



Avd. Matematisk statistik

KTH Matematik

EXAM IN SF2955 COMPUTER INTENSIVE METHODS
WEDNESDAY 18th of August 2010 14.00 – 19.00 in Q17.

Examiner: Timo Koski, tel. 790 71 34. Email: tjtkoski@kth.se

Allowed aids: Formel- och tabellsamling i Matematisk statistik. Formulas in Bayesian Inference. Pocket calculator.

There are six (6) assignments (uppgift). Notation should be explained and defined. Arguments must be detailed enough to make it possible to follow. Numerical answers should be given with the precision of at least 2 significant digits.

Results will be ready before September 1st. You will be informed by email if you specify your email address.

Good luck!

Uppgift 1

The systematic error (i.e., bias) for an estimator $\hat{\theta}$ of the parameter θ is defined as

$$b_F(\theta, \hat{\theta}) = E(\hat{\theta}(X_1, X_2, \dots, X_n)) - \theta.$$

The jackknife estimate of this quantity is

$$\widehat{bias}_{jack} = (n - 1)(\hat{\theta}_{(\cdot)} - \hat{\theta})$$

where

$$\hat{\theta}_{(\cdot)} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(i)}$$

and in addition $\hat{\theta}_{(i)} = \hat{\theta}(X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$, which is the estimate in the i :th jackknife sample when the observation nr i has been removed.

Show that if

$$E(\hat{\theta}) = \theta + \frac{a_1(F)}{n} + \frac{a_2(F)}{n^2} + \dots$$

where F is the underlying distribution for the data, then \widehat{bias}_{jack} has the expectation $a_1(F)/n + O(1/n^2)$. (10 p)

Uppgift 2

a) Assume that a timehomogeneous Markov chain with the transition matrix $P(x, y) = P(X_{n+1} = y | X_n = x)$ and the state space E is *time reversible* with respect to $\pi(x)$, $x \in E$, i.e.,

$$\pi(x)P(x, y) = \pi(y)P(y, x), \quad x, y \in E.$$

Show that the stationary distribution of the chain is π . (5 p)

b) The Metropolis-Hastings algorithm generates a Markov chain in the following manner: If the chain is currently in the state $x \in E$, a proposal y for the next state is given by drawign from the distribution $q(x, y) = P(\text{propose } y | \text{current state is } x)$. This proposal is accepted with probability

$$\alpha(x, y) = \min \left(1, \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \right).$$

If the proposal is not accepted, the chain remains in state x , and a new proposal is drawn.

Prove that the Markov chain so defined has π as a stationary distribution. (5 p)

Uppgift 3

We have 20 observations x_1, x_2, \dots, x_{20} of independent identically distributed variables med expectation m och variance σ^2 . We estimate σ^2 with

$$\hat{\sigma}^2 = \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})^2.$$

We have $\bar{x} = 1.7790$ and $\hat{\sigma}^2 = 22.8095$.

Let us consider bootstrapping of $\hat{\sigma}^2$ with the MATLAB^R code

```
boot=bootstrp(999, 'var', x, 1);
```

which generates 999 bootstrap samples from x_1, x_2, \dots, x_{20} och computes the estimated variance $\hat{\sigma}^2$ in each of the 999 samples. By these 999 estimates we can estimate the percentiles in the bootstrap distribution of $\hat{\sigma}^2$. In particular, the 50:th and 950:th in order of magnitude are 13.5 and 30.5, respectively. You are now expected to compute symmetric 90% confidence intervals for σ with different methods:

a) Compute a simple percentile interval for σ with bootstrap. (2 p)

b) Compute a confidence interval based on the simplified pivotal variable $T = \hat{\sigma} - \sigma$ using bootstrap. (4 p)

c) Consider the pivotal variable

$$T = \frac{\hat{\sigma}^2}{\sigma^2}$$

The rationale behind this variable is that if the data were samples of $N(m, \sigma^2)$, then

$$\frac{\sum_{i=1}^{20} (X_i - \bar{X})^2}{\sigma^2} \text{ would be } \chi^2(19) \text{ distributed.}$$

Compute a 90% symmetric confidence interval for σ based on the chosen pivotal variable by bootstrap. (6 p)

Uppgift 4

Assume that the random variable X is $\text{Bin}(n, \theta)$ distributed and that θ is an outcome of the random variable Θ , which has the distribution $\text{Beta}(a, b)$. Show that the (posterior) distribution of $\Theta \mid X = x, x \in \{0, 1, \dots, n\}$, is $\text{Beta}(x + a, n - x + b)$. (10 p)

Uppgift 5

X_1, \dots, X_n, X_{n+1} are independent and identically distributed with the distribution function $F(x)$, which is continuous, i.e., has a density $f(x)$. We assume that

$$m = \int_{-\infty}^{\infty} x f(x) dx, \quad \sigma^2 = \int_{-\infty}^{\infty} (x - m)^2 f(x) dx$$

exist and are finite.

