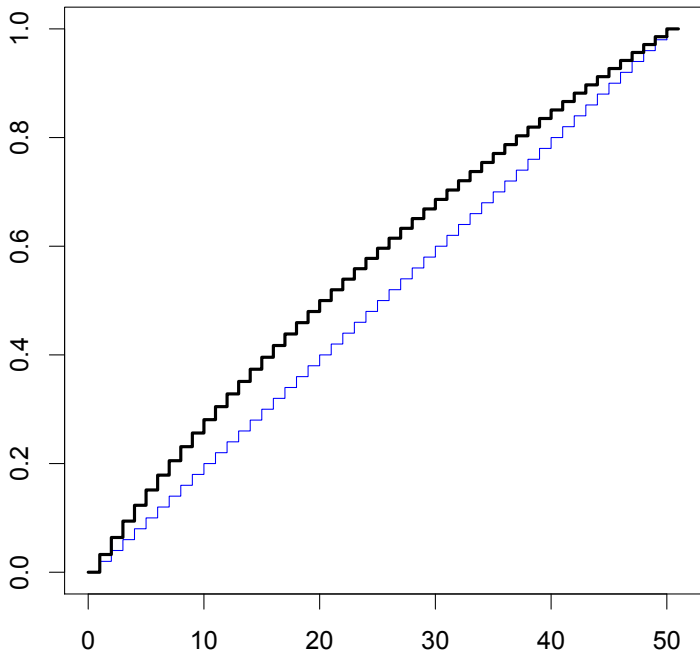
$m = 25$



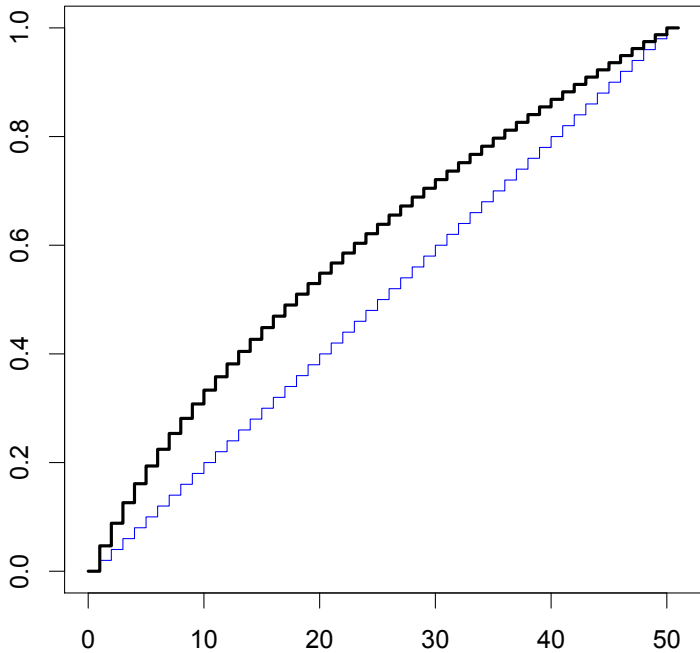
$m = 20$



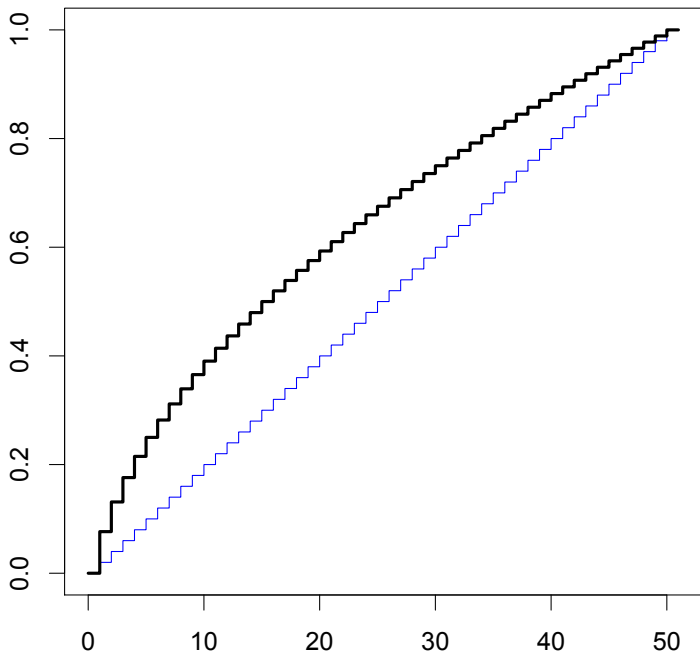
$m = 15$



$m = 10$

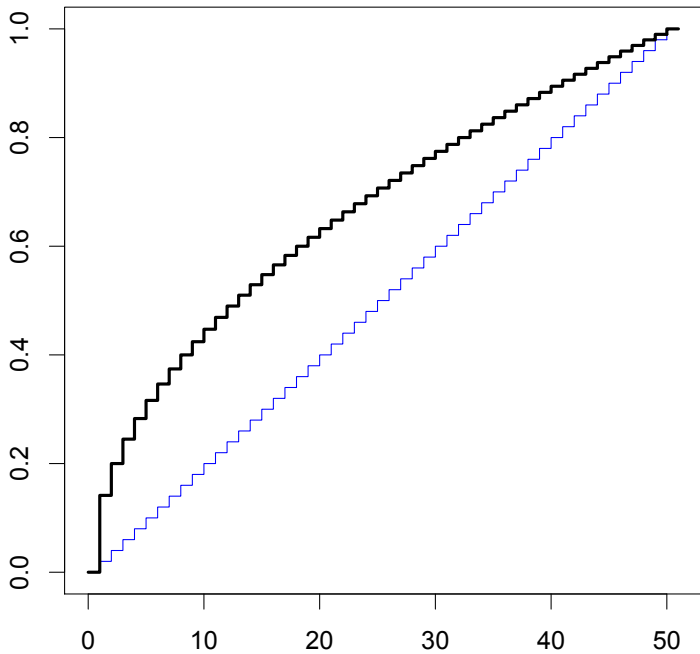


$m = 5$





$m = 0$



### III. Asymptotics

Refine and generalize results by  
Stokes and Sager (1988) and Huang (1997)

Condition on  $\mathbf{R}$  and assume that

$$\frac{N_r}{n} \rightarrow \pi_r \quad \text{for } 1 \leq r \leq k$$

where

$$\begin{cases} \pi_1, \pi_2, \dots, \pi_k > 0 & \text{for } \widehat{F}^S, \\ \pi_1, \pi_k > 0 & \text{for } \widehat{F}^M, \widehat{F}^L. \end{cases}$$

Without loss of generality let

$$F \hat{=} \text{Unif}[0, 1].$$

Empirical processes:

$$\mathbb{V}_r(t) := \sqrt{N_r}(\widehat{F}_r(B_r^{-1}(t)) - t), \quad t \in [0, 1],$$

(Standard empirical process of sample of size  $N_r$  from  $\text{Unif}[0, 1]$ .)

Linear expansion: For  $Z = S, L, M$ ,

$$\sqrt{n}(\widehat{F}^Z - F) \approx \sum_{r=1}^k \gamma_r^Z \mathbb{V}_r \circ B_r$$

