

Likelihood theory: Maximum likelihood estimation

(An overview)

N. Torelli, G. Di Credico, V. Gioia

Fall 2023

University of Trieste

The likelihood function¹

Maximum likelihood estimation: theory

Some numerical aspects

¹Agresti, Kateri: sec 4.2

The likelihood function

The likelihood function

IS A FUNDAMENTAL IDEA IN STATISTICS USED TO ESTIMATE PARAMETERS OF A PROBABILITY DISTRIBUTION BASED ON OBSERVED DATA. IT QUANTIFIES HOW LIKELY IT IS TO OBSERVE THE GIVEN DATA, ASSUMING THE MODEL AND ITS PARAMETERS ARE CORRECT.

Introduced by Sir Ronald Fisher, the **likelihood function** for a certain statistical model $f_\theta(y)$ for the data y is given by the following function of the parameter θ

WE HAVE TO ASSUME A PARAMETRIC MODEL BECAUSE THE LIKELIHOOD IS A FUNCTION OF THE PARAMETER

$$L : \Theta \rightarrow \mathbb{R}^+$$

$$\theta \rightarrow c(y) f_\theta(y), \quad L(\theta; y) = P(X; \theta)$$

IT'S THE PROBABILITY (OR DENSITY) OF THE DATA X , VIEWED AS A FUNCTION OF θ

where $c(y) > 0$ is an arbitrary constant of proportionality.

We may write $L(\theta; y)$ to stress the fact that the data enter the function, though its argument is given by θ .

IF WE PLOT $L(\theta; y)$, THE SHAPE GIVES INSIGHT INTO HOW "PLAUSIBLE" EACH PARAMETER VALUE θ IS.

A SHARP PEAK INDICATES HIGH CONFIDENCE IN $\hat{\theta}_{MLE}$, WHILE A FLAT LIKELIHOOD SUGGESTS UNCERTAINTY.

THE MAIN CONCEPT:

IMAGINE YOU HAVE A STATISTICAL MODEL $P(X; \theta)$ WHERE X REPRESENTS THE OBSERVED DATA AND θ THE PARAMETERS OF THE MODEL
=> THE LIKELIHOOD FUNCTION IS DEFINED AS:

Interpreting the likelihood function

The likelihood function assigns support (credibility) to possible values of θ , meaning that if $L(\theta_1) > L(\theta_2)$ then θ_1 is more supported by the observed data than θ_2 .

THE HIGHER THE VALUE, THE HIGHER IS THE TRUST IN THE PARAMETER

So the *likelihood ratio* $L(\theta_1)/L(\theta_2)$ allows for the comparison between θ_1 and θ_2 ; note that the constant $c(\mathbf{y})$ cancels out.

A mathematical justification for the above interpretation is given by the **Wald inequality**: if θ_t is the **true parameter value**, then

$$E_{\theta_t} \{\log L(\theta_t; \mathbf{Y})\} > E_{\theta_t} \{\log L(\theta; \mathbf{Y})\} \quad \theta \neq \theta_t.$$

The above fact can be proven by straightforward application of the Jensen's inequality.

LET'S SAY θ_2 IS MORE LIKELY TO BE THE TRUE PARAMETER, THEN THE RATIO $L(\theta_1)/L(\theta_2)$ SHOULD ASSUME A VALUE HIGHER THAN 1

The log likelihood function

In the previous slide the **log likelihood function** has been introduced, which is simply the logarithm of $L(\theta)$, namely

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

The log likelihood function carries the same information of the likelihood function, but it is much more manageable. Indeed, for a random sample

$$L(\theta) = \prod_{i=1}^n f_\theta(y_i)$$

$$\ln(L(\theta)) = \ln\left(\prod_{i=1}^n f_\theta(y_i)\right)$$

$$\ell(\theta) = \sum_{i=1}^n \log f_\theta(y_i).$$

OFTEN THE LIKELIHOOD INVOLVES PRODUCTS (FOR MULTIPLE DATA POINTS), WHICH CAN BE COMPUTATIONALLY CHALLENGING AND NUMERICALLY INSTABLE. THE LOG-LIKELIHOOD SIMPLIFIES THIS WHILE PRESERVING THE MAXIMA AND IS OFTEN EASIER TO INTERPRET IN TERMS OF ADDITIVE CONTRIBUTIONS OF INDIVIDUALS DATA POINTS.

but

Notice that $\ell(\theta)$ is defined up to an additive constant, depending only on the data y .

