

# UNIVERSITETET I OSLO

## *Matematisk Institutt*

EXAM IN:                   **STK 4180/9180 – Confidence Distributions**  
                                  **Part I of two parts: The project**  
WITH:                       **Nils Lid Hjort**  
TIME FOR EXAM:           **1.–13.vi.2016**

This is the exam project set for STK 4180/9180, spring semester 2016. It is made available on the course website as of *Wednesday 1 June 12:00*, and candidates must submit their written reports by *Monday 13 June 13:00* (or earlier), to the reception office at the Department of Mathematics, in duplicate. The supplementary four-hour no-book written examination takes place *Thursday June 16* (practical details concerning this are provided elsewhere). Reports may be written in nynorsk, bokmål, riksmål, English or Latin, and should preferably be text-processed (TeX, LaTeX, Word), but may also be hand-processed. Give your ‘student number’ on the first page. Write concisely (in der Beschränkung zeigt sich erst der Meister; brevity is the soul of wit; краткость – сестра таланта). Relevant figures need to be included in the report. Copies of relevant parts of machine programmes used (in R, or matlab, or similar) are also to be included, perhaps as an appendix to the report. Candidates are required to work on their own (i.e. without cooperation with any others). They are graciously allowed not to despair should they not manage to answer all questions well.

Importantly, each student needs to submit *two special extra pages* with her or his report. *The first* (page A) is the ‘erklæring’ (self-declaration form), properly signed; it is available at the webpage as ‘Exam Project, page A, declaration form’. *The second* (page B) is the student’s one-page summary of the exam project report, which should also contain a brief self-assessment of its quality.

This exam set contains four plus one exercises and comprises six pages. The first four exercises are for both the STK 4180 and STK 9180 students, whereas the PhD students taking the STK 9180 version of the course also should do Exercise 5.

### **Exercise 1**

“LIGHT THINKS IT TRAVELS FASTER THAN ANYTHING but it is wrong. No matter how fast light travels, it finds the darkness has always got there first, and is waiting for it.” In this exercise we shall fit Simon Newcomb’s speed of light measurements, from around 1880, to a certain parametric model. You may access the data via `lib.stat.cmu.edu/DASL/Data-files/SpeedofLight.html`, for example. For some discussion of the dataset, see Schweder and Hjort (2016, Appendix A.5). Of the 66 datapoints we discard the two clear outliers at  $-44$  and  $-2$ , so the analyses below are to be based on these 64 datapoints:

28 22 36 26 28 28 26 24 32 30 27 24 33 21 36 32  
 31 25 24 25 28 36 27 32 34 30 25 26 26 25 23 21  
 30 33 29 27 29 28 22 26 27 16 31 29 36 32 28 40  
 19 37 23 32 29 24 25 27 24 16 29 20 28 27 39 23

- (a) Consider the distribution on the real line with cumulative function and density equal to

$$F_0(x) = \frac{\exp(x)}{1 + \exp(x)} \quad \text{and} \quad f_0(x) = \frac{\exp(x)}{\{1 + \exp(x)\}^2}.$$

Fit the model with cumulative distribution function

$$F(y, \xi, \tau) = F_0\left(\frac{y - \xi}{\tau}\right)$$

to the dataset, using maximum likelihood. (I find (27.621, 2.839) for  $(\hat{\xi}, \hat{\tau})$ .) Display a histogram or kernel density estimate of the data, along with the estimated parametric density. Also give associated standard deviation estimates for these estimators.

- (b) Your task now is to construct and display one or more confidence curves for  $p = \Pr_{\xi, \tau}\{Y \leq y_0\}$ , with  $y_0 = 30.5$ .
- (i) Do this using the delta method for  $\hat{p} = F(y_0, \hat{\xi}, \hat{\tau})$ .
  - (ii) Then give such a  $cc(p)$  based on the deviance function  $D(p)$  for  $p$ .
  - (iii) Attempt one of the fine-tuning methods of Chapter 7 of Schweder and Hjort (2016).
- (c) Then consider the three-parameter extension of the model above, with cumulative distribution function

$$F(y, \xi, \tau, \gamma) = \left\{ F_0\left(\frac{y - \xi}{\tau}\right) \right\}^\gamma = \left[ \frac{\exp\{(y - \xi)/\tau\}}{1 + \exp\{(y - \xi)/\tau\}} \right]^\gamma.$$

Estimate the three parameters using maximum likelihood. Obtain a confidence curve for  $\gamma$ . Is the simpler model, with  $\gamma = 1$ , supported by data?

## Exercise 2

QUANTILES MAY BE QUOTED and quartiles quartered. Here we shall look into both parametric and nonparametric confidences for quantiles.