with certain weight functions  $\gamma_1^Z, \gamma_2^Z, \dots, \gamma_k^Z$ .

## Weights for $\widehat{F}^S$

$$\begin{aligned}\sqrt{n}(\widehat{F}^S - F) &= \sum_{r=1}^k \frac{\sqrt{n}(\widehat{F}_r - B_r)}{k} \\ &= \sum_{r=1}^k \frac{\sqrt{n}\mathbb{V}_r \circ B_r}{k\sqrt{N_r}} \\ &\approx \sum_{r=1}^k \gamma_r^S \mathbb{V}_r \circ B_r\end{aligned}$$

with

$$\gamma_r^S := \frac{1}{k\sqrt{\pi_r}}.$$

## Weights for $\widehat{F}^M$

$$\begin{aligned} 0 &\stackrel{!}{=} n\widehat{F}(t) - \sum_{r=1}^k N_r B_r(p) \\ &= \sum_{r=1}^k N_r (\widehat{F}_r(t) - B_r(t) + B_r(t) - B_r(p)) \\ &\approx \sum_{r=1}^k N_r (\widehat{F}_r(t) - B_r(t) - \beta_r(t)(p - t)) \\ &\approx \sqrt{n} \sum_{r=1}^k \sqrt{\pi_r} \mathbb{V}_r(B_r(t)) - n(p - t) \sum_{r=1}^k \pi_r \beta_r(t). \end{aligned}$$

$$0 \stackrel{!}{\approx} \sum_{r=1}^k \sqrt{\pi_r} \mathbb{V}_r(B_r(t)) - \sqrt{n}(\rho - t) \sum_{r=1}^k \pi_r \beta_r(t).$$

