





Chapter 1

Stochastic vectors and matrices

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

- The vector X = (X₁,..., X_p)' is called a random vector or p-dimensional random variable, if the components X₁,..., X_p are one dimensional random variables.
- The vector X is called continuous if there is a function $f(x) = f(x_1, ..., x_p) \ge 0$ such that

$$P(a_1 \leq X_1 \leq b_1, \ldots, a_p \leq X_p \leq b_p) = \int_{a_p}^{b_p} \ldots \int_{a_1}^{b_1} f(x_1, \ldots, x_p) dx_1 \ldots dx_p.$$

The function f is called (joint) probability density function (p.d.f.) of X.

• The random vector X is called discrete, if X has only values in a finite or countable set $\{x_1, x_2, \ldots\} \subset \mathbb{R}^p$. The function f with

$$f(x) = \begin{cases} P(X = x) & x \in \{x_1, x_2, \ldots\} \\ 0 & \text{else} \end{cases}$$

is called probability function or discrete p.d.f. of X.

Continuous random vector

Consider the 2-dimensional continuous random vector $x = (x_1, x_2)'$ with pdf

$$f(x_1, x_2) = \begin{cases} 0.8(x_1 + x_2 + x_1 x_2) & 0 \le x_1 \le 1 \\ 0 \le x_2 \le 1 \\ 0 & \text{else.} \end{cases}$$

Discrete random vector

Consider the two dimensional random vector (X, Y)' with

X = proseminar grade Y = final exam grade

We have the following distribution:

Y/X	1	2	3	4
1	50	40	12	6
	1007	1007	1007	1007
2	<u>55</u>	<u>97</u>	<u>54</u>	<u>10</u>
	1007	1007	1007	1007
3	<u>28</u>	<u>100</u>	<u>68</u>	<u>33</u>
	1007	1007	1007	1007
4	<u>13</u>	79	75	36
	1007	1007	1007	1007
5	9	44	119	79
	1007	1007	1007	1007

- Let the *p*-dimensional random vector X = (X₁,..., X_p)' be partitioned into the p₁-dimensional vector X₁ and the p₂-dimensional Vector X₂, i.e. X = (X'₁, X'₂)'.
- The p_1 -dimensional p.d.f. or probability function $f_{X_1}(x_1)$ of X_1 is then called marginal p.d.f. or marginal probability function of X. It is given by

$$f_{X_1}(x_1) = \int_{-\infty}^{\infty} \ldots \int_{-\infty}^{\infty} f(x_1, x_2) dx_{p_1+1} \ldots dx_p$$

for continuous random vectors, and

$$f_{X_1}(x_1) = \sum_{x_2} f(x_1, x_2)$$

for discrete random vectors.

 The conditional p.d.f. or probability function of X₁ given X₂ = x₂ is defined as

$$f(x_1|x_2) = \begin{cases} \frac{f(x_1, x_2)}{f_{X_2}(x_2)} & \text{for } f_{X_2}(x_2) > 0\\ 0 & \text{else.} \end{cases}$$

The marginal and conditional p.d.f.'s or probability functions for X_2 are defined in complete analogy.

Continuous random vector

Consider the 2-dimensional continuous random vector $x = (x_1, x_2)'$ with pdf

$$f(x_1, x_2) = \begin{cases} 0.8(x_1 + x_2 + x_1 x_2) & 0 \le x_1 \le 1 \\ 0 \le x_2 \le 1 \\ 0 & \text{else.} \end{cases}$$

Discrete random vector								
	Y/X	1	2	3	4			
	1	<u>50</u> 1007	40 1007	<u>12</u> 1007	<u>6</u> 1007			
	2	<u>55</u> 1007	<u>97</u> 1007	<u>54</u> 1007	<u>10</u> 1007			
	3	<u>28</u> 1007	<u>100</u> 1007	<u>68</u> 1007	<u>33</u> 1007			
	4	<u>13</u> 1007	79 1007	75 1007	<u>36</u> 1007			
	5	<u>9</u> 1007	<u>44</u> 1007	<u>119</u> 1007	<u>79</u> 1007			

Expectation or mean vector

Let $X = (X_1, \ldots, X_p)'$ be a *p*-dimensional random vector. Then

$$E(X) = \boldsymbol{\mu} = (\mu_1, \ldots, \mu_p)' = (E(X_1), \ldots, E(X_p))'$$

is called mean vector of X.

Example

Covariance and correlation matrix

The covariance matrix $Cov(X) = \Sigma$ of a *p*-dimensional random vector *X* is defined as

$$Cov(X) = \mathbf{\Sigma} = E(X - \mu)(X - \mu)' = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1p} \\ \vdots & & \vdots \\ \sigma_{p1} & \dots & \sigma_{pp} \end{pmatrix},$$

where $\sigma_{ij} = Cov(X_i, X_j)$, $i \neq j$, is the covariance between X_i and X_j , and $\sigma_{ii} = \sigma_i^2 = Var(X_i)$ is the variance of X_i .

Covariance and correlation matrix

The correlation matrix \boldsymbol{R} of X is defined as

$$\boldsymbol{R} = \begin{pmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \vdots & & \vdots \\ \rho_{p1} & \rho_{p2} & \dots & 1 \end{pmatrix},$$

where

$$\rho_{ij} = \frac{\operatorname{Cov}(X_i, X_j)}{\sqrt{\operatorname{Var}(X_i) \cdot \operatorname{Var}(X_j)}}.$$

Covariance and correlation matrix

Example

Properties of expectations and covariance matrices

Let X and Y be random vectors and **A**, **B**, *a*, *b* matrices and vectors.

1
$$E(X + Y) = E(X) + E(Y)$$

$$E(\mathbf{A}X+b) = \mathbf{A} \cdot E(X) + b$$

3
$$Cov(X) = E(XX') - \mu\mu'$$

• Var
$$(a'X) = a'Cov(X)a = \sum_{i=1}^{k} \sum_{j=1}^{k} a_i a_j \sigma_{ij}$$

The covariance matrix is symmetric and positive semi definite.

•
$$Cov(AX + b) = ACov(X)A'$$

Multivariate Normal Distribution

A continuous *p*-dimensional random vector X = (X₁, X₂,..., X_p)' is said to have a multivariate normal (or Gaussian) distribution if it has p.d.f.

$$f(x) = (2\pi)^{-\frac{p}{2}} |\mathbf{\Sigma}|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(x-\mu)'\mathbf{\Sigma}^{-1}(x-\mu)\right]$$

with $\mu \in {\rm I\!R}^{\rho}$ and positive definite $(\rho \times \rho)$ -matrix ${\boldsymbol \Sigma}$.

- It can be shown that $E(X) = \mu$ and $Cov(X) = \Sigma$.
- We write

 $X \sim N_p(\mu, \mathbf{\Sigma}),$

The special case $\mu = \mathbf{0}$ and $\mathbf{\Sigma} = \mathbf{I}$ is called the (multivariate) standard normal distribution.

Multivariate Normal Distribution

Example $X = (X_1, X_2, X_3, X_4)' \sim N(\mu, \Sigma)$ with $\mu = (1, 2, 3, 4)' \qquad \Sigma = \begin{pmatrix} 3 & 1 & 0 & 1 \\ 1 & 4 & 2 & 0 \\ 0 & 2 & 5 & 2 \\ 1 & 0 & 2 & 4 \end{pmatrix}$

• Let the multivariate normal random variable $X \sim N(\mu, \Sigma)$ be partitioned into the subvectors $Y = (X_1, \dots, X_r)'$ and $Z = (X_{r+1}, \dots, X_p)'$, i.e.

$$X = \begin{pmatrix} Y \\ Z \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu_Y \\ \mu_Z \end{pmatrix}, \quad \mathbf{\Sigma} = \begin{pmatrix} \mathbf{\Sigma}_Y & \mathbf{\Sigma}_{YZ} \\ \mathbf{\Sigma}_{ZY} & \mathbf{\Sigma}_Z \end{pmatrix}.$$

- Then *Y* has an *r*-dimensional normal distribution $Y \sim N(\mu_Y, \Sigma_Y)$.
- The conditional distribution of *Y* given *Z* is again multivariate normal with mean

$$\mu_{Y|Z} = \mu_Y + \boldsymbol{\Sigma}_{YZ} \cdot \boldsymbol{\Sigma}_Z^{-1} (Z - \mu_Z)$$

and covariance matrix

$$\boldsymbol{\Sigma}_{Y|Z} = \boldsymbol{\Sigma}_{Y} - \boldsymbol{\Sigma}_{YZ} \boldsymbol{\Sigma}_{Z}^{-1} \boldsymbol{\Sigma}_{ZY}.$$

- Furthermore, Y and Z are independent if and only if Y and Z are uncorrelated, i.e. if Σ_{YZ} = Σ_{ZY} = 0.
- The equivalence is generally not true for non-normal random vectors: If Y and Z are independent they are also uncorrelated, but in general Σ_{ZY} = 0 does not imply independence.

Example $Y = (X_1, X_2)', Z = (Z_1, Z_2)'$

Linear transformations

Assume $X \sim N_{\rho}(\mu, \mathbf{\Sigma})$ is multivariate normal. Then the linear transformation

 $Y = \mathbf{D}X + d$

with the $m \times p$ matrix **D** and the $m \times 1$ vector **d** is again multivariate normal

 $Y \sim N_m(D\mu + d, D\Sigma D').$







Statistical Inference

Chapter 2

Classical Inference

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Statistical Inference - 2 - Classical Inference - 0 / 108

Classical Inference

Situation

Situation

We consider a **random variable** *X* (discrete or continuous), whose **distribution** P_{θ} depends on an unknown **parameter** θ . Two main goals:

- Estimate the unknown parameter θ using statistical inference.
- Assess the uncertainty of the estimate.

Bernoulli Distribution / Binomial Distribution

Let X be binary with possible values 0 or 1 and probability function

•
$$f(1) = P(X = 1) = \pi$$
,

•
$$f(0) = P(X = 0) = 1 - \pi$$

or more compactly

$$f(x) = egin{cases} \pi^x (1-\pi)^{1-x} & x \in \{0,1\} \ 0 & ext{else.} \end{cases}$$

Here $\theta = \pi = P(X=1) = E(X)$.

Connection to the Binomial distribution

Let $X_1, X_2, ..., X_n$ be i.i.d. **Bernoulli-distributed** random variables with $P(X_i=1) = \pi$, $P(X_i=0) = 1-\pi$; Then

$$X = X_1 + \dots + X_n$$

= 'frequency with which an event E occurs'
 $\sim B(n, \pi)$

The probability function is given by

$$f(x) = \mathsf{P}(X = x) = \begin{cases} \binom{n}{x} \pi^{x} (1 - \pi)^{n - x} & x = 0, 1, \dots, n \\ 0 & \text{else.} \end{cases}$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Some properties of the Binomial distribution

- a) Probability generating function $G_x(t) = E(t^x) = (1 - \pi + t\pi)^n$
- b) Expected Value and Variance

$$E(X) = n\pi$$
$$Var(X) = n\pi(1 - \pi)$$

c) Sum

Let *X*, *Y* be independent with $X \sim B(n,\pi)$ and $Y \sim B(m,\pi)$. Then

$$G_{x+y}(t) = G_x(t)G_y(t) = (1 - \pi + t\pi)^{n+m}$$

This is the probability generating function of the $B(n+m, \pi)$ distribution,

i.e.
$$x+y \sim B(n+m, \pi)$$
.

Proof a)

$$G_X(t) = \mathsf{E}(t^X) = \sum_{x=0}^n t^x \binom{n}{x} \pi^x (1-\pi)^{n-x}$$
$$= \sum_{x=0}^n \binom{n}{x} (t\pi)^x (1-\pi)^{n-x}$$
$$= (1-\pi+t\pi)^n$$

Use here

$$(z+y)^n = \sum_{x=0}^n \binom{n}{x} y^x z^{n-x}$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Proof b)

$$G_X(t) = (1 - \pi + t\pi)^n$$

$$G'_X(t) = n\pi(1 - \pi + t\pi)^{n-1}$$

$$G''_X(t) = n\pi^2(n-1)(1 - \pi + t\pi)^{n-2}$$

$$E(X) = G'_X(1) = n\pi(1 - \pi + \pi)^{n-1} = n\pi$$

$$Var(X) = G''_X(1) + G'_X(1) - (G'_X(1))^2$$

$$= n\pi^2(n-1) + n\pi - n^2\pi^2$$

$$= n\pi(\pi(n-1) + 1 - n\pi)$$

$$= n\pi(1 - \pi)$$

Proof c)

- Let *X*, *Y* be independent with $X \sim B(n, \pi)$ and $Y \sim B(m, \pi)$.
- Then

$$G_{X+Y}(t) = G_X(t)G_Y(t) = (1 - \pi + t\pi)^{n+m}$$

This is the probability generating function of the B(n + m, π) distribution, i.e. X + Y ~ B(n + m, π).

Poisson distribution

Consider the **Poisson distributed** random variable X with probability function

$$P_{\lambda}(X = x) = f_{\lambda}(x) = \begin{cases} rac{\lambda^{x}}{x!} \exp(-\lambda) & x = 0, 1, 2, \dots \\ 0 & \text{else} \end{cases}$$

We have $E(X) = \lambda$, $Var(X) = \lambda$ (see exercise). The goal is to estimate λ , *i.e.* $\theta = \lambda$.

Requirement for *X* **being Poisson distributed**

- Two events can not happen at the same time.
- The probability that an event happens within a small time interval of length Δt is approximately $\lambda \Delta t$.
- The probability that a certain number of events happen within a time interval depends on the length of the interval but not on the specific location at the time axis.
- The number of events in two disjunct time intervals are independent.

Examples

Some examples for Poisson distributions:

Radioactive decay

The duration for an atom nucleus to decay is **Exponentially distributed**. The radioactive half life is the period after that half of the atom nuclei are decayed on average. The number of decays per time unit is **Poisson distributed**.

• Rice grains

Suppose we distribute a number *N* of rice grains randomly to a number of squares. Then the number of rice grains at the squares is (approximately) **Poisson distributed** with $\lambda = \frac{N}{n}$.

Normal distribution

Let $X \sim N(\mu, \sigma^2)$ with parameter μ and σ^2 .

If $\sigma^{\rm 2}$ is known then

$$\theta = \mu$$
,

otherwise we have

$$\theta = \binom{\mu}{\sigma^2}.$$

Arbitrary distribution

Let *X* be an **arbitrary distribution** with existing mean $E(X) = \mu$ and variance $Var(X) = \sigma^2$.

The goal is to estimate μ (and possibly σ^2), i.e. $\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$.

Special case: population

Suppose we are given a population of size *N* with elements y_1, y_2, \ldots, y_N . Aim is to estimate the population mean and variance

$$\mu = \frac{1}{N} \sum_{i=1}^{N} y_i,$$
$$\sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (y_i - \bar{y})^2.$$

Then X ='randomly sampled element of the population' is a random variable with

$$\mathsf{E}(X) = \frac{1}{N}(y_1 + \cdots + y_N) = \bar{y} = \mu$$

and

$$\operatorname{Var}(X) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \bar{y})^2 = \sigma^2.$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Random sample

In order to estimate the unknown parameter θ , we draw an i.i.d. (independent and identically distributed) random sample X_1, X_2, \ldots, X_n , *i.e.* the X_i , $i=1,\ldots,n$, are independent and identically distributed as X.

To estimate θ , we use a so-called **statistic**

$$\hat{\theta} = T(X_1,\ldots,X_N),$$

which is an appropriate function of the sample.
Bernoulli distribution

$$\hat{\pi} =$$
 'relative frequency of ones'
 $= \frac{1}{n}(X_1 + \dots + X_n)$
 $= \bar{X},$
i.e.
 $T(X_1, \dots, X_n) = \bar{X}.$

Poisson distribution

$$\hat{\lambda} = \bar{X}$$

or (because $\lambda = Var(X)$)

$$\hat{\lambda} = \tilde{S}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

or

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

Normal distribution

$$\hat{\mu} = ar{X}$$

 $\hat{\sigma}^2 = ar{S}^2$ or $\hat{\sigma}^2 = S^2$

Arbitrary distribution

$$\hat{\mu} = ar{X}$$

 $\hat{\sigma}^2 = ar{S}^2$ or $\hat{\sigma}^2 = S^2$

Gambler's ruin

- Gambler plays a series of games against the casino,
- Bet 1€,
- Probability of winning the game *p* = 0.52,
- Gambler's stake 5€, casino's stake 50€.

Goal is to find the probability π of ruin! Here we estimate the probability via simulation. Simulate *n* game series and define

$$X_i = \begin{cases} 1 & \text{gambler is ruined in } i\text{-th series} \\ 0 & \text{gambler is not ruined.} \end{cases}$$

Gambler's ruin (continued)

We have

$$X_i \sim B(1,\pi), \quad P(X_i = 1) = \pi, \quad P(X_i = 0) = 1 - \pi.$$

Estimate π through

$$\hat{\pi}=\frac{1}{n}(X_1+\cdots+X_n)=\bar{X}.$$

Question: How large should we choose *n* for the estimate to be precise enough?

Classical Inference

Methods of evaluating estimators

Biased and unbiased estimators

Consider the expected value $E(\hat{\theta})$ of the estimator $\hat{\theta}$ for θ . The **estimator** $\hat{\theta}$ **is unbiased for** θ if

$$E(\hat{ heta}) = heta$$

for <u>all</u> possible parameters θ . Otherwise the **estimator is biased** with

$$Bias(\hat{ heta}) = E(\hat{ heta}) - heta.$$

Frequency interpretation

If we repeat the estimation method several times with realized estimates $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_m$ then

$$\frac{1}{m}(\hat{ heta}_1+\cdots+\hat{ heta}_m)pprox heta_n$$

Frequency interpretation



Figure: Illustration of an unbiased (left) and a biased (right) estimation method.

Arbitrary distribution

We have

$$E(\bar{X}) = E\left(\frac{1}{n}(X_1 + \dots + X_n)\right)$$
$$= \frac{1}{n}(E(X_1) + \dots + E(X_n))$$
$$= \frac{1}{n}(\mu + \dots + \mu)$$
$$= \frac{1}{n}n\mu$$
$$= \mu,$$

i.e. $\hat{\mu} = \bar{X}$ is unbiased for the expected value of a distribution.

