- likelihood
- dichotomous responses (logistic regression) $\rightarrow$ GLM

---

In the previous lectures, we obtain an estimate of $\beta$, by minimizing the least squares criterion

$$\hat{\beta} = \arg\min_{\beta} \sum_{i=1}^{n} (y_i - \bar{X}_i \beta)^2$$

$\hookrightarrow f(x; \beta)$

It really works well $\varepsilon \sim N(0, \sigma^2)$.

We need a more general criterion to estimate the parameters

$\hookrightarrow$ likelihood criterion

$\hookrightarrow$ maximize the likelihood

Define a family of probability distributions,

parametric

Eg., Gaussian $N(\mu, \sigma^2)$

$$\mathcal{F} = \left\{ P_Y(y; \vartheta), \vartheta \in \Theta \right\}$$

$\rightarrow$ parametric space

$\vartheta = (\mu, \sigma^2)$     $\Theta = \mathbb{R} \times \mathbb{R}^+$

$\rightarrow$ parameter

probability density function   $y$ cont
probability function          $y$ discr.

$\vartheta$ depends on a $p$-dimensional parameter, which we need to estimate

in the Gaussian example,   $P_Y(y; \vartheta = (\mu, \sigma^2)) = \dfrac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\dfrac{1}{2\sigma^2}(y - \mu)^2 \right\}$

$X\beta \rightarrow \vartheta = (\beta, \sigma^2)$

The likelihood function is defined as

independent

$$L(\vartheta) = c \, P_Y(y; \vartheta) \quad \rightarrow \text{for } n \text{ obs, } L(\vartheta) = c \prod_{i=1}^{n} P_{Y_i}(y_i; \vartheta)$$

where:
• $c$ is a constant (everything that can be dropped-out in the computations)
• $y$ is fixed (are the data, so given), therefore $L(\vartheta)$ is a function of only $\vartheta$

Since $P_Y(y; \vartheta) \in [0, +\infty)$, it is possible (and more convinient) to work with the log-likelihood

$$\ell(\vartheta) = \log L(\vartheta)$$

if $P_Y(y; \vartheta) = 0$ , $\ell(\vartheta) = -\infty$

Then, to estimate $\theta$, we maximize $L(\theta)$ or $\ell(\theta)$

$$\hat{\theta} = \arg\max_{\theta} L(\theta) = \arg\max_{\theta} \ell(\theta)$$

e.g., in the Gaussian example (given $\sigma^2$)

$$\hat{\beta} = \arg\max_{\beta} \left( \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\frac{1}{2\sigma^2} (y_i - x_i^T\beta)^2 \right\} \right)$$

$$\to c$$

$$= \arg\max_{\beta} \exp\left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - x_i^T\beta)^2 \right\}$$

$$= \arg\max_{\beta} \left( -\sum_{i=1}^{n} (y_i - x_i^T\beta)^2 \right)$$

$$= \arg\min_{\beta} \left( \sum_{i=1}^{n} (y_i - x_i^T\beta)^2 \right)$$

In the special case of Gaussian distribution, the maximum likelihood estimate is the same of least squares one

How to find $\hat{\theta}$ in practice?

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

$\to$ find a stationary point $\quad \dfrac{\partial \ell(\theta)}{\partial \theta} = 0$

$\hat{\theta}$

verify that is a **global** maximum

$$\boxed{\dfrac{\partial^2 \ell(\theta)}{\partial \theta^2} < 0}$$

theoretically more tricky, usually not a problem in practice

$\to$ useful to compute a measure of uncertainty around our estimate (related to the variance of the maximum likelihood estimator)

$$j(\theta) = -\frac{\partial^2 \ell(\theta)}{\partial \theta^2} \qquad \to \quad j(\hat{\theta}) = -\frac{\partial^2 \ell(\theta)}{\partial \theta^2}\Bigg|_{\theta=\hat{\theta}}$$

Fisher observed information, it is the inverse of the variance of $\hat{\theta}$

$$\text{s.e.}(\hat{\theta}) = \text{diag}\left( j(\hat{\theta})^{-1} \right)^{1/2} \qquad i(\theta) = E\left[ j(\theta) \right]$$

$$\hat{\theta} \sim N\left( \theta, j^{-1}(\hat{\theta}) \right) \qquad \text{expected information}$$

approximately distributed

2

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

We can then construct a confidence interval around $\hat{\theta}$, with level $1-\alpha$
For the r-th component of $\hat{\theta}$

$$CI\,(1-\alpha) = \hat{\theta}_r \pm z_{1-\frac{\alpha}{2}} \sqrt{j^{-1}(\theta)_{[r,r]}}$$