We are interested in finding an interval $I_{tol, 1-\alpha}$ of the form

$$I_{tol, 1-\alpha} := [L(X_1, \dots, X_n), U(X_1, \dots, X_n)]$$

such that

$$E[F(U(X_1, \dots, X_n)) - F(L(X_1, \dots, X_n))] = 1 - \alpha,$$

where $0 < \alpha < 1$.

a) Show that

$$E[F(U(X_1, \dots, X_n)) - F(L(X_1, \dots, X_n))] = P(X_{n+1} \in [L(X_1, \dots, X_n), U(X_1, \dots, X_n)]).$$

Then the interval $I_{tol, 1-\alpha}$ is known as a **(mean coverage) tolerance (or prediction) interval for X_{n+1} at confidence level $1 - \alpha$** . What is the interpretation of the interval so constructed? (2 p)

b) Let

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n)$$

be the arithmetic mean of the first n of the variables. Then

$$\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

is a well known estimator of the standard deviation based on the n first variables. We restrict ourselves from now on to intervals $I_{tol, 1-\alpha}$ of the form

$$\begin{aligned} L(X_1, \dots, X_n) &= \bar{X} + t_1 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} \\ U(X_1, \dots, X_n) &= \bar{X} + t_2 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} \end{aligned} \tag{1}$$

If $t_1 = -t_2$, then $I_{tol,1-\alpha}$ is a symmetric tolerance interval. Show now that we must choose t_1 and t_2 so that the equality

$$1 - \alpha = P \left(t_1 \leq \frac{X_{n+1} - \bar{X}}{\hat{\sigma} \sqrt{1 + \frac{1}{n}}} \leq t_2 \right). \quad (2)$$

is satisfied. (2 p)

c) Assume that X_1, \dots, X_n, X_{n+1} are independent and $N(m, \sigma^2)$ -distributed. We assume that $\sigma > 0$ is known, i.e., an estimate is not needed. Show that a symmetric tolerance interval for X_{n+1} with confidence level $1 - \alpha$ is given by

$$L(X_1, \dots, X_n) = \bar{X} - \lambda_{\alpha/2} \sigma \cdot \sqrt{1 + \frac{1}{n}}$$

$$U(X_1, \dots, X_n) = \bar{X} + \lambda_{\alpha/2} \sigma \cdot \sqrt{1 + \frac{1}{n}} \quad (3)$$

where $\alpha/2 = P(Z > \lambda_{\alpha/2})$ for $Z \in N(0, 1)$.

Aid: You can (and should) use the result (2) in b), even if you have not solved b). (4 p)

d) Suggest a symmetric tolerance interval for X_{n+1} with confidence level $1 - \alpha$ of the form (1), when X_1, \dots, X_n, X_{n+1} are independent and $N(m, \sigma^2)$ -distributed and $\sigma > 0$ is not known. (2 p)

Uppgift 6

We continue with construction of mean coverage tolerance (or prediction) intervals for X_{n+1} at confidence level $1 - \alpha$, when X_1, \dots, X_n, X_{n+1} are independent and identically distributed with the distribution function $F(x)$ with density $f(x)$ such that mean and variance exist.

B. Efron and R.J. Tibshirani introduced the notion of nonparametric tolerance (or prediction) intervals for X_{n+1} via the bootstrap t-method. For this we re-consider (2) and regard t_1 and t_2 as percentiles of the statistic $\frac{X_{n+1} - \bar{X}}{\hat{\sigma} \sqrt{1 + \frac{1}{n}}}$.

If the interval is not chosen as symmetric, then we choose $\alpha_1 \in [0, \alpha]$ and t_1 so that

$$P \left(\frac{X_{n+1} - \bar{X}}{\hat{\sigma} \sqrt{1 + \frac{1}{n}}} \leq t_1 \right) = \alpha_1$$

and t_2 so that using $\alpha_2 = 1 - (\alpha - \alpha_1)$ we have

$$P \left(\frac{X_{n+1} - \bar{X}}{\hat{\sigma} \sqrt{1 + \frac{1}{n}}} \leq t_2 \right) = \alpha_2.$$

Give an algorithm (a pseudo-code), which gives a nonparametric tolerance interval for X_{n+1} with the degree of confidence $1 - \alpha$ via the bootstrap t-method. What could be a criterion for selection of the tuning parameter α_1 ? (10 p)

SOLUTIONS TO THE EXAM IN SF2955 COMPUTER INTENSIVE METHODS IN MATHEMATICAL STATISTICS

WEDNESDAY 18th of AUGUST 2009 02.00 p.m. – 07.00 p.m..