Arbitrary distribution (continued)

For the variance we have

$$\mathsf{E}(S^2) = \mathsf{E}\left(\frac{1}{n-1}\sum_{i=1}^n (X_i - \bar{X})^2\right) = \sigma^2,$$

i.e. $\hat{\sigma}^2=S^2$ is unbiased for σ^2 and $\hat{\sigma}^2=\tilde{S}^2$ is biased for σ^2 as

$$E(\tilde{S}^2) = E\left(\frac{1}{n}\sum_{i=1}^n (X_i - \bar{X})^2\right)$$
$$= E\left(\frac{n-1}{n(n-1)}\sum_{i=1}^n (X_i - \bar{X})^2\right)$$
$$= E\left(\frac{n-1}{n}S^2\right) = \frac{n-1}{n}\sigma^2 \neq \sigma^2.$$

Proof $E(S^2)$

We have

$$\begin{split} \sum_{i=1}^{n} (X_i - \mu)^2 &= \sum_{i=1}^{n} \left((X_i - \bar{X}) + (\bar{X} - \mu) \right)^2 \\ &= \sum_{i=1}^{n} (X_i - \bar{X})^2 + 2 \sum_{i=1}^{n} (X_i - \bar{X}) (\bar{X} - \mu) + \sum_{i=1}^{n} (\bar{X} - \mu)^2 \\ &= (n-1)S^2 + n(\bar{X} - \mu)^2 \end{split}$$

Proof continued

Note

$$\sum_{i=1}^{n} (X_i - \bar{X}) (\bar{X} - \mu) = (\bar{X} - \mu) \sum_{i=1}^{n} (X_i - \bar{X})$$
$$= (\bar{X} - \mu) \left(\sum_{i=1}^{n} X_i - \sum_{i=1}^{n} \bar{X} \right)$$
$$= (\bar{X} - \mu) \left(\sum_{i=1}^{n} X_i - n\bar{X} \right)$$
$$= (\bar{X} - \mu) \left(\sum_{i=1}^{n} X_i - \sum_{i=1}^{n} X_i \right)$$
$$= 0$$

Proof continued

Applying expected values on both sides yields

$$\mathsf{E}\left(\sum_{i=1}^{n} (X_i - \mu)^2\right) = (n-1)\mathsf{E}(S^2) + n\underbrace{\mathsf{E}(\bar{X} - \mu)^2}_{\operatorname{Var}(\bar{X})}$$

Because of

$$\mathsf{E}\left(\sum_{i=1}^{n} \left(X_{i}-\mu\right)^{2}\right) = \sum_{i=1}^{n} \mathsf{E}\left(X_{i}-\mu\right)^{2} = n\sigma^{2}$$

we obtain

$$n\sigma^{2} = (n-1)\mathsf{E}(S^{2}) + n\mathsf{Var}(\bar{X}) = (n-1)\mathsf{E}(S^{2}) + \frac{n\sigma^{2}}{n}$$

and therefore

$$\mathsf{E}(S^2) = \sigma^2.$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Remark

Note that $S = \sqrt{S^2}$ is not an unbiased estimator for σ ! Indeed we have

$$\mathsf{E}(S) = \mathsf{E}(\sqrt{S^2}) < \sqrt{\sigma^2} = \sigma$$

because of Jensen's inequality ($\sqrt{\cdots}$ concave).

Comparison of unbiased estimators using the variance

Two (or more) unbiased estimators can be compared with the variance of the estimators. An estimator $\hat{\theta}_1$ is at least as powerful as another estimator $\hat{\theta}_2$ if

$$\operatorname{Var}(\hat{ heta}_1) \leq \operatorname{Var}(\hat{ heta}_2)$$

for <u>all</u> possible values of θ .

Frequency interpretation



Figure: Illustration of the variance of estimators.

Arbitrary distribution

The variance of $\hat{\mu} = \bar{X}$ is given by

V

$$\operatorname{ar}(\bar{X}) = \operatorname{Var}\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}\right)$$
$$= \frac{1}{n^{2}}\left(\sum_{i=1}^{n}\operatorname{Var}X_{i}\right)$$
$$= \frac{1}{n^{2}}n\sigma^{2}$$
$$= \frac{\sigma^{2}}{n} = \frac{\operatorname{Var}(X)}{n}$$

Arbitrary distribution (continued)

In the Bernoulli case we have $Var(X) = \pi(1 - \pi)$ and therefore

$$\operatorname{Var}(\hat{\pi}) = \operatorname{Var}(\bar{X}) = \frac{\operatorname{Var}(X)}{n} = \frac{\pi(1-\pi)}{n}$$



Mean squared error (MSE)

The **mean squared error (MSE)** can be used to compare competing (not necessarily unbiased) estimators. It is defined as

$$\mathsf{MSE}_{ heta}(\hat{ heta}) = \mathsf{E}(\hat{ heta} - heta)^{\mathsf{2}}$$

Justification for the MSE

a) We have the decomposition

$$\mathsf{MSE}(\hat{ heta}) = \mathsf{Var}(\hat{ heta}) + \mathsf{Bias}^2(\hat{ heta})$$

b) For arbitrary $\epsilon > 0$ we have

$$\mathsf{P}\big(\big|\hat{\theta} - \theta\big| < \epsilon\big) > 1 - \frac{1}{\epsilon^2}\mathsf{E}(\hat{\theta} - \theta)^2 = 1 - \frac{1}{\epsilon^2}\mathsf{MSE}(\hat{\theta}),$$

i.e. the 'coverage probability' becomes larger as $MSE(\hat{\theta})$ decreases!

Proof a)

$$MSE_{\theta}(\hat{\theta}) = E(\hat{\theta} - \theta)^{2}$$

$$= E[(\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)]^{2}$$

$$= E(\hat{\theta} - E(\hat{\theta}))^{2} + 2E(\hat{\theta} - E(\hat{\theta}))\underbrace{(E(\hat{\theta}) - \theta)}_{0} + E(E(\hat{\theta}) - \theta)^{2}$$

$$= Var(\hat{\theta}) + Bias^{2}(\hat{\theta})$$

Proof b)

Define the random variable

$$V = \begin{cases} \epsilon^2 & \left| \hat{\theta} - \theta \right| \geq \epsilon \\ 0 & \text{else.} \end{cases}$$

Then

$$\mathsf{E}(\mathsf{V}) = \epsilon^2 \mathsf{P}(\left|\hat{\theta} - \theta\right| \ge \epsilon)$$

and because of V $\leq (\hat{ heta} - heta)^2$

$$\mathsf{E}(V) \leq \mathsf{E}(\hat{\theta} - \theta)^2.$$

Hence

$$\epsilon^2 \, \mathsf{P}(\left| \hat{ heta} - heta
ight) ig| \geq \epsilon) \leq \mathsf{E}(\hat{ heta} - heta)^2 = \mathsf{MSE}(\hat{ heta})$$

and therefore

$$\mathsf{P}(ig|\hat{ heta}- hetaig|)<\epsilon)=\mathsf{1}-\mathsf{P}(ig|\hat{ heta}- hetaig|\geq\epsilon)\geq\mathsf{1}-rac{\mathsf{1}}{\epsilon^2}\mathsf{MSE}(\hat{ heta})$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Frequency interpretation



Figure: Illustration of the MSE of estimators.

Arbitrary distribution / normal distribution

$$\mathsf{MSE}(\hat{\mu}) = \mathsf{Var}(\hat{\mu}) = rac{\sigma^2}{n}$$

Bernoulli distribution

$$\mathsf{MSE}(\hat{\pi}) = \frac{\pi(1-\pi)}{n}$$

Poisson distribution

$$MSE(\hat{\lambda}_1) = \frac{\lambda}{n}$$

An unusual estimation problem

Suppose we are given a lake and we want to estimate the number N of fish in the lake. Estimation approach:

1. Draw a sample of size n_1 , mark the fish and return the marked fish into the lake. Denote by

$$\pi = \frac{n_1}{N} \tag{1}$$

the share of marked fish

An unusual simulation problem (continued)

2. Draw a second sample of size n_2 and estimate the share of marked fish through

$$\hat{\pi} = \frac{X}{n_2},$$

where *X* is the number of marked fish in the sample. Insert $\hat{\pi}$ into (1) and solve for *N* to obtain

$$\hat{N} = \frac{n_1 n_2}{X}$$

(Although $\hat{\pi}$ is unbiased for π , \hat{N} is not unbiased for N)

Case study: Doping - a rare phenomenon?

According to an article published in the "Süddeutsche Zeitung" on 23/12/2006:

Cool, cheerful, devotional, deceptive

A recent study has revealed that nearly every second German top athlete uses unauthorized doping substances.

"They used the well-known questioning technique RRT, which includes special instructions in addition to the normal questions. This method allows further probability calculations."

"It was found out that 48.1% of all German top athletes are using doping."

Initial situation

Sensitive yes-no questions

"Have you ever used unauthorized performance-enhancing drugs?"

"Have you ever made false statements in a tax declaration, for your own benefit?"

. . .

. . .

Objective

Estimate the number of doped athletes.

Estimate the number of taxpayers who have made false statements.

Sample

- Draw a random sample of i = 1, ..., n persons.
- Define for each person the random variables

$$A_{i} = \begin{cases} 1 & i^{\text{th}} \text{ person answers with YES} \\ 0 & i^{\text{th}} \text{ person answers with NO} \end{cases}$$

and

$$R_{i} = \begin{cases} 1 & i^{\text{th}} \text{ person dopes} \\ 0 & i^{\text{th}} \text{ person does not dope} \end{cases}$$

• A_i and R_i follow a Bernoulli-distribution

$$P(A_i = 1) = \pi_A = P(YES-answer)$$
 $P(A_i = 0) = 1 - \pi_A$
and

$$\mathsf{P}(\mathsf{R}_i=\mathsf{1})=\pi_\mathsf{R}=\mathsf{P}(\mathsf{Athlete\ dopes})$$
 $\mathsf{P}(\mathsf{R}_i=\mathsf{0})=\mathsf{1}-\pi_\mathsf{R}.$

Honest answers

- Assuming that we have honest answers, the share π_A of YES-answers and the share π_R of doped athletes are equal, i.e. $\pi_A = \pi_R$.
- The relative frequency

$$\hat{\pi}_R = \bar{A} = \frac{1}{n}(A_1 + A_2 + \ldots + A_n)$$

of YES-answers is an unbiased estimator for the share π_R of doped athletes.

• Expected value and variance of $\hat{\pi}_{R} = \bar{A}$ are given by

$$E(\bar{A}) = \pi_R$$
 $Var(\bar{A}) = \frac{1}{n}\pi_R(1-\pi_R).$

Partially honest answers

• For partially false answers:

$$\begin{split} \mathsf{P}(\mathsf{YES}\text{-answer} \,|\, \mathsf{Athlete \ dopes}) &= \mathsf{P}(A_i = 1 | R_i = 1) = q < 1 \\ \mathsf{P}(\mathsf{YES}\text{-answer} \,|\, \mathsf{Athlete \ does \ not \ dope}) &= \mathsf{P}(A_i = 1 | R_i = 0) = 0. \end{split}$$

hence

 $P(YES-answer) = P(A_i = 1) = P(A_i = 1|R_i = 1)P(R_i = 1) = q \cdot \pi_R.$

• Therefore the relative frequency \overline{A} of YES-answers is a biased estimator for the share π_R of doped athletes with

$$\mathsf{E}(\bar{\mathsf{A}}) = \mathsf{q} \cdot \pi_{\mathsf{R}} < \pi_{\mathsf{R}}.$$

The share of doped athletes is underestimated!

Questioning techniques

- Traditional interviews
- Questionnaire
- Computer-based questionnaire
- Randomized Response

Idea: Guarantee anonymous answers with the help of an additional experiment.

A simple RR-Model

- Every interviewed person tosses or rotates a fair coin, *unobserved* by the interviewer. (Note: a rotated 2 € coin is not fair!)
- If the coin shows "Head" the interviewed person should answer honestly.
- If the coin shows "Tail" the interviewed person must always answer with YES, independent of whether he/she has doped or not.
- The additional experiment guarantees the anonymity, since a YES-answer does not conclude that the person was doped.
- Assumption: The interviewees stick to the rules and answer correctly.

- Even if the answers are only partially honest, i.e. $\pi_A = q\pi_R < \pi_R$, we can make an unbiased estimate of the share π_A of YES-answers through the relative frequency \bar{A} , i.e. $\hat{\pi}_A = \bar{A}$.
- For the probability of YES-answers π_A the following is true

$$\pi_A = P(\text{YES-answer}) = \frac{1}{2}\pi_R + \frac{1}{2} \cdot 1 = \frac{1}{2}\pi_R + \frac{1}{2}$$

and therefore

$$\pi_R = 2\pi_A - 1$$

is true for the probability π_R of doping.

Hence:

$$\hat{\pi}_R = 2\bar{A} - 1$$

is an unbiased estimator for the share of doped athletes.

• The variance is given by

$$Var(\hat{\pi}_R) = \frac{1}{n}\pi_R(1-\pi_R) + \frac{1}{n}(1-\pi_R)$$

= Variance without randomization $+\frac{1}{n}(1-\pi_R)$

 In comparison to the situation of honest answers, the variance of the estimator increases!
The quality of RR in real life - Social insurance fraud

- A study of Heijden et al. (2000) concerning social insurance fraud compares RR, computer-based questionnaires (CASI) and traditional interviews.
- The characteristic is that only persons who evidently committed social insurance fraud were interviewed. The share of social insurance frauds in the sample is thus 100%! Neither the interviewees nor the interviewers were aware of this fact.
- Result: RR results in an estimated share of 43-49%, CASI leads to 19% and traditional interviews lead to 25%.
- RR significantly improves the results, nevertheless there is still a large share of dishonest answers.

Classical Inference

Asymptotic properties of estimators

Random variables as mappings

Consider a random experiment with possible outcome in the **space** Ω and a corresponding **probability measure** *P*.

A random variable X is a mapping that assigns every $\omega \in \Omega$ a real value x, i.e. $X(\omega)=x$, More specifically

$$egin{aligned} X &: \Omega o \mathbb{R} \ & \omega \mapsto X(\omega) = x \end{aligned}$$

Rolling a dice twice

$$\Omega = ig\{(1,1),(1,2),\ldots,(6,6)ig\} \ X((\omega_1,\omega_2)) = \omega_1 + \omega_2$$

e.g.

$$X((3,4)) = 3 + 4 = 7$$

P(X = 3) = P({(1,2), (2,1)}) = $\frac{2}{36}$,

More precise definition of a random variable

Let Ω be the set of possible outcomes of a random experiment. An **event** is a subset of Ω .

In general it is not possible to assign probabilities to every subset of Ω in a consistant way. Instead we assign probabilities to a system of subsets *F*, called σ -field. A σ -field *F* is a set of subsets of Ω such that

i)
$$\Omega \in F$$

ii)
$$A \in F \Rightarrow \overline{A} \in F$$

iii) A_1, A_2, \dots in $F \Rightarrow \bigcup_{i=1}^{\infty} A_i \in F$

We call the **couple** (Ω ,*F*) a measurable space.

More precise definition of a random variable (cont'd.)

A probability measure *P* is a mapping

$$P: F \rightarrow [0, 1]$$

such that

i) $P(\emptyset)=0$

ii) $P(\Omega)=1$

iii) $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$ for sets A; with $A_i \cap A_j = \emptyset$ for all i, j.

The space (Ω , *F*, *P*) is called probability space.

More precise definition of a random variable (cont'd.)

Let (Ω, F, P) be a probability space and (Ω', F') a measurable space. A **random variable** *X* is a *F*-*F'*-measurable function

 $X: \Omega \to \Omega'.$

X is called measurable if

$$X^{-1}(A') \in F$$

for all $A' \in F'$. Measurability is important because for $A' \in F'$:

$$P(X \in A') = P(X^{-1}(A')).$$

Convergence concepts

We consider in the following a sequence of random variables

 X_1, X_2, X_3, \ldots

and investigate the limit behavior for $n \to \infty$.

Arbitrary distribution

Consider the estimator

$$\hat{\mu} = \bar{X}$$

in dependence of *n*:

$$\bar{X}_{1} = X_{1}$$

$$\bar{X}_{2} = \frac{1}{2}(X_{1} + X_{2})$$

$$\bar{X}_{3} = \frac{1}{3}(X_{1} + X_{2} + X_{3})$$

$$\vdots$$

$$\bar{X}_{n} = \frac{1}{n}\sum_{i=1}^{n}X_{i}$$

Almost sure convergence

A sequence X_n of random variables **converges almost surely to a** random variable X if,

$$P(\lim_{n\to\infty}X_n=X)=1.$$

The definition states that $X_n(\omega)$ converges to $X(\omega)$ for all $\omega \in \Omega$, except perhaps for $\omega \in N$ where $N \subset \Omega$ and P(N) = 0.

Notation: $X_n \xrightarrow{\text{a.s.}} X$

 Ω = [0,1] and P the uniform distribution. Define $X_n(\omega) = \omega + \omega^n$ and

$$X(\omega) = \omega.$$

We show that the sequence X_n converges to X almost surely.

Example (continued)

.

• For all $\omega \in [0,1)$ we have $\omega^n o 0$ for $n o \infty$ and therefore

$$X_n(\omega) = \omega + \omega^n \to \omega = X(\omega)$$

- For ω=1 we have X_n(1) = 2 for all n, such that X_n(1) does not converge to X(1) = 1.
- Since P([0, 1)) = 1, X_n converges almost surely to X.

Convergence in probability

 X_n converges to X in probability, if for $\epsilon > 0$

$$\lim_{n\to\infty} P(|X_n-X|\geq\epsilon)=0$$

or

$$\lim_{n\to\infty} P(|X_n-X|<\epsilon)=1.$$

Notation: $X_n \xrightarrow{\mathsf{P}} X$.

Weak law of large numbers

Let X_1, X_2, \ldots be a sequence of i.i.d. random variables with finite expected value μ and finite variance σ^2 .

Then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\mathsf{P}} \mu \quad \text{for} \quad n \to \infty$$

 $(\bar{X}_n \xrightarrow{\text{a.s.}} \mu \text{ can be established as well. Then we speak of the strong law of large numbers).}$

The law says, that for large *n*

$$\left|\bar{X}_n - \mu\right| < \epsilon$$

with high probability.

Example (continued)

Proof: Weak law of large numbers

• Using the Tschebyschov we obtain

$$P(\left|\bar{X}-\mu\right| \geq \epsilon) \geq rac{\operatorname{Var} \bar{X}}{\epsilon^2} = rac{\sigma^2}{n\epsilon^2} o 0 \qquad n o \infty$$

• Tschebyschov inequality: Let X be a random variable with expected value μ and variance σ^2 . For $\epsilon > 0$ the inequality

$$P(|X - \mu| \ge \epsilon) \le \frac{\sigma^2}{\epsilon^2}$$

or equivalently

$$P(|X-\mu|<\epsilon) \ge 1-rac{\sigma^2}{\epsilon^2}$$

holds.

Convergence in probability, not almost sure

Let $\Omega = [0,1]$ and *P* is the uniform distribution. Define $X_1, X_2, ...$ as follows:

$$\begin{aligned} X_{1}(\omega) &= \omega + I_{[0,1]}(\omega) \\ X_{2}(\omega) &= \omega + I_{[0,\frac{1}{2}]}(\omega) \\ X_{3}(\omega) &= \omega + I_{[\frac{1}{2},1]}(\omega) \\ X_{4}(\omega) &= \omega + I_{[0,\frac{1}{3}]}(\omega) \\ X_{5}(\omega) &= \omega + I_{[\frac{1}{3},\frac{2}{3}]}(\omega) \\ X_{6}(\omega) &= \omega + I_{[\frac{2}{3},1]}(\omega) \\ \vdots \end{aligned}$$

Define

$$X(\omega) = \omega.$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Example (continued)

Convergence in probability, not almost sure

• X_n converges to X in probability as

$$P(|X_n - X| \ge \epsilon)$$

is the probability of an interval of ω values whose length is going to 0.

