

Advanced dataanalysis and statistical modelling, week 2

The likelihood principle

Jan K. Møller

DTU Compute
Section for Dynamical Systems
Technical University of Denmark
`jkmo@dtu.dk`

February 2018

An Abstract

- Descriptive statistics, look at the data
- Estimator (A random variable)
- Estimate (a number, the realisation of the estimator)
- The likelihood function

$$L(\theta; x_1, \dots, x_n) = f_{\theta}(x_1, \dots, x_n)$$

- log-likelihood function ($l(\theta; \cdot) = \log L(\theta; \cdot)$)
- Maximum Likelihood Estimate

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

- Score function $S(\theta) = \frac{\partial l}{\partial \theta}$

- Properties of estimators
 - Bias ($E[\hat{\theta} - \theta]$)
 - Variance ($E[(\hat{\theta} - E(\hat{\theta}))^2]$)
 - Mean square error ($E[(\hat{\theta} - \theta)^2]$)
 - Consistency $\hat{\theta}_n \rightarrow \theta$ for $n \rightarrow \infty$
 - MLE's consistent and asymptotically normal
- Hypothesis testing (comparing nested models) and Test statistics
 - Likelihood ratio test
 - Wald test

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood

The beginning of likelihood theory

- Fisher (1922) identified the likelihood function as the key inferential quantity conveying all inferential information in statistical modelling including the uncertainty
- The Fisherian school offers a Bayesian-frequentist compromise

A motivating example

Suppose we toss a thumbtack (used to fasten up documents to a background) 10 times and observe that 3 times it lands point up. Assuming we know nothing prior to the experiment, what is the probability of landing point up, θ ?

- Binomial experiment with $y = 3$ and $n = 10$.
- $P(Y = 3; 10, 0.2) = 0.2013$
- $P(Y = 3; 10, 0.3) = 0.2668$
- $P(Y = 3; 10, 0.4) = 0.2150$

A motivating example

By considering $P_\theta(Y = 3)$ to be a function of the unknown parameter we have the *likelihood function*:

$$L(\theta) = P_\theta(Y = 3)$$

In general, in a Binomial experiment with n trials and y successes, the likelihood function is:

$$L(\theta) = P_\theta(Y = y) = \binom{n}{y} \theta^y (1 - \theta)^{n-y}$$

A motivating example

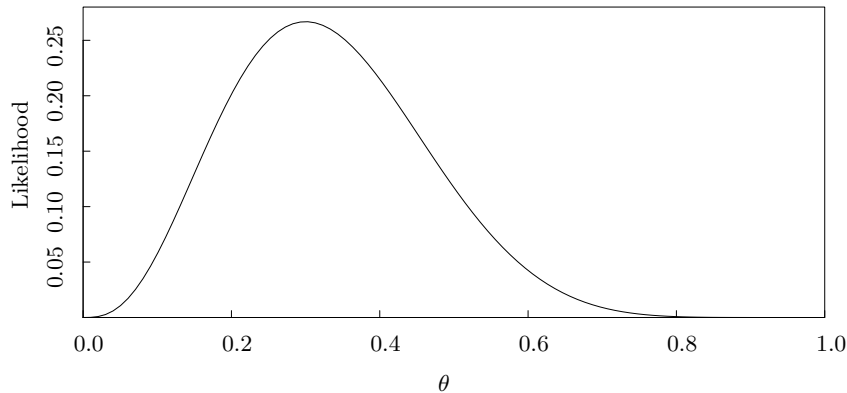


Figure: Likelihood function of the success probability θ in a binomial experiment with $n = 10$ and $y = 3$.

A motivating example

It is often more convenient to consider the log-likelihood function. The log-likelihood function is:

$$\log L(\theta) = y \log \theta + (n - y) \log(1 - \theta) + \text{const}$$

where *const* indicates a term that does not depend on θ .

By solving

$$\frac{\partial \log L(\theta)}{\partial \theta} = 0$$

it is readily seen that the maximum likelihood *estimate* (MLE) for θ is

$$\hat{\theta}(y) = \frac{y}{n} = \frac{3}{10} = 0.3$$

The likelihood principle

- Not just a method for obtaining a point estimate of parameters.
- It is the entire likelihood function that captures all the information in the data about a certain parameter.
- Likelihood based methods are inherently computational. In general numerical methods are needed to find the MLE.
- Today the likelihood principles play a central role in statistical modelling and inference.

