

# Statische Modellen & Data-analyse

## Homework 1

Li Quan

27 april 2011

### Vraag 1

Gegeven de verdelingsfunctie (waarbij  $\lfloor x \rfloor = \text{Floor}(x)$  en  $p \in (0, 1)$ )

$$F_X(x) = \begin{cases} 1 - (1 - p)^{\lfloor x \rfloor} & \text{als } x \geq 1 \\ 0 & \text{elders} \end{cases}$$

kan de dichtheidsfunctie  $f_X(x)$  berekend worden:

$$f_X(x) = \frac{\partial F_X(x)}{\partial x} = \ln(1 - p) \cdot (-(1 - p)^{\lfloor x \rfloor}) \cdot \frac{\partial \lfloor x \rfloor}{\partial x}.$$

De maximum likelihood schatter voor  $p$  (gegeven een steekproef  $x_1, \dots, x_n$  uit  $F_X$ ),

$$\hat{p}_{MLE} = \arg \max_p \sum_{i=1}^n \ln f_X(x_i)$$

kan gevonden worden door volgende vergelijking op te lossen naar  $\hat{p}$ :

$$0 = \sum_{i=1}^n \frac{\partial \ln f_X(x_i)}{\partial p} \Big|_{p=\hat{p}} = \sum_{i=1}^n \left\{ -\frac{1}{\ln(1 - \hat{p})} \frac{1}{1 - \hat{p}} + \frac{\lfloor x_i \rfloor}{1 - \hat{p}} \right\}$$

Uiteindelijk vinden we (na bevestiging dat het om een maximum gaat)

$$\hat{p}_{MLE} = 1 - \exp\left(\frac{-n}{\sum_{i=1}^n \lfloor x_i \rfloor}\right).$$

Om te illustreren dat  $\hat{p}_{MLE}$  een vertekende schatter is van  $p$  werd dit in R gesimuleerd<sup>1</sup> (zie code in Appendix A) voor  $n = 50$ . Tabel 1 toont de resultaten hiervan: hieruit blijkt dat de MLE waarschijnlijk een onderschatter is.

---

<sup>1</sup>Voor het genereren van de steekproefrealisatie werd gebruik gemaakt van de *inverse transformation sampling*-techniek (zie bijvoorbeeld [http://en.wikipedia.org/wiki/Inverse\\_transform\\_sampling](http://en.wikipedia.org/wiki/Inverse_transform_sampling)).

$p$	$\hat{p}_{MLE}^{(k)}$				
	1	2	3	4	5
0.1	0.097	0.104	0.108	0.093	0.090
0.3	0.315	0.323	0.295	0.300	0.300
0.5	0.473	0.482	0.531	0.461	0.553
0.7	0.671	0.655	0.647	0.640	0.654
0.9	0.854	0.770	0.875	0.771	0.780

**Tabel 1:**  $\hat{p}_{MLE}$  voor verschillende waarden van  $p$  met steekproefgrootte  $n = 50$  en  $k = 5$  runs (zie R-code Appendix A).

Om dit te bewijzen definiëren we eerst de bias  $b(\hat{p}) = \mathbb{E}[\hat{p}] - p$ . Dan moeten we aantonen dat  $b(\hat{p}) \neq 0$ . Algemeen werd reeds aangetoond dat MLE schatters onvertekend zijn tot op orde  $n^{-1/2}$ , maar vertekend op orde  $n^{-1}$  [2, Eq. (20)].

Een algemene uitdrukking voor de bias van een MLE schatter  $\hat{\theta}$  is [3]:

$$b(\theta) = \frac{b_1(\theta)}{n} + \frac{b_2(\theta)}{n^2} + \dots$$

Asymptotisch is de MLE dus onvertekend (dit is de asymptotische normaliteitseigenschap van MLE schatters), aangezien  $\lim_{n \rightarrow \infty} b = 0$ .