• X_n does <u>not</u> converge to X almost surely as there is no $\omega \in \Omega$ for which

$$X_n(\omega) \to \omega = X(\omega).$$

Example (continued)

Convergence in probability, not almost sure

- For every ω , $X_n(\omega)$ alternates between ω and $\omega + 1$.
- For example $\omega = \frac{3}{8}$ yields

$$X_1(\frac{3}{8}) = 1\frac{3}{8}$$
$$X_2(\frac{3}{8}) = 1\frac{3}{8}$$
$$X_3(\frac{3}{8}) = \frac{3}{8}$$
$$X_4(\frac{3}{8}) = \frac{3}{8}$$

No pointwise convergence occurs for this sequence.

Convergence in the r-th mean

 X_n converges to X in the r-th mean, if

 $\mathsf{E}(|X_n^r|) < \infty$ for all n

and

$$\lim_{n\to\infty}\mathsf{E}(|X_n-X|^r)=0.$$

For r=2 we say that X_n converges in mean square to X.

Notation: $X_n \xrightarrow{r} X$.

Arbitrary distribution

$$\begin{split} \mathsf{E}(\hat{\mu}) &= \mathsf{E}(\bar{X}) = \mu\\ \mathsf{Var}(\hat{\mu}) &= \frac{\sigma^2}{n}\\ \mathsf{MSE}(\hat{\mu}) &= \mathsf{E}(\hat{\mu} - \mu)^2 = \frac{\sigma^2}{n} \to 0 \qquad \textit{with } n \to \infty, \end{split}$$

i.e. $\hat{\mu}$ converges in mean square to $\mu.$

Convergence in distribution

X_n converges to X in distribution, if

$$\lim_{n\to\infty} P(X_n \le x) = P(X \le x)$$

or

$$\lim_{n\to\infty}F_{X_n}(x)=F_x(x)$$

at all points x where $F_x(x)$ is continuous.

Remark

Note that if

$$f_n \to f$$
 with $n \to \infty$,

where f_n , *f* are the probability functions, then the distributions defined through f_n converge to the distribution defined through *f*.

Reverse is not correct, in general: **Convergence in distribution does not imply that the densities converge.**

Convergence of the Binomial distribution to the Poisson distribution

Let $X_n \sim B(n,\pi)$ with probability function

$$f(x)=P(X=x)=\binom{n}{x}\pi^{x}(1-\pi)^{n-x}.$$

For $n \to \infty$ and $n\pi = \lambda f(x)$ converges to the probability function of the Poisson distribution, i.e.

$$\lim_{n\to\infty} f(x) = \frac{\lambda^x exp(-\lambda)}{x!}$$

Example (continued)

Convergence of the Binomial distribution to the Poisson distribution

$$\lim_{n \to \infty} f(x) = \lim_{n \to \infty} {n \choose x} \pi^x (1 - \pi)^{n - x}$$

$$= \lim_{n \to \infty} \frac{n!}{x! (n - x)!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n - x}$$

$$= \lim_{n \to \infty} \frac{\lambda^x}{x!} \underbrace{\frac{n(n - 1) \cdots (n - x + 1)}{n \cdots n}}_{\to 1} \underbrace{(1 - \frac{\lambda}{n})^n}_{\to \exp(-\lambda)} \underbrace{(1 - \frac{\lambda}{n})^{-x}}_{\to 1}$$
We have used
$$exp(y) = \lim_{n \to \infty} (1 + \frac{y}{n})^n = \sum_{n = 0}^{\infty} \frac{y^n}{n!}$$

Application of the limit theorem

lf

- $X \sim B(n, \pi)$,
- n large,
- π 'small', i.e. $\lambda = n\pi$ 'moderate' (rule of thumb $n > 30, \pi \le 0.05$),

then *X* can be approximated by the Poisson distribution with parameters $\lambda = n\pi$.

Rice grains

Distribute n rice grains randomly to N squares. Let

X ='Number of rice grains in a square'.

We have

$$X \sim B(n, \pi = \frac{1}{N})$$

Because of n > 30, $\pi < 0.05$ we have

$$X \stackrel{a}{\sim} Po(\lambda = n\pi = \frac{n}{N}).$$

The sentence 'at all x for which $F_x(x)$ is continuous' matters! Let

$$X_n \sim N(0, \frac{1}{n})$$

and X a degenerated distribution at 0, i.e. P(X=0) = 1.

Central Limit Theorem

Let X_1, X_2, \ldots be a sequence of i.i.d. random variables with finite mean μ and variance σ^2 .

Then

$$\frac{1}{\sqrt{n}\sigma}\left(\sum_{i=1}^{n}X_{i}-n\mu\right)=\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\frac{X_{i}-\mu}{\sigma}=\frac{\bar{X}-\mu}{\sigma}\sqrt{n}\quad \xrightarrow{\mathsf{D}}\quad N(0,1).$$

Standardizing random variables

Assume X is a random variable mit expected value μ and variance σ^2 . Then

$$Y = \frac{X - \mu}{\sigma} = \frac{X}{\sigma} - \frac{\mu}{\sigma}$$

is the standardized version of X with

$$\mathsf{E}(Y) = \mathsf{E}\left(\frac{X}{\sigma} - \frac{\mu}{\sigma}\right) = \frac{1}{\sigma}\mathsf{E}(X) - \frac{\mu}{\sigma} = \frac{1}{\sigma}\mu - \frac{\mu}{\sigma} = 0$$

and

$$\operatorname{Var}(Y) = \operatorname{Var}\left(\frac{X}{\sigma} - \frac{\mu}{\sigma}\right) = \frac{1}{\sigma^2}\operatorname{Var}(X) = \frac{1}{\sigma^2}\sigma^2 = 1$$

Standardizing random variables

Note that

$$\mathsf{E}\left(\sum_{i=1}^{n} X_{i}\right) = \sum_{i=1}^{n} \mathsf{E}(X_{i}) = \sum_{i=1}^{n} \mu = n\mu$$

and

$$\operatorname{Var}\left(\sum_{i=1}^{n} X_{i}\right) = \sum_{i=1}^{n} \operatorname{Var}(X_{i}) = \sum_{i=1}^{n} \sigma^{2} = n\sigma^{2}$$

such that

$$\frac{1}{\sqrt{n}\sigma}\left(\sum_{i=1}^{n}X_{i}-n\mu\right)$$

is the standardized version of the random variable

$$\sum_{i=1}^n X_i.$$

Utilities for the proof

For a proof of the central limit theorem we need the following facts from analysis:

Exponential limit:

If a sequence a_n converges to a, i.e. $a_n \longrightarrow a$, then

$$\lim_{n \to \infty} \left(1 + \frac{a_n}{n} \right)^n = \exp(a) \tag{1}$$

Utilities for the proof

Taylor series expansion:

Let I be a real valued interval and $f : I \longrightarrow \mathbb{R}$ a function that is r + 1 times continuously differentiable. Then we obtain in a neighborhood around $a \in I$

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots + \frac{f^{(r)}(a)}{r!}(x-a)^r + R(x),$$

where R(x) is a function with

$$lim_{x \longrightarrow a} \frac{R(x)}{(x-a)^r} = 0.$$
 (2)

Proof of the central limit theorem

- Proof in case that the moment generating function $M_{\chi_i}(t)$ exists in a neighborhood around 0. In the general case the proof is similar and based on the so called characteristic function (which is always defined).
- The proof uses property a) (uniqueness) of the moment generating function of random variables (slide 16 of the probability theory slides). We show that the moment generating function of

$$Z_n := \frac{\bar{X}_n - \mu}{\sigma} \sqrt{n}$$

converges for $n \to \infty$ to the moment generating function of the N(0, 1) distribution.

1. Step: Standardizing

- Define the standardized random variable $Y_i = \frac{X_i \mu}{\sigma}$ and let $M_Y(t)$ be the moment generating function of Y_i .
- We now define

$$Z_n := \frac{\bar{X}_n - \mu}{\sigma} \sqrt{n} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$$

and obtain the moment generating function of Z_n :

$$M_{Z_n}(t) = \left(M_Y\left(\frac{t}{\sqrt{n}}\right)\right)^n$$

- Here we used properties c) (with $a = 1/\sqrt{n}$, b = 0) and e) of the moment generating function (slide 16, probability theory slides).
- Note that $M_Y^{(r)}(0) = E(Y^r)$ (property b) page 16 of the slides). This implies $M_Y(0) = 1$, $M_Y'(0) = E(Y) = 0$, $M_Y''(0) = E(Y^2) = Var(Y) = 1$.

2. Step: Taylor series expansion

• A second order Taylor series expansion of $M_Y\left(\frac{t}{\sqrt{n}}\right)$ around 0 yields:

$$M_{Y}\left(\frac{t}{\sqrt{n}}\right) = M_{Y}(0) + \frac{M_{Y}'(0)}{1!}\frac{t}{\sqrt{n}} + \frac{M_{Y}''(0)}{2!}\left(\frac{t}{\sqrt{n}}\right)^{2} + R_{Y}\left(\frac{t}{\sqrt{n}}\right)$$
$$= 1 + \frac{0}{1!}\frac{t}{\sqrt{n}} + \frac{1}{2!}\left(\frac{t}{\sqrt{n}}\right)^{2} + R_{Y}\left(\frac{t}{\sqrt{n}}\right)$$
$$= 1 + \frac{1}{2}\left(\frac{t}{\sqrt{n}}\right)^{2} + R_{Y}\left(\frac{t}{\sqrt{n}}\right)$$

2. Step: Taylor series expansion

• For fixed $t \neq 0$ we have $\frac{t}{\sqrt{n}} \rightarrow 0$ and because of (2)

$$\lim_{n\to\infty}\frac{R_{Y}\left(\frac{t}{\sqrt{n}}\right)}{\left(\frac{t}{\sqrt{n}}\right)^{2}}=0$$

• Since $t \neq 0$ is fixed we also obtain

$$\lim_{n \to \infty} \frac{R_Y\left(\frac{t}{\sqrt{n}}\right)}{\left(\frac{1}{\sqrt{n}}\right)^2} = \lim_{n \to \infty} n R_Y\left(\frac{t}{\sqrt{n}}\right) = 0$$

• This statement is also valid for t = 0 because $R_Y\left(\frac{0}{\sqrt{n}}\right) = 0$.
3. Step: Limit

• For fixed t we finally obtain

$$\lim_{n \to \infty} \left(M_Y\left(\frac{t}{\sqrt{n}}\right) \right)^n = \lim_{n \to \infty} \left[1 + \frac{\left(t/\sqrt{n}\right)^2}{2!} + R_Y\left(\frac{t}{\sqrt{n}}\right) \right]^n$$
$$= \lim_{n \to \infty} \left[1 + \frac{1}{n} \left(\frac{t^2}{2} + nR_Y\left(\frac{t}{\sqrt{n}}\right)\right) \right]^n$$
$$= \exp\left(\frac{t^2}{2}\right)$$

• Thereby we used the exponential limit (1) from above with

$$a_n = \frac{t^2}{2} + n R_Y \left(\frac{t}{\sqrt{n}}\right) \rightarrow \frac{t^2}{2}$$

• Since $\exp\left(\frac{t^2}{2}\right)$ is the moment generating function of the N(0, 1) distribution the theorem is proven.

Relationship among modes of convergence



Construction of confidence intervals I

Let X_1, \dots, X_n be an i.i.d sample with $X_i \sim N(\mu, \sigma^2)$ and σ^2 known.

Then

$$\left[\bar{X} - z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}}\right]$$

is a $1 - \alpha$ confidence interval for μ .

Proof

We have

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

respectively

$$\frac{\bar{X}-\mu}{\sigma}\sqrt{n}\sim N(0,1). \tag{3}$$

Hence for 0 $<\alpha<$ 1 we obtain

$$P(-z_{1-\frac{\alpha}{2}} \leq \frac{\bar{X}-\mu}{\sigma}\sqrt{n} \leq z_{1-\frac{\alpha}{2}}) = 1-\alpha,$$

where $z_{1-\frac{\alpha}{2}}$ is the $1-\frac{\alpha}{2}$ Quantile of the *N*(0,1) distribution.

Proof (continued)

Rearranging terms yields

$$P(\bar{X}-z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}} \le \mu \le \bar{X}+z_{1-\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}})=1-\alpha,$$

such that

$$\left[\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right]$$

is a 1- α confidence interval for μ .

Because of the CLT(3) and with it the CI is valid for arbitrary distributions provided that *n* is large enough (rule of thumb $n \ge 30$).

Special Case: Approximation of the Binomial Distribution

Let $X \sim B(n, \pi)$. X can be represented as $X = X_1 + \cdots + X_n$, where

$$X_i \in \{0,1\}, \quad P(X_i = 1) = \pi, \quad P(X_i = 0) = 1 - \pi$$

Applying the CLT we obtain

$$\frac{X - nE(X_i)}{\sqrt{n \operatorname{Var}(X_i)}} = \frac{X - n\pi}{\sqrt{n\pi(1 - \pi)}} \stackrel{\text{D}}{\sim} N(0, 1)$$

Thus approximately we have

$$X \stackrel{a}{\sim} N(n\pi, n\pi(1-\pi))$$
$$P(X \le x) \approx \phi\left(\frac{x - n\pi}{\sqrt{n\pi(1-\pi)}}\right).$$

Unknown variance

So far we have only considered the case where the variance σ^2 is known.

To construct confidence intervals for unknown σ^2 we would like to replace

$$rac{\mathbf{X}-\mu}{\sigma}\sqrt{n} \quad \stackrel{\mathsf{D}}{\sim} \quad N(0,1)$$

by

$$rac{ar{X}-\mu}{S}\sqrt{n}$$
 $\stackrel{ ext{D}}{\sim}$ $N(0,1).$

Continuous mapping theorem

If $g:\mathbb{R}
ightarrow \mathbb{R}$ is a continuous function, then

(i)
$$X_n \xrightarrow{a.s.} X \Rightarrow g(X_n) \xrightarrow{a.s.} g(X)$$

(ii) $X_n \xrightarrow{p} X \Rightarrow g(X_n) \xrightarrow{p} g(X)$
(iii) $X_n \xrightarrow{d} X \Rightarrow g(X_n) \xrightarrow{d} g(X)$

Sums and products

(a)
$$X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \Rightarrow X_n + Y_n \xrightarrow{p} X + Y$$

(b) $X_n \xrightarrow{r} X, Y_n \xrightarrow{r} Y \Rightarrow X_n + Y_n \xrightarrow{r} X + Y$
(c) $X_n \xrightarrow{d} X, Y_n \xrightarrow{d} a \Rightarrow X_n + Y_n \xrightarrow{d} X + a$ (Slutzky)
(d) $X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \Rightarrow X_n \cdot Y_n \xrightarrow{p} X \cdot Y$
(e) $X_n \xrightarrow{d} X, Y_n \xrightarrow{d} a \Rightarrow X_n \cdot Y_n \xrightarrow{d} a \cdot X$ (Slutzky)

Example

Sampling variance



Proof

We have

$$\frac{1}{n}\sum_{i=1}^{n}X_{i}^{2}\xrightarrow{\mathsf{P}}E(X^{2})$$

because of the weak law of large numbers.

- As $\bar{X} \stackrel{\mathsf{P}}{\longrightarrow} \mu$ the continuous mapping theorem implies $\bar{X}^2 \stackrel{\mathsf{P}}{\longrightarrow} \mu^2$.
- Then

$$\frac{1}{n}\sum X_i^2 - \bar{X}^2 \xrightarrow{\mathsf{P}} E(X^2) - \mu^2 = \sigma^2$$

applying slide 96 a).

• Finally as $n/(n-1) \rightarrow 1$ we obtain

$$S^{2} = \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^{n} X_{i}^{2} - \bar{X}^{2} \right] \stackrel{\mathsf{P}}{\longrightarrow} \sigma^{2}.$$

Construction of confidence intervals II

Let X_1, X_2, \ldots be an i.i.d. sample with finite mean μ and variance σ^2 . We show that

$$\frac{\bar{X}-\mu}{S}\sqrt{n} \sim N(0,1).$$

Proof

- Since $S^2 \xrightarrow{P} \sigma^2$ we have $S \xrightarrow{P} \sigma$ according to the continuous mapping theorem.
- This implies (using slide 96 d)

$$\frac{1}{S} \xrightarrow{\mathsf{P}} \frac{1}{\sigma} \quad \text{and} \quad \frac{\sigma}{S} \xrightarrow{\mathsf{P}} 1.$$

According to Slutsky we obtain

$$\frac{\bar{X} - \mu}{S} \sqrt{n} = \underbrace{\frac{\sigma}{S}}_{\stackrel{P}{\longrightarrow} 1} \underbrace{\frac{\bar{X} - \mu}{\sigma} \sqrt{n}}_{\stackrel{D}{\longrightarrow} N(0, 1)} \xrightarrow{D} N(0, 1)$$

• This result implies the following CI for μ and large *n*:

$$\left[\bar{X} - z_{1-\frac{\alpha}{2}}\frac{S}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}}\frac{S}{\sqrt{n}}\right]$$

Asymptotically unbiased estimators

An estimator $\hat{\theta}$ for θ is asymptotically unbiased, if for all θ

$$\lim_{n\to\infty} E(\hat{\theta}) = \theta$$

Obviously every unbiased estimator $\hat{\theta}$ is also asymptotically unbiased.

Example

Arbitrary distribution

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

is unbiased for σ^2 .

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

is however biased. Since

$$E(\tilde{S}^2) = E\left(\frac{n-1}{n}S^2\right) = \frac{n-1}{n}\sigma^2$$

 \tilde{S}^2 is at least asymptotically unbiased.

Consistency of an estimator

An estimator $\hat{\theta}$ for θ is called (weakly) consistent , if $\hat{\theta}$ converges on probability to θ , i.e.

$$\hat{\theta} \stackrel{\mathsf{P}}{\longrightarrow} \theta \qquad \qquad n \to \infty$$

(strong consistency means convergence almost surely).

An **estimator** $\hat{\theta}$ **is consistent in square mean**, if $\hat{\theta}$ converges to θ in square mean, i.e.

$$\hat{\theta} \xrightarrow{r=2} \theta.$$

Example

Consistency of \bar{X} for μ

 X_1, X_2, \dots i.i.d. with finite expected value μ and variance σ^2 . As $\overline{X} \xrightarrow{P} \mu$ (*law of large numbers*), $\hat{\mu} = \overline{X}$ is consistent for μ . It follows that

$$\hat{\pi} = \bar{X}$$

and

$$\hat{\lambda}=\bar{X}$$

are consistent for π in the Binomial case and λ in the Poisson case.

Example

Consistency of the variance estimator

We have $S^2 \xrightarrow{\mathsf{P}} \sigma^2$. It follows that S^2 is consistent for σ^2 .

Because of the continuous mapping theorem $S \xrightarrow{P} \sigma$, so that *S* is consistent for σ .

Relationships between the various consistency concepts



If the *MSE* converges to 0, then $\hat{\theta} \xrightarrow{r=2} \theta$, because $MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2$. This implies consistency in square mean and in turn weak consistency. Note that **every consistent estimator is asymptotically unbiased.**

Note also, that **the consistency of** $\hat{\theta}$ **for** θ **implies the consistency of** $g(\hat{\theta})$ **for** $g(\theta)$, if *g* is continuous (because of the continuous mapping theorem).