Overview

- 1 The likelihood principle
- 2 Point estimation theory**
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood

Point estimation theory

We will assume that the statistical model for \mathbf{y} is given by a parametric family of joint densities:

$$\{f_{\mathbf{Y}}(y_1, y_2, \dots, y_n; \boldsymbol{\theta})\}_{\boldsymbol{\theta} \in \Theta^k}$$

Remember that when the n random variables are independent, the joint probability density equals the product of the corresponding marginal densities or:

$$f(y_1, y_2, \dots, y_n) = f_1(y_1) \cdot f_2(y_2) \cdot \dots \cdot f_n(y_n)$$

Point estimation theory

Definition (Unbiased estimator)

Any estimator $\hat{\theta} = \hat{\theta}(\mathbf{Y})$ is said to be *unbiased* if

$$\mathbb{E}[\hat{\theta}] = \theta$$

for all $\theta \in \Theta^k$.

Definition (Minimum mean square error)

An estimator $\hat{\theta} = \hat{\theta}(\mathbf{Y})$ is said to be *uniformly minimum mean square error* if

$$\mathbb{E} \left[(\hat{\theta}(\mathbf{Y}) - \theta)(\hat{\theta}(\mathbf{Y}) - \theta)^T \right] \leq \mathbb{E} \left[(\tilde{\theta}(\mathbf{Y}) - \theta)(\tilde{\theta}(\mathbf{Y}) - \theta)^T \right]$$

for all $\theta \in \Theta^k$ and all other estimators $\tilde{\theta}(\mathbf{Y})$.

Point estimation theory

Dispersion matrix

The matrix $\text{Var} \left[\hat{\boldsymbol{\theta}}(\mathbf{Y}) \right]$ is often called a variance covariance matrix since it contains variances in the diagonal and covariances outside the diagonal. This important matrix is often termed the *Dispersion matrix*.

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function**
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood

The likelihood function

- The likelihood function is built on an assumed parameterized statistical model as specified by a parametric family of joint densities for the observations $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$.
- The *likelihood* of any specific value $\boldsymbol{\theta}$ of the parameters in a model is (proportional to) the probability of the actual outcome, $Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n$, calculated for the specific value $\boldsymbol{\theta}$.
- The likelihood function is simply obtained by considering the likelihood as a function of $\boldsymbol{\theta} \in \Theta^k$.

The likelihood function

Definition (Likelihood function)

Given the parametric density $f_Y(\mathbf{y}, \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \Theta^P$, for the observations $\mathbf{y} = (y_1, y_2, \dots, y_n)$ the *likelihood function* for $\boldsymbol{\theta}$ is the function

$$L(\boldsymbol{\theta}; \mathbf{y}) = c(y_1, y_2, \dots, y_n) f_Y(y_1, y_2, \dots, y_n; \boldsymbol{\theta})$$

where $c(y_1, y_2, \dots, y_n)$ is a constant.

- Very often it is more convenient to consider the *log-likelihood* function defined as

$$l(\boldsymbol{\theta}; \mathbf{y}) = \log(L(\boldsymbol{\theta}; \mathbf{y})).$$

Example: Likelihood function for mean of normal distribution

An automatic production of a bottled liquid is considered to be stable. A sample of three bottles were selected at random from the production and the volume of the content volume was measured. The deviation from the nominal volume of 700.0 ml was recorded.

The deviations (in ml) were $y_1 = 4.6$; $y_2 = 6.3$; and $y_3 = 5.0$.

Example: Likelihood function for mean of normal distribution

First a *model* is formulated

- i Model: C+E (center plus error) model, $Y = \mu + \epsilon$
- ii Data: $Y_i = \mu + \epsilon_i$
- iii Assumptions:
 - Y_1, Y_2, Y_3 are independent
 - $Y_i \sim N(\mu, \sigma^2)$
 - σ^2 is known, $\sigma^2 = 1$,

Thus, there is only one unknown model parameter, $\mu_Y = \mu$.