Hence we expect that

$$\sqrt{n}(\hat{F}^M - F) \approx \sum_{r=1}^k \gamma_r^M \mathbb{V}_r \circ B_r$$

with

$$\gamma_r^M := \frac{\sqrt{\pi_r}}{\sum_{s=1}^k \pi_s \beta_s}.$$

## Weights for $\widehat{F}^L$

With

$$w_r := \frac{\beta_r}{B_r(1 - B_r)}$$

we write

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{\partial}{\partial p} L(t, p) = \sum_{r=1}^k N_r w_r(p) [\widehat{F}_r(t) - B_r(p)] \\ &\approx \sum_{r=1}^k N_r w_r(t) (\widehat{F}_r(t) - B_r(t) - \beta_r(t)(p - t)) \\ &\approx \sqrt{n} \sum_{r=1}^k \sqrt{\pi_r} w_r(t) \mathbb{V}_r(B_r(t)) - n(p - t) \sum_{r=1}^k \pi_r w_r(t) \beta_r(t). \end{aligned}$$

$$0 \stackrel{!}{\approx} \sum_{r=1}^k \sqrt{\pi_r} w_r(t) \mathbb{V}_r(B_r(t)) - \sqrt{n} (p - t) \sum_{r=1}^k \pi_r w_r(t) \beta_r(t).$$

Hence we expect that

$$\sqrt{n} (\hat{F}^L - F) \approx \sum_{r=1}^k \gamma_r^L \mathbb{V}_r \circ B_r$$

with

$$\gamma_r^L := \frac{\sqrt{\pi_r} w_r}{\sum_{s=1}^k \pi_s w_s \beta_s}.$$



**Theorem 1** (Linear expansions). For  $Z = S, L, M$  let

$$\mathbb{V}^Z := \sum_{r=1}^k \gamma_r^Z \mathbb{V}_r \circ B_r.$$

For any  $\delta \in [0, 1/2)$ ,

$$\sup_{t \in (0,1)} \frac{|\sqrt{n}(\widehat{F}^Z(t) - t) - \mathbb{V}^Z(t)|}{t^\delta(1-t)^\delta} \rightarrow_p 0 \quad \text{as } n \rightarrow \infty.$$

$$\sup_{t \in (0,c) \cup (1-c,1)} \frac{\sqrt{n}|\mathbb{V}^Z(t)|}{t^\delta(1-t)^\delta} \rightarrow_p 0 \quad \text{as } n \rightarrow \infty, c \downarrow 0.$$

All three estimators asymptotically equivalent in the tail regions!

**Theorem 2** (Tail behaviour). Define

$$\mathbb{V}^{(\ell)} := \frac{\mathbb{V}_1 \circ B_1}{k\sqrt{N_1/n}} \quad \text{and} \quad \mathbb{V}^{(r)} := \frac{\mathbb{V}_k \circ B_k}{k\sqrt{N_k/n}}.$$

For  $Z = S, L, M$  and any fixed  $\kappa \in [1/2, 1)$ ,

$$\sup_{t \in (0, c)} \frac{|\sqrt{n}(\widehat{F}^Z(t) - t) - \mathbb{V}^{(\ell)}(t)|}{t^\kappa} \rightarrow_p 0$$

$$\sup_{t \in (1-c, 1)} \frac{|\sqrt{n}(\widehat{F}^Z(t) - t) - \mathbb{V}^{(r)}(t)|}{(1-t)^\kappa} \rightarrow_p 0$$

as  $n \rightarrow \infty$  and  $c \downarrow 0$ .

**Theorem 3** (Limiting distributions). For  $Z = S, L, M$ ,

$$(\sqrt{n}(\widehat{F}^Z(t) - t))_{t \in [0,1]} \rightarrow_{\mathcal{L}} \mathbb{G}^Z = \sum_{r=1}^k \gamma_r^Z \mathbb{G}_r \circ B_r$$

with independent **Brownian bridges**  $\mathbb{G}_1, \mathbb{G}_2, \dots, \mathbb{G}_k$ .