Uppgift 1Notera att jackknife-stickproven består av $n - 1$ observationer vardera så

$$E(\widehat{\theta}_{(i)}) = \theta + \frac{a_1(F)}{n-1} + \frac{a_2(F)}{(n-1)^2} + \dots$$

Om vi nu studerar $E(\widehat{\theta}_{(\cdot)})$ erhålls

$$\begin{aligned} E(\widehat{\theta}_{(\cdot)}) &= E\left(\frac{1}{n} \sum_{i=1}^n \widehat{\theta}_{(i)}\right) = \frac{1}{n} \sum_{i=1}^n E(\widehat{\theta}_{(i)}) = \\ &= \theta + \frac{a_1(F)}{n-1} + \frac{a_2(F)}{(n-1)^2} + \dots \end{aligned}$$

som ger

$$E\left((n-1)(\widehat{\theta}_{(\cdot)} - \widehat{\theta})\right) = \frac{a_1(F)}{n} + O\left(\frac{1}{n^2}\right).$$

Vi får alltså

$$\widehat{bias}_{jack} = (n-1)(\widehat{\theta}_{(\cdot)} - \widehat{\theta}) \approx \frac{a_1(F)}{n} \approx bias_F(\widehat{\theta}, \theta)$$

Uppgift 2

a)

Resultatet inses om man summerar över \mathbf{x} eftersom man då erhåller för vänsterledet

$$\sum_{\mathbf{x}} \pi(\mathbf{x})P(\mathbf{x}, \mathbf{y})$$

som är fördelningen efter ett tidssteg i Markovkedjan då man startat den med fördelningen π (Jämför $\mathbf{p}^{(n+1)} = \mathbf{p}^{(n)}P$). För högerledet ger en summering över \mathbf{x}

$$\sum_{\mathbf{x}} \pi(\mathbf{y})P(\mathbf{y}, \mathbf{x}) = \pi(\mathbf{y}) \sum_{\mathbf{x}} P(\mathbf{y}, \mathbf{x}) = \pi(\mathbf{y})$$

eftersom P är en övergångsmatrix och alltså

$$\sum_{\mathbf{x}} P(\mathbf{y}, \mathbf{x}) = 1$$

vilket alltså innebär att π är stationär fördelning till kedjan.

b)

Vi har för $\mathbf{y} \neq \mathbf{x}$

$$P(\mathbf{x}, \mathbf{y}) = P(\text{att föreslå } \mathbf{y} | \text{start i } \mathbf{x}) P(\text{acceptera förslaget}) = q(\mathbf{x}, \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y})$$

Vi får då att för $\mathbf{y} \neq \mathbf{x}$

$$\begin{aligned} \pi(\mathbf{x}) P(\mathbf{x}, \mathbf{y}) &= \pi(\mathbf{x}) q(\mathbf{x}, \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{x}) q(\mathbf{x}, \mathbf{y}) \min \left(\frac{\pi(\mathbf{y}) q(\mathbf{y}, \mathbf{x})}{\pi(\mathbf{x}) q(\mathbf{x}, \mathbf{y})}, 1 \right) = \\ &= \min (\pi(\mathbf{y}) q(\mathbf{y}, \mathbf{x}), \pi(\mathbf{x}) q(\mathbf{x}, \mathbf{y})) \end{aligned}$$

och på samma sätt (byt \mathbf{x} och \mathbf{y})

$$\pi(\mathbf{y}) P(\mathbf{y}, \mathbf{x}) = \min (\pi(\mathbf{x}) q(\mathbf{x}, \mathbf{y}), \pi(\mathbf{y}) q(\mathbf{y}, \mathbf{x}))$$

och dessa är lika! Om $\mathbf{x} = \mathbf{y}$ är å andra sidan villkoret för tidsreversibilitet trivialt uppfyllt! Eftersom kedjan är tidsreversibel med avseende på π så har den också π till stationär fördelning enligt a-delen.

Uppgift 3

Vi uppskattar 5%- respektive 95%-percentilerna med det 50:de respektive 950:de i storleksordning som ger 13.5 respektive 30.5.

a) Det enkla percentilintervallet för σ^2 blir alltså (13.5, 30.5) dvs för σ får vi $(\sqrt{13.5}, \sqrt{30.5}) \approx (3.67, 5.52)$

b) Om vi väljer tal a respektive b så att

$$1 - \alpha = P(a \leq \hat{\sigma} - \sigma \leq b) = P(\hat{\sigma} - b \leq \sigma \leq \hat{\sigma} - a)$$

så vi skulle alltså få konfidensintervallet $(\hat{\sigma} - b, \hat{\sigma} - a)$ för σ .