Example 1: the Poisson model

For a random sample y_1, \dots, y_n , with $Y_i \sim \mathcal{P}(\lambda)$ i.i.d., we readily get

$$L(\lambda) = \frac{\lambda^{\sum_{i=1}^n y_i} \exp\{-n\lambda\}}{\prod_{i=1}^n y_i!},$$

so that

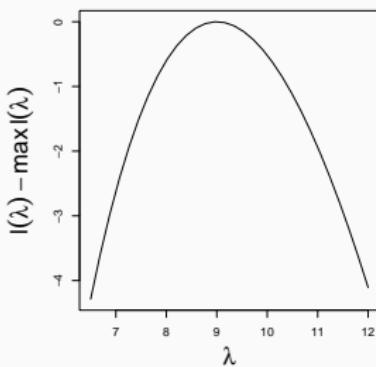
$$\ell(\lambda) = \log(\lambda) \sum_{i=1}^n y_i - n\lambda,$$

neglecting the term which does not depend on λ .

R lab: the Poisson log likelihood

Assume that for a sample $n = 10$ we observe $\sum_i y_i = 90$.

```
lik_pois <- function(lam, n, sumy) log(lam) * sumy - n * lam
xx <- seq(6.5, 12, l = 30)
ll <- sapply(xx, lik_pois, sumy = 90, n = 10)
par(pty = "s")
plot(xx, ll - max(ll), type = "l", xlab = expression(lambda),
     ylab = expression(l(lambda)-max(l(lambda))), cex.lab = 2)
```



Example 2: the normal model

For a random sample y_1, \dots, y_n , with $Y_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d.

$$L(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \mu)^2}{2\sigma^2}\right\},$$

and then with some simple algebra

$$\ell(\mu, \sigma^2) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2.$$

Sufficient statistics

The definition of sufficient statistic, given in the probability part, can be re-interpreted for the log likelihood function: $t(\mathbf{y})$ is **sufficient** for θ if $L(\theta)$ can be written as

$$L(\theta) = h(\mathbf{y}) g_{\theta}\{t(\mathbf{y})\}.$$

FRACTION THAT ONLY DEPENDS ON THE DATA FRACTION THAT DEPENDS ON THE PARAMETER AND ON THE DATA THROUGH THE SUFFICIENT STATISTIC (sum y_i, sum y_i², etc...)

The **minimal sufficient statistic** allows for the maximal reduction of dimensionality, in the sense that a minimal sufficient statistic is a function of every other sufficient statistic.

For the Poisson model, the $\sum_i y_i$ (or, equivalently, the sample mean \bar{y}) is sufficient for λ , whereas for the normal model the sufficient statistic is given by the pair $(\sum_i y_i, \sum_i y_i^2)$ (or, equivalently, by the pair (\bar{y}, s^2)).

These two statistical models are an instance of an **exponential family**, an important model class that includes also other important elements, such as the binomial distribution. They play an important role in the theory of *generalized linear models*. AND HAVE A LOT OF USEFUL DISTRIBUTION THEY CAN BE WRITTEN IN THE FORM ABOVE

WHAT ARE SUFFICIENT STATISTICS?

A SUFFICIENT STATISTIC IS A SUMMARY OF THE DATA THAT CONTAINS ALL THE INFORMATION NEEDED TO COMPUTE THE LIKELIHOOD OF A PARAMETER θ .

IN SIMPLE TERMS:

ONCE YOU HAVE THE SUFFICIENT STATISTIC,
YOU DON'T NEED RAW DATA ANYMORE TO
ESTIMATE θ .

FORMALLY, A STATISTIC $T(y)$ IS SUFFICIENT FOR θ IF THE CONDITIONAL DISTRIBUTION OF THE DATA y , GIVEN $T(y)$, DOES NOT DEPEND ON θ .
MATHEMATICALLY:

$$P(y | T(y), \theta) = P(y | T(y))$$

EXAMPLE: FOR A NORMAL DISTRIBUTION $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$:

THE SAMPLE MEAN \bar{X} AND THE SAMPLE VARIANCE S^2 ARE SUFFICIENT FOR μ AND σ^2 .

FOR A BINOMIAL $X \sim \text{Bin}(n, p)$:

THE NUMBER OF SUCCESSES $\sum X_i$ IS SUFFICIENT FOR p .