## Vraag 2

We beschouwen eerst  $Y$  en  $Z$ . Dan is  $\begin{pmatrix} Y \\ Z \end{pmatrix} \sim \mathcal{N}_2\left(\begin{pmatrix} \mu_Y \\ \mu_Z \end{pmatrix}, \begin{pmatrix} \sigma_{YY} & \sigma_{YZ} \\ \sigma_{ZY} & \sigma_{ZZ} \end{pmatrix}\right)$ , waarbij we volgende vergelijkingen moeten oplossen ( $\mu_Z = 0$  en  $\sigma_Z^2 = \sigma_{ZZ} = 1$ ):

$$\begin{aligned} \mathbb{E}[Y|Z = z] &= 1 + z \\ &= \mu_{YY} + \sigma_{YZ}\sigma_Z^{-1}(z - \mu_Z) = \mu_Y + \sigma_{YZ}z \\ \text{Cov}[Y|Z = z] &= 1 \\ &= \sigma_{YY} - \sigma_{YZ}\sigma_Z^{-1}\sigma_{ZY} = \sigma_{YY} - \sigma_{YZ}\sigma_{ZY} \end{aligned}$$

De oplossing hiervoor is  $\mu_Y = 1$ ,  $\sigma_{YZ} = \sigma_{ZY} = 1$  en  $\sigma_{YY} = 2$ .

Dus dan is<sup>2</sup>  $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \sim \mathcal{N}_3\left(\begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{pmatrix}\right)$ , met bijkomende vergelijkingen:

$$\mathbb{E}[X|Y = y, Z = z] = 1 - y$$

$$= \mu_1 + \begin{pmatrix} \sigma_{12} & \sigma_{13} \end{pmatrix} \begin{pmatrix} \sigma_{22} & \sigma_{23} \\ \sigma_{32} & \sigma_{33} \end{pmatrix}^{-1} \begin{pmatrix} y - \mu_2 \\ z - \mu_3 \end{pmatrix}$$

$$\text{Cov}[X|Y = y, Z = z] = 1$$

$$= \sigma_{11} - \begin{pmatrix} \sigma_{12} & \sigma_{13} \end{pmatrix} \begin{pmatrix} \sigma_{22} & \sigma_{23} \\ \sigma_{32} & \sigma_{33} \end{pmatrix}^{-1} \begin{pmatrix} \sigma_{21} \\ \sigma_{31} \end{pmatrix}$$

De oplossingen van het volledige stelsel zijn dan:  $\mu_1 = 0$ ,  $\mu_2 = 1$ ,  $\mu_3 = 0$ ,  $\sigma_{11} = 3$ ,  $\sigma_{22} = 2$ ,  $\sigma_{33} = 1$ ,  $\sigma_{12} = \sigma_{21} = 2$ ,  $\sigma_{13} = \sigma_{31} = 1$  en  $\sigma_{23} = \sigma_{32} = 1$ . Dus is

$$\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \sim \mathcal{N}_3\left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 3 & 2 & 1 \\ 2 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix}\right).$$

Na herschikken tot  $\begin{pmatrix} Y \\ X \\ Z \end{pmatrix} \sim \mathcal{N}_3\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 1 & 1 \end{pmatrix}\right)$  kunnen we gemakkelijk  $Y|X, Z \sim \mathcal{N}(\mu_Y, \sigma_{YY})$  berekenen waarbij

$$\mu_Y = 1 + \begin{pmatrix} 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 \\ 1 & 1 \end{pmatrix}^{-1} \begin{pmatrix} X \\ Z \end{pmatrix} = \frac{X + Z}{2}$$

$$\sigma_{YY} = 2 - \begin{pmatrix} 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 1 \\ 1 & 1 \end{pmatrix}^{-1} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = \frac{1}{2}$$

Voor de verdeling van  $(U, V)^T$  zoeken we eerst de verdeling van  $U$  en  $V$ . Hiervoor gebruiken we volgende eigenschap: als  $X \sim \mathcal{N}(\mu, \sigma^2)$ , dan is een lineaire transformatie  $aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$ . Uit  $U = 1 + Z$  en  $V = 1 - Y$  volgt dan dat  $U \sim \mathcal{N}(1, 1)$  en  $V \sim \mathcal{N}(0, 2)$ . Verder volgt uit de lineariteitseigenschap van de covariantie dat  $\text{cov}(U, V) = \text{cov}(1 + Z, 1 - Y) = \text{cov}(Z, -Y) = -\text{cov}(Z, Y) = -1$  zodat

$$\begin{pmatrix} U \\ V \end{pmatrix} \sim \mathcal{N}_2\left(\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & -1 \\ -1 & 2 \end{pmatrix}\right).$$

---

<sup>2</sup>Hierbij worden cijferindices gebruikt:  $X \equiv 1$ ,  $Y \equiv 2$  en  $Z \equiv 3$ .