Standard errors

Consider an i.i.d. random sample X_1, \ldots, X_n . Let $\hat{\theta} = T(X_1, \ldots, X_n)$ be an estimator for θ . Let *V* be a consistant estimator for Var(T).

Then \sqrt{V} is consistent for $\sqrt{Var(T)}$.

The quantity

$$se(\hat{ heta}) = \sqrt{V}$$

is called standard error of $\hat{\theta}$.

Example

Arbitrary distribution

 X_1, \ldots, X_n i.i.d. random sample with finite expected value μ and variance σ^2 . Then

$$Var(\bar{X}) = \frac{\sigma^2}{n}$$

and

$$V = \frac{S^2}{n}$$

is consistent for $Var(\bar{X})$. Therefore

$$se(\hat{\mu}) = rac{S}{\sqrt{n}}$$







Chapter 3

Likelihood based inference

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Statistical Inference - 3 - Likelihood based inference - 0 / 62

Likelihood based inference

Likelihood and log-likelihood

Likelihood and log-likelihood

We throw a coin n = 10 times and observe y = 8 times tail.

Which value for

 $\pi =$ 'Probability of tail'

is most likely?

Let X='Number of tails', then

 $X \sim B(10, \pi).$

Likelihood and log-likelihood (continued)

The following table provides $P_{\pi}(X=8)$ in dependence of π :

π	<i>Ρ</i> _π (<i>X</i> =8)
0.1	0
0.2	0.00007373
÷	÷
0.7	0.23347444
0.8	0.3019898
0.9	0.1937
1	0

Likelihood and log-likelihood (continued)

The most plausible value for the parameter π is $\pi = 0.8!$

Likelihood theory uses the estimator $\hat{\pi}$ that maximizes the function

$$L(\pi) = P_{\pi}(X = 8)$$

with respect to π .

 $L(\pi)$ is denoted as likelihood. In likelihood theory, all inferential conclusions are based on the likelihood.

Definition: Likelihood

Let *x* be the realized value of a random variable *X*. Let $f(x, \theta)$ be the corresponding probability function (if *X* is discrete) respectively density (if *X* is continuous).

Then

$$L(\theta) = f(x,\theta)$$

is called likelihood (function).

Definition: Maximum likelihood estimator

The **maximum likelihood estimator** $\hat{\theta}_{ML}$ of a parameter θ is given by maximizing the likelihood:

$$\hat{\theta}_{\mathsf{ML}} = \operatorname*{argmax}_{\theta} L(\theta).$$

Usually it is more convenient to maximize the log-likelihood

$$\ell(\theta) = \log L(\theta).$$

Because of the monotonicity of the logarithm we have

$$\hat{ heta}_{\mathsf{ML}} = rgmax_{ heta} \, {oldsymbol{\mathcal{L}}}(heta) = rgmax_{ heta} \, {oldsymbol{\ell}}(heta).$$

Example

Binomial distribution

Let $X \sim B(n, \pi)$ and x be the realized value. Then

$$L(\pi) = \binom{n}{k} \pi^{x} (1 - \pi)^{n - x} \propto \pi^{x} (1 - \pi)^{n - x}$$
$$l(\pi) = x \log(\pi) + (n - x) \log(1 - \pi)$$
$$l'(\pi) = \frac{x}{\pi} + \frac{(n - x)}{1 - \pi} (-1) = \frac{x}{\pi} - \frac{n - x}{1 - \pi}$$

Example (continued)

Binomial distribution		ì
Setting to zero yields		l
	$\frac{x}{\pi} - \frac{n-x}{1-\pi} = 0$	l
and therefore	Y	l
	$\hat{\pi}_{ML} = \frac{\lambda}{n}.$	

Likelihood and ML-estimator for a random sample

Consider an i.i.d. sample $X_1, ..., X_n$ with probability function or density of X_i given by $f(x_i, \theta)$. If

$$x = (x_1, \ldots, x_n)'$$

is the realized sample of

$$X=(X_1,\ldots,X_n)',$$

then because of the independence we have

$$L(\theta) = f(x_1, \theta) \cdot f(x_2, \theta) \cdot \cdots \cdot f(x_n, \theta) = \prod_{i=1}^n f(x_i, \theta)$$

and

$$\ell(\theta) = \log f(x_1, \theta) + \cdots + \log f(x_n, \theta) = \sum_{i=1}^n \log f(x_i, \theta).$$

Example

Poisson distribution

Let X_1, \ldots, X_n be a random sample with $X_i \sim Po(\lambda)$. The realizations are given by x_1, x_2, \ldots, x_n . The probability function for X_i is given by

$$f(x_i) = P(X_i = x_i) = \frac{\lambda^{x_i}}{x_i!} \exp(-\lambda).$$

Example (continued)

Poisson distribution

$$L_{i}(\lambda) = \frac{\lambda^{x_{i}}}{x_{i}!} \exp(-\lambda) \propto \lambda^{x_{i}} \exp(-\lambda)$$

$$l_{i}(\lambda) = \log(\lambda^{x_{i}} \exp(-\lambda)) = x_{i}\log(\lambda) - \lambda$$

$$l_{i}'(\lambda) = \frac{x_{i}}{\lambda} - 1$$

$$l'(\lambda) = \frac{1}{\lambda} \sum_{i=1}^{n} x_{i} - n$$

Setting to zero and solving for λ yields

$$\hat{\lambda}_{\mathsf{ML}} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Likelihood based inference

Score function and Fisher information

Definition: Score function

The first derivative of the log-likelihood

$$\mathcal{S}(heta) = rac{\partial \ell(heta)}{\partial heta} = \ell'(heta)$$

is called **score function**.

Definition: Fisher information

The negative second derivative of the log-likelihood

$$I(\theta) = -\frac{\partial^2 \ell(\theta)}{\partial \theta^2} = -\frac{\partial S(\theta)}{\partial \theta} = -\ell''(\theta)$$

is called Fisher information.

Evaluating the Fisher information at the maximum $\hat{\theta}_{ML}$ yields the observed Fisher information.
First interpretation of the (observed) Fisher information

At the maximum $\hat{\theta}_{\rm ML}$ the second derivative of the log-likelihood is negative and therefore I($\hat{\theta}_{\rm ML}$) positive.

The larger $I(\hat{\theta}_{ML})$, the larger is the curvature of the log-likelihood, i.e. the more *'information'* to locate the maximum.

Example

Binomial distribution

Let
$$X \sim B(n, \pi)$$
.
 $l(\pi) = x \log(\pi) + (n - x) \log(1 - \pi)$
 $S(\pi) = l'(\pi) = \frac{x}{\pi} - \frac{n - x}{1 - \pi}$
 $l(\pi) = -l''(\pi) = -\left(-\frac{x}{\pi^2} + \frac{n - x}{(1 - \pi)^2}(-1)\right) = \frac{x}{\pi^2} + \frac{n - x}{(1 - \pi)^2}$
 $l(\hat{\pi}) = \frac{x}{(\frac{x}{n})^2} + \frac{n - x}{(1 - (\frac{x}{n}))^2} = \dots = \frac{n}{\hat{\pi}(1 - \hat{\pi})}$

Example

Poisson distribution

 X_1, \ldots, X_n i.i.d. Poisson distributed with $X_i \sim Po(\lambda)$ and realizations x_1, \ldots, x_n .

$$S_{i}(\lambda) = l'_{i}(\lambda) = \frac{x_{i}}{\lambda} - 1$$

$$S(\lambda) = \frac{1}{\lambda} \sum_{i=1}^{n} x_{i} - n\lambda$$

$$I_{i}(\lambda) = -l''_{i}(\lambda) = \frac{x_{i}}{\lambda^{2}}$$

$$I(\lambda) = \frac{1}{\lambda^{2}} \sum_{i=1}^{n} x_{i}$$

$$I(\hat{\lambda}) = \frac{n^{2}}{\left(\sum_{i=1}^{n} x_{i}\right)^{2}} \sum_{i=1}^{n} x_{i} = \frac{n^{2}}{\sum_{i=1}^{n} x_{i}} = \frac{n}{\bar{x}}$$

Definition: Expected Fisher information

The expected value of the Fisher information $I(\theta)$, considered as a function of the sample variables $X = (X_1, ..., X_n)$ ', is called **expected Fisher information**:

$$J(\theta) = \mathsf{E}\big(\mathsf{I}(\theta)\big).$$

Note that $J(\theta)$ is additive for an i.i.d. sample with density $f(x,\theta)$. Let $J_1(\theta)$ be the expected Fisher information of x_i . Then $J(\theta) = nJ_1(\theta)$.

Example

Binomial distribution

Let $X \sim B(n, \pi)$. Then

$$I(\pi) = \frac{X}{\pi^2} + \frac{n - X}{(1 - \pi)^2}$$

$$J(\pi) = E\left(\frac{X}{\pi^2} + \frac{n - X}{(1 - \pi)^2}\right)$$

$$= \frac{n\pi}{\pi^2} + \frac{n - n\pi}{(1 - \pi)^2} = \frac{n}{\pi} + \frac{n}{1 - \pi}$$

$$= \frac{n(1 - \pi) + n\pi}{\pi(1 - \pi)} = \frac{n}{\pi(1 - \pi)}$$

$$= \frac{1}{Var(\hat{\pi})},$$

Likelihood based inference

Expected value and variance of the score function

Expected value and variance

Under regularity conditions we show

$$E(S(\theta)) = 0$$
 and $Var(S(\theta)) = J(\theta)$.

Regularity conditions

Exchangeability of Differentiation and Integration (or Differentiation and Summation).

For example fulfilled for fixed integral limits, f and $\frac{\partial f}{\partial t}$ continuous, then

$$\frac{\partial}{\partial t} \int_{a}^{b} f(x,t) \, \mathrm{d}x = \int_{a}^{b} \frac{\partial f(x,t)}{\partial t} \, \mathrm{d}x$$

Proof (for continuous distributions)

Because

$$S(\theta) = \sum_{i=1}^{n} S_i(\theta),$$

it is sufficient to show $E(S_i(\theta)) = 0$ and we can restrict ourselves to n = 1, then $L(\theta) = f(x, \theta)$.

• We have the following equivalent representations of the score function:

$$S(\theta) = \frac{\partial I(\theta)}{\partial \theta} = \frac{\partial \log L(\theta)}{\partial \theta} = \frac{1}{L(\theta)} \frac{\partial L(\theta)}{\partial \theta}$$

$$E(S(\theta)) = \int S(\theta)f(x,\theta) dx$$

= $\int \frac{1}{L(\theta)} \frac{\partial L(\theta)}{\partial \theta} f(x,\theta) dx$
= $\int \frac{\partial L(\theta)}{\partial \theta} dx$
= $\frac{\partial}{\partial \theta} \left(\int L(\theta) dx \right)$
= $\frac{\partial}{\partial \theta} (1)$
= 0

Since
$$E(S(\theta)) = 0$$
 we have $Var(S(\theta)) = E(S(\theta)^2)$. Then

$$J(\theta) = \mathsf{E}\left(-\frac{\partial^{2} \mathsf{logL}(\theta)}{\partial \theta^{2}}\right) = \mathsf{E}\left(-\frac{\partial}{\partial \theta}\left(\underbrace{\frac{\partial \mathsf{L}(\theta)}{\partial \theta}}_{=f}, \underbrace{\frac{1}{\mathsf{L}(\theta)}}_{=g}\right)\right)$$
$$= \mathsf{E}\left(-\frac{\frac{\partial^{2}\mathsf{L}(\theta)}{\partial \theta^{2}}\mathsf{L}(\theta) - \left(\frac{\partial\mathsf{L}(\theta)}{\partial \theta}\right)^{2}}{\mathsf{L}(\theta)^{2}}\right) = -\mathsf{E}\left(\frac{\frac{\partial^{2}\mathsf{L}(\theta)}{\partial \theta^{2}}}{\mathsf{L}(\theta)}\right) + \mathsf{E}\left(\frac{\left(\frac{\partial\mathsf{L}(\theta)}{\partial \theta}\right)^{2}}{\mathsf{L}(\theta)^{2}}\right)$$
$$= -\int \frac{\frac{\partial^{2}\mathsf{L}(\theta)}{\partial \theta^{2}}f(x,\theta)\,\mathsf{d}x + \int \frac{\left(\frac{\partial\mathsf{L}(\theta)}{\partial \theta}\right)^{2}}{\mathsf{L}(\theta)^{2}}f(x,\theta)\,\mathsf{d}x$$
$$= -\frac{\partial^{2}}{\partial \theta^{2}}\underbrace{\int \mathsf{L}(\theta)\,\mathsf{d}x}_{=1} + \int \left(\frac{\partial}{\partial \theta}\mathsf{logL}(\theta)\right)^{2}f(x,\theta)\,\mathsf{d}x$$
$$= \mathsf{E}\left(\mathsf{S}(\theta)^{2}\right) = \mathsf{Var}\left(\mathsf{S}(\theta)\right)$$

We thereby used the following rules regarding derivatives

$$\big(\frac{f}{g}\big)' = \frac{f'g - fg'}{g^2}$$

and

$$rac{\partial}{\partial heta} \mathsf{logL}(heta) = rac{rac{\partial \mathsf{L}(heta)}{\partial heta}}{\mathsf{L}(heta)}$$

Likelihood based inference

Cramer-Rao bound

Cramer-Rao bound

Consider an unbiased estimator $\hat{\theta}$ for θ , i.e.

$$E(\hat{ heta}) = heta.$$

Let $J(\theta)$ be the expected Fisher information. Then we have

$$Var(\hat{ heta}) \geq rac{1}{J(heta)}$$
 (Cramer-Rao bound)

under regularity conditions, i.e. exchangeability of integration and differentiation.

Application of the Cramer-Rao bound

If we can show in applications, that

$$/ar(\hat{ heta}) = rac{1}{J(heta)},$$

i.e. the bound is reached, then $\hat{\theta}$ is the estimator with lowest variance among all unbiased estimators.

Then $\hat{\theta}$ is called **efficient**.

Proof

Let

$$f(\vec{x},\theta) = \prod_{i=1}^{n} f(x_i,\theta).$$

Then the squared correlation between $\mathcal{S}(\theta)$ and $\hat{\theta}$ is

$$\rho^{2}(\mathsf{S}(\theta), \hat{\theta}) = \frac{\mathsf{Cov}(\mathsf{S}(\theta), \hat{\theta})^{2}}{\mathsf{Var}(\mathsf{S}(\theta))\mathsf{Var}(\hat{\theta})} \leq 1$$

Because of $Var(S(\theta)) = J(\theta)$ we have:

$$\frac{\operatorname{Cov}(\mathsf{S}(\theta),\hat{\theta})^2}{J(\theta)\operatorname{Var}(\hat{\theta})} \leq 1.$$

Rearranging yields

$$\mathsf{Var}(\hat{ heta}) \geq rac{\mathsf{Cov}ig(\mathsf{S}(heta), \hat{ heta}ig)^2}{J(heta)}$$

It remains to show that $Cov(S(\theta), \hat{\theta}) = 1$.

$$Cov(S(\theta), \hat{\theta}) = E(S(\theta)\hat{\theta}) - \underbrace{E(S(\theta))}_{=0} E(\hat{\theta})$$

$$= \int S(\theta)\hat{\theta}f(\vec{x}, \theta) d\vec{x} = \int \frac{\partial f(\vec{x}, \theta)}{\partial \theta} \frac{1}{f(\vec{x}, \theta)} \hat{\theta}f(\vec{x}, \theta) d\vec{x}$$

$$= \int \frac{\partial f(\vec{x}, \theta)}{\partial \theta} \hat{\theta} d\vec{x}$$

$$= \frac{\partial}{\partial \theta} \int f(\vec{x}, \theta)\hat{\theta} d\vec{x}$$

$$= \frac{\partial}{\partial \theta} \left(E(\hat{\theta})\right)$$

$$= \frac{\partial}{\partial \theta}(\theta)$$

$$= 1$$

Example

Binomial distribution

 $X \sim B(n, \pi).$

Then

$$\hat{\pi} = rac{X}{n},$$
 $J(\pi) = rac{n}{\pi(1-\pi)}.$

$$\operatorname{Var}(\hat{\pi}) = \frac{n\pi(1-\pi)}{n^2} = \frac{\pi(1-\pi)}{n} = \frac{1}{J(\pi)},$$

i.e. $\hat{\pi}$ reaches the Cramer-Rao lower bound and is therefore efficient.

Likelihood based inference

Distribution of the score statistic

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Statistical Inference - 3 - Likelihood based inference - 33 / 62

Distribution of the score statistic

Let X_1, \ldots, X_n be a random sample with density $f(x, \theta)$. Then

$$\frac{\mathsf{S}(\theta)}{\sqrt{J(\theta)}} \xrightarrow{\mathsf{D}} \mathsf{N}(0,1) \tag{1}$$

i.e. approximately we have

$$S(\theta) \stackrel{a}{\sim} N(0, J(\theta)).$$

The proof is a direct application of the CLT

$$\frac{\mathcal{S}(\theta)}{\sqrt{J(\theta)}} = \frac{\sum \mathcal{S}_i(\theta) - 0}{\sqrt{J(\theta)}} \stackrel{\mathsf{D}}{\longrightarrow} N(0, 1)$$

Remark

In (1) the expected Fisher information can be replaced by the (observed) Fisher information, i.e.

$$rac{\mathsf{S}(heta)}{\sqrt{\mathrm{I}(heta)}}, \qquad rac{\mathsf{S}(heta)}{\sqrt{J(\hat{ heta}_{\mathsf{ML}})}}, \qquad rac{\mathsf{S}(heta)}{\sqrt{\mathrm{I}(\hat{ heta}_{\mathsf{ML}})}}.$$

Example

Poisson distribution

$$\mathsf{T} = rac{\mathsf{S}(\lambda)}{\sqrt{J(\lambda)}} = rac{\sum X_i}{\lambda} - n \quad \stackrel{\mathsf{a}}{\sim} \quad \mathsf{N}(0,1).$$

Application: Score Test

Aim is to test the hypothesis

$$H_0: \theta = \theta_0$$
 versus $H_1: \theta \neq \theta_0$

Use the test statistic

$$\frac{\mathsf{S}(\theta_0)}{\sqrt{J(\theta_0)}} \quad \stackrel{a}{\sim} \quad \textit{N}(0,1) \qquad (\text{under } \mathsf{H}_0).$$

Application: Score Test (continued)

If $|S(\theta_0)|$ large, then the difference between θ_0 and the ML-estimator $\hat{\theta}_{ML}$ is large (note that $S(\hat{\theta}_{ML})=0$), i.e. we reject H₀.

More precisely we reject H₀ if

$$\frac{\left|\mathsf{S}(\theta_0)\right|}{\sqrt{J(\theta_0)}} > z_{1-\frac{\alpha}{2}}.$$

Note that the computation of the ML-estimator is <u>not</u> required to use the score test.

Example

Poisson distribution

 $H_0: \lambda = \lambda_0$ versus $H_1: \lambda \neq \lambda_0$

Use

$$\mathsf{T} = \frac{\frac{\sum X_i}{\lambda_0} - n}{\sqrt{\frac{n}{\lambda_0}}} \quad \stackrel{\mathsf{a}}{\sim} \quad \mathsf{N}(0, 1)$$

as the test statistic.