SUFFICIENT STATISTICS ARE CRITICAL TO MLE BECAUSE THEY REDUCE THE COMPLEXITY OF THE ESTIMATION PROBLEM:

- INSTEAD OF WORKING WITH THE FULL DATA X , WE CAN FOCUS ONLY ON $T(X)$, THE SUFFICIENT STATISTICS,

THIS IS ESPECIALLY USEFUL WHEN:

- THE DATA IS LARGE, AND SUMMARIZING IT REDUCES COMPUTATIONAL OVERHEAD
- THE SUFFICIENT STATISTIC MAKES IT EASIER TO WRITE DOWN AND SOLVE THE LIKELIHOOD EQUATIONS.

Maximum likelihood estimation

Given the interpretation of the (log) likelihood, the maximum of $\ell(\theta)$ is the value of the parameter which is most supported by the data.

A natural step is to take it as the point estimate, the **maximum likelihood estimate** (MLE) of θ

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \ell(\theta)$$

IDENTIFY THE PARAMETER θ THAT MAXIMISE THE LIKELIHOOD

Notice that since $\ell(\theta)$ is also a function of y , the MLE is a statistic.

The MLE in the two examples

For the Poisson model, simple calculus gives

$$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^y y_i = \bar{y}.$$

For the normal model, we need to maximize a function of two variables, and we get

$$\begin{cases} \hat{\mu} = \bar{y} \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2. \end{cases}$$

MLE: comments

Maximum likelihood estimation has a **central role** in modern statistics (and machine learning). There are several reasons for this:

1. The MLE algorithm is **automatic**: given a parametric statistical model for the data, the MLE follows from the chosen model.
2. The MLE of a function of a parameter $\psi = g(\theta)$ is defined by the simple plug-in rule $\hat{\psi} = g(\hat{\theta})$, which is very convenient for applications.
3. The MLE has **excellent properties**, which we illustrate in what follows.

RECAP

WE HAVE DEFINED THE LIKELIHOOD FUNCTION FOR A PARAMETER θ , ASSUMING AN i.i.d SAMPLE, AS THE PRODUCT OF UNIVARIATE PROBABILITY DENSITY FUNCTION:

$$L(\theta) = \prod_{i=0}^n f_\theta(y_i)$$

AND THIS COME FROM THE JOINT PROBABILITY DENSITY FUNCTION.

WHEN WE OBSERVED OUR SAMPLE AND INSERT THEM INTO THE LIKELIHOOD FUNCTION, WE CONSIDER IT LIKE A FUNCTION OF THE PARAMETER AND WE CAN SEE WHICH VALUES OF THE PARAMETER SPACE ARE THE MORE "SUPPORTED".

THINK OF THE LIKELIHOOD AS TRYING TO LOCATE A "SWEET SPOT" FOR A MACHINE SETTING (θ) THAT MAKES IT PRODUCE THE OBSERVED OUTCOMES (X) MOST CONSISTENTLY.

THE SUFFICIENT STATISTIC IS LIKE A DASHBOARD SUMMARY THAT CONSENSUS ALL RELEVANT DATA / SO YOU DON'T HAVE TO INSPECT EVERY PIECE OF THE MACHINE'S HISTORY TO FIGURE OUT THE OPTIMAL SETTING.

Maximum likelihood estimation: theory

Likelihood quantities

THE FIRST DERIVATIVE
IDENTIFY THE MAX
THE SECOND DERIVATIVE
IDENTIFY THE RATE AT
WHICH WE OBSERVE INFORMATION

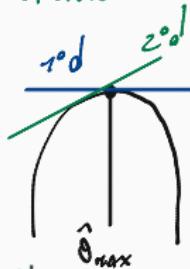
The first two derivatives of $\ell(\theta)$ play an important role.

The vector of first derivatives is called the score function

↳ THE GRADIENT

$$U(\theta) = U(\theta; \mathbf{y}) = \frac{\partial \ell(\theta)}{\partial \theta}$$

$$U(\hat{\theta}_{\text{max}}) = 0$$



The matrix of second derivatives, with negative sign, is called the **observed information matrix**:

A STEEPER CURVATURE INDICATES MORE CONFIDENCE IN THE PARAMETER ESTIMATE

$$J(\theta) = J(\theta; \mathbf{y}) = -\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top}$$

$$\theta = (\mu, \sigma^2)$$

$$U(\theta) = \begin{pmatrix} \frac{\partial \ell(\theta)}{\partial \mu} \\ \frac{\partial \ell(\theta)}{\partial \sigma^2} \end{pmatrix}$$

$$\left\{ \begin{array}{l} \frac{\partial \ell(\theta)}{\partial \mu} = 0 \\ \frac{\partial \ell(\theta)}{\partial \sigma^2} = 0 \end{array} \right.$$