$\hookrightarrow$ quantile of level $1-\frac{\alpha}{2}$ of a standard normal distribution $\mathcal{N}(0;1)$

$$j(\theta) = \frac{\partial^2 \ell(\theta)}{\partial \theta^2} = \begin{bmatrix} \frac{\partial^2 \ell(\theta)}{\partial \theta_{(1)} \partial \theta_{(1)}} & \cdots & \frac{\partial^2 \ell(\theta)}{\partial \theta_{(1)} \partial \theta_{(\cdot)}} \\ \vdots & & \\ \frac{\partial^2 \ell(\theta)}{\partial \theta_{(\cdot)} \partial \theta_{(1)}} & \cdots & \frac{\partial^2 \ell(\theta)}{\partial \theta_{(\cdot)} \partial \theta_{(\cdot)}} \end{bmatrix}$$

dimension of $\theta$

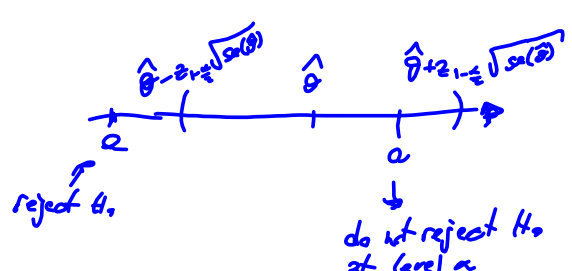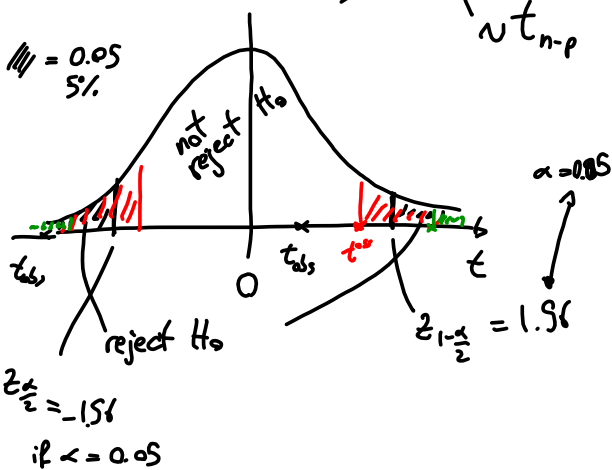$$\frac{\partial^2 \ell(\theta)}{\partial \theta_{(r)} \partial \theta_{(r)}}$$

Confidence intervals and hypothesis testing are closely connected

$$H_0 : \theta_r = a$$

For a fixed statistical significance level $\alpha$, the Wald test is used to testing the hypothesis, and it is based on the t-statistic

$$t = \frac{\theta_r - a}{\hat{s.e.}(\hat{\theta}_r)} \overset{H_0}{\sim} \mathcal{N}(0;1)$$

$\sim t_{n-p}$ in the case of Gaussian distribution
$n > 30$, $t$ is very close to $\mathcal{N}(0;1)$

$\alpha = 0.05$
$5\%$

not reject $H_0$

$t_{obs}$

$t_{obs}$

$t^{**}$

$O$

$t$

reject $H_0$

$z_{1-\frac{\alpha}{2}} = 1.96$

$z_{\frac{\alpha}{2}} = -1.96$
if $\alpha = 0.05$

$\alpha = 0.05$

$\hat{\theta} - z_{1-\frac{\alpha}{2}}\sqrt{se(\theta)}$    $\hat{\theta}$    $\hat{\theta} + z_{1-\frac{\alpha}{2}}\sqrt{se(\theta)}$

$a$    $a$

reject $H_0$

do not reject $H_0$ at level $\alpha$

Alternatively, one can compute the p-value : $2\min\left(\Phi(t), 1-\Phi(t)\right)$

$$2\Phi(-|t|) \longrightarrow \text{p-value that we contrast with } \alpha(0.05)$$

3

More generally, for

$$\theta_{(p-q+1:p)} = a_{[1:q]}$$

we can use the likelihood ratio test

$$w = 2\left(\ell(\hat{\theta}) - \ell(\theta_0)\right)$$

where $\theta_0 = (\hat{\theta}_{1:q}, a_{0,1q})$

Since $w \sim \chi^2_q$, the p-value is computed as

$$\text{p-value} = \Pr(\chi^2 > w) \qquad \chi^2 \sim \chi^2_q$$

$q = 2 \qquad p = 4$

$\theta_1 \; \theta_2 \; \theta_3 \; \theta_4$

$\hat{\theta} = (\hat{\theta}_1, \; \hat{\theta}_2, \; \hat{\theta}_3, \; \hat{\theta}_4)$

$\theta_0 = (\hat{\theta}_1, \; \hat{\theta}_2, \; a_1, \; a_2)$



---

Gaussian example (with only one variable $x$)     $\beta$ - 1-dimensional

$$Y \sim N(x\beta, \sigma^2)$$

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

$$= c \, (\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - x_i\beta)^2\right\}$$

$$\ell(\beta, \sigma^2) = \log L(\theta) = -\frac{n}{2}\log\sigma^2 - \frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - x_i\beta)^2$$

$$\ell_\beta(\beta, \sigma^2) = \frac{\partial \log L(\theta)}{\partial \beta} = +\frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - x_i\beta)x_i$$

$$\ell_{\sigma^2}(\beta, \sigma^2) = \frac{\partial \log L(\theta)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2}\sum_{i=1}^{n}(y_i - x_i\beta)^2$$

$$\ell_\beta(\beta, \sigma^2) = 0 \;\rightarrow\; \frac{1}{2\sigma^2}\sum_{i=1}^{n}(x_i y_i - x_i^2\beta) = 0 \qquad \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

if $x_i$ is $p$-dimensional

$$= (x^T x)^{-1} x^T y$$

$$\ell_{\sigma^2}(\hat{\beta}, \sigma^2) = 0 \;\rightarrow\; -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2}\sum_{i=1}^{n}(y_i - x_i\hat{\beta})^2 = 0$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n}(y_i - x_i\hat{\beta})^2}{n} = \frac{D(\hat{\beta})}{n}$$