Bootstraphypotesen innebär att vi tror att $a^* \approx a$ och $b^* \approx b$ där a^* och b^* uppfyller

$$1 - \alpha = P(a^* \leq \hat{\sigma}^* - \hat{\sigma} \leq b^*) = P(\hat{\sigma} + a^* \leq \hat{\sigma}^* \leq \hat{\sigma} + b^*)$$

och vi får alltså att $\hat{\sigma} + a^* = \sqrt{13.5}$ och $\hat{\sigma} + b^* = \sqrt{30.5}$ om vi bortser från att vi bara gjort ett ändligt antal bootstrap-stickprov.

Alltså får vi $a^* = \sqrt{13.5} - \sqrt{22.8095} \approx -1.10$ och $b^* = \sqrt{30.5} - \sqrt{22.8095} \approx 0.75$. Detta ger konfidensintervallet $(\sqrt{22.8095} - 0.75, \sqrt{22.8095} - (-1.10)) \approx (4.03, 5.88)$.

c) Med a och b uppfyllande

$$1 - \alpha = P \left(a \leq \frac{\hat{\sigma}^2}{\sigma^2} \leq b \right) = P \left(\frac{\hat{\sigma}^2}{b} \leq \sigma^2 \leq \hat{\sigma}^2 a \right) = P \left(\frac{\hat{\sigma}}{\sqrt{b}} \leq \sigma \leq \frac{\hat{\sigma}}{\sqrt{a}} \right)$$

fås alltså ett konfidensintervall för σ till

$$\left(\frac{\hat{\sigma}}{\sqrt{b}}, \frac{\hat{\sigma}}{\sqrt{a}} \right).$$

Vi skattar a och b med $a^* = 13.5/\hat{\sigma}^2 = 13.5/22.8095 \approx 0.59$ resp. $b^* = 30.5/\hat{\sigma}^2 = 30.5/22.8095 \approx 1.34$ och vi får konfidensintervallet

$$\left(\frac{\sqrt{22.8095}}{\sqrt{1.34}}, \frac{\sqrt{22.8095}}{\sqrt{0.59}} \right) \approx (4.13, 6.22)$$

Uppgift 4

Se kurskompendium sid. 187.

Uppgift 5

a) We start from the right hand side by conditioning

$$\begin{aligned} & P(X_{n+1} \in [L(X_1, \dots, X_n), U(X_1, \dots, X_n)]) \\ &= \int_{R^n} P\left(X_{n+1} \in [L(x_1, \dots, x_n), U(x_1, \dots, x_n)] \mid x_1, \dots, x_n\right) \prod_{i=1}^n f(x_i) dx_1 \dots dx_n \end{aligned}$$

and since X_1, \dots, X_n, X_{n+1} are independent

$$= \int_{R^n} P(X_{n+1} \in [L(x_1, \dots, x_n), U(x_1, \dots, x_n)]) \prod_{i=1}^n f(x_i) dx_1 \dots dx_n$$

and by definition of $F(x)$

$$= \int_{R^n} \{F(U(x_1, \dots, x_n)) - F(L(x_1, \dots, x_n))\} \prod_{i=1}^n f(x_i) dx_1 \dots dx_n$$

which is the expectation of a function of the random variables X_1, \dots, X_n , or

$$= E[F(U(X_1, \dots, X_n)) - F(L(X_1, \dots, X_n))].$$

The interpretation of the interval is that with probability $1-\alpha$ we obtain a sample X_1, \dots, X_n, X_{n+1} such that $X_{n+1} \in [L(X_1, \dots, X_n), U(X_1, \dots, X_n)]$. This is an interval for prediction of X_{n+1} .

b) From a) we get

$$\begin{aligned} P(X_{n+1} \in [L(X_1, \dots, X_n), U(X_1, \dots, X_n)]) &= P\left(\bar{X} + t_1 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} < X_{n+1} < \bar{X} + t_2 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}}\right) \\ &= P\left(t_1 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} \leq X_{n+1} - \bar{X} \leq t_2 \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}}\right) \\ &= P\left(t_1 \leq \frac{X_{n+1} - \bar{X}}{\hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}}} \leq t_2\right), \end{aligned}$$

as was to be shown.

