

VE401 RC Week9

Wang Yangyang

UM-SJTU JI

2022 Spring

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Categorical Random Variables

A random variable X that can take on the values $1, \dots, k$ with respective probabilities p_1, \dots, p_k as above. A random sample of size n from X is collected and the results are expressed as a *random vector*

$$(X_1, X_2, \dots, X_k) \quad \text{with} \quad X_1 + X_2 + \dots + X_k = n$$

The Multinomial Distribution: A random vector $((X_1, \dots, X_k), f_{X_1 X_2 \dots X_k})$ where

$$f_{X_1 X_2 \dots X_k}(x_1, \dots, x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$$

$p_1, \dots, p_k \in (0, 1)$, $n \in \mathbb{N} \setminus \{0\}$ is said to have a multinomial distribution with parameters n and p_1, \dots, p_k

- $E[X_i] = np_i, \quad i = 1, \dots, k$
- $\text{Var}[X_i] = np_i(1 - p_i), i = 1, \dots, k$
- $\text{Cov}[X_i, X_j] = -np_i p_j, 1 \leq i < j \leq k$

The Pearson Statistics

Let $((X_1, \dots, X_k), f_{X_1 X_2 \dots X_k})$ be a multinomial random variable with parameters n and p_1, \dots, p_k . For large n the *Pearson statistic*

$$\sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} = \sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i}$$

follows an approximate chi-squared distribution with $k - 1$ degrees of freedom (because we have $k - 1$ independent cells).

Cochran's Rule: This tell us how large n needs to be for the chi-squared distribution to be a good approximation to the true distribution of the Pearson statistic when

$$\begin{aligned} E[X_i] = np_i &\geq 1, & \text{for all } i = 1, \dots, k \\ E[X_i] = np_i &\geq 5, & \text{for 80\% of all } i = 1, \dots, k \end{aligned}$$

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Test for Multinomial Distribution

Let (X_1, \dots, X_k) be a sample of size n from a categorical random variable with parameters (p_1, \dots, p_k) . We perform the chi-squared goodness-of-fit test.

Note: In this test, we directly make assumptions on parameters p_i *without estimation based on samples*. This may happen when we already have some prior knowledge of the distribution (e.g. PRNG).

Procedures

- ① Set

$$H_0 : p_i = p_{i0}, \quad i = 1, \dots, k$$

- ② Calculate the expected values

$$E_i = np_{i0}$$

Then test whether the *Cochran's rule* is satisfied.

Test for Multinomial Distribution

Procedures

- ③ If satisfied, calculate the Pearson statistic.

$$\chi^2_{k-1} = \sum_{i=1}^k \frac{(X_i - np_{i0})^2}{np_{i0}}$$

which follows a chi-squared distribution with

degrees of freedom: independent cells $-m = k - 1$
independent cells $= k - 1, m=0$

- ④ We reject H_0 at significance level α if $\chi^2_{k-1} > \chi^2_{\alpha, k-1}$.

Goodness-of-Fit Test for Discrete Distribution

Now, we calculate the *estimates for parameters*, to make assumptions indirectly

Procedures

- ① Suppose we guess that data *follow some distribution*, so we set the hypothesis as

H_0 : A specific distribution with unknown parameter p_i
e.g., H_0 : A Poisson distribution with parameter k .

- ② Estimate parameters p_i from the sample based on your hypothesis.
e.g., for *Poisson distribution*, estimate k by $\hat{k} = \bar{X}$.
Then we can calculate p_i . Suppose we have three categories $x = 0, x = 1, x \geq 2$, then

$$p_0 = P[X = 0] = \frac{e^{-\hat{k}} \hat{k}^0}{0!}, \quad p_1 = P[X = 1] = \frac{e^{-\hat{k}} \hat{k}^1}{1!},$$

$$p_2 = P[X \geq 2] = 1 - P[X = 0] - P[X = 1]$$

Goodness-of-Fit Test for Discrete Distribution

Procedures

- ③ Calculate the expected values $E_i = np_i$. Then make a table by yourself as below,

Category i	Exp. Frequency E_i	Obs. Frequency O_i
0	np_0	x_0
1	np_1	x_1
2	np_2	x_2
...

- ④ Test whether the *Cochran's Rule* is satisfied. If not satisfied, then go back to procedure (ii) and (iii) and change your number of categories.

Goodness-of-Fit Test for Discrete Distribution

Procedures

- ⑤ If (iv) satisfied, calculate the *Pearson statistic*.

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

now follows a chi-squared distribution with

$$\# \text{ independent cells} - m = k - 1 - m$$

degrees of freedom, where m is the *number of parameters that we estimate*. e.g., for the previous Poisson distribution test with 3 categories,

$$k = 3, m = 1$$

- ⑥ We reject H_0 at significance level α if χ^2 exceeds the critical value.