$\leftarrow n-p$  because $\hat{\sigma}^2$ is biased, but they are equal for $n \to \infty$

---

To evaluate the hypothesis

$$\theta_{p-q+1} = 0$$

we contrast the deviance of the constrained model and the full model
When the error is Gaussian, we use the F

$$F = \frac{[D(\theta_0) - D(\hat{\theta})]/q}{D(\hat{\theta})/(n-p)} \sim F_{q, n-p}$$

$\rightarrow$ it is an F distribution with
$q$ degrees of freedom at the numerator and
$n-p$ at the denominator

Binomial distribution

$Y_i \in \{0, 1\}$     $Y_i \sim Bi(1, \tilde{\pi})$    Bernoulli distribution

                       $\tilde{\pi}$ : probability of success

$Y = \sum_{i=1}^{n} Y_i \sim Bi(n, \tilde{\pi})$    number of successes in # trials

$$Pr(y; \tilde{\pi}) = \binom{n}{y} \tilde{\pi}^y (1 - \tilde{\pi})^{n-y}$$

$$L(\tilde{\pi}) = \binom{n}{y} \tilde{\pi}^y (1 - \tilde{\pi})^{n-y}$$

$$\ell(\tilde{\pi}) = \overset{c}{=} y \log \tilde{\pi} + (n-y) \log(1 - \tilde{\pi})$$

$$\ell_{\pi}(\tilde{\pi}) = \frac{y}{\tilde{\pi}} + \frac{n-y}{1-\tilde{\pi}}(-1) \quad \to \quad \frac{y}{\tilde{\pi}} - \frac{n-y}{1-\tilde{\pi}} = 0 \quad \frac{y - \tilde{\pi}y - n\tilde{\pi} + \tilde{\pi}y}{} = 0$$

$$\hat{\pi} = \frac{y}{n}$$

$$\ell_{\pi\pi}(\tilde{\pi}) = -\frac{y}{\tilde{\pi}^2} - \frac{n-y}{(1-\tilde{\pi})^2}$$

$$j(\hat{\pi}) = \frac{y}{\tilde{\pi}^2} + \frac{n-y}{(1-\tilde{\pi})^2} \qquad j(\hat{\pi}) = \frac{y \, n^2}{y^2} + \frac{n-y}{(1 - \frac{y}{n})^2} = \frac{n}{\hat{\pi}} + \frac{(n-y)n^2}{(n-y)^2}$$

$$= \frac{n}{\hat{\pi}} + \frac{n}{1 - \hat{\pi}}$$

$$se(\hat{\pi}) = \sqrt{\left(\frac{n}{\hat{\pi}} + \frac{n}{1-\hat{\pi}}\right)^{-1}} = \sqrt{\left(\frac{n - n\hat{\pi} + n\hat{\pi}}{\hat{\pi}(1-\hat{\pi})}\right)^{-1}} = \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}$$

## Comparison of two groups

Brasilian bank example    $Y = $ satisfaction $\in \{low, high\}$

                        $X = $ age $\in \{\underset{1}{'young'}, \underset{2}{'old'}\}$    young $< 55$

$$\ell(\tilde{\pi}_1, \tilde{\pi}_2) = c + y_1 \log \hat{\pi}_1 + (n_1 - y_1) \log(1 - \hat{\pi}_1) + y_2 \log \hat{\pi}_2 + (n_2 - y_2) \log(1 - \hat{\pi}_2)$$

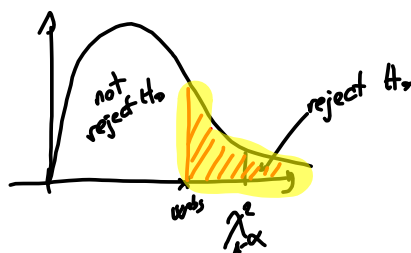$H_0 : \tilde{\pi}_1 = \tilde{\pi}_2$    $H_0: \tilde{\pi}_1 - \tilde{\pi}_2 = 0$

where: $\hat{\pi}_1 = \frac{y_1}{n_1} \simeq 0.73$

$W = 2\left(\underset{2}{\ell(\hat{\pi}_1, \hat{\pi}_2)} - \underset{1}{\ell(\hat{\pi}, \hat{\pi})}\right)$

$\hat{\pi}_2 = \frac{y_2}{n_2} \simeq 0.82$

$\hat{\pi} = \frac{y_1 + y_2}{n_1 + n_2} \simeq 0.76$

$W \sim \chi_1^2$



$Pr(X < w^{obs}) = $ pchisq($w_{obs}$, 1)