Example: Likelihood function for mean of normal distribution

The joint probability density function for Y_1, Y_2, Y_3 is given by ¹

$$f_{Y_1, Y_2, Y_3}(y_1, y_2, y_3; \mu) = \frac{1}{(\sqrt{2\pi})^3} \exp \left[-\sum_{i=1}^3 \frac{(y_i - \mu)^2}{2} \right]$$

which for every value of μ is a function of the three variables y_1, y_2, y_3 . Now, we can establish the likelihood function

$$\begin{aligned} L_{4.6, 6.3, 5.0}(\mu) &= f_{Y_1, Y_2, Y_3}(4.6, 6.3, 5.0; \mu) \\ &= \frac{1}{(\sqrt{2\pi})^3} \exp \left[-\frac{(4.6 - \mu)^2}{2} - \frac{(6.3 - \mu)^2}{2} - \frac{(5.0 - \mu)^2}{2} \right] \end{aligned}$$

The function depends only on μ .

¹Remember that the normal probability density is: $f(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$

Example: Likelihood function for mean of normal distribution

Reducing the expression one finds

$$\begin{aligned} L_{4.6,6.3,5.0}(\mu) &= \frac{1}{(\sqrt{2\pi})^3} \exp\left[-\frac{1.58}{2}\right] \exp\left[-\frac{3(5.3 - \mu)^2}{2}\right] \\ &= \frac{1}{(\sqrt{2\pi})^3} \exp\left[-\frac{1.58}{2}\right] \exp\left[-\frac{3(\bar{y} - \mu)^2}{2}\right] \end{aligned}$$

which shows that (except for a factor not depending on μ), the likelihood function does only depend on the observations (y_1, y_2, y_3) through the average $\bar{y} = \sum y_i / 3$.

Example: Likelihood function for mean of normal distribution

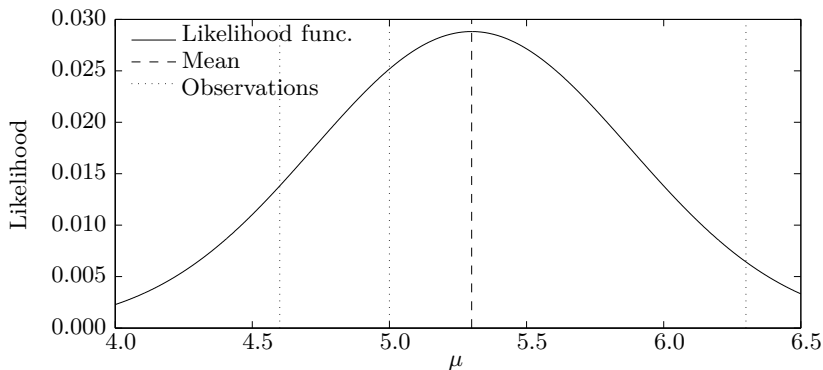


Figure: The likelihood function for μ given the observations $y_1 = 4.6$; $y_2 = 6.3$ and $y_3 = 5.0$.

Sufficient statistic

- The primary goal in analysing observations is to characterise the information in the observations by a few numbers.
- A *statistics* $t(Y_1, Y_2, \dots, Y_n)$ is a function of the observations.
- In estimation a sufficient statistic is a statistic that contains all the information in the observations. Formally

$$f_{\mathbf{Y}}(\mathbf{Y}; \theta) = h(\mathbf{y})g(t(\mathbf{y}); \theta)$$

with the factor $h(\mathbf{y})$ not depending on the parameter θ , and the factor $g(t(\mathbf{y}); \theta)$ only depending on \mathbf{y} through the (vector valued) function $t()$.

