

# Linear regression and dummy variable encoding

Miaoyan Wang

Department of Statistics  
UW Madison

Reading: Chapter 6.1-6.3, 6.6 in RC; Chapter 3 in JF.

## Another view of T-test

- Recall the simple linear regression (SLR) model

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad \varepsilon_i \sim i.i.d. N(0, \sigma^2),$$

for all  $i = 1, \dots, n$ .

- Equivalently

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{MVN}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

where  $\mathbf{X}_{n \times 2} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$  denote the  $n \times 2$  design matrix.

- One-sample test is a special case of SLR.
- Two-sample test is also a special case of SLR.

# Equivalence to one-sample test

- Let

$$Y_i = \beta_0 + \varepsilon_i, \quad \varepsilon_i \sim i.i.d. \ N(0, \sigma^2),$$

for all  $i = 1, \dots, n$ .

- Equivalently

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{MVN}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

where  $\mathbf{X}_{n \times 1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$  denote the  $n \times 1$  design matrix, and  $\boldsymbol{\beta} = \beta_0$ .

- The one-sample mean test is equivalent to

$$H_0 : \beta_0 = \mu \text{ vs. } H_A : \beta_0 \neq \mu$$

# Equivalence to two-sample test

- Let

$$Y_i = \beta_0 \mathbb{1}_{i \text{ is in group 1}} + \beta_1 \mathbb{1}_{i \text{ is in group 2}} + \varepsilon_i, \quad \varepsilon_i \sim i.i.d. N(0, \sigma^2),$$

for all  $i = 1, \dots, n$ .

- Equivalently

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{MVN}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

where  $\mathbf{X}_{n \times 2} = \begin{bmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}$  denote the  $n \times 2$  design matrix, and

$$\boldsymbol{\beta} = (\beta_0, \beta_1)'.$$

- The unpaired two sample mean test is equivalent to

$$H_0 : \beta_0 - \beta_1 = 0 \text{ vs. } H_A : \beta_0 - \beta_1 \neq 0$$

# Dummy variable

- The predictors in the linear model can be either continuous (e.g., age, height) or categorical (e.g., gender, group)
- For a categorical predictor that has  $p$  categories, define  $p - 1$  **dummy variables**:

$$X_{ik} = \begin{cases} 1 & \text{observation } i \text{ is in category } k \\ 0 & \text{otherwise} \end{cases}$$

where  $k = 1, \dots, p - 1$ .

- Include dummy variables as predictors in the linear model.
- Example. Consider  $n$  i.i.d. observations from the following model:

$$Y = \beta_0 + \beta_1 \text{Age} + \beta_2 X + \varepsilon, \quad \text{where } \varepsilon \sim i.i.d. N(0, \sigma^2),$$

with  $X = 1$  if male,  $X = 0$  if female.

- **What is the interpretation for  $\beta_0$ ,  $\beta_1$ , and  $\beta_2$ ?**

## Example with categorical variables

Consider the effect of education on hourly wages ( $Y$ ). The education is classified into three categories:

Option in Survey ( $O$ )	Meaning ( $M$ )
1	College dropout
2	College
3	MS and above

Which model makes more sense?

- $Y = \beta_0 + \beta_1 O + \varepsilon$ ?
- $Y = \beta_0 + \beta_1 \mathbf{1}_{\text{college}} + \beta_2 \mathbf{1}_{\text{MS and above}} + \varepsilon$ ?
- $Y = \beta_0 + \beta_1 \mathbf{1}_{\text{college dropout}} + \beta_2 \mathbf{1}_{\text{college}} + \varepsilon$ ?

(In all cases, assume  $\varepsilon \sim i.i.d.N(0, \sigma^2)$ )

## Example (Cont.)

- To include the education as predictor in a regression model, define 2 dummy variables  $X_1$  and  $X_2$ :

Option in Survey ( $O$ )	Meaning ( $M$ )	$X_1$	$X_2$
1	College dropout	0	0
2	College	1	0
3	MS and above	0	1

- Baseline (all dummies 0): college dropout;
- $X_1 = 1$ , if the highest degree is college, 0 otherwise;
- $X_2 = 1$ , if degree with MS and above, 0 otherwise.

Include  $X_1$  and  $X_2$  as dummy variables in a regression model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \underbrace{\beta_3 X_3 + \dots + \beta_p X_p}_{\text{other predictors, e.g., age}} + \varepsilon, \quad \varepsilon \sim i.i.d. N(0, \sigma^2).$$

## Inference on the linear contrast

Recall the study that investigates the effect of education on hourly salary ( $Y$ ):

Education	$X_1$	$X_2$
College dropout	0	0
College	1	0
MS and above	0	1

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon, \quad \text{where } \varepsilon \sim i.i.d. N(0, \sigma^2).$$

Suppose we are interested in testing:

- The mean salary for “MS and above” is the same as for “College”:  
 $H_0 : \beta_1 = \beta_2 \longleftrightarrow H_0 : 0 * \beta_0 + 1 * \beta_1 - 1 * \beta_2 = 0$
- The mean salary for “College” is the same as for “College dropout”:  
 $H_0 : \beta_1 = 0 \longleftrightarrow H_0 : 0 * \beta_0 + 1 * \beta_1 + 0 * \beta_2 = 0$
- Compared to college dropout, the mean salary increase for “MS and above” is twice as that for “College”:  
 $H_0 : \beta_2 = 2\beta_1 \longleftrightarrow H_0 : 0 * \beta_0 + 2 * \beta_1 - 1 * \beta_2 = 0$



# Inference on the linear contrast

- All these hypothesis tests could be expressed as a linear contrast:

$$H_0 : c_0\beta_0 + c_1\beta_1 + c_2\beta_2 = 0 \quad \text{v.s.} \quad H_\alpha : c_0\beta_0 + c_1\beta_1 + c_2\beta_2 \neq 0,$$

for a given vector  $\mathbf{c} = (c_0, c_1, c_2)$ . Let  $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2)'$ .

- What is the distribution of  $\mathbf{c}'\hat{\boldsymbol{\beta}}$  under the null? Multivariate normal with

$$\mathbb{E}(\mathbf{c}'\hat{\boldsymbol{\beta}}) = \mathbf{c}'\boldsymbol{\beta}, \quad \text{Var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) = \underline{\hspace{2cm}} = \sigma^2 \mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}$$

- In case  $\sigma^2$  is unknown, plug in the estimator  $\hat{\sigma}^2$ . (what is the form of  $\hat{\sigma}^2$ ?)

$$\frac{\mathbf{c}'\hat{\boldsymbol{\beta}} - \mathbf{c}'\boldsymbol{\beta}}{\sqrt{\widehat{\text{Var}}(\mathbf{c}'\hat{\boldsymbol{\beta}})}} \sim T_{n-3}$$