- (a) Suppose observations  $Y_1, \dots, Y_n$  are independent from the normal  $N(\mu, \sigma^2)$  distribution. Show that  $p$ -quantile may be expressed as  $\psi_p = \mu + z_p\sigma$ , with  $z_p = \Phi^{-1}(p)$ . With  $\hat{\psi}_p = \hat{\mu} + z_p\hat{\sigma}$ , in terms of the standard estimators for the two normal parameters, show that

$$\sqrt{n}(\hat{\psi}_p - \psi_p) \rightarrow_d N\left(0, \left(1 + \frac{1}{2}z_p^2\right)\sigma^2\right)$$

as  $n$  increases. Use this to put up an approximate confidence distribution for  $\psi_p$ .

(b) Show that

$$t = \frac{\sqrt{n}(\psi_p - \hat{\psi}_p)}{\hat{\sigma}}$$

is a pivot, and explain how this may be used to construct an exact confidence distribution for  $\psi_p$  for the normal model.

(c) Then we switch gears and aim at constructing confidence distributions for quantiles without parametric assumptions. We take  $Y_1, \dots, Y_n$  to be independent from a continuous and strictly increasing cumulative distribution function  $F(y)$ , for concreteness on the half-line  $[0, \infty)$ . Let  $\psi_p = F^{-1}(p)$ . With  $Y_{(1)} < \dots < Y_{(n)}$  the ordered data, show that

$$\Pr_F\{\psi_p \leq Y_{(j)}\} = \Pr\{p \leq U_{(j)}\},$$

where  $U_{(1)} < \dots < U_{(n)}$  are the ordered observations from a sample of independent  $U_1, \dots, U_n$  from the uniform distribution on the unit interval.

(d) It is well-known that the  $U_{(j)}$  has a beta distribution, with parameters  $(j, n + 1 - j)$ . Explain how this leads to a nonparametric confidence distribution  $C(\cdot)$  for the  $p$ -quantile, with

$$C(Y_{(j)}) = 1 - G(p, j, n + 1 - j) = 1 - \int_0^p g(u, j, n + 1 - j) du,$$

in terms of the cumulative Beta  $(a, b)$  distribution function  $G(u, a, b)$  with density  $g(u, a, b)$  (these are available in R as `pbeta` and `dbeta`).

(e) To see these methods in action, go to the booksite [feb.kuleuven.be/public/ndbaf-45/modelselection/](http://feb.kuleuven.be/public/ndbaf-45/modelselection/) of Claeskens's and Hjort's *Model Selection and Model Averaging*, and access the dataset on low birthweights for 189 newborns; see also 'small babies' in Schweder and Hjort (2016, p. 437). Then get hold of the 189 weights of the 189 mothers (prior to pregnancy), where you also should convert these from pounds to kilograms. Apply the normal-based as well as the nonparametric method to compute and display confidence curves for the 0.90-quantile of this weight distribution. Comment on your findings. You may also try to display confidence curves for 0.10, 0.50, 0.90 quantiles, in the same diagram.

### Exercise 3

"POWER IS NOT A MEANS; IT IS an end. The object of power is power." Orwell did not necessarily have confidence in mind, but here we are concerned with confidence power.

(a) Let  $X$  and  $Y$  be independent and exponentially distributed with positive parameters  $a$  and  $a + \delta$ , i.e. with joint density

$$a \exp(-ax)(a + \delta) \exp\{-(a + \delta)y\} \quad \text{for } x > 0, y > 0.$$

Show that the conditional density of  $Y$  given  $Z = X + Y$  is

$$g(y|z) = \frac{\delta \exp(-\delta y)}{\int_0^z \delta \exp(-\delta y') dy'} = \frac{\delta \exp(-\delta y)}{1 - \exp(-\delta z)} \quad \text{for } 0 \leq y \leq z.$$

- (b) Find the optimal confidence distribution  $C(\delta) = C(\delta, x, y)$  for  $\delta$ , after having observed  $(x, y)$ .
- (c) Assume now that we observe three independent pairs of exponentials, with

$$X_j \sim \text{Expo}(a_j) \text{ and } Y_j \sim \text{Expo}(a_j + \delta) \quad \text{for } j = 1, 2, 3,$$

with the data pairs being

$$(2.864, 0.156), (1.354, 0.314), (1.438, 0.377).$$

Compute and display the three corresponding confidence curves  $cc_j(\delta) = |1 - 2C_j(\delta)|$ , for  $j = 1, 2, 3$ , in the same diagram.

- (d) Above we have three confidence distributions, or three confidence curves, for the same parameter  $\delta$ . These may be combined in different ways. Show that

$$C^*(\delta) = \Phi\left(\sum_{j=1}^3 \frac{1}{\sqrt{3}} \Phi^{-1}(C_j(\delta))\right)$$

is a confidence distribution, and display the consequent confidence curve  $cc^*(\delta)$  alongside the three component confidence curves.

- (e) Attempt to find, compute and display a better confidence curve than  $cc^*(\delta)$ .