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix**
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood

The Score function

Definition (Score function)

Consider $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k) \in \Theta^k$, and assume that Θ^k is an open subspace of \mathbb{R}^k , and that the log-likelihood is continuously differentiable. Then consider the first order partial derivative (gradient) of the log-likelihood function:

$$l'_{\boldsymbol{\theta}}(\boldsymbol{\theta}; \mathbf{y}) = \frac{\partial}{\partial \boldsymbol{\theta}} l(\boldsymbol{\theta}; \mathbf{y}) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} l(\boldsymbol{\theta}; \mathbf{y}) \\ \vdots \\ \frac{\partial}{\partial \theta_k} l(\boldsymbol{\theta}; \mathbf{y}) \end{pmatrix}$$

The function $l'_{\boldsymbol{\theta}}(\boldsymbol{\theta}; \mathbf{y})$ is called the *score function* often written as $S(\boldsymbol{\theta}; \mathbf{y})$.

The information matrix

Definition (Observed information)

The matrix

$$\mathbf{j}(\boldsymbol{\theta}; \mathbf{y}) = - \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} l(\boldsymbol{\theta}; \mathbf{y})$$

with the elements

$$\mathbf{j}(\boldsymbol{\theta}; \mathbf{y})_{ij} = - \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\boldsymbol{\theta}; \mathbf{y})$$

is called the *observed information* corresponding to the observation \mathbf{y} , evaluated in $\hat{\boldsymbol{\theta}}$.

The observed information is thus equal to the Hessian (with opposite sign) of the log-likelihood function evaluated at $\boldsymbol{\theta}$. The Hessian matrix is simply (with opposite sign) the *curvature* of the log-likelihood function.

The information matrix

Definition (Expected information)

The expectation of the observed information

$$i(\theta) = E[j(\theta; \mathbf{Y})],$$

where the expectation is determined under the distribution corresponding to θ , is called the *expected information*, or the *information matrix* corresponding to the parameter θ . The expected information is also known as the *Fisher information matrix*

Example: Score function, Observed and Expected Information

In order to determine the expected information it is necessary to perform analogous calculations substituting the data by the corresponding random variables Y_1, Y_2, Y_3 .

The likelihood function can be written as

$$L_{y_1, y_2, y_3}(\mu) = \frac{1}{(\sqrt{2\pi})^3} \exp \left[-\frac{\sum (y_i - \bar{y})^2}{2} \right] \exp \left[-\frac{3(\bar{y} - \mu)^2}{2} \right].$$

Example: Score function, Observed and Expected Information

Introducing the random variables (Y_1, Y_2, Y_3) instead of (y_1, y_2, y_3) and taking logarithms one finds

$$l(\mu; Y_1, Y_2, Y_3) = -\frac{3(\bar{Y} - \mu)^2}{2} - 3\ln(\sqrt{2\pi}) - \frac{\sum(Y_i - \bar{Y})^2}{2},$$

and hence the score function is

$$l'_\mu(\mu; Y_1, Y_2, Y_3) = 3(\bar{Y} - \mu),$$

and the observed information

$$j(\mu; Y_1, Y_2, Y_3) = 3.$$

Example: Score function, Observed and Expected Information

It is seen in this (Gaussian) case that the observed information (curvature of log likelihood function) does not depend on the observations Y_1, Y_2, Y_3 , and hence the expected information is

$$i(\mu) = E[j(\mu; Y_1, Y_2, Y_3)] = 3.$$

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)**
- 6 Model selection
- 7 Profile likelihood

The Maximum Likelihood Estimate (MLE)

Definition (Maximum Likelihood Estimate (MLE))

Given the observation $\mathbf{y} = (y_1, y_2, \dots, y_n)$ the *Maximum Likelihood Estimate (MLE)* is a function $\hat{\boldsymbol{\theta}}(\mathbf{y})$ such that

$$L(\hat{\boldsymbol{\theta}}; \mathbf{y}) = \sup_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}; \mathbf{y})$$

The function $\hat{\boldsymbol{\theta}}(\mathbf{Y})$ over the sample space of observations \mathbf{Y} is called an *ML estimator*.

In practice it is convenient to work with the log-likelihood function $l(\boldsymbol{\theta}; \mathbf{y})$.

The Maximum Likelihood Estimate (MLE)

The *score function* can be used to obtain the estimate, since the MLE can be found as the solution to

$$l'_{\theta}(\boldsymbol{\theta}; \mathbf{y}) = 0$$

which are called the *estimation equations for the ML-estimator*, or, just the ML equations.

- It is common practice, especially when plotting, to normalize the likelihood function to have unit maximum and the log-likelihood to have zero maximum.