IF θ IS A VECTOR,
WE OBTAIN A VECTOR
OF DERIVATIVES.
IF SO, WE HAVE TO
COMPUTE A SYSTEM
TO OBTAIN THE MAX

IF WE WANT THE 2nd DERIVATIVES
OF A VECTOR, WE OBTAIN A
MATRIX THAT DEPENDS ON THE SIZE 13

Some properties

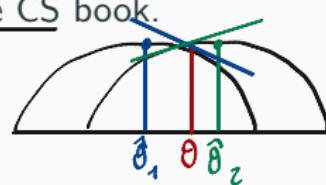
THIS PROPERTIES ARE ASSUMED IF WE CAN REPRODUCE THE TEST A LARGE AMOUNT OF TIME

The derivatives of the log likelihood function satisfy some important properties, provided that some **regularity conditions** hold (we shall return on them later on).

The proofs are simple, and they are reported in the CS book.

1. Zero expected score

$$E_{\theta} \{ U(\theta; \mathbf{Y}) \} = 0$$



2. 2nd Bartlett identity

ASSUMING θ THE TRUE VALUE, IF WE REPEAT SAMPLING AND EVALUATE THE SCORE FUNCTION, ON AVERAGE THE EXPECTATION IS EQUAL TO 0.

$$\text{cov}_{\theta} \{ U(\theta; \mathbf{Y}) \} = E_{\theta} \{ J(\theta; \mathbf{Y}) \} = \mathcal{I}(\theta)$$

The expected value $\mathcal{I}(\theta)$ of the observed information matrix is called the *Fisher information matrix* (or just the *expected information matrix*).

$$\mathcal{I}(\theta) = E \left[\left(\frac{\partial \ln f_{\theta}(\mathbf{y})}{\partial \theta} \right)^2 \right]$$

The Cramér-Rao lower bound

The third property is important, and we first state it for a one-parameter model (scalar θ).

3. *The Cramér-Rao lower bound:* the variance of *any unbiased estimator* $\tilde{\theta}$ cannot be smaller than the reciprocal of the expected information:

$$\text{var}_\theta\{\tilde{\theta}(\mathbf{Y})\} \geq \frac{1}{\mathcal{I}(\theta)}.$$

IT PROVIDES A THEORETICAL LOWER BOUND FOR THE VARIANCE OF ANY UNBIASED ESTIMATOR $\hat{\theta}$

Actually, by differentiation of the unbiasedness condition with respect to θ it follows that $\text{cov}_\theta\{\tilde{\theta}, U(\theta; \mathbf{Y})\} = 1$, which readily implies the Cramér-lower bound.

THE MLE IS EFFICIENT IN LARGE SAMPLES, MEANING IT ACHIEVES THIS LOWER BOUND ASYMPTOTICALLY

The extension to multiparameter models is given by the condition that the matrix $\text{cov}(\tilde{\theta}) = \mathcal{I}(\theta)^{-1}$ is positive semi-definite.

LINEAR PROPERTIES 3

$$|\text{Cov}(\tilde{\theta}(y), U(\theta))| = \frac{|\text{Cov}(\tilde{\theta}(y), U(\theta))|}{\sqrt{\text{Var}(\tilde{\theta}(y)) \text{Var}(U(\theta))}} \leq 1 \quad \text{and} \quad \text{Var}(\tilde{\theta}(y)) \geq \frac{\text{Cov}(\tilde{\theta}(y), U(\theta))^2}{\text{Var}(U(\theta))}$$

$$\begin{aligned} \text{Cov}(\tilde{\theta}(y), U(\theta))^2 &= E(\tilde{\theta}(y) U(\theta)) - E(\tilde{\theta}(y)) E(U(\theta)) \\ &= E(\tilde{\theta}(y) U(\theta)) = \int \tilde{\theta}(y) U(\theta) \cdot f_{\theta}(y) dy \\ &= \int \tilde{\theta}(y) \frac{\partial \log f_{\theta}(y)}{\partial \theta} f_{\theta}(y) dy \\ &= \int \tilde{\theta}(y) \frac{1}{f_{\theta}(y)} \frac{\partial f_{\theta}(y)}{\partial \theta} \cancel{f_{\theta}(y)} dy \\ &= \frac{\partial}{\partial \theta} \int \frac{\tilde{\theta}(y) f_{\theta}(y)}{f_{\theta}(y)} dy = \frac{\partial E(\tilde{\theta}(y))}{\partial \theta} = \frac{\partial \theta}{\partial \theta} = 1 \end{aligned}$$

$$\text{Var}(\tilde{\theta}(y)) \geq \frac{1}{\chi(\theta)}$$

Consistency of MLE

EX OF NOT CONSISTENT:
THE NUMBER OF
PARAMETERS INCREASES WITH
THE SAMPLE SIZE

We are ready to state the first crucial property of the MLE:

Maximum likelihood estimators are **usually** consistent, that is if the sample size tends to infinity $\hat{\theta}$ tends to θ_t , the **true** parameter value.

A justification for the result is given by the fact that in regular situations $\ell(\theta)/n \rightarrow E_{\theta}\{\ell(\theta)\}/n$ as $n \rightarrow \infty$, so that eventually the maximum of $\ell(\theta)$ and $E\{\ell(\theta)\}$ must coincide at θ_t by the Wald inequality.

The formal proof (typically) employs the law of large numbers.

Consistency can fail if the number of parameters increases with the sample size.

Large-sample distribution of MLE

We establish it by a Taylor expansion for the score function:

$$U(\hat{\theta}) \doteq U(\theta_t) - (\hat{\theta} - \theta_t) J(\theta_t),$$

with equality when $n \rightarrow \infty$ since $\hat{\theta} - \theta_t \rightarrow \mathbf{0}$. (due to consistency)

From the definition of $\hat{\theta}$, we get $U(\hat{\theta}) = \mathbf{0}$. Under mild assumptions

$$\frac{J(\theta_t)}{n} \rightarrow \frac{\mathcal{I}(\theta_t)}{n},$$

whereas $U(\theta_t)$ is a random vector with mean vector $\mathbf{0}$ and covariance matrix $\mathcal{I}(\theta_t)$.

In the large sample limit

$$\hat{\theta} - \theta_t \stackrel{\text{d}}{\sim} \mathcal{I}(\theta_t)^{-1} U(\theta_t; \mathbf{y}),$$

implying that $E(\hat{\theta} - \theta_t) = \mathbf{0}$ and $\text{cov}(\hat{\theta} - \theta_t) = \mathcal{I}(\theta_t)^{-1}$.

Large-sample normality of MLE

ON THE AGRESTI BOOK
SECTION 8.1

In the case when the sample is formed by independent observations, it follows that the log likelihood is the sum of independent contributions: under mild conditions the central limit theorem applies, and in the large sample limit

$$\hat{\theta} \sim \mathcal{N}\{\theta_t, \mathcal{I}(\theta_t)^{-1}\}.$$

Notice that whenever this holds, it would be possible (and recommendable, in some sense) to use $J(\theta_t)$ in place of $\mathcal{I}(\theta_t)$.

Again, since θ_t is unknown, we replace it by $\hat{\theta}$, obtaining the following estimated standard error for the k -th component of θ

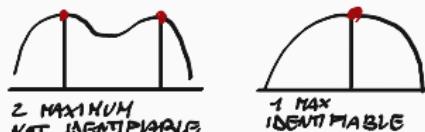
$$SE(\hat{\theta}_k) = \sqrt{\left[J(\hat{\theta})^{-1} \right]_{kk}} \xrightarrow{\text{VALUE ON THE DIAGONAL}}$$

Note: for *regular models* (see next slide), the observed information is positive definite at $\hat{\theta}$, so that the SE above is well defined.

Regularity conditions

We end the summary of the theory by mentioning the regularity conditions, which are some assumptions on the statistical model, required for the previous results to be valid.

The CS book lists the following ones:



1. The pdf of \mathbf{y} defined by different values of θ are distinct, namely the model is *identifiable*. (**IDENTIFIABILITY**)
2. The true parameter value θ_t is interior to Θ . (**COMPACTNESS**)
3. Within some neighbourhood of θ_t , the first three derivatives of $\ell(\theta)$ exist and are bounded, while the expected information satisfies the 2nd Bartlett identity, is positive definite and finite. (**SMOOTHNESS**)

These are mild conditions, which are generally valid in most cases.