Recap Hypotheses testing

• Aim is to test the hypothesis

 $H_0: \theta = \theta_0$ versus $H_1: \theta \neq \theta_0$.

- Two sources of error:
 - H_0 is true, but we decide for H_1 (type I error)
 - H_1 is true, but we decide to keep H_0 (type II error)
- Test is constructed asymmetrically such that type 1 error is under our control, e.g. $\alpha = 0.05$. Type II error is not controlled directly, but can be reduced usually through sample size.
- Test decision is based on a test statistics whose distribution under *H*₀ must be known. Often distribution is not known under *H*₁, then type II error is not available/computable.

Recap Hypotheses testing

- Caution I:
 - A nonsignificant result, i.e. *H*₀ is kept, does not necessarily imply that *H*₀ is true!
 - Suppose we toss a coin 4 times and obtain 3 times head. The aim is to test H_0 : $\pi = 0.5$ versus H_1 : $\pi \neq 0.5$ where π is the probability of head.
 - Obviously any reliable test will not reject *H*₀ on basis of the given data.
 - However, already our intuition says, that the nonsignificant results does not imply that *H*₀ is true.
- Caution II:

Think of the power of the test. Is there a real chance for a significant result with the test at hand?

Likelihood based inference

Distribution of the ML-estimator

Distribution of the ML-estimator

Let X_1, \ldots, X_n be an i.i.d random sample with density (or probability function) $f(x, \theta)$ and $\hat{\theta}_{ML}$ the ML-estimator for θ . Let θ_0 be the true parameter.

Then under regularity conditions the ML-estimator is consistent, i.e.

$$\hat{\theta}_{ML} \stackrel{\mathsf{P}}{\longrightarrow} \theta_0$$

and

$$\sqrt{n}(\hat{\theta}_{ML}-\theta_0) \stackrel{\mathsf{d}}{\longrightarrow} N\!\left(0, \frac{1}{J_1(\theta_0)}\right)$$

or

$$(\hat{ heta}_{ML} - heta_0)\sqrt{nJ_1(heta_0)} = (\hat{ heta}_{ML} - heta_0)\sqrt{J(heta_0)} \stackrel{\mathsf{d}}{\longrightarrow} N(0, 1),$$

i.e. we obtain approximately

$$\hat{\theta}_{ML} \stackrel{\mathsf{a}}{\sim} N(\theta_0, J(\theta_0)^{-1})$$

Remarks I

The asymptotic results suggest the following statistical properties of the maximum likelihood estimator:

- $\hat{\theta}_{ML}$ is weakly consistent, i.e. $\hat{\theta}_{ML} \xrightarrow{\mathsf{P}} \theta_0$.
- Consistency implies that $\hat{\theta}_{ML}$ is asymptotically unbiased. Note however, that $\hat{\theta}_{ML}$ is generally not unbiased in finite samples!
- For large *n* the approximate variance of $\hat{\theta}_{ML}$ is

$$ext{Var}(\hat{ heta}_{ML}) pprox rac{1}{J(heta_0)},$$

i.e. the ML estimator reaches the Cramer- Rao bound as *n* tends to infinity.

• This in turn implies that $\hat{\theta}_{ML}$ is asymptotically efficient, i.e. among all asymptotically unbiased estimators it is the estimator with the lowest variance.

Remarks II

- In $\hat{\theta}_{ML} \stackrel{a}{\sim} N(\theta_0, J(\theta_0)^{-1})$ we can replace $J(\theta_0)$ by $I(\theta_0)$.
- Because $\hat{\theta}_{ML}$ is consistent for θ_0 we also have

$$\hat{\theta}_{ML} \stackrel{a}{\sim} N(\theta_0, J(\hat{\theta}_{ML})^{-1}) \quad \text{and} \quad \hat{\theta}_{ML} \stackrel{a}{\sim} N(\theta_0, I(\hat{\theta}_{ML})^{-1}).$$

For the construction of tests and CI's this is important because $J(\hat{\theta}_{ML})$ does not depend on the unknown parameter.

• The second remark yields approximate standard errors

$$se(\hat{ heta}_{ML}) pprox \sqrt{rac{1}{J(\hat{ heta}_{ML})}}$$

or

$$se(\hat{ heta}_{ML}) pprox \sqrt{rac{1}{I(\hat{ heta}_{ML})}}$$

Regularity conditions for consistency

- **R1:** The parameter is identifiable, i.e. if $\theta \neq \theta'$ then $f(x,\theta) \neq f(x,\theta')$.
- **R2:** The densities $f(x,\theta)$ have common support and $f(x,\theta)$ is differentiable in θ .
- **R3:** The parameter space Θ contains an open set of which the true parameter is an interior point.

Additional regularity conditions for asymptotic normality

- **R4:** The density $f(x,\theta)$ is three times differentiable with respect to θ , the third derivative is continuous in θ and $\int f(x,\theta) dx$ can be differentiated three times under the integral sign.
- **R5:** For any θ_0 , there exists a positive number *c* and a function M(*x*) (both of which may depend on θ_0) such that

$$\left|rac{\partial^3 \mathsf{log} f(x, heta)}{\partial heta^3}
ight| \leq \mathsf{M}(x)$$

for all x and $\theta_0 - c < \theta < \theta_0 + c$ with $E(M(x)) < \infty$.

Proof of asymptotic normality

• A second order Taylor series expansion of the score function at the ML-estimator $\hat{\theta}_{ML}$ about the true value θ_0 yields

$$s\left(\hat{\theta}_{ML}\right) = s(\theta_0) + \left(\hat{\theta}_{ML} - \theta_0\right) I''(\theta_0) + \frac{1}{2} \left(\hat{\theta}_{ML} - \theta_0\right)^2 I'''(\theta^*)$$
(3)

with θ^* between θ_0 and $\hat{\theta}_{ML}$.

• Using $s(\hat{\theta}_{ML}) = 0$ and rearranging (3) yields

$$(\hat{\theta}_{ML} - \theta_0)\sqrt{n} = \frac{\frac{s(\theta_0)}{\sqrt{n}}}{-\frac{1}{n}I''(\theta_0) - \frac{1}{2n}(\hat{\theta}_{ML} - \theta_0)I'''(\theta^*)}.$$
 (4)

Now we can derive the following three properties:

Because of the asymptotic properties of the score function we have

$$\frac{s(\theta_0)}{\sqrt{n}} \stackrel{d}{\longrightarrow} N(0, J_1(\theta_0)).$$

• For $-l''(\theta_0)/n$ the central limit theorem yields

$$-\frac{I''(\theta_0)}{n} = \frac{I(\theta_0)}{n} = \frac{1}{n} \sum_{i=1}^n I_i(\theta_0) \stackrel{d}{\longrightarrow} J_1(\theta_0).$$
(5)

) The second summand in the denominator tends in probability to 0 as $n \rightarrow \infty$, i.e.

$$\frac{1}{2n}(\hat{\theta}_{ML}-\theta_0)I^{\prime\prime\prime}(\theta^*) \xrightarrow{\rho} 0.$$
(6)

This is true because $\hat{\theta}_{ML} - \theta_0 \xrightarrow{p} 0$ (consistency of the ML-estimator) and $1/nI'''(\theta^*)$ is bounded in probability. The latter is due to regularity condition R5. More specifically,

$$\left|\frac{1}{n}I'''(\theta^*)\right| = \left|\frac{1}{n}\sum_{i=1}^n I''_i(\theta^*)\right| \leq \frac{1}{n}\sum_{i=1}^n M(X_i) \stackrel{d}{\longrightarrow} E(M(X)) < \infty.$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

- Properties 2 and 3 imply that the denominator in (4) converges to $J_1(\theta_0)$ in probability, while property 1 shows that the numerator converges in distribution to $N(0, J_1(\theta_0))$.
- Slutzky's theorem shows that the ratio converges in distribution to $N(0, 1/J_1(\theta_0))$, which completes the proof.
Likelihood based inference

Approximate Tests and confidence intervals

Wald-statistic

• We investigate the hypotheses

$$H_0: \theta = \theta_0$$
 versus $H_1: \theta \neq \theta_0$

• For testing we use

$$\mathsf{T}_{1} = \left(\hat{\theta}_{\mathsf{ML}} - \theta_{0}\right) \sqrt{\mathsf{I}(\hat{\theta}_{\mathsf{ML}})}$$

or

$$T_2 = \left(\hat{\theta}_{ML} - \theta_0\right) \sqrt{J(\hat{\theta}_{ML})}.$$

Wald-statistic

• Under H₀ we have

$$\begin{array}{rcl} T_1 & \stackrel{a}{\sim} & N(0,1) \\ T_2 & \stackrel{a}{\sim} & N(0,1) \end{array}$$

- Rejection:
 - We reject H₀ if

$$||\mathbf{T}_{1}|| > z_{1-\frac{\alpha}{2}}$$

respectively

$$\left|\mathsf{T}_{2}\right| > z_{1-\frac{\alpha}{2}}.$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Poisson distribution

 $H_0: \lambda = \lambda_0$ versus $H_1: \lambda \neq \lambda_0$

The various variants of the Fisher information are given by

$$I(\lambda) = \frac{n\bar{X}}{\lambda^2} \qquad J(\lambda) = \frac{n}{\lambda}$$
$$I(\hat{\lambda}_{ML}) = \frac{n\bar{X}}{\bar{X}^2} = \frac{n}{\bar{X}} \qquad J(\hat{\lambda}_{ML}) = \frac{n}{\bar{X}}.$$

Then

$$\mathsf{T}_1 = \mathsf{T}_2 = ig(ar{X} - \lambda_0ig)\sqrt{rac{n}{ar{X}}} \quad \stackrel{\mathrm{a}}{\sim} \quad N(0,1).$$

Binomial distribution

$$H_0: \pi = \pi_0 \quad \text{versus} \quad H_1: \pi \neq \pi_0$$
$$J(\pi) = \frac{n}{\pi(1 - \pi)}, \quad J(\hat{\pi}) = \frac{n}{\bar{X}(1 - \bar{X})}, \quad \bar{X} = \frac{X}{n}$$
$$T_2 = (\hat{\pi} - \pi_0) \sqrt{\frac{n}{\hat{\pi}(1 - \hat{\pi})}}$$

Wald confidence intervals

Because of the duality between tests and confidence intervals we obtain the following approximate confidence intervals:

$$\left[\hat{\theta}_{\mathsf{ML}} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{\mathsf{I}(\hat{\theta}_{\mathsf{ML}})}}\right] = \left[\hat{\theta}_{\mathsf{ML}} \pm z_{1-\frac{\alpha}{2}} \mathsf{se}(\hat{\theta}_{\mathsf{ML}})\right]$$

respectively

$$\left[\hat{\theta}_{\mathsf{ML}} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{J(\hat{\theta}_{\mathsf{ML}})}}\right]$$

Poisson distribution

An approximate confidence interval is given by

$$\left[\bar{X} \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{X}}{n}}\right]$$

Binomial distribution

An approximate confidence interval is given by

$$\left[\hat{\pi} \pm z_{1-\frac{\alpha}{2}} \frac{\hat{\pi}(1-\hat{\pi})}{n}\right]$$

as

$$J(\hat{\pi})=\frac{n}{\bar{X}(1-\bar{X})}.$$

Likelihood based inference

Likelihood ratio test

Likelihood ratio test

• Our goal is again to investigate the hypotheses

$$H_0: \theta = \theta_0$$
 versus $H_1: \theta \neq \theta_0$

- The idea of the likelihood ratio test is to compare $\ell(\hat{\theta}_{ML})$ with $\ell(\theta_0)$. If the difference $\ell(\hat{\theta}_{ML}) \ell(\theta_0)$ is "large", then θ_0 is not plausible and H_0 is rejected.
- The exact test statistic is

$$\mathsf{W} = \mathsf{2}\left(\ell(\hat{ heta}_{\mathsf{ML}}) - \ell(heta_0)
ight) = \mathsf{2}\lograc{\mathsf{L}(\hat{ heta}_{\mathsf{ML}})}{\mathsf{L}(heta_0)}.$$

• Under H₀ we have

$$W \stackrel{a}{\sim} \chi_1^2$$

and we reject H_0 if $W > \chi_1^2(1 - \alpha)$.

Poisson distribution

 $H_0: \lambda = \lambda_0$ versus $H_1: \lambda \neq \lambda_0$

$$I(\lambda) = \sum_{i=1}^{n} X_i \log(\lambda) - n\lambda = n\bar{X}\log(\lambda) - n\lambda$$
$$W = 2(I(\hat{\lambda}) - I(\lambda_0)) = 2(n\bar{X}\log\bar{X} - n\bar{X} - n\bar{X}\log(\lambda_0) + n\lambda_0)$$

Derivation of the asymptotic distribution of W

• To derive the asymptotic distribution of *W* under H₀ we compute the first order Taylor series expansion of the log-likelihood about θ_0 at $\hat{\theta}_{ML}$

$$I(\hat{\theta}_{ML}) = (\theta_0) + \left(\hat{\theta}_{ML} - \theta_0\right)I'(\theta_0) + \frac{1}{2}\left(\hat{\theta}_{ML} - \theta_0\right)^2I''(\theta^{**}) \quad (7)$$

with θ^{**} between $\hat{\theta}_{ML}$ and θ_0 .

• For the derivation of the asymptotic distribution of the ML-estimator we also computed the second order Taylor series expansion of $l'(\hat{\theta}_{ML}) = s(\hat{\theta}_{ML})$ given by (3). Because of $l'(\hat{\theta}_{ML}) = 0$ we can rearrange (3) to obtain

$$-I'(heta_0) = \left(\hat{ heta}_{ML} - heta_0
ight)I''(heta_0) + rac{1}{2}\left(\hat{ heta}_{ML} - heta_0
ight)^2I'''(heta^*).$$

• Inserting this into (7) yields

$$I(\hat{\theta}_{ML}) - I(\theta_{0}) = -(\hat{\theta}_{ML} - \theta_{0}) \left[(\hat{\theta}_{ML} - \theta_{0}) I''(\theta_{0}) + \frac{1}{2} (\hat{\theta}_{ML} - \theta_{0})^{2} I'''(\theta^{*}) \right] \\ + \frac{1}{2} (\hat{\theta}_{ML} - \theta_{0}) I''(\theta^{**}) \\ = -n (\hat{\theta}_{ML} - \theta_{0})^{2} \left[\frac{1}{n} I''(\theta_{0}) + \frac{1}{2n} (\hat{\theta}_{ML} - \theta_{0}) I'''(\theta^{*}) - \frac{1}{2n} I''(\theta^{**}) \right]$$

- We already know from (6) that the second summand in brackets tends to 0 in probability.
- The other two summands in brackets both tend to $-J_1(\theta_0)$ (see the derivation of (5)). Hence $2(I(\hat{\theta}_{ML}) I(\theta_0))$ has the same limit distribution as $(\hat{\theta}_{ML} \theta_0)^2 n J_1(\theta_0)$. The asymptotic normality of the ML-estimator provides

$$\left(\hat{\theta}_{ML}-\theta_{0}
ight)\sqrt{nJ_{1}(\theta_{0})}\overset{d}{\longrightarrow}N(0,1).$$

Invoking the continuous mapping theorem we finally arrive at

$$\left(\hat{\theta}_{ML}-\theta_{0}\right)^{2}$$
 nJ₁(θ_{0}) $\stackrel{d}{\longrightarrow}\chi_{1}^{2}$.

(Note that if X i.i.d. N(0, 1) then $X^2 \sim \chi_1^2$.)







Chapter 4

Likelihood inference for vector valued parameters

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck) Statistical Inference – 4 – Likelihood inference for vector valued parameters – 0/38

Likelihood inference for vector valued parameters

Situation

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck) Statistical Inference – 4 – Likelihood inference for vector valued parameters – 1/38

Situation

- Let $X_1,...,X_n$ be a random sample with probability function or density $f_i(x_i,\theta)$.
- The X_i 's are still independent but no longer identically distributed.
- Our goal is to estimate θ .

(Normal distribution)

The random sample $X_1, ..., X_n$ is assumed to be i.i.d. with $X_i \sim N(\mu, \sigma^2)$, i.e.

$$\boldsymbol{\theta} = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$$

Linear Model

• $Y_1,...,Y_n$ independent with $Y_i \sim N(\mu_i,\sigma^2)$ and

$$\mu_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik}.$$

• Define
$$\theta = (\beta_0, \beta_1, \dots, \beta_k, \sigma^2)'$$
.

• We assume $X_{i1},...,X_{ik}$ non stochastic.

Binary regression models

- We assume an independent random sample $Y_1, ..., Y_n$ with $Y_i \sim B(1, \pi_i)$ being Bernoulli distributed, i.e. $P(Y_i = 1) = \pi_i$ and $P(Y_i = 0) = 1 \pi_i$.
- Logit model:

$$\pi_{i} = \frac{\exp(\beta_{0} + \beta_{1}x_{i1} + \dots + \beta_{k}x_{ik})}{1 + \exp(\beta_{0} + \beta_{1}x_{i1} + \dots + \beta_{k}x_{ik})}$$
$$= \frac{\exp(\mathbf{x}_{i}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}_{i}'\boldsymbol{\beta})}$$
$$= \frac{\exp(\eta_{i})}{1 + \exp(\eta_{i})}$$

Binary regression models

Probit model

$$\pi_i = \Phi(\beta_0 + \cdots + \beta_k x_{ik}) = \Phi(\mathbf{x}'_i \beta) = \Phi(\eta_i),$$

where $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ik})'$, $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)'$ and $\eta_i = \mathbf{x}'_i \boldsymbol{\beta}$. • Here $\boldsymbol{\theta} = (\beta_0, \dots, \beta_k)'$.

• For the logit model we further obtain

$$g(\pi_i) = \log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik},$$

where g is called link function (here the logit-link).