Invariance property

Theorem (Invariance property)

Assume that $\hat{\theta}$ is a maximum likelihood estimator for θ , and let $\psi = \psi(\theta)$ denote a one-to-one mapping of $\Omega \subset \mathbb{R}^k$ onto $\Psi \subset \mathbb{R}^k$. Then the estimator $\psi(\hat{\theta})$ is a maximum likelihood estimator for the parameter $\psi(\theta)$.

Example: Invariance property

Consider the binomial example

$$l(p) = y \log p + (n - y) \log(1 - p) + \text{const}$$

The optimisation needs to be bounded ($p \in [0, 1]$), an alternative parametrisation would be

$$p(\theta) = \frac{1}{1 + e^{-\theta}}$$

here $p(-\infty) = 0$ and $p(\infty) = 1$, hence θ need not to be bounded.

Distribution of the ML estimator

Theorem (Distribution of the ML estimator)

We assume that $\hat{\boldsymbol{\theta}}$ is consistent. Then, under some regularity conditions,

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \rightarrow N(0, \mathbf{i}(\boldsymbol{\theta})^{-1})$$

where $\mathbf{i}(\boldsymbol{\theta})$ is the expected information or the information matrix.

- The results can be used for inference under very general conditions. As the price for the generality, the results are only asymptotically valid.
- The practical significance of this result is that the MLE makes efficient use of the available data for large data sets.

Distribution of the ML estimator

In practice, we would use

$$\hat{\theta} \sim N(\theta, j^{-1}(\hat{\theta}))$$

where $j(\hat{\theta})$ is the observed (Fisher) information.

This means that asymptotically

- i) $E[\hat{\theta}] = \theta$
- ii) $D[\hat{\theta}] = j^{-1}(\hat{\theta})$

Distribution of the ML estimator

- The standard error of $\hat{\theta}_i$ is given by

$$\hat{\sigma}_{\hat{\theta}_i} = \sqrt{\text{Var}_{ii}[\hat{\theta}]}$$

where $\text{Var}_{ii}[\hat{\theta}]$ is the i 'th diagonal term of $\mathbf{j}^{-1}(\hat{\theta})$

- Hence we have that an estimate of the dispersion (variance-covariance matrix) of the estimator is

$$D[\hat{\theta}] = \mathbf{j}^{-1}(\hat{\theta})$$

- An estimate of the uncertainty of the individual parameter estimates are obtained by decomposing the dispersion matrix as follows:

$$D[\hat{\theta}] = \hat{\sigma}_{\hat{\theta}} \mathbf{R} \hat{\sigma}_{\hat{\theta}}$$

into $\hat{\sigma}_{\hat{\theta}}$, which is a diagonal matrix of the standard deviations of the individual parameter estimates, and \mathbf{R} , which is the corresponding correlation matrix.

The Wald Statistic

A test of an individual parameter

$$\mathcal{H}_0 : \theta_i = \theta_{i,0}$$

is given by the *Wald statistic*:

$$Z_i = \frac{\hat{\theta}_i - \theta_{i,0}}{\hat{\sigma}_{\hat{\theta}_i}}$$

which under \mathcal{H}_0 is approximately $N(0, 1)$ -distributed.

The $(1 - \alpha)$ Wald confidence interval is therefore

$$CI_{\theta_i} = \hat{\theta}_i \pm z_{1-\alpha/2} \hat{\sigma}_{\theta_i}$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile in the standard normal distribution.

Quadratic approximation of the log-likelihood

- A second-order Taylor expansion around $\hat{\theta}$ provides us with a quadratic approximation of the normalized log-likelihood around the MLE.
- A second-order Taylor's expansion around $\hat{\theta}$ we get

$$l(\theta) \approx l(\hat{\theta}) + l'(\hat{\theta})(\theta - \hat{\theta}) - \frac{1}{2}j(\hat{\theta})(\theta - \hat{\theta})^2$$

and then

$$\log \frac{L(\theta)}{L(\hat{\theta})} \approx -\frac{1}{2}j(\hat{\theta})(\theta - \hat{\theta})^2$$

- In the case of normality the approximation is exact which means that a quadratic approximation of the log-likelihood corresponds to normal approximation of the $\hat{\theta}(\mathbf{Y})$ estimator.

