

Lecture 2: *t*-Test, The Wald Test and Likelihood Ratio Test

Lecturer: Prof. Jingyi Jessica Li

Subscriber: Kexin Li

**Disclaimer:** These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

## 2.1 Multivariate Linear Model

The vector form of the model is:

$$Y_{n \times 1} = \mathbf{X}_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1},$$

where the *design matrix*

$$\mathbf{X} = [1_{p \times 1} \quad \mathbf{x}_1 \quad \dots \quad \mathbf{x}_{p-1}] , \mathbf{x}_j \in \mathbb{R}^n \text{ are predictors/features/covariates.}$$

### Fixed Design vs. Random Design

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \mathbf{x}_i \in \mathbb{R}^p \text{ are observations.}$$

For random design, we assume that  $\mathbf{x}_i$ 's are random and thus  $(\mathbf{x}_i, Y_i)$  are jointly random. Then we have to check whether  $\{\mathbf{x}_i, Y_i\}_{i=1}^n$  can be reasonably assumed to be independent and identically distributed (*i.i.d.*). Random design is widely used in econometrics. A related concept is the measurement error model, which assumes that the marginal distribution of  $\mathbf{x}_i$  takes a specific form.

For fixed design, the only randomness comes from  $\epsilon_i$ 's. The assumptions that all the observations follow the same linear model and the condition that  $\epsilon_i$ 's are *i.i.d* is easier to satisfy. Fixed design is widely used statistics.

If the linear model is reasonably assumed, it is sufficient to consider fixed design in most cases.

For a detailed discussion about random design and fixed design, please refer to

Buja, A., et al. "A conspiracy of random predictors and model violations against classical inference in regression." Arxiv preprint. Available at <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.753.718&rep=rep1&type=pdf> (2014).

## 2.2 Inference of $\beta$ and $\sigma^2$

- $\beta$  indicates the effect of predictors on the response

- $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$ 
  - exact if  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$
  - the distribution of  $\hat{\beta}$  totally depends on our assumptions
  - asymptotically holds when  $\epsilon_i$ 's are *i.i.d* but not Gaussian (Central Limit Theorem)  
For the asymptotic result, please see <http://cameron.econ.ucdavis.edu/e240a/asymptotic.pdf>.
- the distribution of  $\hat{\beta}$  can be viewed as from infinitely repeated sampling (Frequentist View)

### 2.2.1 *t*-test

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \beta_1 : \text{the intercept}$$

Hypothesis testing:

$$H_0 : \beta_j = 0; \quad H_1 : \beta_j \neq 0$$

The *t*-statistic:

$$T = \frac{\hat{\beta}_j}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}} = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2 \times j^{\text{th}} \text{ diagonal entry of } (\mathbf{X}^T \mathbf{X})^{-1}}}$$

For the denominator, we use the *plug-in estimator* (replacing the true value by the estimator).

In order to do a valid statistical inference, we need to know the distribution of the *t*-statistic.

### 2.2.2 Hat Matrix

The fitted value:

$$\hat{Y} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$$

We define the *Hat Matrix* as:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

The residuals can be expressed as:

$$e = Y - \hat{Y} = Y - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y = (\mathbf{I}_n - \mathbf{H})Y$$

#### Properties of the Hat Matrix

- The Hat Matrix is symmetric, i.e.

$$\mathbf{H}^T = \mathbf{H}$$

- The Hat Matrix is idempotent, i.e.

$$\mathbf{H}^2 = \mathbf{H}$$

- $\text{trace}(\mathbf{H}) = \text{trace}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{trace}(\mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}) = \text{trace}(\mathbf{I}_p) = p$ ,  $\text{trace}(\mathbf{I}_n - \mathbf{H}) = n - p$ .

### 2.2.3 Estimators for $\sigma^2$

Given the Gaussian assumption, the MLE estimator:

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$$

We can show that

$$\frac{n\hat{\sigma}_{MLE}^2}{\sigma^2} = \frac{\sum_{i=1}^n e_i^2}{\sigma^2} \sim \chi_{n-p}^2 \quad \text{given that } \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

and thus

$$\mathbb{E}\left[\frac{n\hat{\sigma}_{MLE}^2}{\sigma^2}\right] = n - p, \quad \mathbb{E}[\hat{\sigma}_{MLE}^2] = \frac{n-p}{n}\sigma^2 \text{ (biased)}$$

An easy modification to generate an unbiased estimator is:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{n-p} \sum_{i=1}^n e_i^2 = \frac{RSS}{n-p} = \frac{SSE}{n-p}$$

where *RSS* (Residual Sum of Squares) and *SSE* (Sum of Squared Errors) are different names for the same item.

### 2.2.4 *t*-statistic and *t*-distribution

*t*-distribution:

- (i) The numerator follows the standard Gaussian distribution.
- (ii) The denominator is the square root of a variable following the  $\chi^2$  distribution divided by the degree of freedom.
- (iii) The numerator and the denominator are independent.

We are to show that the *t*-statistic  $T$  follows the *t*-distribution under the null hypothesis only conditional on normality assumption on the error term:  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ .

**Proof.** For (iii), an important property about  $\mathcal{N}(\cdot, \cdot)$  is:

If  $\mathbb{R}^p \ni U \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbf{A} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{B} \in \mathbb{R}^{m \times p}$ , then

$$\mathbf{AU} \text{ and } \mathbf{BU} \text{ are independent} \iff \mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^T = \text{Cov}(\mathbf{AU}, \mathbf{BU}) = \mathbf{0}_{n \times m}$$

Here in our case,

$$\begin{aligned} \hat{\beta} &= \mathbf{AY}, \quad \mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ e &= \mathbf{BY}, \quad \mathbf{B} = \mathbf{I}_n - \mathbf{H} \\ Y &\sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}), \quad \boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_n \end{aligned} \left. \right\} \Rightarrow \mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^T = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{I}_n - \mathbf{H}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T - \mathbf{X}^T \mathbf{H}) = \mathbf{0}$$

$$\Rightarrow \hat{\beta} \perp e$$

$$\Rightarrow \hat{\beta}_j \perp \hat{\sigma}_{unbiased}^2 = \frac{1}{n-p} \sum_{i=1}^n e_i^2$$

$$\Rightarrow \hat{\beta}_j \perp \sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}$$

Together with

$$\hat{\beta}_j \sim \mathcal{N}(\beta_j, \sigma^2 \cdot j^{th} \text{ diagonal entry of } (\mathbf{X}^T \mathbf{X})^{-1})$$

$$\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)} = \sqrt{\hat{\sigma}_{unbiased}^2 \cdot j^{th} \text{ diagonal entry of } (\mathbf{X}^T \mathbf{X})^{-1}}$$

and

$$\frac{\hat{\sigma}_{unbiased}^2}{\sigma^2} \sim \frac{\chi_{n-p}^2}{n-p}$$

Under the null hypothesis ( $H_0$ ), we can say that the *t*-statistic follows the *t*-distribution with  $(n - p)$  degree of freedom:

$$T = \frac{\hat{\beta}_j}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}} \stackrel{H_0}{\sim} t_{n-p}$$

Asymptotically, the *t*-distribution converges to the standard normal distribution if viewing  $\hat{\sigma}^2$  as the true value of  $\sigma^2$  as  $n \rightarrow \infty$ :

$$t_{n-p} \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, 1)$$

## 2.3 The Wald Test

$$\beta = \begin{pmatrix} \beta_{1(p_1 \times 1)} \\ \beta_{2(p_2 \times 1)} \end{pmatrix}, \quad \widehat{\text{Var}}(\hat{\beta}) = \begin{pmatrix} \widehat{\text{Var}}(\hat{\beta}_1) & \widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) \\ \widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{Var}}(\hat{\beta}_2) \end{pmatrix}$$

Hypothesis testing:

$$H_0 : \beta_2 = \mathbf{0}_{p_2 \times 1}; \quad H_1 : \beta_2 \neq \mathbf{0}_{p_2 \times 1}$$

The Wald statistic:

$$W = \hat{\beta}_2^T \widehat{\text{Var}}(\hat{\beta}_2)^{-1} \hat{\beta}_2$$

Note that if  $p_2 = 1$ ,  $W = T^2$ .

Given the normality assumptions that  $\hat{\beta}_2 \sim \mathcal{N}(\cdot, \cdot)$ ,  $e \sim \mathcal{N}(\cdot, \cdot)$ , the Wald statistic follows the exact *F*-distribution:

$$W/p_2 \stackrel{H_0}{\sim} F(p_2, n - p)$$

Asymptotically, the normality assumptions can be relaxed with large sample size (Central Limit Theorem):

$$W \underset{n \rightarrow \infty}{\stackrel{H_0}{\sim}} \chi_{p_2}^2$$

## 2.4 Likelihood Ratio Test

The test-statistic  $\Lambda$ :

$$\begin{aligned}\Lambda &= \frac{\max_{\beta \in H_0} L(\beta)}{\max_{\beta} L(\beta)} \\ \log \Lambda &= \max_{\beta \in H_0} \log L(\beta) - \max_{\beta} \log L(\beta) \\ &= -\frac{1}{2\sigma^2} [RSS_{H_0} - RSS] \quad (\text{in linear model case with Gaussian assumptions}) \\ -2 \log \Lambda &= \frac{RSS_{H_0} - RSS}{\sigma^2} \\ -2 \log \hat{\Lambda} &= \frac{RSS_{H_0} - RSS}{RSS/(n-p)}\end{aligned}$$

Question: why do we use  $RSS/(n-p)$  instead of  $RSS_{H_0}/(n-p+p_2)$  as the estimator of  $\sigma^2$ ?

Answer: the reason is that we assume the full model with a  $p$ -dim  $\beta$  as the true model.

Asymptotically,

$$-2 \log \hat{\Lambda} \underset{n \rightarrow \infty}{\overset{H_0}{\sim}} \chi_{p_2}^2$$

**Notes:**

- (i) LRT: the most powerful test
- (ii) Wald  $\iff$  LRT under linear model assumptions.
- (iii) **Practical tips:**

- do  $t$ -test when only 1 coefficient interested
- do LRT(calculate RSS) or the Wald test when more coefficients interested

	Tests	Decisions to make	Description
(iv)	$H_0^{(j)} : \beta_j = 0, H_1^{(j)} : \beta_j \neq 0$ $j = 1, 2, \dots, p$	$p$	Multiple testing Type I error is hard to control
	$H_0 : \beta = 0, H_1 : \beta \neq 0$	1	One test for $p$ parameters jointly

## 2.5 ANOVA (Analysis of Variance)

$$\begin{aligned}\sum_{i=1}^n (Y_i - \bar{Y})^2 &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \\ SST &= RSS(\text{or } SSE) + SSR\end{aligned}$$

- $\hat{Y}_i = \bar{Y}_i$