The covariance function  $K^Z$  of  $\mathbb{G}^Z$  equals

$$K^Z(s, t) = \sum_{r=1}^k \gamma_r^Z(s) \gamma_r^Z(t) K(B_r(s), B_r(t))$$

with

$$K(s, t) := \min\{s, t\} - st.$$

In the **balanced case**, i.e.  $\pi_1 = \dots = \pi_k = 1/k$ ,

$$\mathbb{G}^S \equiv \mathbb{G}^M.$$

**Theorem 4** (Asymptotic relative efficiencies). For any  $t \in (0, 1)$ ,

$$K^L(t, t) \leq K^S(t, t) \quad (\text{equality for at most one } t)$$

$$K^L(t, t) \leq K^M(t, t) \quad (\text{equality iff } k = 2 \text{ and } t = 1/2).$$

On the other hand,

$$\sup_{\pi} \frac{K^S(t, t)}{K^L(t, t)} = \infty$$

$$\sup_{\pi} \frac{K^M(t, t)}{K^L(t, t)} = \frac{\rho(t) + \rho(t)^{-1} + 2}{4} < \infty$$

with

$$\rho(t) := \max_r w_r(t) / \min_r w_r(t).$$

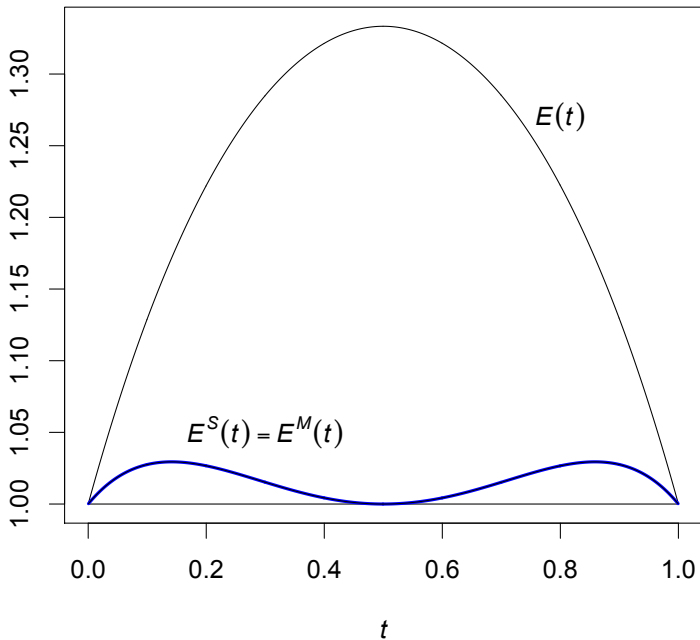
**The case  $k = 2$**

$$\frac{K^M(t, t)}{K^L(t, t)} \leq \frac{1}{1 - (2t - 1)^2/9} \leq 1.125.$$

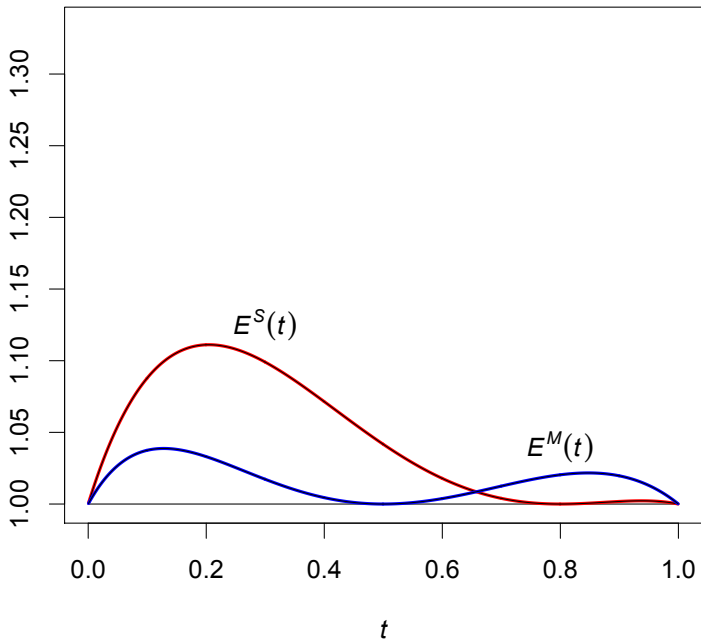
Plots of asymptotic relative (in)efficiencies

$$(0, 1) \ni t \mapsto \begin{cases} E^S(t) & := & K^S(t, t)/K^L(t, t) \\ E^M(t) & := & K^M(t, t)/K^L(t, t) \\ E(t) & := & K(t, t)/K^L(t, t) \end{cases}$$

