

# A Cantelli-type inequality for constructing non-parametric p-boxes based on exchangeability

Matthias C. M. Troffaes    Tathagata Basu

Department of Mathematical Sciences  
Durham University, UK

July 2019



# Outline

- 1 Cantelli's inequality and p-boxes
- 2 Contributions
- 3 Conclusions

# Cantelli's inequality & induced p-box

- random variable  $X$ , known mean  $\mu$  and variance  $\sigma^2$
- **Cantelli's inequality** [1]:

$$0 \leq P\left(\frac{X - \mu}{\sigma} \leq \lambda\right) \leq \frac{1}{1 + \lambda^2} \quad \text{if } \lambda \leq 0, \quad (1a)$$

$$\frac{\lambda^2}{1 + \lambda^2} \leq P\left(\frac{X - \mu}{\sigma} \leq \lambda\right) \leq 1 \quad \text{if } \lambda \geq 0. \quad (1b)$$

- induces a p-box (lower & upper cdf, bounding a set of probability measures)

## Issue

What if only sample mean and sample standard deviation are known?

# Contributions: Problem Statement

## Assumptions

$X_1, \dots, X_n, X_{n+1}$  is a finite sequence of discrete exchangeable random variables.

## Notation

$$\bullet \bar{X} := \frac{1}{n} \sum_{j=1}^n X_j \quad \bullet S^2 := \frac{\sum_{j=1}^n (X_j - \bar{X})^2}{(n-1)}$$

## Objective

Find functions  $\underline{f}$  and  $\bar{f}$  such that

$$\underline{f}(\lambda, n) \leq P\left(\frac{X_{n+1} - \bar{X}}{S} \leq \lambda\right) \leq \bar{f}(\lambda, n) \quad (2)$$

# Contributions: Cantelli-type inequality

## Theorem

For every  $\lambda \geq 0$ ,

$$\frac{1}{n+1} \left\lceil \frac{(n+1)\lambda_n^2}{\lambda_n^2 + 1} \right\rceil \leq P \left( \frac{X_{n+1} - \bar{X}}{S + \frac{\Delta_n}{\sqrt{n}}} < \lambda \right) \leq 1 \quad (3)$$

where  $\lambda_n := \frac{n}{\sqrt{n^2-1}}\lambda$  and  $\Delta_n := \sqrt{\frac{n+1}{n-1}}(\max X_j - \min X_j)$ .

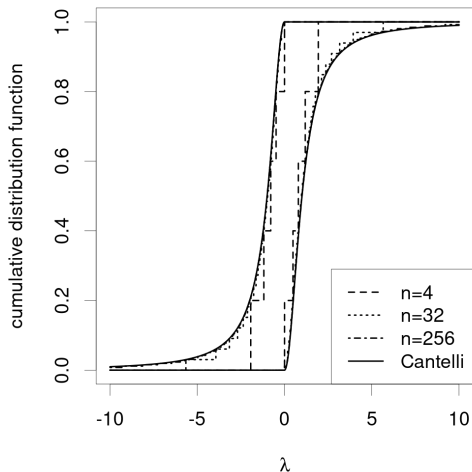
Similarly, for every  $\lambda \leq 0$ ,

$$0 \leq P \left( \frac{X_{n+1} - \bar{X}}{S + \frac{\Delta_n}{\sqrt{n}}} \leq \lambda \right) \leq \frac{1}{n+1} \left\lfloor \frac{n+1}{\lambda_n^2 + 1} \right\rfloor. \quad (4)$$

Here,  $\lfloor x \rfloor := \max\{n \in \mathbb{Z} : n \leq x\}$  and  $\lceil x \rceil := -\lfloor -x \rfloor$ .

# Contributions: Cantelli-type inequality

Plotting left and right hand sides from inequalities in theorem:



# Contributions: P-box and Prediction Interval

Inequalities from theorem induce the following:

## Non-parametric p-box

... on the random variable  $Z_{n+1} = \frac{X_{n+1} - \bar{X}}{S + \frac{\Delta_n}{\sqrt{n}}}$ .

Not a p-box on  $X_{n+1}$  directly!

## Prediction interval

... on the random variable  $X_{n+1}$ .

For any  $l_1$  and  $l_2$ , we can calculate  $\alpha_1$  and  $\alpha_2$  such that

$$\alpha_1 \leq P(\bar{X} - l_1 S_n < X_{n+1} \leq \bar{X} + l_2 S_n) \leq \alpha_2 \quad (5)$$

where,  $S_n := S + \frac{\Delta_n}{\sqrt{n}}$ .

