

last change 01.11.12@18.40

## 1.1 Maximum Likelihood Estimation

**Example:**  $n$  Bernoulli trials 0–1-sequence  $\mathbf{y} = (y_1, \dots, y_n)$ . success probability  $p$ .

$$L(p | \mathbf{y}) = \prod_1^n p^{y_i} (1-p)^{1-y_i}$$

ml estimator

$$l(p | \mathbf{y}) = \log L(p | \mathbf{y}) = \sum y \log p + (n - \sum y) \log(1-p)$$

MLE<sup>1</sup>

$$\hat{p} = \arg \max_p l(p | \mathbf{y}) = \frac{\sum y}{n}$$

**properties:**

- unbiasedness:  $\mathbb{E}\hat{p} = p$
- variance:  $\mathbb{V}\hat{p} = p(1-p)/n$
- consistency:  $\mathbb{P}(|\hat{p} - p| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$  (law of large numbers)
- asymptotic distribution: (CLT) as  $n \rightarrow \infty$

$$\frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \rightarrow \mathcal{N}(0, 1)$$

### properties of the mle

density:

$$y_1, \dots, y_n \text{ iid } \sim f_\theta$$

in the discrete case, replace the density by probability  $p_\theta$

likelihood

$$L(\theta | y_1, \dots, y_n) = f_\theta(y_1, \dots, y_n) = \prod f_\theta(y_i) = \prod L(\theta | y_i)$$

take logs for numerical and mathematical (CLT) reasons

---

<sup>1</sup>**note:** strictly, maximization has to take place under the restrictions  $0 \leq p \leq 1$ . By introducing the reparametrization

$$p = \frac{e^\beta}{1 + e^\beta}$$

this maximization in  $\beta$  is without restriction. this will be used in the logit model.

$$\begin{aligned}
l(\theta | y_1, \dots, y_n) &= \log L(\theta | y_1, \dots, y_n) = \log f_\theta(y_1, \dots, y_n) \\
&= \sum \log f_\theta(y_i) = \sum \log L(\theta | y_i) = \sum l(\theta | y_i)
\end{aligned}$$

maximize

$$\hat{\theta} = \arg \max_{\theta} l(\theta | y_1, \dots, y_n)$$

notes

1. as the likelihood is a function of the unknown parameter given the observations
2. as the observations are random variables the likelihood function in each possible value  $\theta$  is a random variable of the observations
3. the likelihood in every possible value  $\theta$  has expectation and variance w.r. to the random observations
4. the same holds for  $\hat{\theta}$

define **Fisher information**:

$$\mathbb{I}_n(\theta) = -\mathbb{E}\left(\sum_i^n \frac{d^2}{d\theta^2} l(\theta | y_i)\right)$$

under *iid assumption*

$$\mathbb{I}_n(\theta) = -n\mathbb{E}\left(\frac{d^2}{d\theta^2} l(\theta | y)\right) = -n \int \frac{d^2}{d\theta^2} l(\theta | y) f_\theta(y) dy = n\mathbb{I}_1$$

Notes

1. the likelihood function is supposed to be **log-concave** in order to have a unique maximum.
2. the (Fisher) information increases linearly with the sample size
3. the (Fisher) information tends to  $\infty$  for i.i.d. observations
4. the Fisher information is a misnomer and should not be confounded with the Shannon information

**asymptotic normality of mle:**

$$\sqrt{\mathbb{I}_n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}(0, 1)$$

sometimes this is written loosely as  $\hat{\theta} \rightarrow \mathcal{N}(\theta, \mathbb{I}_n^{-1})$

Note

$$\mathbb{V}\hat{\theta} \rightarrow \mathbb{I}_n^{-1}$$

Furthermore, it is *asymptotically unbiased*:  $\mathbb{E}\hat{\theta} \rightarrow \theta$

**Example:** for the binomial probability ( $n$  Bernoulli trials)

$$l(p | \mathbf{y}) = \log L(p | \mathbf{y}) = \sum y \log p + (n - \sum y) \log(1 - p)$$

$\Rightarrow$

$$\begin{aligned} \frac{d}{dp} l(p | \mathbf{y}) &= \frac{\sum y}{p} - \frac{n - \sum y}{1 - p} \\ \frac{d^2}{dp^2} l(p | \mathbf{y}) &= -\frac{\sum y}{p^2} - \frac{n - \sum y}{(1 - p)^2} \\ \mathbb{E} \frac{d^2}{dp^2} l(p | \mathbf{y}) &= -\left( \frac{np}{p^2} + \frac{n - np}{(1 - p)^2} \right) \\ &= -n \left( \frac{1}{p} + \frac{1}{1 - p} \right) = \frac{-n}{p(1 - p)} \end{aligned}$$

$\Rightarrow$

$$\sqrt{\mathbb{I}_n}(\hat{p} - p) = \sqrt{\frac{n}{p(1 - p)}}(\hat{p} - p) \rightarrow \mathcal{N}(0, 1)$$