Aangezien  $\text{cov}(U, Y) = \text{cov}(1 + Z, Y) = \text{cov}(Z, Y) = 1$ , is  $\begin{pmatrix} U \\ Y \end{pmatrix} \sim \mathcal{N}_2\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}\right)$ .

Dan is  $\mathbb{E}(Y|U = u) = \mu_Y + \sigma_{YU}\sigma_{UU}^{-1}(u - \mu_U)$  en dus

$$\mathbb{E}(Y|U = 2) = 1 + \frac{1}{1}(2 - 1) = 2.$$

### Vraag 3

Definieer  $\mathbf{X}'_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ X_{i3} \end{pmatrix} - \begin{pmatrix} X_{i1} \\ X_{i1} \\ X_{i1} \end{pmatrix}$ . Dan is  $\mathbb{E}[\mathbf{X}'_i] = \boldsymbol{\mu}' = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_1 \\ \mu_1 \end{pmatrix} = \begin{pmatrix} 0 \\ \mu_2 - \mu_1 \\ \mu_3 - \mu_1 \end{pmatrix}$ .

Dan is de hypothese equivalent met  $\mathbf{A}\boldsymbol{\mu}' = \boldsymbol{\mu}_0 = \mathbf{0}$ , waarbij  $\mathbf{A} = \mathbf{I}_3$  (of eender welke niet-singuliere matrix, bijvoorbeeld  $\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix}$ ). Dus kunnen we eenvoudig de

(affien invariante)  $T^2$ -test van Hotelling gebruiken waarbij:

$$\begin{aligned} \bar{\mathbf{X}}' &= \sum_{i=1}^n \mathbf{X}'_i \\ \mathbf{S} &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}'_i - \bar{\mathbf{X}}')(\mathbf{X}'_i - \bar{\mathbf{X}}')^\tau \\ T^2 &= n(\bar{\mathbf{X}}' - \boldsymbol{\mu}_0)^\tau \mathbf{S}^{-1}(\bar{\mathbf{X}}' - \boldsymbol{\mu}_0) = n\bar{\mathbf{X}}'^\tau \mathbf{S}^{-1}\bar{\mathbf{X}}' \end{aligned}$$

### Vraag 4

Definieer  $\mathbf{D}_i = \mathbf{X}_{i1} - \mathbf{X}_{i2}$ ,  $\mathbf{S}_D = \frac{(n_1-1)\mathbf{S}_1 + (n_2-1)\mathbf{S}_2}{n_1+n_2-2} = \frac{\mathbf{S}_1 + \mathbf{S}_2}{2}$  (aangezien  $n = n_1 = n_2$ ) en  $\boldsymbol{\mu}_D = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ . (Alle hypothesetesten gebeuren op het significantieniveau  $\alpha = 0.05$ .)

De nulhypothese  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$  is dan equivalent met  $\boldsymbol{\mu}_D = \mathbf{0}$ . We gebruiken dus de multivariate *gepaarde*  $T^2$ -test van Hotelling [1] waarbij  $T^2 = n\bar{\mathbf{d}}^\tau \mathbf{S}_D^{-1}\bar{\mathbf{d}}$ . Dan kan volgende  $F$ -statistiek gebruikt worden:  $F = \frac{n-p}{p(n-1)}T^2 \sim F_{p,n-p}$ .

Voor de gegeven realisatie ( $n = 10$  en  $p = 2$ ) wordt dit dus:

$$\begin{aligned} \bar{\mathbf{d}} &= \begin{pmatrix} 3 & 1 \end{pmatrix}^\tau - \begin{pmatrix} 1 & 1 \end{pmatrix}^\tau = \begin{pmatrix} 2 & 0 \end{pmatrix}^\tau \\ \mathbf{s}_D &= \frac{1}{2} \left\{ \begin{pmatrix} 4 & -1 \\ -1 & 2 \end{pmatrix} + \begin{pmatrix} 2 & -2 \\ -2 & 4 \end{pmatrix} \right\} = \begin{pmatrix} 3 & -1.5 \\ -1.5 & 3 \end{pmatrix} \\ t^2 &= 10 \begin{pmatrix} 2 & 0 \end{pmatrix} \begin{pmatrix} 0.444 & 0.222 \\ 0.222 & 0.444 \end{pmatrix} \begin{pmatrix} 2 & 0 \end{pmatrix}^\tau = 17.778 \\ f &= \frac{8}{18}t^2 = 7.902 \end{aligned}$$

De  $p$ -waarde is  $p = 2P_{H_0}(F \geq |f|) = 2(1 - 0.9872) = 0.0255 < 0.05$ . Dus we verwerpen de nulhypothese.

