

# STAT 231 Online Spring 2020

Linxuntian

July 2020

## Contents

<b>1</b>	<b>Week 7</b>	<b>2</b>
1.1	Test of Hypothesis Part I . . . . .	2
1.2	Test of Hypothesis Part II . . . . .	2
<b>2</b>	<b>Week 8</b>	<b>3</b>
<b>3</b>	<b>Week 9</b>	<b>4</b>
3.1	Gaussian Response Models Part I . . . . .	4
3.2	Gaussian Response Models Part II . . . . .	4
<b>4</b>	<b>Week 10</b>	<b>4</b>
4.1	Gaussian Response Models Part III . . . . .	4
4.2	Comparing Means of Two Populations . . . . .	5
<b>5</b>	<b>Week 11</b>	<b>6</b>
5.1	Gaussian Response Models Part IV . . . . .	6
5.2	Multinomial Models and Goodness of Fit . . . . .	6
<b>6</b>	<b>Week 12</b>	<b>7</b>
6.1	Two-Way Tables and Tests of Hypothesis . . . . .	7
6.2	Cause and Effect . . . . .	7

# 1 Week 7

## 1.1 Test of Hypothesis Part I

### Definition

A **test statistic** or **discrepancy measure** is a function of the data  $D = g(Y)$  that is constructed to measure the degree of “agreement” between the data  $Y$  and the null hypothesis  $H_0$ .

### Summary (Steps of a Statistical Test of Hypothesis)

1. Assume that the null hypothesis  $H_0$  will be tested using data  $Y$
2. Adopt a test statistic or discrepancy measure  $D(Y)$  for which large value of  $D$  are less consistent with  $H_0$ . Let  $d = D(y)$  be the corresponding observed value of  $D$
3. Calculate

$$\begin{aligned} p\text{-value} &= P(D \geq d; \text{assuming } H_0 \text{ is true}) \\ &= P(D \geq d; H_0) \end{aligned}$$

4. Draw a conclusion based on the  $p$ -value

## 1.2 Test of Hypothesis Part II

### Summary (Hypothesis Test for Binomial Model)

1. Test  $H_0: \theta = \theta_0$  using  $Y \sim \text{Bi}(n, \theta)$
2. Test statistic  $D(Y) = |Y - n\theta_0|$ . Let  $d = |y - n\theta_0|$  be the observed value of  $D$
3. Calculate

$$\begin{aligned} p\text{-value} &= P(D \geq d; \theta = \theta_0) \\ &= P(|Y - n\theta_0| \geq d) \end{aligned}$$

where  $Y \sim \text{Bi}(n, \theta)$

If  $n$  is large, we can approximate

$$\begin{aligned} p\text{-value} &\approx P\left(|Z| \geq \frac{d}{\sqrt{n\theta_0(1-\theta_0)}}\right) \quad \text{where } Z \sim N(0, 1) \\ &= 2 \left[ 1 - P\left(Z \leq \frac{d}{\sqrt{n\theta_0(1-\theta_0)}}\right) \right] \end{aligned}$$

4. Draw a conclusion

### Summary (Hypothesis Test for Gaussian Model)

When the **variance** is unknown

1. 123

When the **mean** is unknown

## 2 Week 8

### Summary (Likelihood Ratio Test of Hypothesis)

1. To test  $H_0: \theta = \theta_0$

2. Likelihood ratio test statistic  $\Lambda(\theta_0) = -2 \log \left[ \frac{L(\theta_0)}{L(\hat{\theta})} \right],$

with the observed value  $\lambda(\theta_0) = -2 \log \left[ \frac{L(\theta_0)}{L(\hat{\theta})} \right]$

3. Calculate

$$\begin{aligned} p\text{-value} &\approx P(W \leq \lambda(\theta_0)) \quad \text{where } W \sim \chi^2(1) \\ &= 2[1 - P(Z \leq \sqrt{\lambda(\theta_0)})] \quad \text{where } Z \sim G(0, 1) \end{aligned}$$