Notes

1. asymptotic normality  $\mathbb{I}(\theta)$  depends on the unknown value replace by  $\mathbb{I}(\hat{\theta})$
2. if the expectation required in the computation of  $\mathbb{I}(\theta)$  may be hard to find, estimate  $\mathbb{I}(\theta)$  by  $\hat{\mathbb{I}}_n = \sum_i^n \frac{d^2}{d\theta^2} l(\hat{\theta} | y_i)$
3. loosely it could be stated

$$\hat{\theta} \rightarrow \mathcal{N}(\theta, \hat{\mathbb{I}}_n^{-1})$$

as  $n \rightarrow \infty$ , where

$$\hat{\mathbb{I}}_n = -\sum \frac{d^2}{d\theta^2} l(\hat{\theta} | y_i)$$

**higher dimensional parameters**

sample  $Y_1, \dots, Y_n$  iid  $f(y | \theta)$ ,  $\theta' = (\theta_1, \dots, \theta_J) \Rightarrow$  log likelihood function:

$$l(\theta | y_1, \dots, y_n) = l(\theta | \mathbf{y}) = \sum_i^n \log f(y_i | \theta) = \sum_i^n l(\theta | y_i)$$

mle

$$\hat{\theta} = \arg \min_{\theta} l(\theta | \mathbf{y})$$

**score function** ( $J$ -vector):

$$\frac{\partial}{\partial \theta} l(\theta | \mathbf{y}) = \left[ \frac{\partial}{\partial \theta_j} l(\theta | \mathbf{y}) \right] = \sum_i^n \frac{\partial}{\partial \theta} l(\theta | y_i)$$

mle ( $J$ -vector) solution to score equation:

$$\frac{\partial}{\partial \theta} l(\hat{\theta} | \mathbf{y}) = \left[ \frac{\partial}{\partial \theta_j} l(\hat{\theta} | \mathbf{y}) \right]_{j=1}^J = \mathbf{0}$$

**Fisher information** ( $J \times J$ -matrix) of  $n$  observations ( $\mathbb{I}_n$ ) resp of 1 observation ( $\mathbb{I}_1$ ).

$$\mathbb{I}_n(\theta) = -\mathbb{E} \frac{\partial^2}{\partial \theta^2} l(\theta | \mathbf{y}) = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta_{j_1} \partial \theta_{j_2}} l(\theta | \mathbf{y}) \right]_{j_1, j_2}^{J, J} = - \left[ \sum_i^n \mathbb{E} \frac{\partial^2}{\partial \theta_{j_1} \partial \theta_{j_2}} l(\theta | y_i) \right]$$

Notes:

1. In the iid case  $\mathbb{I}_n = n\mathbb{I}_1$
2. since  $\mathbb{I}_n(\theta)$  is a positive definite symmetric matrix it can be factored into the product of two matrices<sup>2</sup>  $(\sqrt{\mathbb{I}_n})' \sqrt{\mathbb{I}_n} = \mathbb{I}_n$ .

**asymptotic normality of mle:**

$$\sqrt{\mathbb{I}_n} (\hat{\theta} - \theta) \rightarrow \mathcal{N}_J(0, I)$$

this holds under certain under regularity conditions (existence of derivatives of log likelihoods etc, *true value parameter not boundary value* etc)

**LR-test**

$$\begin{aligned} H_0 &: \theta \in \Theta_0 \\ H_1 &: \theta \in \Theta - \Theta_0 \end{aligned}$$

where  $\Theta_0 \subset \Theta$   
compute

$$\begin{aligned} \hat{\theta}_0 &= \arg \min_{\theta \in \Theta_0} l(\theta | \mathbf{y}) \\ \hat{\theta} &= \arg \min_{\theta \in \Theta} l(\theta | \mathbf{y}) \end{aligned}$$

---

<sup>2</sup>this is related to the **principal axis theorem**, (see wikipedia and elsewhere)

It can be shown that

$$2 \sum_i^n \left( l(\hat{\theta} | y_i) - l(\hat{\theta}_0 | y_i) \right) \rightarrow \chi^2(df)$$

where  $df = \dim(\Theta) - \dim(\Theta_0)$  number of free parameters

**Wald test**

(case  $J = 1$ ) under  $H_0 : \theta = \theta_0$

$$\sqrt{\mathbb{I}_n(\theta_0)} (\hat{\theta} - \theta_0) \rightarrow \mathcal{N}_1(0, 1) \implies \mathbb{I}_n(\theta_0) (\hat{\theta} - \theta_0)^2 \rightarrow \chi^2(1)$$

This follows directly from the asymptotic normality of the ML-estimator (Taylor developing the score function). For case  $p \geq 1$

$$(\hat{\theta} - \theta_0)' \mathbb{I}_n(\theta_0) (\hat{\theta} - \theta_0) \rightarrow \chi^2(p)$$