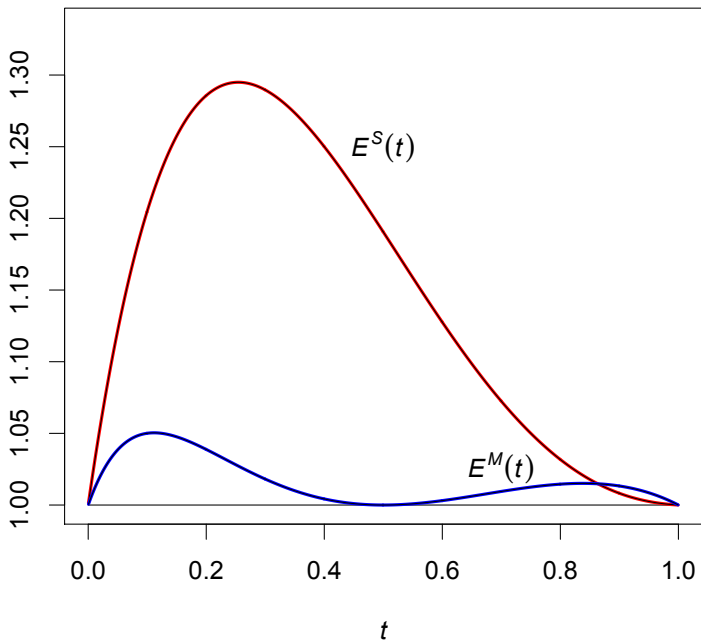
$$\pi_1 = 0.5$$



$$\pi_1 = 0.4$$

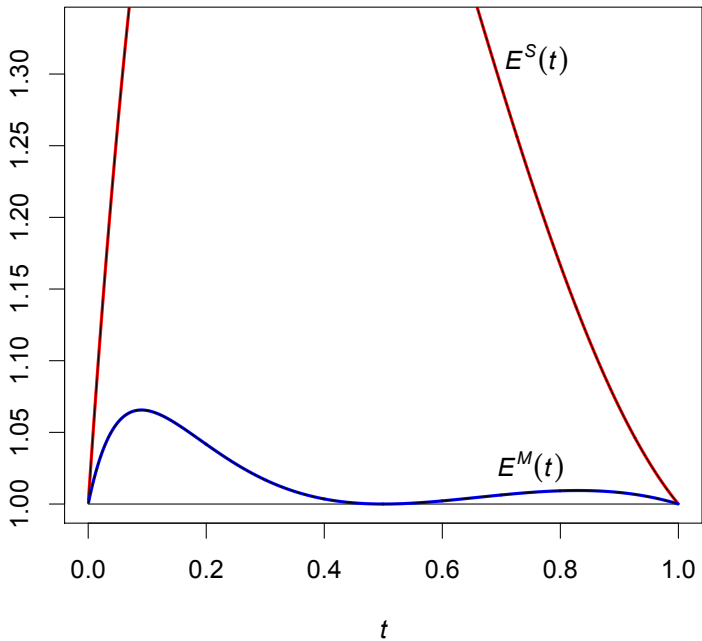


$$\pi_1 = 0.3$$

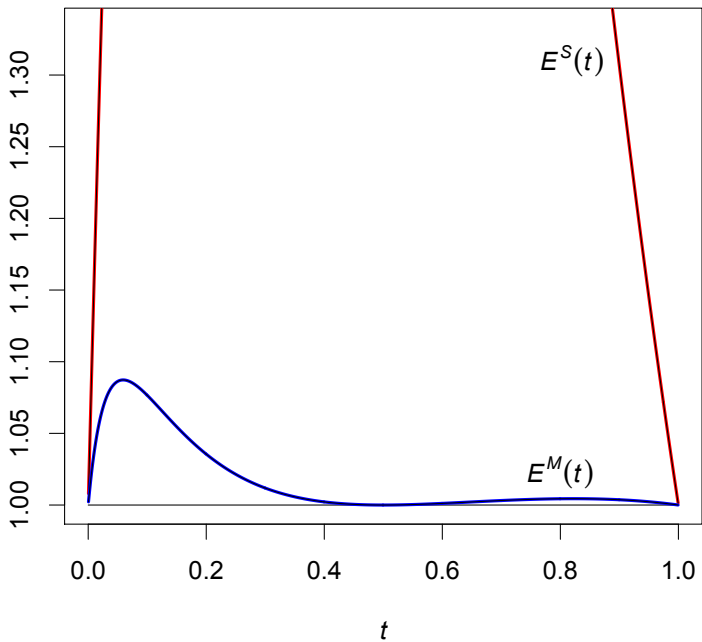




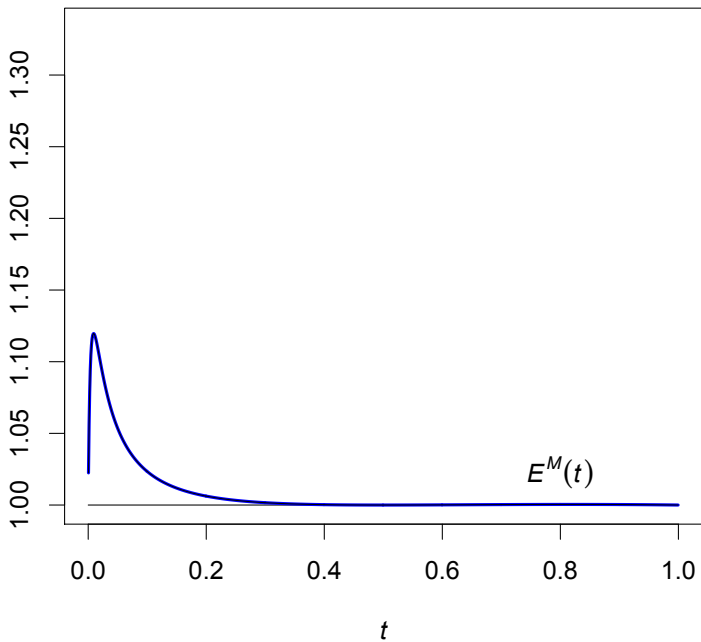
$$\pi_1 = 0.2$$

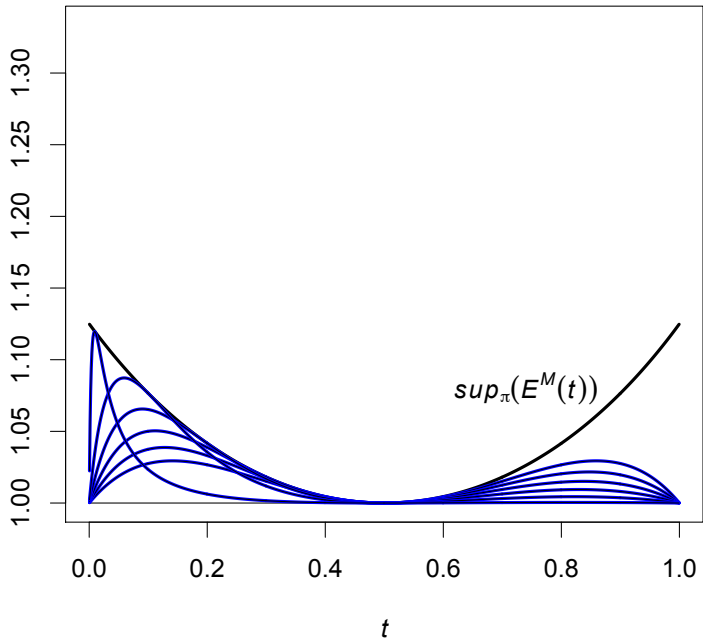


$$\pi_1 = 0.1$$



$$\pi_1 = 0.01$$





## IV. Exact Inference for Finite Samples

### Pointwise confidence intervals

Adaptation of a method by [Terpstra and Miller \(2006\)](#) ...

If  $F(x) = p \in [0, 1]$ , then conditional on  $\mathbf{N} = (N_r)_{r=1}^k$ ,

$$n\hat{F}(x) \sim \sum_{r=1}^k Y_{r,p}$$

with

$Y_{1,p}, Y_{2,p}, \dots, Y_{k,p}$  independent,  
 $Y_{r,p} \sim \text{Bin}(N_r, B_r(p))$ .