Example: Quadratic approximation of the log-likelihood

Consider again the thumbtack example.

The log-likelihood function is:

$$l(\theta) = y \log \theta + (n - y) \log(1 - \theta) + \text{const}$$

The score function is:

$$l'(\theta) = \frac{y}{\theta} - \frac{n - y}{1 - \theta},$$

and the observed information:

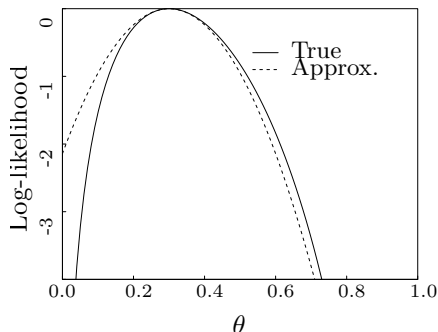
$$j(\theta) = \frac{y}{\theta^2} + \frac{n - y}{(1 - \theta)^2}.$$

For $n = 10$, $y = 3$ and $\hat{\theta} = 0.3$ we obtain

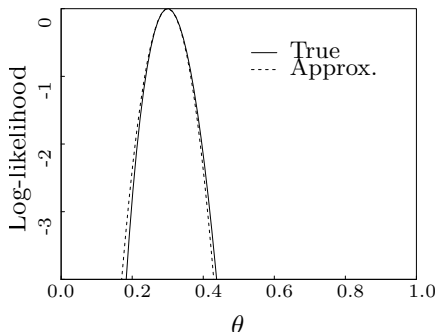
$$j(\hat{\theta}) = 47.6$$

Example: Quadratic approximation of the log-likelihood

The quadratic approximation is poor in this case. By increasing the sample size to $n = 100$, but still with $\hat{\theta} = 0.3$ the approximation is much better.



(a) $n = 10, y = 3$



(b) $n = 100, y = 30$

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection**
- 7 Profile likelihood

Example: Poisson regression

Consider the data

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|
| x | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| y | 3 | 0 | 4 | 5 | 6 | 4 | 9 | 7 | 4 | 10 |

the following regression model is proposed

$$Y_i \sim \text{Pois}(\lambda_i); \quad \lambda_i = \exp\{\theta_0 + \theta_1 x_i\}$$

Fit the parameters and estimate the uncertainty of the parameters.

Likelihood ratio tests

- Methods for testing hypotheses using the likelihood function.
- The basic idea: determine the maximum likelihood estimates under both a null and alternative hypothesis.
- It is assumed that a sufficient model with $\theta \in \Omega$ exists.
- Then consider some theory or assumption about the parameters $\mathcal{H}_0 : \theta \in \Omega_0$ where $\Omega_0 \subset \Omega$; $\dim(\Omega_0) = r$ and $\dim(\Omega) = k$
- The purpose of the testing is to analyze whether the observations provide sufficient evidence to reject this theory or assumption. If not we accept the null hypothesis.

Likelihood ratio tests

Definition (Likelihood ratio)

Consider the hypothesis $\mathcal{H}_0 : \boldsymbol{\theta} \in \Omega_0$ against the alternative $\mathcal{H}_1 : \boldsymbol{\theta} \in \Omega \setminus \Omega_0$ ($\Omega_0 \subseteq \Omega$), where $\dim(\Omega_0) = r$ and $\dim(\Omega) = k$.

For given observations y_1, y_2, \dots, y_n the *likelihood ratio* is defined as

$$\lambda(\mathbf{y}) = \frac{\sup_{\boldsymbol{\theta} \in \Omega_0} L(\boldsymbol{\theta}; \mathbf{y})}{\sup_{\boldsymbol{\theta} \in \Omega} L(\boldsymbol{\theta}; \mathbf{y})}$$

- If λ is small, then the data are seen to be more plausible under the alternative hypothesis than under the null hypothesis.
- Hence the hypothesis (\mathcal{H}_0) is rejected for small values of λ .