Winding up

The previous results have illustrated that

1. The MLE is a **consistent estimator**.
2. The MLE is **asymptotic efficient**, since its asymptotic variance attains the Cramér-Rao lower bound.
3. The large sample distribution (aka the approximate distribution) of the MLE is **multivariate normal**, with standard error that can be estimated by the observed information evaluated at the parameter estimate.

Example 1: Poisson model

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Pos}(\lambda) \quad \hat{\lambda} = \frac{\sum y_i}{n} \Rightarrow \sum y_i \sim \text{Pos}(n\lambda)$$

Here $\hat{\lambda} = \bar{y}$, and consistency follows from the law of large numbers, in agreement with likelihood theory.

$$\hat{\lambda}_{ML} = \frac{\sum y_i}{n}$$

WE HAVE TO CHECK IF WE OBTAIN THE SAME RESULT AS THE ONE OBTAINED FROM THE CENTRAL LIMIT THEOREM USING THE RESULTS FROM THE LIKELIHOOD THEORY

Furthermore, the CLT states that for large n

$$\hat{\lambda} \sim \mathcal{N}(\lambda, \lambda/n). \quad (1)$$

$$\Rightarrow E[\hat{\lambda}_{ML}] = \lambda \quad (2)$$

$$\text{Var}(\hat{\lambda}_{ML}) = E[J(\lambda)]^{-1}$$

This result can be obtained also from likelihood theory. Indeed, we get

$$J(\lambda) = \frac{\sum y_i}{\lambda^2}$$

$$U(\lambda) = \frac{\sum y_i}{\lambda} - n$$

$$J'(\lambda) = -\frac{\partial^2 U(\lambda)}{\partial \lambda^2}$$

$$= -\frac{2}{\lambda}$$

$$= \frac{2}{\lambda} \left[\frac{\sum y_i}{\lambda} - n \right]$$

$$= -\left[-\frac{\sum y_i}{\lambda^2} \right]$$

$$= \frac{\sum y_i}{\lambda^2}$$

so that $I(\lambda) = n/\lambda$ and $I(\lambda)^{-1} = \lambda/n$.

NOW, WE WANT THE FISHER INFORMATION MATRIX:

$$\begin{aligned} X(\lambda) &= E[J(\lambda)] \\ &= E\left[\frac{\sum y_i}{\lambda^2}\right] \\ &= \frac{n\lambda}{\lambda^2} = \frac{n}{\lambda} \end{aligned}$$

FROM THIS $\text{Var}(\hat{\lambda}_{ML}) = \frac{\lambda}{n}$ (3)
 IF WE COMBINE (2) AND (3) WE GET:
 $\hat{\lambda}_{ML} \sim \mathcal{N}(\lambda, \frac{\lambda}{n})$ OBTAINING THE SAME RESULT
 GIVEN BY (1)

Example 2: normal example

Here we get

$$J(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & \frac{n}{\sigma^4} (\bar{y} - \mu) \\ \frac{n}{\sigma^4} (\bar{y} - \mu) & -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_{i=1}^n (y_i - \mu)^2 \end{pmatrix}$$

and therefore

$$\mathcal{I}(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}$$

The implication is that $\hat{\mu}$ and $\hat{\sigma}^2$ are (asymptotically) uncorrelated, and the two estimated standard errors are *(due to the uncorrelation)*

$$\text{SE}(\hat{\mu}) = \frac{\hat{\sigma}}{\sqrt{n}}, \quad \text{SE}(\hat{\sigma}^2) = \frac{\sqrt{2}\hat{\sigma}^2}{\sqrt{n}}.$$

Some numerical aspects

SOME TIMES, IN COMPLEX MODEL, YOU CAN SEE THE MAXIMISATION PROBLEM AS AN OPTIMIZATION PROBLEM AND APPLY SOME ALGORITHM

Numerical optimisation

The algorithmic nature of the MLE estimation method translates the statistical model into an optimisation problem: once a (sensible) statistical model has been specified for the data, we obtain parameter estimates with excellent properties by maximizing the log likelihood.

In some simple settings, like in the examples above, it is possible to find the analytical expression for the MLE, but in general we must resort to **numerical optimisation** of the log likelihood.

There are indeed several methods available for the task. Some knowledge of the most important issues related to it turns out particularly useful even for the application of off-the-shelf routines in R (or other environments).

Newton's method

Newton's method for optimisation is commonly used for minimization, in this case of the objective function $f(\theta) = -\ell(\theta)$.