c) We derive the distribution of $X_{n+1} - \bar{X}$. We have

$$E(X_{n+1} - \bar{X}) = E(X_{n+1}) - E(\bar{X}) = m - m = 0,$$

since $E(\bar{X}) = \frac{1}{n}(E(X_1) + \dots + E(X_n)) = \frac{1}{n} \cdot n \cdot m = m$. As X_1, \dots, X_n are independent $N(m, \sigma^2)$ variables, it holds that

$$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n) \in N(m, \sigma^2/n)$$

Since X_1, \dots, X_n, X_{n+1} are independent $N(m, \sigma^2)$, \bar{X} and X_{n+1} are independent random variables and we obtain

$$\text{Var}(X_{n+1} - \bar{X}) = \text{Var}(X_{n+1}) + \text{Var}(\bar{X}) = \sigma^2 + \frac{\sigma^2}{n} = \sigma^2 \left(1 + \frac{1}{n}\right).$$

Thus

$$\frac{X_{n+1} - \bar{X}}{\sigma \sqrt{1 + \frac{1}{n}}} \in N(0, 1).$$

Therefore by the case b) above we consider the equation

$$1 - \alpha = P\left(t_1 \leq \frac{X_{n+1} - \bar{X}}{\sigma \cdot \sqrt{1 + \frac{1}{n}}} \leq t_2\right)$$

and this equation is solved by $t_1 = -t_2, t_2 = \lambda_{\alpha/2}$, where $\alpha/2 = P(Z > \lambda_{\alpha/2})$ for $Z \in N(0, 1)$. This yields the tolerance interval as desired in (3).

d) We recall a basic theorem of mathematical statistics. If $X \in N(0, 1)$ and $Y \in \chi^2(f)$, and X and Y are independent, then

$$\frac{X}{\sqrt{Y/f}} \in t(n-1),$$

i.e., the ratio is a Student t-distributed random variable with $n-1$ degrees of freedom. Since

$$\frac{X_{n+1} - \bar{X}}{\sigma \sqrt{1 + \frac{1}{n}}} \in N(0, 1).$$

and, as is well known,

$$(n-1) \frac{\hat{\sigma}^2}{\sigma^2} \in \chi^2(n-1),$$

and as these are independent, we get that

$$\frac{X_{n+1} - \bar{X}}{\hat{\sigma} \sqrt{1 + \frac{1}{n}}} = \frac{X_{n+1} - \bar{X}}{\sigma \sqrt{1 + \frac{1}{n}} \sqrt{(n-1)\hat{\sigma}^2/(\sigma^2 \cdot (n-1))}}$$

is $t(n-1)$ distributed. Hence we take

$$L(X_1, \dots, X_n) = \bar{X} - t_{\alpha/2}(n-1)\hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}}$$

$$U(X_1, \dots, X_n) = \bar{X} + t_{\alpha/2}(n-1)\hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}}$$

(4)

where $\alpha/2 = P(Z > t_{\alpha/2}(n-1))$, where $Z \in t(n-1)$.

Uppgift 6

We have an I.I.D sample x_1, x_2, \dots, x_n from F . An algorithm is as follows

Step 1. Choose B , the number of bootstrap replicates.

Step 2. Generate B independent bootstrap samples $x_i^*(b)$, $i = 1, \dots, n$, $b = 1, 2, \dots, B$, by drawing with replacement from x_1, x_2, \dots, x_n .

Step 3. Generate B independent bootstrap samples $x_{n+1}^*(b)$, $b = 1, 2, \dots, B$ by drawing with replacement from x_1, x_2, \dots, x_n .

Step 4. For $b = 1, 2, \dots, B$ compute

$$\bar{x}^*(b) = \frac{1}{n} \sum_{i=1}^n x_i^*(b)$$

and

$$\hat{\sigma}^2(b) = \frac{1}{n-1} \sum_{i=1}^n (x_i^*(b) - \bar{x}^*(b))^2.$$

Step 5. Evaluate for $b = 1, 2, \dots, B$ the bootstrap t -statistic

$$T^*(b) = \frac{x_{n+1}^*(b) - \bar{x}^*(b)}{\hat{\sigma}(b) \cdot \sqrt{1 + \frac{1}{n}}}$$

Step 6. Derive from the histogram of $T^*(b)$ the quantiles t_1^* and t_2^* and select α_1 , for example, so that the width $t_2^* - t_1^*$ is minimized.

Step 7. The bootstrapped tolerance interval is

$$\begin{aligned} L(X_1, \dots, X_n) &= \bar{X} + t_1^* \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} \\ U(X_1, \dots, X_n) &= \bar{X} + t_2^* \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n}} \end{aligned} \tag{5}$$