Binary regression models

- The ratio log $\frac{\pi_i}{1-\pi_i}$ is called log-odds, which can be regarded as a linear combination of the covariates.
- For the odds ratio $\frac{\pi_i}{1-\pi_i}$ we have a multiplicative model, i.e.

$$\frac{\pi_i}{1-\pi_i} = \exp(\beta_0) \exp(\beta_1 x_{i1}) \cdots \exp(\beta_k x_{ik}).$$

Likelihood inference for vector valued parameters

Likelihood, score function and Fisher information

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck) Statistical Inference – 4 – Likelihood inference for vector valued parameters – 8/38

Likelihood

• The likelihood of the sample $X_1,...,X_n$ is given by

$$L(\boldsymbol{\theta}) = \prod_{i=1}^{n} f_i(X_i, \boldsymbol{\theta}) = \prod_{i=1}^{n} L_i(\boldsymbol{\theta})$$

with

$$L_i(\boldsymbol{\theta}) = f_i(X_i, \boldsymbol{\theta}).$$

• The log-likelihood is given by

$$I(\boldsymbol{\theta}) = \sum_{i=1}^{n} I_i(\boldsymbol{\theta}) = \sum_{i=1}^{n} \log(f_i(X_i, \boldsymbol{\theta})),$$

where

$$I_i(\boldsymbol{\theta}) = \log(f_i(X_i, \boldsymbol{\theta})).$$

Normal distribution

$$L_{i}(\mu, \sigma^{2}) = \frac{1}{\sqrt{2\pi\sigma^{2}}} \exp\left(-\frac{1}{2\sigma^{2}}(x_{i} - \mu)^{2}\right)$$

$$\propto \frac{1}{\sqrt{\sigma^{2}}} \exp\left(-\frac{1}{2\sigma^{2}}(x_{i} - \mu)^{2}\right)$$

$$l_{i}(\mu, \sigma^{2}) = \log(1) - \frac{1}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(x_{i} - \mu)^{2}$$

$$= -\frac{1}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(x_{i} - \mu)^{2}$$

$$l(\mu, \sigma^{2}) = \sum_{i=1}^{n} l_{i}(\mu, \sigma^{2})$$

Linear regression

$$L_{i}(\boldsymbol{\beta}, \sigma^{2}) = \frac{1}{\sqrt{2\pi\sigma^{2}}} \exp\left(-\frac{1}{2\sigma^{2}}(y_{i} - \boldsymbol{x}_{i}^{\prime}\boldsymbol{\beta})^{2}\right)$$
$$\propto \frac{1}{\sqrt{\sigma^{2}}} \exp\left(-\frac{1}{2\sigma^{2}}(y_{i} - \boldsymbol{x}_{i}^{\prime}\boldsymbol{\beta})^{2}\right)$$
$$l_{i}(\boldsymbol{\beta}, \sigma^{2}) = -\frac{1}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y_{i} - \boldsymbol{x}_{i}^{\prime}\boldsymbol{\beta})^{2}$$
$$l(\boldsymbol{\beta}, \sigma^{2}) = -\frac{1}{2}n\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}\sum_{i=1}^{n}(y_{i} - \boldsymbol{x}_{i}^{\prime}\boldsymbol{\beta})^{2}$$

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck) Statistical Inference – 4 – Likelihood inference for vector valued parameters – 11/38

Binary regression

We obtain

$$L_i(\beta) = \pi_i^{y_i} (1 - \pi_i)^{1 - y_i}$$

$$l_i(\beta) = y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i)$$

for the likelihood and log-likelihood.

Example (continued)

Binary regression

Because of

$$\pi_i = \frac{\exp(\mathbf{x}_i'\beta)}{1 + \exp(\mathbf{x}_i'\beta)}, \quad 1 - \pi_i = \frac{1}{1 + \exp(\mathbf{x}_i'\beta)}$$

and

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \mathbf{x}_i' \boldsymbol{\beta}.$$

we further obtain

$$\begin{split} l_i(\boldsymbol{\beta}) &= y_i \boldsymbol{x}_i' \boldsymbol{\beta} + \log \left(\frac{1}{1 + \exp(\boldsymbol{x}_i' \boldsymbol{\beta})} \right) \\ &= y_i \boldsymbol{x}_i' \boldsymbol{\beta} - \log(1 + \exp(\boldsymbol{x}_i' \boldsymbol{\beta})) \end{split}$$

in case of the logit model.

Score Function

$$S_{i}(\theta) = \left(\frac{\partial l_{i}(\theta)}{\partial \theta_{1}}, \dots, \frac{\partial l_{i}(\theta)}{\partial \theta_{p}}\right)'$$
$$S(\theta) = \sum_{i=1}^{n} S_{i}(\theta)$$

The ML-estimator is the solution to the following system of equations

 $S(\theta) = \mathbf{0}.$

^{© 2020} Stefan Lang (Dept. of Statistics, Universität Innsbruck) Statistical Inference - 4 - Likelihood inference for vector valued parameters - 14/38

Fisher Information

$$I_{i}(\boldsymbol{\theta}) = -\begin{pmatrix} \frac{\partial^{2} I_{i}(\boldsymbol{\theta})}{\partial \theta_{1} \partial \theta_{1}} & \cdots & \frac{\partial^{2} I_{i}(\boldsymbol{\theta})}{\partial \theta_{1} \partial \theta_{p}} \\ \vdots & \vdots \\ \frac{\partial^{2} I_{i}(\boldsymbol{\theta})}{\partial \theta_{p} \partial \theta_{1}} & \cdots & \frac{\partial^{2} I_{i}(\boldsymbol{\theta})}{\partial \theta_{p} \partial \theta_{p}} \end{pmatrix}$$
$$I(\boldsymbol{\theta}) = \sum_{i=1}^{n} I_{i}(\boldsymbol{\theta})$$

The observed Fisher-Information $I(\hat{\theta}_{ML})$ can be calculated by inserting the ML-estimator $\hat{\theta}_{ML}$. Expected Fisher-Information $J(\theta)$ is the expected value of $I(\theta)$.

i=1

Normal distribution

$$l_{i}(\mu, \sigma^{2}) = -\frac{1}{2} \log(\sigma^{2}) - \frac{1}{2\sigma^{2}} (x_{i} - \mu)^{2}$$
$$\frac{\partial l_{i}(\mu, \sigma^{2})}{\partial \mu} = -\frac{1}{2\sigma^{2}} 2(-1)(x_{i} - \mu) = \frac{1}{\sigma^{2}} (x_{i} - \mu)$$
$$\frac{\partial l_{i}(\mu, \sigma^{2})}{\partial \sigma^{2}} = -\frac{1}{2\sigma^{2}} + \frac{1}{2(\sigma^{2})^{2}} (x_{i} - \mu)^{2}$$
$$S\binom{\mu}{\sigma^{2}} = \left(\frac{\frac{1}{\sigma^{2}} \sum_{i=1}^{n} (x_{i} - \mu)}{-\frac{n}{2\sigma^{2}} + \frac{1}{2(\sigma^{2})^{2}} \sum_{i=1}^{n} (x_{i} - \mu)}\right)$$

Example (continued)

Normal distribution

The ML-estimator is given as the solution to

I:
$$\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) = 0$$

II: $-\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0$

Example (continued)

Normal distribution

From I we have $\hat{\mu} = \bar{x}$.

Inserting $\hat{\mu} = \bar{x}$ in II yields

$$-\frac{n}{\sigma^2} + \frac{1}{(\sigma^2)^2} \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

We obtain

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

Remark

- The ML-estimator for σ^2 in the previous example is biased.
- Alternatively to the ML-estimator for variance parameters, the so called Restricted ML-estimator (REML) is often used.
- It maximizes the marginal likelihood

$$\mathsf{RL}(\sigma^2) = \int_{-\infty}^{\infty} L(\mu, \sigma^2) \, \mathrm{d}\mu$$
$$= \int_{-\infty}^{\infty} \prod_{i=1}^{n} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x_i - \mu)^2\right) \right) \mathrm{d}\mu.$$

 REML estimators are used regularly for variance parameters, e.g. in linear models or linear mixed models, in the context of likelihood based inference.

Linear regression

$$\begin{split} l_i(\beta, \sigma^2) &= -\frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (y_i - \mathbf{x}'_i \beta)^2 \\ l(\beta, \sigma^2) &= -\frac{1}{2} n \log(\sigma^2) - \frac{1}{2\sigma^2} (Y - X\beta)' (Y - X\beta) \\ &= -\frac{1}{2} n \log(\sigma^2) - \frac{1}{2\sigma^2} (Y'Y - 2Y'X\beta + \beta'(X'X)\beta) \\ \frac{\partial l(\beta, \sigma^2)}{\partial \beta} &= -\frac{1}{2\sigma^2} (-2X'Y + 2X'X\beta) \\ &= -\frac{1}{\sigma^2} (X'X\beta - X'Y) \\ \frac{\partial l(\beta, \sigma^2)}{\partial \sigma^2} &= -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2(\sigma^2)^2} (Y - X\beta)' (Y - X\beta) \end{split}$$

Example (continued)

Linear regression

Hence the score function is given by

$$\mathsf{S}\begin{pmatrix}\boldsymbol{\beta}\\\sigma^2\end{pmatrix} = \begin{pmatrix} -\frac{1}{\sigma^2} (X' X \beta - X' Y) \\ -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (X - X \beta)' (Y - X \beta) \end{pmatrix}$$

To obtain the ML-estimator we solve

$$\mathsf{S}\begin{pmatrix}\boldsymbol{\beta}\\\sigma^2\end{pmatrix} = \begin{pmatrix}\mathbf{0}\\\mathbf{0}\end{pmatrix}$$
Example (continued)

Linear regression

We immediately obtain

$$\hat{\beta}_{\rm ML} = (X'X)^{-1}X'Y$$
$$\hat{\sigma}_{\rm ML}^2 = \frac{1}{n}(Y - X\hat{\beta})'(Y - X\hat{\beta}).$$

Again $\hat{\sigma}_{\rm ML}^2$ is biased. The REML estimator

$$\hat{\sigma}_{\text{REML}}^2 = \frac{1}{n-k-1} (Y - X\hat{\beta})' (Y - X\hat{\beta})$$

is unbiased and maximizes the marginal likelihood

$$\mathsf{RL}(\sigma^2) = \int \mathsf{L}\binom{\boldsymbol{\beta}}{\sigma^2} \, \mathsf{d}\boldsymbol{\beta}.$$

Example

Binary regression: Logit model

$$I_i(\boldsymbol{\beta}) = y_i \boldsymbol{x}_i' \boldsymbol{\beta} - \log(1 + \exp(\boldsymbol{x}_i' \boldsymbol{\beta})) = y_i \eta_i - \log(1 + \exp(\eta_i))$$

The individual score function is given by

$$s_i(\boldsymbol{\beta}) = \frac{\partial l_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{\partial l_i(\boldsymbol{\beta})}{\partial \eta_i} \frac{\partial \eta_i}{\partial \boldsymbol{\beta}}$$
$$= \left(y_i - \frac{1}{1 + \exp(\eta_i)} \exp(\eta_i) \right) \boldsymbol{x}_i$$
$$= (y_i - \pi_i) \boldsymbol{x}_i.$$

Example (continued)

Binary regression: Logit model

This implies

$$s(\boldsymbol{\beta}) = \sum_{i=1}^{n} (y_i - \pi_i) \boldsymbol{x}_i = \boldsymbol{X}'(\boldsymbol{y} - \boldsymbol{\pi}),$$

where

$$\boldsymbol{X} = \left(\begin{array}{ccccc} 1 & x_{11} & \ldots & x_{1k} \\ 1 & x_{21} & \ldots & x_{2k} \\ \vdots & \vdots & & \vdots \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \ldots & x_{nk} \end{array} \right)$$

is a design matrix and $\mathbf{y} = (y_1, \dots, y_n)'$ and $\pi = (\pi_1, \dots, \pi_n)'$. The ML-estimator solves the nonlinear equation system

$$s(oldsymbol{eta}) = oldsymbol{X}'ig(oldsymbol{y} - oldsymbol{\pi}ig) = oldsymbol{0}$$

numerically.

Example (continued)

Binary regression: Logit model

$$J(\boldsymbol{\beta}) = \mathsf{E}\left(-\frac{\partial^2 l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'}\right) = \mathsf{Cov}(\boldsymbol{S}(\boldsymbol{\beta})) = \mathsf{E}(\boldsymbol{S}(\boldsymbol{\beta})\boldsymbol{S}'(\boldsymbol{\beta})).$$

After some tedious calculations we finally obtain

$$J(\boldsymbol{\beta}) = \sum_{i=1}^{n} \boldsymbol{x}_{i} \boldsymbol{x}_{i}' \pi_{i} (1 - \pi_{i}).$$

Numerical computation of the ML-estimator

• In most practical situations

$$s(heta) = \mathbf{0}$$

can not be solved analytically. In this case we need numerical algorithms to determine the ML estimator.

• Here: Newton method and Fisher scoring.

Iterative solution: θ one dimensional

• Set
$$oldsymbol{ heta}=oldsymbol{ heta}^{(0)}$$

2 Fit at $\theta = \theta^{(0)}$ a tangent to $s(\theta)$. The tangent is given by

$$y = s\left(\theta^{(0)}\right) + \frac{\partial s\left(\theta^{(0)}\right)}{\partial \theta} \cdot \left(\theta - \theta^{(0)}\right)$$

Set $\theta = \theta^{(1)}$ as the root of the tangent, i.e.

$$\boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}^{(0)} - \left(\frac{\partial \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right)}{\partial \boldsymbol{\theta}}\right)^{-1} \cdot \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right) = \boldsymbol{\theta}^{(0)} + \mathsf{I}^{-1}\left(\boldsymbol{\theta}^{(0)}\right) \cdot \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right)$$

Proceed with step 2 thereby replacing $\theta^{(0)}$ by $\theta^{(1)}$. Iterate until the new solution $\theta^{(1)}$ and the previous solution do not change.









General multiparameter case $\theta = (\theta_1, \ldots, \theta_p)'$:

$$\boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}^{(0)} - \left(\frac{\partial \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right)}{\partial \boldsymbol{\theta}'}\right)^{-1} \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right) = \boldsymbol{\theta}^{(0)} + \left(\mathsf{I}\left(\boldsymbol{\theta}^{(0)}\right)\right)^{-1} \boldsymbol{s}\left(\boldsymbol{\theta}^{(0)}\right)$$

with

$$s(\theta) = \left(\frac{\partial l(\theta)}{\partial \theta_1}, \dots, \frac{\partial l(\theta)}{\partial \theta_p}\right)'$$
$$l(\theta) = -\frac{\partial s(\theta)}{\partial \theta'} = -\left(\begin{array}{ccc} \frac{\partial^2 l(\theta)}{\partial \theta_1 \partial \theta_1} & \dots & \frac{\partial^2 l(\theta)}{\partial \theta_1 \partial \theta_p} \\ \vdots & \vdots \\ \frac{\partial^2 l(\theta)}{\partial \theta_p \partial \theta_1} & \dots & \frac{\partial^2 l(\theta)}{\partial \theta_p \partial \theta_p} \end{array}\right)$$

Fisher scoring

Replace the Fisher information $I(\theta)$ by the expected Fisher Information $J(\theta) = E(I(\theta))$, i.e.

$$\boldsymbol{ heta}^{(1)} = \boldsymbol{ heta}^{(0)} + \left(J\left(\boldsymbol{ heta}^{(0)}
ight)
ight)^{-1} \boldsymbol{s}\left(\boldsymbol{ heta}^{(0)}
ight).$$

Advantage: In many cases numerically more favorable.

Properties of the ML-estimator

- a) The ML-estimator is asymptotically unbiased and consistent.
- b) Distributional properties

$$\hat{\boldsymbol{ heta}}_{ML} \stackrel{a}{\sim} N_{\mathcal{P}}\left(\boldsymbol{ heta}, J(\boldsymbol{ heta})^{-1}
ight)$$

 $J(\theta)$ can be substituted by $J(\hat{\theta}_{ML})$, $I(\theta)$ or $I(\hat{\theta}_{ML})$.

c) For $n \rightarrow \infty$ the ML-estimator converges to the Cramér - Rao bound:

$$Cov(\hat{ heta}_{ML}) = J(heta)^{-1}$$

Hence, the ML-estimator is efficient for $n \to \infty$.

Statistical tests

Linear hypotheses

$$H_0: \mathbf{C}\boldsymbol{\theta} = \boldsymbol{d}$$
 versus $H_1: \mathbf{C}\boldsymbol{\theta} \neq \boldsymbol{d}$,

where **C** is a $r \times p$ with full column rank.

Some special cases

• j-th parameter zero

$$H_0: heta_j = 0$$
 versus $H_1: heta_j \neq 0$
 $\mathbf{C} = (0 \dots \underbrace{1}_j \dots 0)$ $\mathbf{d} = 0$

• Parameter θ_l equal to θ_j

$$H_0: \theta_I = \theta_j$$
 versus $H_1: \theta_I \neq \theta_j$
 $\mathbf{C} = (0 \dots \underbrace{1}_{I} \dots \underbrace{-1}_{j} \dots 0)$ $\mathbf{d} = 0.$

Test principles

• Likelihood-ratio-test

 $ilde{ heta}$ the ML-estimator for heta under $extbf{C} heta = extbf{d}$.

$$W := 2(l(\hat{\theta}_{ML}) - l(\tilde{\theta})) \stackrel{a}{\sim} \chi_r^2$$

Wald–test

$$\mathcal{T} := (\mathbf{C} \hat{m{ heta}}_{ML} - m{d})' (\mathbf{C} \ \widehat{\mathit{Cov}(\hat{m{ heta}}_{ML})} \, \mathbf{C}')^{-1} (\mathbf{C} \hat{m{ heta}}_{ML} - m{d}) \stackrel{a}{\sim} \chi^2_r$$

Score–test

$$U := s(\tilde{\theta})' \widehat{Cov(\hat{\theta}_{ML})}^{-1} s(\tilde{\theta}) \stackrel{a}{\sim} \chi_r^2$$

Wald-Test (special case)

• Test the hypotheses

$$H_0: \theta_j = \theta_{j0}$$
 against $H_1: \theta_j \neq \theta_{j0}$

Possible test statistic

$$au = rac{\hat{ heta}_j - heta_{j0}}{oldsymbol{se}\left(\hat{ heta}_j
ight)} \stackrel{ extsf{a}}{\sim} extsf{N}(0,1)$$

Reject
$$H_0$$
 if $|T| > z_{1-\frac{\alpha}{2}}$.

Alternative test statistic

$$T^2 = rac{(\hat{ heta}_j - heta_{j0})^2}{se\left(\hat{ heta}_j
ight)^2} \stackrel{a}{\sim} \chi^2(1)$$

Reject H_0 if $T^2 > \chi_1^2(1 - \alpha)$.







Chapter 5

Model choice and variable selecion

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Model choice and variable selecion

Introduction

- In many applications, a large (potentially enormous) number of candidate predictor variables is available, and we face the challenge and decision as to which of these variables to include in the regression model.
- The following are two naive (but often practiced) approaches to the model selection problem:
 - Strategy 1: Estimate the most complex model which includes all available covariates.
 - Strategy 2: First estimate a model with all variables. Then, remove all insignificant variables from the model.

- We investigate the simulated data illustrated in Figure 8 a).
- The true model used for simulation is $y_i = -1 + 0.3x_i + 0.4x_i^2 0.8x_i^3 + \varepsilon_i$ with $\varepsilon_i \sim N(0, 0.07^2)$.
- The scatter plot suggests polynomial modeling of the relationship between *y* and *x* resulting in the regression model

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_l x_i^l + \varepsilon_i.$$

Figure 8 c) - e) show the estimated relationship for *l* = 1 (regression line), *l* = 2, and *l* = 5.

• Figure 8 f) additionally displays the mean squared error

$$\mathsf{MSE}(I) = \frac{1}{50} \sum_{i=1}^{50} (y_i - \hat{y}_i(I))^2$$

of the fitted models depending on the order of the polynomial (continuous line). Clearly, MSE(I) decreases monotonically with increased *I*.

• Figure 8 b) shows additionally simulated observations for every design point x_i , i = 1, ..., 50. We refer to this data set as the validation sample, whereas we refer to the first data set (used for estimation) as the training set.

- Figure 8 f) shows the mean squared error of ŷ^{*}_i for the data y^{*}_i (dashed line) in the validation set.
- Apparently, the fit to the new data is initially getting better with an increase of the polynomial order. However, from the polynomial order *I* = 3 onward, the fit is getting worse.
- The more complex the model, the better is the fit to the data that were used for estimation. However with new data resulting from the same data generating process, models that are too complex can cause a poorer fit.