4. Draw a conclusion

### 3 Week 9

#### 3.1 Gaussian Response Models Part I

##### Confidence Interval for beta $\beta$

A 100p% Confidence Interval for  $\beta$  is given by

$$\hat{\beta} \pm a \frac{s_e}{\sqrt{S_{xx}}}$$

where  $P(T \leq a) = \frac{(1+p)}{2}$  and  $T \sim t(n-2)$

#### 3.2 Gaussian Response Models Part II

##### Summary (Hypothesis testing of No Relationship for beta $\beta$ )

1. Test  $H_0: \beta = 0$

2. Test statistic  $D = \frac{|\tilde{\beta} - 0|}{s_e/\sqrt{S_{xx}}}$ , with observed value  $d = \frac{|\hat{\beta} - 0|}{s_e/\sqrt{S_{xx}}}$

3. Calculate

$$\begin{aligned} p\text{-value} &= P\left(|T| \geq \frac{|\hat{\beta} - 0|}{s_e/\sqrt{S_{xx}}}\right) \\ &= 2 \left[ 1 - P\left(T \leq \frac{|\hat{\beta} - 0|}{s_e/\sqrt{S_{xx}}}\right) \right] \end{aligned}$$

where  $T \sim t(n-2)$

4. Draw Conclusion

##### Confidence Interval for Mean $\mu(x) = \alpha + \beta x$

A 100p% Confidence Interval for  $\mu(x) = \alpha + \beta x$  is given by

$$\begin{aligned} \hat{\mu}(x) \pm as_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \\ = \hat{\alpha} + \hat{\beta}x \pm as_e \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \end{aligned}$$

where  $P(T \leq a) = \frac{(1+p)}{2}$  and  $T \sim t(n-2)$

### 4 Week 10

#### 4.1 Gaussian Response Models Part III

##### Prediction Interval For Future Response $Y$

A 100p% Prediction Interval for a future response  $Y$  is given by

$$\begin{aligned} \hat{\mu}(x) \pm as_e \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \\ = \hat{\alpha} + \hat{\beta}x \pm as_e \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \end{aligned}$$

where  $P(T \leq a) = \frac{(1+p)}{2}$  and  $T \sim t(n-2)$

### Definition

For the simple linear regression model, let

$$\hat{\mu}_i = \hat{\alpha} + \hat{\beta}x_i$$

and let

$$\hat{r}_i = y_i - \hat{\mu}_i \quad i = 1, 2, \dots, n$$

The  $\hat{r}_i$ 's are called **residuals** since  $\hat{r}_i$  represents what is “left” after the model has been “fitted” to the data

### Definition

Define the **standardized residuals** as

$$\hat{r}_i = \frac{\hat{r}_i}{s_e} = \frac{y_i - \hat{\mu}_i}{s_e} \quad i = 1, 2, \dots, n$$

## 4.2 Comparing Means of Two Populations

### Definition

**Pooled estimate of variance** is defined as

$$\begin{aligned} s_p^2 &= \frac{1}{n_1 + n_2 - 2} \left[ \sum_{i=1}^{n_1} (y_{1i} - \bar{y}_1)^2 + \sum_{i=1}^{n_2} (y_{2i} - \bar{y}_2)^2 \right] \\ &= \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \end{aligned}$$

### C.I. for Difference in Means $\mu_1 - \mu_2$

A 100p% C.I. for  $\mu_1 - \mu_2$  is given by

$$\bar{y}_1 - \bar{y}_2 \pm as_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

where  $P(T \leq a) = \frac{(1+p)}{2}$  and  $T \sim t(n_1 + n_2 - 2)$

### Summary (Hypothesis testing of Difference in Means $\mu_1 - \mu_2$ )

1. Test  $H_0: \mu_1 - \mu_2 = 0$

2. Test statistic  $D = \frac{|\bar{Y}_1 - \bar{Y}_2 - 0|}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$  and observed value  $d = \frac{|\bar{y}_1 - \bar{y}_2 - 0|}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

3. Calculate

$$\begin{aligned} p\text{-value} &= P(|T| \geq d) \\ &= 2[1 - P(T \leq d)] \end{aligned}$$

where  $T \sim t(n_1 + n_2 - 2)$

4. Draw a conclusion

### C.I. for $\sigma$

A 100p% Confidence Interval for  $\sigma$  is given by

## 5 Week 11

### 5.1 Gaussian Response Models Part IV

#### Summary (Hypothesis testing of Difference in Means $\mu = \mu_1 - \mu_2$ with-in-pair)

1. Test  $H_0: \mu = 0$  using  $Y_i = Y_{1i} - Y_{2i} \sim G(\mu_1 - \mu_2, \sigma)$ ,  $i = 1, 2, \dots, n$ , independently