Likelihood ratio tests

- It is sometimes possible to transform the likelihood ratio into a statistic, the exact distribution of which is known under \mathcal{H}_0 . This is for instance the case for the General Linear Model for Gaussian data.
- In most cases, however, we must use the following important result regarding the asymptotic behavior:

Theorem (Wilk's Likelihood Ratio test)

For $\lambda(\mathbf{y})$ as above, then under the null hypothesis \mathcal{H}_0 , the random variable $-2 \log \lambda(\mathbf{Y})$ converges in law to a χ^2 random variable with $(k - r)$ degrees of freedom, i.e.,

$$-2 \log \lambda(\mathbf{Y}) \rightarrow \chi^2(k - r)$$

under \mathcal{H}_0 .

$-2 \log \lambda(\mathbf{Y})$ is also called the deviance and is denoted D

Example: Linear regression

Assume we want to the the following hypothesis

$$\mathcal{H}_0 : Y_i = \beta_0 + \epsilon_i; \quad \epsilon_i \sim N(0, \sigma^2)$$

against the alternative

$$\mathcal{H}_1 : Y_i = \beta_0 + \beta_1 x_i + \epsilon_i; \quad \epsilon_i \sim N(0, \sigma^2)$$

Also assume that σ^2 is known. The log-likelihood in each of the two cases is

$$\begin{aligned} l_0(\hat{\beta}_0) &= -\frac{1}{2\sigma^2} \sum (y_i - \hat{\beta}_0)^2 &= -\frac{1}{2\sigma^2} \sum e_{0,i}^2 \\ l_1(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{1}{2\sigma^2} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 &= -\frac{1}{2\sigma^2} \sum e_{1,i}^2 \end{aligned}$$

Example: Linear regression, cont.

Assuming \mathcal{H}_0 is true, then

$$2l_0 \sim \chi^2(n-1); \quad 2l_1 \sim \chi^2(n-2)$$

and the deviance between the models is

$$D = 2(l_1 - l_0) \sim \chi^2(1)$$

Note that

- In this case the test is exact (normal errors and known variance)
- \mathcal{H}_0 is rejected if D is large, i.e. the error in \mathcal{H}_0 are small compared to \mathcal{H}_1 .
- When σ^2 is not known, exact distribution are also available (see chap. 3)

Hypothesis chains

Consider a *chain* of hypotheses specified by a sequence of parameter spaces

$$\mathbb{R} \subseteq \Omega_M \dots \subset \Omega_2 \subset \Omega_1 \subset \mathbb{R}^n.$$

For each parameter space Ω_i we define a hypothesis

$$\mathcal{H}_i : \boldsymbol{\theta} \in \Omega_i$$

with $\dim(\Omega_i) < \dim(\Omega_{i-1})$.

Partial likelihood ratio test

Definition (Partial likelihood ratio test)

Assume that the hypothesis \mathcal{H}_i allows the sub hypothesis $\mathcal{H}_{i+1} \subset \mathcal{H}_i$. The *partial likelihood ratio test* for \mathcal{H}_{i+1} under \mathcal{H}_i is the likelihood ratio test for the hypothesis \mathcal{H}_{i+1} under the assumption that the hypothesis \mathcal{H}_i holds. The likelihood ratio test statistic for this partial test is

$$\lambda_{\mathcal{H}_{i+1}|\mathcal{H}_i}(y) = \frac{\sup_{\boldsymbol{\theta} \in \Omega_{i+1}} L(\boldsymbol{\theta}; \mathbf{y})}{\sup_{\boldsymbol{\theta} \in \Omega_i} L(\boldsymbol{\theta}; \mathbf{y})}$$

When \mathcal{H}_{i+1} holds, the distribution of $\lambda_{\mathcal{H}_{i+1}|\mathcal{H}_i}(Y)$ approaches a $\chi^2(f)$ distribution with $f = \dim(\Omega_i) - \dim(\Omega_{i+1})$.

Partial tests

Theorem (Partitioning into a sequence of partial tests)

Consider a chain of hypotheses.