Figure: Simulated training data y_i (panel a) and validation data y_i^* (panel b). Panels *c*)-*e*) show estimated polynomials based on the training set. Panel f) displays MSE(*I*) in relation to the polynomial degree (solid line). The dashed line shows MSE(*I*), if the estimated polynomials are used to predict the validation data y_i^* .

- Consider the n = 150 observations $(y_i, x_{i1}, x_{i2}, x_{i3}), i = 1, \dots, 150$, in the scatter plot matrix in Figure 9.
- The variables x_1 and x_3 are independent and uniformly distributed on [0,1]. The variable x_2 is defined as $x_2 = x_1 + u$, where u is also uniformly distributed on [0,1].
- The response variable y is simulated according to the model

$$y \mid x_1, x_2, x_3 \sim N(-1 + 0.3x_1 + 0.2x_3, 0.2^2).$$

The conditional mean of y is thus dependent on x_1 and x_3 , but not on x_2 .

- We first estimate a regression model with all available covariates x_1 , x_2 , and x_3 , see Table 1.
- Clearly, x_1 and x_2 are nonsignificant. If we followed strategy 2, we would not only eliminate the nonrelevant covariate x_2 , but also the relevant variable x_1 .
- If we estimate a correctly specified model with true predictor variables x₁ and x₃, we obtain the results shown in Table 2.
- When having a correct model specification, not only is *x*₃ significant, but so is the previously insignificant variable *x*₁.



Figure: Scatter plot matrix for the variables y, x_1 , x_2 and x_3 .

		Standard-			95% Confidence-
Variable	Coefficient	error	t-value	p-value	interval
intercept	-0.970	0.047	-20.46	< 0.001	-1.064 -0.877
<i>x</i> ₁	0.146	0.187	0.78	0.436	-0.224 0.516
<i>X</i> ₂	0.027	0.177	0.15	0.880	-0.323 0.377
<i>x</i> ₃	0.227	0.052	4.32	< 0.001	0.123 0.331

Table: Results for the model based on covariates x_1 , x_2 and x_3 .

		Standard-			95% Confidence-
Variable	Coefficient	error	t-value	p-value	interval
intercept	-0.967	0.039	-24.91	< 0.001	-1.042 -0.889
<i>x</i> ₁	0.173	0.055	3.17	0.002	0.065 0.281
<i>x</i> ₃	0.226	0.052	4.33	< 0.001	0.123 0.330

Table: Results for the correctly specified model based on covariates x_1 and x_3 .

Model choice and variable selecion

Theoretical insights

Theoretical insights

We focus on the following questions:

- Irrelevant Variables: What can be said about the bias and the variance of the least squares estimator, in the case that we include irrelevant variables in the model?
- Missing Variables: What can be said about the bias and the variance of the least squares estimator, if we omit relevant variables in the model?
- Prediction Quality: What effect does the model specification, more specifically the selected variables in the model, have on prediction?

- Consider a partition of the available explanatory variables $\mathbf{x} = (x_0, x_1, \dots, x_k)'$ with $x_0 \equiv 1$ into the subsets $\mathbf{x}_1 = (x_0, x_1, \dots, x_{k_1})'$ and $\mathbf{x}_2 = (x_{k_1+1}, \dots, x_k)'$.
- We look at the two models

$$m{y} = m{X}m{eta} + m{arepsilon} = m{X}_1m{eta}_1 + m{X}_2m{eta}_2 + m{arepsilon}$$

and

$$\boldsymbol{y} = \boldsymbol{X}_1 \boldsymbol{\beta}_1 + \boldsymbol{u}.$$

• The first model uses all available variables. The second model uses only the subset *x*₁.

• For the submodel we obtain the least squares estimators

$$\hat{oldsymbol{eta}} = (oldsymbol{X}'oldsymbol{X})^{-1}oldsymbol{X}'oldsymbol{y}$$
 and $ilde{oldsymbol{eta}}_1 = (oldsymbol{X}_1'oldsymbol{X}_1)^{-1}oldsymbol{X}_1'oldsymbol{y}$

respectively.

• For the estimator $ilde{oldsymbol{eta}}_1$ of the submodel, we obtain

$$\mathsf{E}(ilde{oldsymbol{eta}}_1)=oldsymbol{eta}_1+(oldsymbol{X}_1'oldsymbol{X}_1)^{-1}oldsymbol{X}_1'oldsymbol{X}_2oldsymbol{eta}_2$$

and

$$\operatorname{Cov}(\tilde{\boldsymbol{\beta}}_1) = \sigma^2 (\boldsymbol{X}_1' \boldsymbol{X}_1)^{-1}.$$

• For the estimator $\hat{oldsymbol{eta}}$ in the full model we obtain

$$\mathsf{E}(\hat{eta}) = eta \qquad \mathsf{Cov}(\hat{eta}) = \sigma^2 \, (\mathbf{X}'\mathbf{X})^{-1}.$$

We now investigate the following two situations:

- Missing variables: Even though the complete model y = Xβ + ε is correct, we mistakenly estimate the reduced model y = X₁β₁ + u. In this case we neglect the relevant variables x₂.
- *Irrelevant Variables:* Even though the reduced model $y = X_1\beta_1 + u$ is correct, we mistakenly estimate the full model $y = X\beta + \varepsilon$. In this case, we included irrelevant variables in the model. The variables in x_2 are redundant.

In the first case of missing variables the following applies:

- β_1 is biased. An exception is the case when $X'_1X_2 = 0$, i.e. every variable in X_1 is uncorrelated to every variable in X_2 .
- It can be shown that the difference $\operatorname{Cov}(\hat{\beta}_1) \operatorname{Cov}(\tilde{\beta}_1)$ of covariance matrices is positive semi-definite. This implies that the components of the estimator $\tilde{\beta}_1$ based on the submodel $\boldsymbol{y} = \boldsymbol{X}_1 \beta_1 + \boldsymbol{u}$ show a smaller variance than the corresponding components of the estimator $\hat{\beta}_1$ based on the correct model $\boldsymbol{y} = \boldsymbol{X}\beta + \boldsymbol{\varepsilon}$. Thus we have $\operatorname{Var}(\hat{\beta}_j) \geq \operatorname{Var}(\tilde{\beta}_j)$.
- It can be shown that situations exist, in which the components in $\tilde{\beta}_1$ based on the misspecified submodel actually show a smaller MSE than the components in $\hat{\beta}_1$, which are based on the full model, i.e. $MSE(\hat{\beta}_j) \ge MSE(\tilde{\beta}_j)$.

In the second case of irrelevant variables, we have:

- Even though irrelevant variables were considered, $\hat{\beta}$ is unbiased. Of course, the estimator $\tilde{\beta}_1$ based on the true model is also unbiased.
- It can be shown that the estimators for the components in β₁ based on β̂ have larger variance than based on β̃₁. Thus we have Var(β̂_j) ≥ Var(β̃_j). If the estimated model contains irrelevant variables, then the precision of the estimators decreases.

We can reach the following conclusion: Preferably, the specified model should not contain irrelevant covariates. Moreover, we should aim for a sparse model so that bias and variance, and thus MSE, are small.
- Next we take a look at prediction quality in linear models. Thereby, we do not necessarily assume that the model is correctly specified.
- We assume independent observations y_i, i = 1,..., n, with expectation E(y_i) = μ_i and variance Var(y_i) = σ².
- The variables $x_0 = 1, x_1, ..., x_k$ are available as potential regressors.
- We assume that a subset *M* ⊂ {0, 1, 2, ..., *k*} of the available variables will be used for estimation.
- The specified model is defined by the subset of included covariates with corresponding design matrix *X*_M.

• For the least squares estimator we obtain

$$\hat{\boldsymbol{\beta}}_{M} = (\boldsymbol{X}_{M}^{\prime}\boldsymbol{X}_{M})^{-1}\boldsymbol{X}_{M}^{\prime}\boldsymbol{y}.$$

• An estimator \hat{y}_M for the vector μ of means $\mu_i = E(y_i)$ is given by

$$\hat{\boldsymbol{y}}_M = \boldsymbol{X}_M \hat{\boldsymbol{\beta}}_M.$$

- We can view the estimator ŷ_{iM} also as a prediction for future observations y_{n+i} = μ_i + ε_{n+i}, i = 1,..., n, with given covariates x_{i1},..., x_{ik}.
- In the following, we derive a formula for the sum of the expected squared prediction errors $\sum E(y_{n+i} \hat{y}_{iM})^2$.

To do so, we need the following, easily verifiable, properties of \hat{y}_M :

• Expectation:

$$\mathsf{E}(\hat{\boldsymbol{y}}_M) = \boldsymbol{X}_M (\boldsymbol{X}_M' \boldsymbol{X}_M)^{-1} \boldsymbol{X}_M' \mathsf{E}(\boldsymbol{y}).$$

• Covariance matrix:

$$\operatorname{Cov}(\hat{\boldsymbol{y}}_M) = \sigma^2 \boldsymbol{X}_M (\boldsymbol{X}'_M \boldsymbol{X}_M)^{-1} \boldsymbol{X}'_M.$$

• Sum of the variances:

$$\sum_{i=1}^{n} \operatorname{Var}(\hat{y}_{iM}) = \sigma^{2} \operatorname{tr}(\boldsymbol{X}_{M}(\boldsymbol{X}_{M}^{\prime}\boldsymbol{X}_{M})^{-1}\boldsymbol{X}_{M}^{\prime}) = |\boldsymbol{M}| \sigma^{2},$$

where |M| represents the cardinal number of M, i.e. the number of the covariates included in the model. The sum of the variances increases as more covariates are included in the model.

• Sum of the mean squared errors (SMSE):

SMSE =
$$\sum_{i=1}^{n} E(\hat{y}_{iM} - \mu_i)^2$$

= $\sum_{i=1}^{n} E((\hat{y}_{iM} - \mu_{iM}) + (\mu_{iM} - \mu_i))^2$
= $\sum_{i=1}^{n} Var(\hat{y}_{iM}) + 2\sum_{i=1}^{n} E((\hat{y}_{iM} - \mu_{iM})(\mu_{iM} - \mu_i)) + \sum_{i=1}^{n} (\mu_{iM} - \mu_i)^2$
= $|M| \sigma^2 + \sum_{i=1}^{n} (\mu_{iM} - \mu_i)^2$.

Here we used $\mu_{iM} = E(\hat{y}_{iM})$ as an abbreviation for the expectation of the estimator \hat{y}_{iM} .

These properties provide us with the expected squared prediction error:

SPSE =
$$\sum_{i=1}^{n} E(y_{n+i} - \hat{y}_{iM})^2$$

= $\sum_{i=1}^{n} E((y_{n+i} - \mu_i) - (\hat{y}_{iM} - \mu_i))^2$
= $\sum_{i=1}^{n} (E(y_{n+i} - \mu_i)^2 - 2E((y_{n+i} - \mu_i)(\hat{y}_{iM} - \mu_i)) + E(\hat{y}_{iM} - \mu_i)^2)$
= $\sum_{i=1}^{n} E(y_{n+i} - \mu_i)^2 + \sum_{i=1}^{n} E(\hat{y}_{iM} - \mu_i)^2$
= $n\sigma^2 + SMSE$
= $n\sigma^2 + |M|\sigma^2 + \sum_{i=1}^{n} (\mu_{iM} - \mu_i)^2$

Note that in line 3 of the above derivation, the expectation for the cross product term can be written as the product of expectations due to the independence of \hat{y}_{iM} and y_{n+i} . This way the entire term becomes zero.

Thus, the expected squared prediction error can be decomposed into three additive terms:

- *Irreducible Prediction Error:* The first term $n \sigma^2$ depends on the error variance. Hence, it cannot be reduced, even by sophisticated inference techniques. This term is therefore referred to as the irreducible prediction error.
- *Variance:* The second term consists of the sum of variances $Var(\hat{y}_{iM})$ of the estimators \hat{y}_{iM} . This term can be manipulated through model choice. It becomes smaller as fewer variables are included in the model.
- Squared Bias: The last term $\sum (\mu_{iM} \mu_i)^2$ can be seen as a bias term. It consists of the squared bias of the estimator \hat{y}_{iM} for the expectation μ_i . This term can also be manipulated through model choice and becomes smaller as more variables are included in the model.

- The decomposition of the expected prediction error into an irreducible error, a variance term, and a squared bias term is not limited to linear models, but rather a *fundamental property of prediction* in all statistical models.
- The formula for SPSE shows a classical bias-variance trade-off.
- The more complex the model, the smaller the squared bias and the greater the variance. On the contrary, simpler models show a greater squared bias and in return for that a smaller variance.
- This bias-variance trade-off is not only characteristic for linear models, but for all statistical models.

Model choice and variable selecion

Model choice

Model choice: Minimizing SPSE

Estimate SPSE using new and independent data

- If in fact additional observations y_{n+i} are available, we are able to estimate SPSE = $\sum_{i=1}^{n} E(y_{n+i} \hat{y}_{iM})^2$ simply by $\widehat{\text{SPSE}} = \sum_{i=1}^{n} (y_{n+i} \hat{y}_{iM})^2$.
- In practice, it is usually not possible to use this approach, as it rarely happens that additional observations are collected.
- An alternative procedure is the following:
 - Randomly split the data into two parts, i.e. a test and a validation sample.
 - Use the test data set to estimate the specified model.
 - Use the validation set to assess the goodness-of-fit, i.e. for the estimation of SPSE.

Model choice: Minimizing SPSE

Estimate SPSE using existing data

- A naive estimator for SPSE would be the use of the squared sum of residuals $\sum (y_i \hat{y}_{iM})^2$.
- Note that this sum *underestimates* on average the expected prediction error, as it can be shown that

$$\mathsf{E}\left(\sum_{i=1}^{n}(y_i-\hat{y}_{iM})^2
ight)=\mathsf{SPSE}-2|M|\sigma^2.$$

• Thus a better estimate for SPSE is given by

$$\widehat{\mathsf{SPSE}} = \sum_{i=1}^n (y_i - \hat{y}_{iM})^2 + 2|M|\hat{\sigma}^2.$$

Model choice: Minimizing SPSE

- Accordingly, we choose a model that minimizes \widehat{SPSE} . In doing so we have to keep in mind that we always use the same estimator for $\hat{\sigma}^2$. Preferably, this estimator should be based on the full model with all available variables, in order to keep the bias in $\hat{\sigma}^2$ small.
- The criterion SPSE has the typical structure of many model choice criteria. It consists of two terms: The first term, the sum of squared residuals, measures the goodness-of-fit and becomes smaller the more complex the model becomes. The second term $2|M|\hat{\sigma}^2$ measures model complexity and becomes smaller as models become simpler.

Model choice: Corrected coefficient of determination

- When comparing different models, the use of the coefficient of determination is limited, since the coefficient of determination will always increase (never decrease) with the addition of a new covariate into the model.
- The *corrected coefficient of determination* adjusts for this problem, by including a correction term for the number of parameters.
- The corrected coefficient of determination is defined by

$$\bar{R}^2 = 1 - \frac{n-1}{n-p}(1-R^2).$$

We advice against its usage, since the "penalty" for newly included covariates appears to be too small. It can be shown that \overline{R}^2 already increases when a variable with a t-value greater than 1 is included in the model, implying we would include variables with a p-value of about 0.3.

Model choice: Mallow's C_p

• Mallow's C_p ("Complexity parameter") is defined by

$$C_p = rac{\sum_{i=1}^{n} (y_i - \hat{y}_{iM})^2}{\hat{\sigma}^2} - n + 2|M|.$$

- C_{ρ} can be understood as an estimate of SMSE/ σ^2 .
- Thus minimizing C_p produces the same optimal model as minimizing SPSE.

Model choice: Akaike information criterion

• AIC is defined by

$$AIC = -2 \cdot I(\hat{\beta}_M, \hat{\sigma}^2) + 2(|M| + 1),$$

where $I(\hat{\beta}_M, \hat{\sigma}^2)$ is the maximum value of the log-likelihood.

- Smaller values of the AIC correspond to a better model fit.
- In a linear model with Gaussian errors, we obtain

$$AIC = n \cdot \log(\hat{\sigma}^2) + 2(|M| + 1).$$

Model choice: Akaike information criterion

- Figure 10 plots AIC for the simulated data from the Example on page 237 as a function of the polynomial degree.
- AIC obtains a minimum for l = 2 resulting in a reasonable model, even though we do not obtain the polynomial order of the true model with l = 3.

Model choice: Akaike information criterion



Figure: AIC as a function of the polynomial degree for the simulated data of the example on page 237.

Model choice: Cross Validation

Cross validation is based on the following general principle:

- Partition the data set into r subsets $1, \ldots, r$, of similar size.
- Start with the first data set as validation set and use the combined remaining r - 1 data sets for parameter estimation. Based on the estimates, obtain predictions for the validation set and determine the sum of the squared prediction errors.
- Cycle through the partitions using the second, third, up to the *r*th data set as validation sample, and all other data sets for estimation. Determine the sum of squared prediction errors.
- Use the model with the smallest sum of squared prediction errors, where the final parameter estimates reflect all data. The partition into test and validation samples serves only to estimate the expected squared prediction error.

Model choice: "leave-one-out" Cross Validation

- "leave-one-out" cross validation uses all observations with the exception of *one* for the estimation of model parameters.
- We use this "leave-one-out" estimator to predict the deleted observation and to determine the squared prediction error.
- If we denote the "leave-one-out" estimator with \hat{y}_{iM}^{-i} , we obtain the cross validation score

$$CV = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_{iM}^{-i})^2.$$

It can be shown

$$\mathrm{CV} = \frac{1}{n} \sum_{i=1}^{n} \left(\frac{y_i - \hat{y}_{iM}}{1 - h_{iiM}} \right)^2,$$

where h_{iiM} are the diagonal elements of the hat matrix $\mathbf{H}_{\mathbf{M}} = \mathbf{X}_{\mathbf{M}} (\mathbf{X}'_{\mathbf{M}} \mathbf{X}_{\mathbf{M}})^{-1} \mathbf{X}'_{\mathbf{M}}.$

Model choice and variable selecion

Practical Use of Model Choice Criteria

We can use the various model choice criteria to select the most promising models from candidate models. We recommend the following approach:

- On the basis of scientific knowledge, perhaps gained from previous research, we obtain a preselection of potential models. The models can differ in the number of variables but also in model type (e.g. linear vs. nonlinear). The total number of potential models should be as small as possible.
- All potential models can now be assessed with the aid of one of the various model choice criteria. The summary of the results should not be restricted to the "best" model. As a rule, there are a number of competitive models having approximately equal model fit, differing only in small aspects from each other. These differences cause some uncertainty regarding the conclusions.

If the number of potential models is large, we can use the following partially heuristic methods:

- Complete Model Selection (All-Subset-Selection): In case that the number of covariates is smaller than about 40, we can determine the best model (in the sense of a model choice criterion) with the "leaps and bounds" algorithm. An implementation can be found e.g. in the R package leaps.
- *Forward Selection:* Based on a starting model, forward-selection includes one additional variable in every iteration of the algorithm. The variable which offers the greatest reduction of a preselected model choice criteria (*C*_p, AIC, CV) is chosen. The algorithm terminates if no further reduction is possible.