Let

$$G_{\mathbf{N},p}(y) := \mathbb{P}\left(\sum_{r=1}^k Y_{r,p} \leq y\right).$$

Exact p-values:

$$\text{Null hypothesis } "F(x) \geq p" : G_{\mathbf{N},p}(n\hat{F}(x)),$$

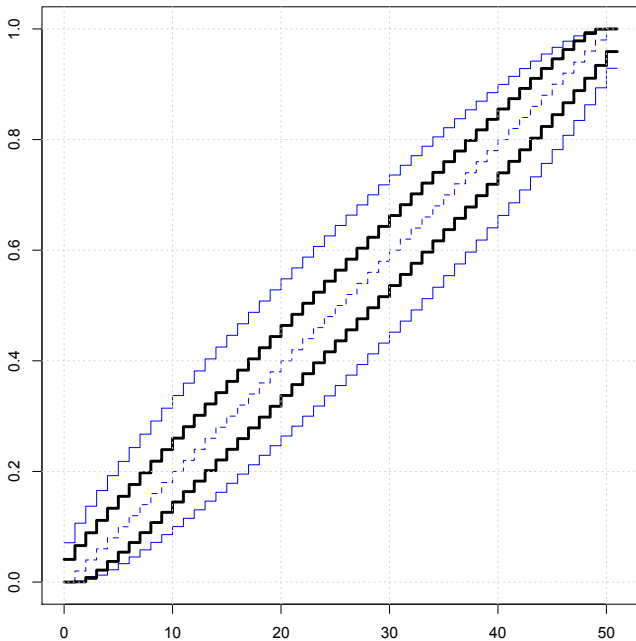
$$\text{Null hypothesis } "F(x) \leq p" : 1 - G_{\mathbf{N},p}(n\hat{F}(x) - 1).$$

Resulting **one-sided confidence bounds**:

$$\{p \in [0, 1] : G_{\mathbf{N},p}(n\hat{F}(x)) \geq \alpha\} = [0, b_{\alpha}(\mathbf{N}, n\hat{F}(x))],$$

$$\{p \in [0, 1] : G_{\mathbf{N},p}(n\hat{F}(x) - 1) \leq 1 - \alpha\} = [a_{\alpha}(\mathbf{N}, n\hat{F}(x)), 1].$$

95%-confidence intervals in RSS setting,  $n = k = 50$

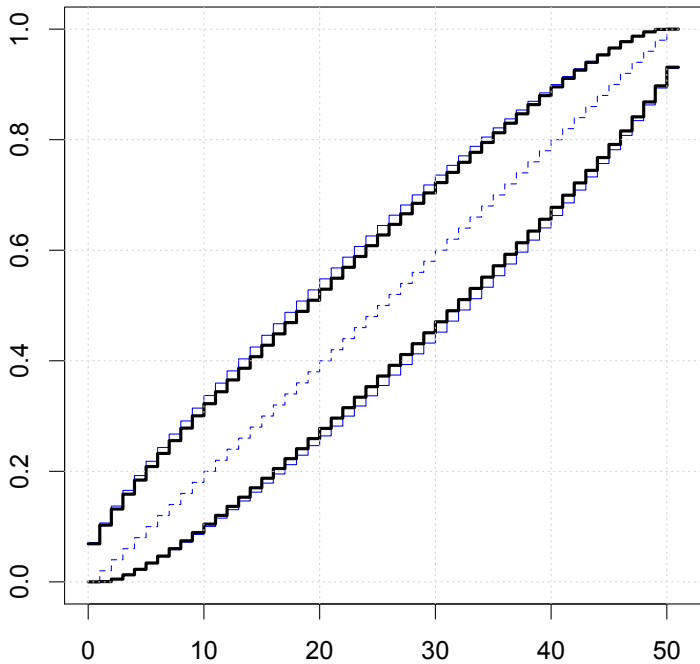


95%-confidence intervals when

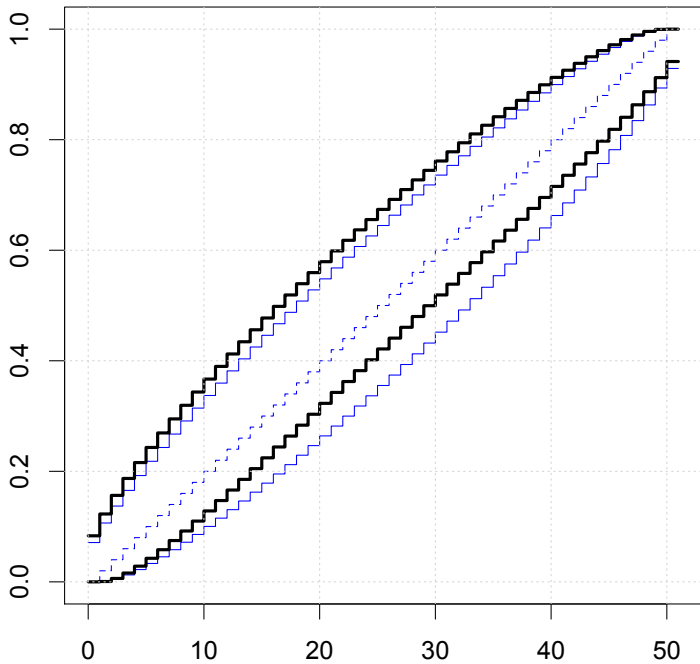
$$k = 2 \quad \text{and} \quad \mathbf{N} = (m, 50 - m)^{\top}$$