The theory is well described in the CS book, here we mention the most important aspects. The idea is to locally approximate $f(\theta)$ as a quadratic function, which is repeatedly minimised.

$$\theta_{m+1} = \theta_m - \frac{U(\theta_m)}{H(\theta_m)}$$

The resulting method consists in an **iterative algorithm**, which is started with $k = 0$ and a *guesstimate* $\theta^{[0]}$, and iterates the following steps:

1. Evaluate $\ell(\theta^{[k]}), U(\theta^{[k]})$ and $J(\theta^{[k]})$. $k=0 \dots n$
2. If $U(\theta^{[k]}) \doteq \mathbf{0}$ and $J(\theta^{[k]})$ is positive definite then stop.
3. If $H = J(\theta^{[k]})$ is not positive definite, perturb it so that it is.
4. Solve $H\delta = U(\theta^{[k]})$ for the search direction δ .
5. If $\ell(\theta^{[k]} + \delta) > \ell(\theta^{[k]})$, repeatedly halve δ until it is (*this is the step-length control*).
6. Set $\theta^{[k+1]} = \theta^{[k]} + \delta$, increment k by one and return to step 1.

Fisher scoring and Quasi-Newton.

THIS MEASURES THE CURVATURE OF THE ACTUAL OBSERVED $\ell(\theta)$ IS THE NEGATIVE OF THE HESSIAN $I_{obs}(\theta) = -H(\theta)$

Whenever available, it is always a good idea to replace the observed information with the expected information $I(\theta^{[k]})$ in the Newton's method.

The resulting algorithm has a long successful tradition in statistics, it is called **Fisher scoring** and, indeed, it has better convergence properties.

Another variant avoids the computation of either $J(\theta^{[k]})$ or $I(\theta^{[k]})$, by building an approximation to the second derivative of $\ell(\theta)$ as the optimization proceeds. This is the approach of the **Quasi-Newton** methods, such as the widely used BFGS algorithm.

Quasi-Newton methods are implemented in several R functions and packages; see the CRAN Task View for *Optimisation* (<https://cran.r-project.org/web/views/Optimization.html>).

METHOD	USES OBSERVED INFORMATION	USES EXPECTED INFORMATION	COMPUTATIONAL COMPLEXITY
Newton's Method	YES	No	High (requires Hessian)
Fisher Scoring	No	Yes	Moderate
Quasi-Newton	No	Approximation ONLY	Low (Hessian's approx.)

STABILITY

LESS STABLE (HESSIAN CAN BE ILL-CONDITIONED)

More STABLE

STABLE, EFFICIENT

An example: logistic regression

We follow the MASS book for a simple example on a dose-response model.

Namely, we assume that y_i is the number of dead budworms (out of 20) for a dose of insecticide x_i^* . In particular, the statistical model is

$$Y_i \sim \mathcal{B}_i(20, \pi_i) \quad i = 1, \dots, 12, \text{ independent}$$

with

$$\pi_i(\alpha, \beta) = \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}}$$

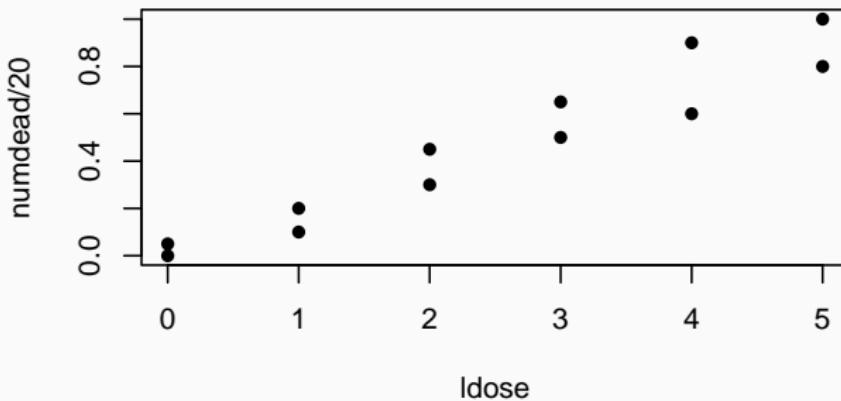
with $x_i = \log(x_i^*)$.

This is a simple instance of a *logistic regression model*.