#### Exercise 4

“THERE IS NO CORRELATION BETWEEN A CHILDHOOD SUCCESS and a professional athlete” (says a man with nine Olympic gold medals). This exercise concerns a certain intraskater correlation parameter, used in the speedskating model described and employed in the CLP book, Section 14.5, to analyse a list of Sprint Speedskating World Championships regarding what influences the 1000-metre results. The intraskater correlation relates to the degree of similarity between a top skater’s performance on Saturday and on Sunday.

You might wish to have a brief look at that model, and its use and its context, but for the most part you should be able to handle the present exercise without knowing the many details involved in the speedskating analyses. Essentially, it suffices to know that the intraskater correlation estimates, say  $\hat{\rho}_j$  for the world championships of year  $j$ , are independent, with  $\hat{\rho}_j \sim N(\rho_j, \sigma_j^2)$ , where the standard errors (estimated standard deviations) are well estimated and taken as known. The table below displays these  $\hat{\rho}_j$  and  $\sigma_j$ , for the world championships of Nagano 2014, Salt Lake City 2013, Calgary 2012, Heerenveen 2011, Obihiro 2010:

year	$\hat{\rho}_j$	$\sigma_j$
2014	0.736	0.098
2013	0.801	0.076
2012	0.584	0.134
2011	0.737	0.093
2010	0.902	0.039

- (a) Make a figure of the  $k = 5$  (approximate) 90% confidence intervals  $\hat{\rho}_j \pm 1.645 \sigma_j$ .
- (b) Assume now that the  $\rho_j$  are not constant from championship to championship, but may vary with a certain degree of variability, modelled as  $\rho_j \sim N(\rho_0, \tau^2)$ . Show that

$$\hat{\rho}_j \sim N(\rho_0, \sigma_j^2 + \tau^2) \quad \text{for } j = 1, \dots, k.$$

- (c) Define

$$\hat{\rho}_0(\tau) = \frac{\sum_{j=1}^k \hat{\rho}_j}{\sum_{j=1}^k \frac{1}{\sigma_j^2 + \tau^2}}.$$

Write down the log-likelihood function  $\ell(\rho_0, \tau)$ , and from this derive an expression for the profile log-likelihood function  $\ell_{\text{prof}}(\tau)$ . Compute and display this profile function in a diagram, and find the maximum likelihood estimates  $(\hat{\rho}_0, \hat{\tau})$ .

- (d) Construct a confidence distribution for  $\tau$  based on

$$Q(\tau) = \sum_{j=1}^k \frac{\{\hat{\rho}_j - \hat{\rho}_0(\tau)\}^2}{\sigma_j^2 + \tau^2},$$

using the fact that  $Q(\tau)$ , under the true  $\tau$ , has a  $\chi_{k-1}^2$  distribution. (This is shown in an exercise for Chapter 13 in CLP, and does not need to be demonstrated here.) Compute its point-mass at zero.

- (e) Also construct a second confidence distribution for  $\tau$ , using

$$C^*(\tau) = \Pr_{\tau}\{R(\tau) \geq R_{\text{obs}}(\tau)\},$$

where

$$R(\tau) = \sum_j \frac{|\hat{\rho}_j - \hat{\rho}_0(\tau)|}{(\sigma_j^2 + \tau^2)^{1/2}}.$$

Since  $R(\tau)$  does not have any clear distribution, you need to compute  $C^*(\tau)$  by simulation. Display the two confidence distribution for  $\tau$  in the same diagram.

- (f) Then construct a confidence curve for the overall mean intraskater correlation parameter  $\rho_0$ , perhaps using t-bootstrapping.

### Exercise 5: for the PhD students taking the STK 9180 exam

The article ‘Confidence Distribution, the Frequentist Distribution Estimator of a Parameter: A Review’, by Min-ge Xie and Kesar Singh, was published in *International Statistical Review*, in 2013. The article may be found on Min-ge Xie’s webpage

[www.stat.rutgers.edu/home/mxie/RCPapers/insr.12000.pdf](http://www.stat.rutgers.edu/home/mxie/RCPapers/insr.12000.pdf)

along with various discussion contributions, by the famous statisticians Sir David Cox, Brad Efron, Don Fraser, Emmanuel Parzen, Christian Robert, Tore Schweder and Nils Lid

Hjort. These discussion contributions, with a rejoinder from Xie and Singh, can also be found at Xie's webpage:

[www.stat.rutgers.edu/home/mxie/RCPapers/insr.12001.discussion.pdf](http://www.stat.rutgers.edu/home/mxie/RCPapers/insr.12001.discussion.pdf)

Your task is to single out *two* of the discussion contributions, though not the one by Schweder and Hjort, and then (a) provide a brief summary of the points being made, (b) provide your own views concerning these points, perhaps along with further comments from your side.