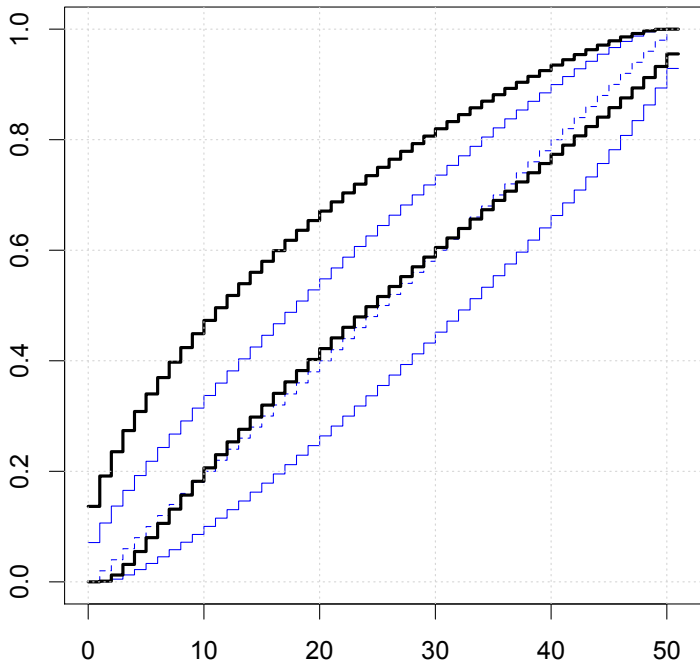
$m = 25$



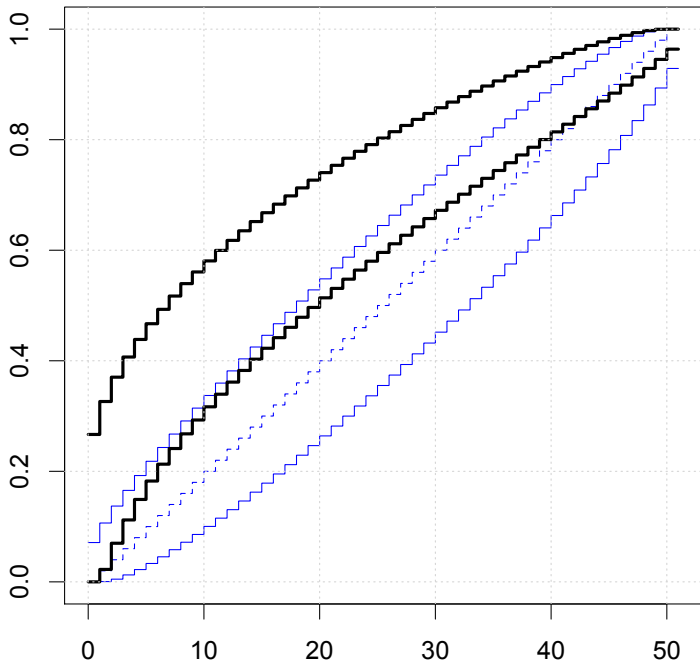
$m = 20$



$m = 10$



$m = 0$



## Confidence bands

With confidence  $1 - \alpha$ ,

$$F(x) \in [\widehat{F}^Z(x) \pm \kappa^Z(\mathbf{N}, \alpha)] \quad \text{for all } x \in \mathbb{R}$$

where

$$\kappa^Z(\mathbf{N}, \alpha) := (1 - \alpha)\text{-quantile of } \mathcal{L}\left(\sup_{t \in [0,1]} |\widehat{B}^Z(t) - t| \mid \mathbf{N}\right)$$

$$\widehat{B}^Z : \quad \text{distributed as } \widehat{F}^Z \text{ in case of } F \hat{=} \text{Unif}[0, 1]$$

Estimation of  $\kappa^Z(\mathbf{N}, \alpha)$  via Monte Carlo simulations,  
particularly easy for  $Z = \mathbb{M} \dots$

$$\kappa^M(m) := \kappa^M((m, 50 - m), 5\%):$$