- Backward Elimination: Backward elimination starts with the full model containing all potential covariates. Subsequently, in every iteration, the covariate which provides the greatest reduction of the model choice criteria (C_p, AIC, CV) is eliminated from the model. The algorithm terminates if no further reduction is possible.
- Stepwise Selection: Stepwise selection is a combination of forward selection and backward elimination. In every iteration of the algorithm, the inclusion and the deletion of a variable are both possible.

- The listed procedures should not be confounded with an algorithm proposed by Efroymson in the 1960s, even though the approach is similar.
- In contrast to what has been proposed above, the Efroymson algorithm includes or excludes those variables in/from the model, which have the highest or lowest t-value.
- The procedure terminates when no variable that potentially needs to be included has a p-value of less than a previously fixed maximal p-value (e.g. 0.05) and when no variable that needs to be excluded has a p-value greater than a minimal p-value (e.g. 0.1).

This automatic procedure, which is implemented in all major statistical software packages, is often viewed as obscure among statisticians due to the following two reasons:

- Forward, backward, and stepwise selection usually provide different results. This also happens when using a global model choice criterion such as AIC. We can, however, compare the different selected models with the help of the global model choice criterion. When using the Efroymson approach, discrimination between the different models is impossible.
- The repetitive use of the *t*-test statistic, to assess whether or not a regression coefficient is different from zero, suggests exact tests. However, the *t*-test statistic does not follow a t-distribution under the null hypothesis, since during the selection process we do not test an arbitrary variable, but rather the variable with the *maximal* t-value.

- We illustrate the approaches for model choice using data from the sales price of pre-owned VW Golf automobiles.
- Our goal is to model the relationship between the sales price in 1000 Euro (variable *price*) and the five explanatory variables "age of the car in months" (*age*), "kilometer reading in 1000 km" (*kilometer*), "number of months until the next appointment with the Technical Inspection Agency" (*TIA*), "ABS brake yes/no" (*extras1*), and "sunroof yes/no" (*extras2*).

The plots of the Figure on page 281 suggest the following effects:

- We can assume a linear or monotonically decreasing nonlinear effect for the variables *age* and *kilometer*, which could be appropriately modeled using (orthogonal) polynomials of degree three or less.
- The variable *TIA* appears to either have no effect or a very weak linear effect on the average sales price.
- Cars with ABS (*extras1*) seem to be slightly more expensive than cars without ABS; the effect, however, remains arguable.
- We can attest to no difference in the average sales price for models with or without sunroof (*extras2*).
- All in all, there seems to be a relationship with age and the kilometer reading. The effects of the remaining variables appear doubtful.

- We first examine eight regression models (see Table 3), which do not differ in the modeling of the variables *age* and *kilometer*. For the remaining three regressors, all possible model combinations will be tested.
- Using the AIC criterion, we obtain the first model in Table 3 as the preliminary best model. Figure 11 displays the AIC values for the eight models under consideration. In addition, the AIC for a ninth model based on automatic variable selection is provided; see below.

- Since only five explanatory variables are available, we can even determine the AIC best model with the help of the "Leaps and Bounds" algorithm. This model attains an AIC value of 389.35. It differs from the current "best model", in that it only makes use of polynomials of second degree (not third) for the variables *age* and *kilometer*.
- Figure 11 shows that the obtained AIC for this ninth model is considerably smaller than the best AIC value of all the models that we examined so far.



Model	kilometer	age	extras1	extras2	TIA	AIC
	degree 3	degree 3		linear		
1	х	х				393.234 (1)
2	х	х	х			394.566 (2)
3	х	х		х		395.119 (4)
4	х	х			х	394.973 (3)
5	х	х	х	х		396.481 (6)
6	х	х	х		х	396.143 (5)
7	х	х		х	х	396.881 (7)
8	х	х	х	x	х	398.085 (8)

Table: *Prices of used cars: potential models. The values in brackets indicate the rank of the models according to AIC.*



Figure: Prices of used cars: AIC values for the potential models.

		Standard-	95% Confidence-				
Variable	Coefficient	error	t-value	p-value	interval		
intercept	3.397	0.056	60.220	< 0.001	3.285	3.508	
ageop1	-0.705	0.061	-11.470	< 0.001	-0.826	-0.584	
ageop2	0.187	0.057	3.270	0.001	0.074	0.300	
kilometerop1	-0.439	0.061	-7.170	< 0.001	-0.560	-0.318	
kilometerop2	0.141	0.057	2.460	0.015	0.028	0.254	

Table: Prices of used cars: estimation results for the best model according to AIC.



Figure: Prices of used cars: model 9 based on all-subset selection, effects of age and kilometer reading including partial residuals.







Statistical Inference

Chapter 6

Bayesian Inference

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Statistical Inference - 6 - Bayesian Inference - 0 / 33

Bayesian Inference

Basic Concepts

Basic Concepts

- The fundamental difference to likelihood-based inference is that the unknown parameters $\theta = (\theta_1, \dots, \theta_p)'$ are not considered as fixed, deterministic quantities but as random variables with a *prior distribution*.
- *Prior distribution:* Any (subjective) information about the unknown parameter θ is expressed by specifying a probability distribution $p(\theta)$ for θ . The prior describes the *degree of uncertainty* about the unknown parameters prior to the statistical analysis.
- Observation model: The observation model specifies the conditional distribution of observable quantities, that is the random sample variables $\mathbf{Y} = (Y_1, \dots, Y_n)'$, given the parameters. The p.d.f. or probability function of this conditional distribution is proportional to the likelihood $L(\theta)$ and will be denoted by $p(\mathbf{y} | \theta)$.
Basic Concepts

- Based on the prior and the observation model, Bayes' theorem determines the distribution of θ after the data are known through the statistical experiment, that is the conditional distribution of θ given the observations $\mathbf{y} = (y_1, \dots, y_n)'$.
- We obtain

$$p(\theta \mid \boldsymbol{y}) = \frac{p(\boldsymbol{y} \mid \theta) p(\theta)}{\int p(\boldsymbol{y} \mid \theta) p(\theta) d\theta} = c \cdot p(\boldsymbol{y} \mid \theta) p(\theta),$$

with the normalizing constant $c = [\int p(\mathbf{y} | \theta) p(\theta) d\theta]^{-1}$. This conditional distribution is called *posterior (distribution)*.

Example

Poisson Distribution

- Consider an i.i.d. sample Y_1, \ldots, Y_n from a Poisson distribution, i.e. $Y_i \sim Po(\lambda)$.
- The joint probability for the observed sample $\mathbf{y} = (y_1, \dots, y_n)'$ is

$$p(\mathbf{y} \mid \lambda) = \frac{1}{y_1! \cdots y_n!} \lambda^{\sum_{i=1}^n y_i} \exp(-n\lambda).$$

• We specify a gamma distribution with parameters *a* and *b* for λ , i.e. $\lambda \sim Ga(a, b)$. It follows that λ has p.d.f.

$$p(\lambda) = k \,\lambda^{a-1} \exp(-b\lambda)$$

with
$$k = \frac{b^a}{\Gamma(a)}$$
.

Poisson Distribution

• The posterior is obtained as

$$p(\lambda | \mathbf{y}) = \frac{p(\mathbf{y} | \lambda) p(\lambda)}{\int p(\mathbf{y} | \lambda) p(\lambda) d\lambda}$$

= $c \frac{1}{y_1! \cdots y_n!} \lambda \sum_{i=1}^n y_i \exp(-n\lambda) k \lambda^{a-1} \exp(-b\lambda).$

• To determine the type of this distribution, we can ignore all factors that do not depend on λ . This gives

$$p(\lambda \mid \mathbf{y}) \propto \lambda^{\sum_{i=1}^{n} y_i} \exp(-n\lambda) \lambda^{a-1} \exp(-b\lambda)$$
$$= \lambda^{a+\sum_{i=1}^{n} y_i - 1} \exp(-(b+n)\lambda).$$

Poisson Distribution

• This has the form of a gamma distribution with parameters $a' = a + \sum_{i=1}^{n} y_i$ and b' = b + n, i.e.

$$\lambda \mid \mathbf{y} \sim Ga\left(a + \sum_{i=1}^{n} y_i, b + n\right),$$

and the posterior has the same type of distribution as the prior.

• We call the prior as conjugate to the Poisson model because the posterior is of the same type as the prior.

• We consider a logit model with a single covariate *x*:

$$Y_i = B(1, \pi_i), \quad \pi_i = rac{\exp(\eta_i)}{1 + \exp(\eta_i)}, \quad \eta_i = \beta_0 + \beta_1 x_i, \quad i = 1, \dots, n.$$

• Assuming, as usual, (conditionally) independent response variables, the observation model is given by

$$p(\mathbf{y} | \boldsymbol{\beta}) \propto L(\boldsymbol{\beta}) = \prod_{i=1}^{n} \pi_i^{y_i} (1 - \pi_i)^{1-y_i},$$

where $\beta = (\beta_0, \beta_1)'$ is the vector of regression coefficients.

• Since estimated regression coefficients are often approximately normally distributed, it is reasonable to assume a two-dimensional normal prior, i.e.

 $p(m{eta}) \sim N_2(m{m},m{M})$

with prior mean m and prior covariance matrix M.

- If results from a previous statistical analysis are available, we could choose the previous point estimate as *m* and its estimated covariance matrix as *M*.
- If the previous analysis has been carried out some time ago, we may also multiply *M* with a factor *a* > 1 to express increased uncertainty.

- Increasing the variances in **M**, the normal prior becomes very flat and approximates a uniform distribution.
- In the limiting case the prior becomes proportional to a constant, i.e.

 $p(oldsymbol{eta}) \propto ext{const.}$

We also write $p(\beta) \propto 1$.

• The integral of this flat prior over \mathbb{R}^2 is not finite, so that $p(\beta)$ is not a density in the usual sense. Such a prior is called improper or diffuse.

- Such diffuse priors are admissible as long as the posterior, resulting from Bayes' theorem, is a proper distribution. i.e. its integral over ${\rm I\!R}^2$ is finite. In a Bayesian logit model this is the case if a finite MLE exists.
- With a flat, diffuse prior the posterior density is

$$p(\boldsymbol{\beta} \mid \boldsymbol{y}) = \frac{p(\boldsymbol{\beta})p(\boldsymbol{y} \mid \boldsymbol{\beta})}{\int p(\boldsymbol{\beta})p(\boldsymbol{y} \mid \boldsymbol{\beta}) d \boldsymbol{\beta}} \propto p(\boldsymbol{y} \mid \boldsymbol{\beta}) = \prod_{i=1}^{n} \pi_{i}^{y_{i}} (1 - \pi_{i})^{1 - y_{i}}.$$

• Although the posterior is proper, it has no known distributional type.

Bayesian Point Estimates

• The posterior mean is given by

$$\hat{\boldsymbol{ heta}} = EW(\boldsymbol{ heta} \mid \boldsymbol{y}) = \int \boldsymbol{ heta} \, p(\boldsymbol{ heta} \mid \boldsymbol{y}) \, d\boldsymbol{ heta} = c \cdot \int \boldsymbol{ heta} \, p(\boldsymbol{y} \mid \boldsymbol{ heta}) \, p(\boldsymbol{ heta}) \, d\boldsymbol{ heta}.$$

• The posterior mode is the value $\hat{\theta}$ that (globally) maximizes the posterior density, i.e.

$$\hat{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \boldsymbol{y}) = \operatorname*{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{y} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}).$$

The second expression shows that no integration is necessary to compute the posterior mode, because the normalizing constant is not needed.

• The posterior median, that is the median of the posterior distribution, is sometimes preferred to the posterior mean because it is more robust against outliers.

Poisson Distribution

• The posterior mean is

$$EW(\lambda \mid \mathbf{y}) = rac{a + \sum_{i=1}^{n} y_i}{b + n}.$$

- The smaller *a* (in relation to ∑ y_i) and *b* (in relation to *n*), the closer the posterior mean is to the usual MLE λ̂ = ȳ.
- The larger the prior information, i.e. the larger *a* and *b* are, the more the posterior mean and the MLE differ from each other.

Bayesian Interval Estimates

- For the posterior mean, a natural measure is the posterior variance.
- For the posterior median, the interquartile distance seems to be appropriate to measure its variability.
- In case of the posterior mode, the curvature of the posterior at the mode, i.e. the observed Fisher information, is a natural choice.
- Another way of assessing uncertainty are Bayesian confidence intervals or *credible intervals* or, more generally, *credible regions*: A region C ⊂ Θ of the parameter space is said to be a (1 − α)-credible region for θ if

$$\mathsf{P}(\boldsymbol{\theta} \in \boldsymbol{C} \,|\, \boldsymbol{y}) = 1 - \alpha.$$

If $\mathcal{C} \subseteq \mathbb{R}$ is an interval it is called credible interval.

• A credible region contains (at least) a probability mass $1 - \alpha$ of the posterior.

Bayesian Inference

Markov Chain Monte Carlo Methods

© 2020 Stefan Lang (Dept. of Statistics, Universität Innsbruck)

Statistical Inference - 6 - Bayesian Inference - 14 / 33

Basic Idea

- MCMC methods allow to draw samples from posterior distributions (and, in principle, from any distribution) that are usually not available analytically and to estimate characteristics of the posterior such as the mean, the variance or quantiles, or the posterior density itself.
- The most important advantage compared to more traditional methods of drawing a sample from a distribution, for example importance or rejection sampling, is that samples can be drawn from high-dimensional densities, even for dimensions in the thousands.
- Another advantage is that the normalizing constant, often a high-dimensional integral that cannot be computed with traditional numerical methods, does not have to be known.

MCMC Methods - Basis Idea

- Let θ be the unknown vector of parameters in a Bayesian model and $p(\theta \mid \mathbf{y})$ the posterior density (we assume here that θ is continuous).
- Instead of directly drawing an i.i.d. sample from $p(\theta | \mathbf{y})$, a Markov chain is generated such that the iterations of the transition kernel converge to the posterior of interest.
- In this way random numbers are generated that can be considered as a (correlated) sample from the posterior after some time of convergence, the *burn-in phase*.

Metropolis-Hastings Algorithm

To draw random numbers from the density $p(\theta | \mathbf{y})$, the Metropolis-Hastings algorithm proceeds as follows:

- Initialize $\theta^{(0)}$ and the number *T* of iterations. Set t = 1.
- 2 Draw a random number θ^* from the proposal density $q(\theta^* | \theta^{(t-1)})$ and accept it as the new state $\theta^{(t)}$ with probability $\alpha(\theta^* | \theta^{(t-1)})$, otherwise set $\theta^{(t)} = \theta^{(t-1)}$.
- Stop if t = T, otherwise set t = t + 1 and go to 2.

After a *burn-in phase* t_0 , the random numbers $\theta^{(t_0+1)}, \ldots, \theta^{(T)}$ can be considered as a (correlated) sample from $p(\theta | \mathbf{y})$.

• We consider the following simulated logit model with two covariates x_1 and x_2 :

$$Y_i = B(1, \pi_i) \quad i = 1, \dots, 500,$$

$$\pi_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)},$$

$$\eta_i = -0.5 + 0.6 x_{i1} - 0.3 x_{i2}.$$

- The covariates x_1 and x_2 are drawn independently from a standard normal distribution.
- We want to construct a Metropolis-Hastings algorithm to estimate the parameter $\beta = (-0.5, 0.6, -0.3)'$ given this simulated data.

- We specify independent diffuse priors $p(\beta_j) \propto \text{const.}$
- The posterior is then proportional to the likelihood:

$$p(\boldsymbol{\beta} \mid \boldsymbol{y}) \propto \prod_{i=1}^{500} \pi_i^{y_i} (1-\pi_i)^{1-y_i}.$$

- As a proposal density for the Metropolis-Hastings algorithm we choose a 3-dimensional normal distribution, with the current state $\beta^{(t-1)}$ as its mean.
- For its covariance matrix, we start with the diagonal matrix $\Sigma = diag(0.4^2, 0.4^2, 0.4^2)$.
- Figure 13 (first row) shows the first 2000 random numbers for β_0 and β_1 drawn with this proposal density.
- Since we have specified diffuse priors, Bayes estimates for the regression coefficients should not differ too much from the MLEs. Therefore, the MLEs are displayed as horizontal lines in the plots.



Figure: Sampling paths for β_0 and β_1 for different MH algorithms. The third column shows the respective autocorrelation functions for β_1 .



Figure: Sampling paths for β_0 and β_1 for different MH algorithms. The third column shows the respective autocorrelation functions for β_1 .

- Clearly, only a few of the proposed random numbers are accepted with this first algorithm, sometimes the state remains unchanged for more than 100 iterations.
- Thus, the acceptance probabilities are far too small.
- We obtain larger acceptance probabilities if the variances of the proposal density are decreased to $\Sigma = diag(0.1^2, 0.1^2, 0.1^2)$.
- The second row in Figure 13 shows the first 2000 random numbers for β_0 and β_1 resulting from this second MH algorithm.
- We recognize a short burn-in phase of about 50 iterations, followed by reasonable iterations with relatively large acceptance rates.

- If we further decrease the variance to $\Sigma = diag(0.02^2, 0.02^2, 0.02^2)$, acceptance rates are further increased, but successive draws remain almost in the same state, see the first row in Figure 14.
- A useful and important tool for assessing the quality of MCMC algorithms is the autocorrelation function of the sample.
- Ideally, autocorrelations should rapidly converge to zero with increasing lags. The smaller the autocorrelation of successive parameters, the better the characteristics of the posterior can be estimated, based on the same length *T* of the sample.

- For practical work, 'thinning' is carried out for the original sample, i.e. only every *k*th random number is kept in the sample, so that the remaining random numbers are almost uncorrelated. In this way, memory space can be saved without worsening estimation results.
- To generate an approximately uncorrelated sample of size 1000 with our second MH algorithm, we would have to generate about 20000 random numbers after a short burn-in phase and then keep only every 20th random number in the thinned sample.

We can conclude the following:

- Small variances of the proposal density lead to high acceptance rates.
- In contrast, acceptance rates become small for large variances.
- For very large or very small variances autocorrelations of successive random numbers are high.
- The art of designing good MH algorithms is therefore the choice of appropriate proposal densities that combine high acceptance rates with low autocorrelations.
- Furthermore, automated methods are desirable that do not require subjective tuning of parameters of the proposal density.

- An algorithm with these desirable properties is the MH algorithm based on IWLS proposals, see the last column of Fig. 14.
- Using this algorithm a Markov chain was generated and, after the burn-in phase, 20000 random numbers were drawn. Saving every 20th random number led to a thinned sample of size 1000. Based on this thinned sample all characteristics of interest of the posterior can be approximated.
- To approximate the posterior mean we compute the arithmetic means for the sample, resulting in $\hat{\beta} = (-0.64, 0.65, -0.38)'$.
- Estimation of credible intervals can be based on the quantiles of the sampled random numbers. For example, we obtain 95% credible intervals by choosing the 2.5% quantiles as lower and 97.5% quantiles as upper bounds. This results in the intervals [-0.87, -0.42], [0.52 0.78] and [-0.52, -0.26] for the sample generated in our example.