Now, assume that \mathcal{H}_1 holds, and consider the minimal hypotheses $\mathcal{H}_M : \theta \in \Omega_M$ with the alternative $\mathcal{H}_1 : \theta \in \Omega_1 \setminus \Omega_M$. The likelihood ratio test statistic $\lambda_{\mathcal{H}_M|\mathcal{H}_1}(y)$ for this hypothesis may be factorized into a chain of partial likelihood ratio test statistics $\lambda_{\mathcal{H}_{i+1}|\mathcal{H}_i}(y)$ for \mathcal{H}_{i+1} given \mathcal{H}_i , $i = 1, \dots, M$.

- The partial tests "corrects" for the effect of the parameters that are in the model at that particular stage
- When interpreting the test statistic corresponding to a particular stage in the hierarchy of models, one often says that there is "controlled for", or "corrected for" the effect of the parameters that are in the model at that stage.

Example: likelihood ratio test

Consider the data

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|
| x | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| y | 3 | 0 | 4 | 5 | 6 | 4 | 9 | 7 | 4 | 10 |

the following regression model is proposed

$$Y_i \sim \text{Pois}(\lambda_i); \quad \lambda_i = \theta_0 + \theta_1 x_i$$

Test the hypothesis $\mathcal{H}_1 : \theta_1 \neq 0$

Variable selection in hypothesis chains

In-sample methods for model selection

The model complexity is evaluated using the same observations as those used for estimating the parameters of the model.

- The *training data* is used for evaluating the performance of the model.
- Any extra parameter will lead to a reduction of the loss function.
- In the in-sample case statistical tests are used to assess the significance of extra parameters, and when the improvement is small in some sense the parameters are considered to be non-significant.
- The classical statistical approach.
- *test set*: used for assessing the generalized performance, i.e. the performance on new data

Variable selection in hypothesis chains

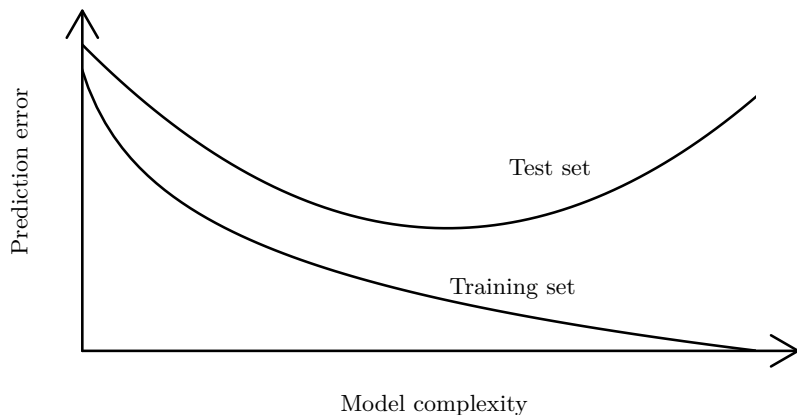


Figure: A typical behavior of the (possibly generalized) training and test prediction error as a function of the model complexity.

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood**

Profile likelihood

Definition (Profile likelihood)

Assume that the statistical model for the observations Y_1, Y_2, \dots, Y_n is given by the family of joint densities, $\theta = (\tau, \zeta)$ and τ denoting the parameter of interest. Then the *profile likelihood function* for τ is the function

$$L_P(\tau; \mathbf{y}) = \sup_{\zeta} L((\tau, \zeta); \mathbf{y})$$

where the maximization is performed at a fixed value of τ .

Profile likelihood - Confidence intervals

Consider again $\theta = (\tau, \zeta)$. From the likelihood ratio test it is seen that

$$\left\{ \tau; \frac{L_P(\tau; y)}{L(\hat{\theta}, y)} > \exp\left(-\frac{1}{2}\chi_{1-\alpha}^2(p)\right) \right\} \quad (1)$$

defines a set of values of τ (p -dimensional) that constitutes a $100(1 - \alpha)\%$ confidence region for τ .

- For normal distributions the confidence is exact and else it is an approximation.
- Likelihood based confidence intervals have the advantage of possibly being asymmetric when this is relevant.

Overview

- 1 The likelihood principle
- 2 Point estimation theory
- 3 The likelihood function
- 4 The information matrix
- 5 The maximum likelihood estimate (MLE)
- 6 Model selection
- 7 Profile likelihood