# Conclusions I

- novel Cantelli-type inequality
- induces non-parametric p-box and prediction interval
- only assumes exchangeability (rather than conditional independence)
- only uses sample mean and sample standard deviation
- similar to Saw [2] (but Saw does not induce a p-box)
- useful for modelling when only sample mean and sample standard deviation are known (e.g. measurement problems)



## Conclusions II

- p-box only on  $Z_{n+1}$  and not on  $X_{n+1}$
- use of prediction interval not entirely clear
- p-box on  $X_{n+1}$ , conditional on  $\bar{X}$  and  $S$ , using exchangeability, remains an open problem (might be impossible as pointed out by a kind reviewer)

# Thank you for your attention!

## A Cantelli-type inequality for constructing non-parametric p-boxes based on exchangeability

Matthias C. M. Troffaes and Tathagata Basu

Department of Mathematical Sciences, Durham University, UK



### Introduction

- We derive a Cantelli-type inequality to produce a non-parametric p-box and prediction interval.
- Based on sample mean and sample standard deviation only (i.e. no parametric assumptions).
- We assume exchangeability (rather than conditional independence).
- Useful for modelling when only sample mean and sample standard deviation are known (e.g. measurement problems).

### Exchangeability [3]

We say a finite sequence  $X_1, X_2, \dots, X_n$  of discrete random variables is **exchangeable** if, for all  $\sigma \in \mathcal{S}_n$  and all  $x_1, x_2, \dots, x_n \in \mathcal{X}$

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_{\sigma(1)} = x_1, \dots, X_{\sigma(n)} = x_n)$$

Independent  $\implies$  Exchangeable

### Cantelli's inequality [1] and p-boxes [2]

Let  $X$  be a random variable with known mean  $\mu$  and known variance  $\sigma^2$ . Then the **Cantelli's inequality** is given by:

$$0 \leq P\left(\frac{X - \mu}{\sigma} \leq -\lambda\right) \leq \frac{1}{1 + \lambda^2} \quad \forall \lambda > 0$$
$$\frac{\lambda}{1 + \lambda^2} \leq P\left(\frac{X - \mu}{\sigma} \leq \lambda\right) \leq 1 \quad \forall \lambda < 0$$

A **p-box** is specified by two cumulative distribution functions,  $F$  and  $\bar{F}$  and represents the set of all cumulative distribution functions bounded by  $F$  and  $\bar{F}$ :

$$\{F \leq \bar{F}, \bar{F} \leq P(x) \leq F(x), \forall x \in \mathcal{X}\}$$

Cantelli's inequality gives distribution free p-box for known  $\mu$  and  $\sigma^2$ .

Cantelli's inequality can also be written as a p-box on  $X$ .

$$0 \leq P(X \leq \mu - \lambda\sigma) \leq \frac{\lambda^2}{1 + \lambda^2} \quad \forall \lambda > 0$$
$$\frac{\lambda}{1 + \lambda^2} \leq P(X \leq \mu + \lambda\sigma) \leq 1 \quad \forall \lambda < 0$$

### Sav's inequality (Chebyshev) [4]

Let  $X_1, \dots, X_n, Y_n$  be a sequence of discrete exchangeable random variables. Define  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$  and  $\bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2$  and  $\bar{Y}_n^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2$ . Then for every  $\lambda \geq 1$ :

$$P(|\bar{X}_n - \bar{Y}_n| > \lambda) \leq \frac{1}{\lambda^2} \frac{|\bar{X}_n^2 - \bar{Y}_n^2|}{|\bar{X}_n - \bar{Y}_n|}$$

where  $\lambda_n = \sqrt{\frac{|\bar{X}_n^2 - \bar{Y}_n^2|}{|\bar{X}_n - \bar{Y}_n|}}$  and  $\lambda(x) = \max\{\lambda \in \mathbb{R} : \lambda \leq x\}$ .