2. Test statistic  $D = \frac{|\bar{Y} - 0|}{S/\sqrt{n}}$ , with observed value  $d = \frac{|\bar{y} - 0|}{s/\sqrt{n}}$

3. Calculate

$$p\text{-value} = 2[1 - P(T \leq d)]$$

where  $T \sim t(n-1)$

4. Draw a conclusion

### 5.2 Multinomial Models and Goodness of Fit

#### Summary (Hypothesis testing of $\theta$ )

1. Test  $H_0: \theta = \theta_0 = \left(\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k}\right)$  using  $f(y_1, y_2, \dots, y_k; \theta_1, \theta_2, \dots, \theta_k)$

2. Likelihood ratio test statistic  $\Lambda = 2 \sum_{j=1}^k Y_j \log \left( \frac{Y_j}{E_j} \right)$

and observed value  $\lambda = 2 \sum_{j=1}^k y_j \log \left( \frac{y_j}{e_j} \right)$

3. Calculate

$$p\text{-value} = P(\Lambda \geq \lambda; H_0) \approx P(W \geq \lambda)$$

where  $W \sim \chi^2(k-1-p)$

4. Draw a conclusion

## 6 Week 12

### 6.1 Two-Way Tables and Tests of Hypothesis

#### Definition

An example of **Two-way table**...

Let  $\theta_{ij}$  be the probability a randomly selected individual is combined type  $(A_i, B_j)$  and note that  $\sum_{i=1}^a \sum_{j=1}^b \theta_{ij} = 1$ . The  $a \times b$  frequencies  $(Y_{11}, Y_{12}, \dots, Y_{ab})$  follow a Multinomial distribution with  $k = ab$  classes.

#### Summary (Hypothesis testing of Dependence of $A$ and $B$ class)

1. Test  $H_0: \theta_{ij} = \alpha_i \beta_j$

Note that the expected frequencies under the hypothesis  $H_0$  are given by  $e_{ij} = \frac{r_i c_j}{n}$

2. Likelihood ratio statistic  $\Lambda = 2 \sum_{i=1}^a \sum_{j=1}^b Y_{ij} \log \left( \frac{Y_{ij}}{E_{ij}} \right)$ ,

with observed value  $\lambda = 2 \sum_{i=1}^a \sum_{j=1}^b y_{ij} \log \left( \frac{y_{ij}}{e_{ij}} \right)$

3. Calculate

$$p\text{-value} \approx P(W \geq \lambda)$$

where  $W \sim \chi^2((a-1) \times (b-1))$

**Note:**

- (a) If  $((a-1) \times (b-1)) = 1$ , then

$$p\text{-value} = 2[1 - P(Z \leq \sqrt{\lambda})], \quad \text{where } Z \sim G(0, 1)$$

- (b) If  $((a-1) \times (b-1)) = 2$ , then

$$p\text{-value} = P(W \geq \lambda) = e^{-\lambda/2}, \quad \text{where } W \sim \chi^2(2) = \text{Exponential}(2)$$

4. Draw a conclusion

### 6.2 Cause and Effect

#### Definition

Let be  $y$  a response variate and let be  $x$  an explanatory variate associated with units in a population or process. Then, if all other factors that affect are held constant, let us change  $x$  (or observe different values of  $x$ ) and see if  $y$  changes. If it does, we say that  $x$  has a **causal effect** on  $y$ .

#### Definition

$x$  has a **causal effect** on  $Y$  if, when all other factors that affect are held constant, a change in  $x$  (or observing different values of  $x$ ) induces a change in a property of the distribution of  $Y$ .

**Reasons two variates could be related include the following:**

1. The explanatory variate is the direct cause of the response variate
2. The response variate is causing a change in the explanatory variate
3. The explanatory variate is a contributing but not sole cause of the response variate

4. Confounding variates may exist
5. Both variates may result from a common cause
6. Both variates are changing with time
7. The association may be due to coincidence