R lab: budworm data

There are two observations at each dose (M/F budworms), but here for the sake of simplicity we ignore the different sex.

```
ldose <- rep(0:5, 2)  
numdead <- c(1, 4, 9, 13, 18, 20, 0, 2, 6, 10, 12, 16)  
plot(ldose, numdead / 20, pch=16)
```



Logistic regression: likelihood quantities

With some simple algebra we get:

$$\ell(\alpha, \beta) = \sum_i \{y_i (\alpha + \beta x_i) - 20 \log(1 + e^{\alpha + \beta x_i})\}$$

$$U(\alpha, \beta) = \begin{pmatrix} \sum_i \{y_i - 20 \pi_i(\alpha, \beta)\} \\ \sum_i \{y_i - 20 \pi_i(\alpha, \beta)\} x_i \end{pmatrix}$$

$$\mathcal{I}(\alpha, \beta) = \begin{pmatrix} \sum_i 20 \pi_i(\alpha, \beta) \{1 - \pi_i(\alpha, \beta)\} & \sum_i 20 \pi_i(\alpha, \beta) \{1 - \pi_i(\alpha, \beta)\} x_i \\ \sum_i 20 \pi_i(\alpha, \beta) \{1 - \pi_i(\alpha, \beta)\} x_i & \sum_i 20 \pi_i(\alpha, \beta) \{1 - \pi_i(\alpha, \beta)\} x_i^2 \end{pmatrix}$$

Notice that for this model $J(\alpha, \beta) = \mathcal{I}(\alpha, \beta)$.

R lab: likelihood and score functions

```
loglik <- function(theta, data){  
    eta <- theta[1] + theta[2] * data$x  
    out <- sum(data$y * eta - 20 * log(1+exp(eta)))  
    return(out)  
}  
  
score <- function(theta, data){  
    prob <- plogis(theta[1] + theta[2] * data$x)  
    out <- c(sum(data$y - prob * 20),  
            sum((data$y - prob * 20) * data$x))  
    return(out)  
}
```

R lab: information function

```
info <- function(theta, data){  
    prob <- plogis(theta[1] + theta[2] * data$x)  
    info11 <- sum(20 * prob * (1-prob))  
    info12 <- sum(20 * prob * (1-prob) * data$x)  
    info22 <- sum(20 * prob * (1-prob) * data$x^2)  
    out <- matrix(c(info11, info12, info12, info22), 2, 2)  
    return(out)  
}
```

R lab: starting point

Let's start from $\alpha = \beta = 0$: we obtain

```
theta0 <- c(0, 0); budw <- data.frame(y = numdead, x = ldose)

loglik(theta0, budw)

## [1] -166.3553

score(theta0, budw)

## [1] -9 105

info(theta0, budw)

##      [,1] [,2]
## [1,]    60   150
## [2,]   150   550
```

R lab: first step

```
H <- info(theta0, budw)
u0 <- score(theta0, budw)
delta <- solve(H, u0)
theta1 <- theta0 + delta

theta1
## [1] -1.9714286  0.7285714
loglik(theta1, budw)
## [1] -114.7219
```

which is clearly an improvement.

R lab: first 10 steps

```
## k = 1 theta= -1.971429 0.7285714 loglik= -114.7219
## k = 2 theta= -2.621436 0.9572079 loglik= -111.8192
## k = 3 theta= -2.760585 1.004947 loglik= -111.734
## k = 4 theta= -2.766079 1.006804 loglik= -111.7339
## k = 5 theta= -2.766087 1.006807 loglik= -111.7339
## k = 6 theta= -2.766087 1.006807 loglik= -111.7339
## k = 7 theta= -2.766087 1.006807 loglik= -111.7339
## k = 8 theta= -2.766087 1.006807 loglik= -111.7339
## k = 9 theta= -2.766087 1.006807 loglik= -111.7339
## k = 10 theta= -2.766087 1.006807 loglik= -111.7339
```

The algorithm converges quickly, and actually after 10 iterations

```
cat(score(theta10, budw), det(info(theta10, budw)),
    sqrt(diag(solve(info(theta10, budw)))))
```

```
## 1.776357e-15 5.329071e-15 2361.462 0.3701342 0.1235889
```

R lab: glm analysis

```
budworm.lg0 <- glm(cbind(y, 20-y) ~ x, binomial, budw)
summary(budworm.lg0, cor = FALSE)

## 
## Call:
## glm(formula = cbind(y, 20 - y) ~ x, family = binomial, data =
## 
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.7661   ^ 0.3701  -7.473 7.82e-14 ***
## x            1.0068   ^ 0.1236   8.147 3.74e-16 ***
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
## Null deviance: 124.876  on 11  degrees of freedom
## Residual deviance: 16.984  on 10  degrees of freedom
```