Voor de hypotheses testen voor de aparte componenten gebruiken we volgende formule voor het  $100(1 - \alpha)\%$  betrouwbaarheidsinterval (BI):

$$\mathbf{u}^\tau \bar{\mathbf{D}} \pm t_{n-1, 1-\alpha/2} \sqrt{\frac{\mathbf{u}^\tau \mathbf{S}_D \mathbf{u}}{n}}.$$

Voor de nulhypothese  $H_0 : \mu_{11} = \mu_{21}$  kiezen we  $\mathbf{u} = \begin{pmatrix} 1 & 0 \end{pmatrix}^\tau$ . We krijgen dan als BI

$$\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 0 \end{pmatrix} \pm t_{9, 0.975} \sqrt{\frac{\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} 3 & -1.5 \\ -1.5 & 3 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix}^\tau}{10}}.$$

Met  $t_{9, 0.975} = 2.262$  geeft dit dan als BI =  $[0.761; 3.239]$ . We verwerpen dus  $H_0$  aangezien het BI niet 0 bevat.

Analoog voor de nulhypothese  $H_0 : \mu_{12} = \mu_{22}$  (waarbij  $\mathbf{u} = \begin{pmatrix} 0 & 1 \end{pmatrix}^\tau$ ), verkrijgen we uiteindelijk het BI =  $[-1.239; 1.239]$ . Hier aanvaarden we dus wel  $H_0$ .

## Referenties

- [1] T. W. Anderson. *An Introduction to Multivariate Statistical Analysis*. Wiley-Interscience, second edition, Sept. 1984.
- [2] D. R. Cox and E. J. Snell. A general definition of residuals. *Journal of the Royal Statistical Society. Series B (Methodological)*, 30(2):pp. 248–275, 1968.
- [3] D. Firth. Bias reduction of maximum likelihood estimates. *Biometrika*, 80(1):27–38, Mar. 1993.

## A. R-code

```
#Li Quan - April 27, 2011
#homework 1: question 1, part 2

#the given cdf function
mycdf <- function(p, x) {
  ifelse(x >= 1, 1 - (1-p)**(floor(x)), 0)
}

#set the p in (0,1)
p <- 0.9
cdf <- function(x) mycdf(p, x)

#inverts the cdf (using numeric solver)
# http://en.wikipedia.org/wiki/Inverse_transform_sampling
# http://stackoverflow.com/questions/1594121/r-how-do-i-best-simulate-an-arbitrary
#   -univariate-random-variate-using-its-probab
inversefun.cdf<-function(x, cdf, starting.value=0, precision=1e-6){
  #we first search a lower bound
  low.found<-FALSE
  low<-starting.value
  while(!low.found){
    if(cdf(low)>=(x-precision))
      low<-low-(low-starting.value)^2-1
    else
      low.found<-TRUE
  }
  #... and an upper bound
  up.found<-FALSE
  up<-starting.value
  while(!up.found){
    if(cdf(up)<=(x+precision))
      up<-up+(up-starting.value)^2+1
    else
      up.found<-TRUE
  }
  # solve this equation
  uniroot(function(y) cdf(y)-x, c(low, up))$root
}

#sample size
n <- 50
#generates n random variables of the distribution using inverse transform sampling
vars<- sapply(runif(n), function(x) inversefun.cdf(x,cdf))
hist(vars)

#calculates the estimator for p
#we immediately use the explicit formulation for p_MLE
calcEstimator<-function(vars) {
  n <- length(vars)
  vars_floors <- sapply(vars, 'floor') #we apply the floor function on each
    element of the list
  result <- 1 - exp(-n/sum(vars_floors));result
}

pMLE <- calcEstimator(vars)
#print some information, we could use hypothesis tests but we leave it up to the
  reader
cat("Sample size n was", n, "\n\nThe original p was", p, "\n\nThe mle for p is", pMLE)
```