### Main Result

Let  $X_1, \dots, X_n, Y_n$  be a finite sequence of discrete exchangeable random variables. Let  $\mathcal{X} \subseteq \mathbb{R}$  denote the range of the  $X_i$  (i.e.  $\mathcal{X} = \{x \in \mathbb{R} : \exists i, X_i = x\}$ ), where  $\lambda_n$  is the maximum value that can be attained by  $X_n$  and  $\mu_n$  is the minimum value. Let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $\bar{Y}_n$  and  $\bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2$  and  $\bar{Y}_n^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2$ . Then for every  $\lambda \geq 1$ :

$$\frac{1}{n+1} \frac{|\bar{X}_n^2 - \bar{Y}_n^2|}{|\bar{X}_n - \bar{Y}_n|} \leq P\left(\frac{\bar{X}_n - \bar{Y}_n}{\lambda_n} \leq -\lambda\right) \leq \frac{1}{1 + \lambda^2} \quad (1)$$

where  $\lambda_n = \frac{\mu_n - \bar{X}_n}{\lambda_n}$  and  $\lambda_n = \frac{\bar{X}_n - \mu_n}{\lambda_n}$ . Similarly for  $\lambda < 0$ ,

$$\frac{\lambda}{1 + \lambda^2} \leq P\left(\frac{\bar{X}_n - \bar{Y}_n}{\lambda_n} \leq \lambda\right) \leq \frac{1}{1 + \lambda^2} \quad (2)$$

Here,  $\lambda(x) = \max\{\lambda \in \mathbb{R} : \lambda \leq x\}$  and  $\lambda(x) = \min\{-\lambda, x\}$ .

### P-box and Imprecise Prediction Interval

We use our main result to construct a p-box for the random variable:

$$\bar{X}_n - \bar{Y}_n = \frac{\bar{X}_n - \bar{Y}_n}{\lambda_n}$$

However, this cannot be used as a p-box for  $X_{n+1}$ . We cannot substitute the observed values for  $\bar{X}_n$  and  $\bar{Y}_n$  in the equation. Our result allows to construct an asymmetric **prediction interval**. Our result gives, for every  $\lambda$  and  $\alpha$ , the following bounds:

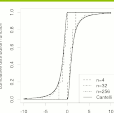
$$0 \leq P(X_{n+1} \leq \mu_n - \lambda\sigma_n) \leq \frac{\alpha}{1 + \lambda^2}$$
$$0 \leq P(X_{n+1} \leq \mu_n + \lambda\sigma_n) \leq 1 - \alpha$$

This gives us the following bounds on the coverage probability:

$$0 \leq \beta \leq P(X_{n+1} \leq \mu_n + \lambda\sigma_n) \leq 1 - \beta \leq 1 - \alpha$$

where,  $\lambda_n = \frac{\mu_n - \bar{X}_n}{\lambda_n}$ .

### Illustration



We compare our bounds with Cantelli's bounds. The correction term is larger for smaller sample sizes, which results in tighter bounds for smaller samples. For instance, for  $\alpha = 0$  this value is approximately  $0.015 \times \Delta$ , whereas for  $\alpha = 0.05$  it is only  $0.003 \times \Delta$ .

### Discussion

- We provide a Cantelli-type inequality based on exchangeability and construct a non-parametric p-box.
- We propose a prediction interval based on exchangeability.
- We can construct p-box only for the random variable  $\bar{X}_n$ .
- Finding bounds on the cumulative distribution of  $X_{n+1}$ , conditional on  $\bar{X}_n$  and  $\bar{Y}_n$  using exchangeability remains an open problem.



### References

- [1] P. Cantelli. Su certi delti probabilisti. In: *Atti del Congresso Internazionale di Matematica* 6(1926). Bologna, pp. 47–59.
- [2] Scott Ferson et al. Constructing Probability Boxes and Distributions. *Stochastic Structures, Inst. Syst. SAND98-0015*. Unpublished version, Sandia National Laboratories, Jan. 2003.
- [3] J. C. Kingman. Sites of Exchangeability. In: *The Annals of Probability* 5:2 (Apr. 1978), pp. 185–197. doi:10.1214/aop/1176935316.
- [4] John G. Sun, Mark C. K. Yang and Tai-Chiu Mei. Chebyshev Inequality with Estimated Mean and Variance. In: *The American Statistician* 38:2 (1984), pp. 130–132. URL: <https://www.jstor.org/stable/2332243>.



# We look forward to seeing you at our poster!

# References

-  F. P. Cantelli. “Sui confini della probabilità”. In: *Atti del Congresso Internazionale dei Matematici* 6 (1928). Bologna, pp. 47–59.
-  John G. Saw, Mark C. K. Yang, and Tse Chin Mo. “Chebyshev Inequality with Estimated Mean and Variance”. In: *The American Statistician* 38.2 (1984), pp. 130–132.