Outline

① Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

② Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

③ Supplementary Materials

- Prepare MMA File

Test for Independence

We define the marginal row and column sums

$$n_{i\cdot} = \sum_{j=1}^c n_{ij}, \quad n_{\cdot j} = \sum_{i=1}^r n_{ij}$$

For a contingency table as below,

	column 1	column 2	column 3	
row 1	n_{11}	n_{12}	n_{13}	$n_{1\cdot}$
row 2	n_{21}	n_{22}	n_{23}	$n_{2\cdot}$
row 3	n_{31}	n_{32}	n_{33}	$n_{3\cdot}$
	$n_{\cdot 1}$	$n_{\cdot 2}$	$n_{\cdot 3}$	n

Test for Independence

Procedures

- ① If the hypothesis is that row and column categorizations are independent, then it should be the case that

$$H_0 : p_{ij} = p_i \cdot p_j$$

- ② Estimates for the row and column probabilities are $\widehat{p}_{i\cdot} = \frac{n_{i\cdot}}{n}$, $\widehat{p}_{\cdot j} = \frac{n_{\cdot j}}{n}$, so if H_0 is assumed,

$$\widehat{p}_{ij} = \widehat{p}_{i\cdot} \cdot \widehat{p}_{\cdot j} = \frac{n_{i\cdot} \cdot n_{\cdot j}}{n^2}$$

- ③ Calculate the expected values

$$E_{ij} = n \cdot \widehat{p}_{ij} = \frac{n_{i\cdot} \cdot n_{\cdot j}}{n}$$

Then test whether the *Cochran's rule* is satisfied.

Test for Independence

Procedures

- ④ If satisfied, calculate the Pearson statistic.

$$\chi^2_{(r-1)(c-1)} = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

which follows a *chi-squared distribution* with degrees of freedom:

$$\text{\#independent cells} - m = rc - 1 - (r - 1 + c - 1) = (r - 1)(c - 1)$$

- #independent cells = $rc - 1$ because we have rc categories.
- $m = r - 1 + c - 1$ because we estimate p_i with $i \leq r - 1$ and p_j with $j \leq c - 1$.

- ⑤ We reject H_0 if the value of $\chi^2_{(r-1)(c-1)}$ exceeds the critical value.

Test for Comparing Proportions

Now the row totals are fixed, rewrite the table in terms of proportions:

	column 1	column 2	column 3	
row 1	p_{11}	p_{12}	p_{13}	$p_{1.} = 1$ (fixed)
row 2	p_{21}	p_{22}	p_{23}	$p_{2.} = 1$ (fixed)
row 3	p_{31}	p_{32}	p_{33}	$p_{3.} = 1$ (fixed)
row 4	p_{41}	p_{42}	p_{43}	$p_{4.} = 1$ (fixed)

We want to compare proportions from each row, so

$$H_0 : \begin{cases} p_{11} = p_{21} = p_{31} = p_{41} \\ p_{12} = p_{22} = p_{32} = p_{42} \\ p_{13} = p_{23} = p_{33} = p_{43} \end{cases}$$

Test for Comparing Proportions

Procedure

- ① Supposing that H_0 is true,

$$p_j := p_{1j} = p_{2j} = p_{3j} = p_{4j}$$

where p_j is the proportion of all objects following into the j th column. If H_0 is assumed, estimates for the column proportions are

$$\hat{p}_j = \frac{n_{.j}}{n}$$

- ② Calculate the expected values

$$E_{ij} = n_i \cdot \hat{p}_{ij} = n_i \cdot \hat{p}_j = \frac{n_i \cdot n_{.j}}{n}$$

Then test whether the *Cochran's rule* is satisfied.

Test for Comparing Proportions

Procedure

- ① If satisfied, calculate the Pearson statistic.

$$X^2_{(r-1)(c-1)} = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

which follows a *chi-squared* distribution with

independent cells $- m = r(c - 1) - (c - 1) = (r - 1)(c - 1)$

degrees of freedom.

- #independent cells $= r(c - 1)$ because only the number of objects in the first $c - 1$ columns can be independently chosen, so we have a total of $r \cdot (c - 1)$ independent cells.
- $m = (c - 1)$ because we estimate p_j for $c - 1$ times.

- ② We reject H_0 if the value of $X^2_{(r-1)(c-1)}$ exceeds the critical value.

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Simple Linear Regression Model

We assume that

$$Y \mid x = \beta_0 + \beta_1 x + E,$$

where $E[E] = 0$. We want to find estimators

$$\begin{aligned} B_0 &:= \widehat{\beta_0} = \text{estimator for } \beta_0, & b_0 &= \text{estimate for } \beta_0, \\ B_1 &:= \widehat{\beta_1} = \text{estimator for } \beta_1, & b_1 &= \text{estimate for } \beta_1. \end{aligned}$$

Assumptions

- For each value of x , the random variable follows a normal distribution with variance σ^2 and mean $\mu_{Y|x} = \beta_0 + \beta_1 x$.
- The random variables $Y \mid x_1$ and $Y \mid x_2$ are independent if $x_1 \neq x_2$.

Least Square Estimator

We have the error *sum of squares*

$$SS_E := \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

To *minimize* it, we take

$$\frac{\partial SS_E}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0$$

$$\frac{\partial SS_E}{\partial b_1} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) x_i = 0$$

which gives

$$b_1 = \frac{S_{xy}}{S_{xx}}, \quad b_0 = \bar{y} - b_1 \bar{x}$$

Properties

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x}) x_i = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

$$S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \bar{y}) y_i = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x}) y_i = \sum_{i=1}^n (y_i - \bar{y}) x_i =$$

$$\sum_{i=1}^n x_i y_i - n\bar{x} \cdot \bar{y} = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right).$$

$$b_1 = \frac{S_{xy}}{S_{xx}}, \quad b_0 = \bar{y} - b_1 \bar{x}, \quad SS_E = S_{yy} - b_1 S_{xy}.$$

Calculation Procedure

- ① Find $\sum x_i$, $\sum y_i$, $\sum x_i^2$, $\sum y_i^2$, $\sum x_i y_i$ and calculate

$$S_{xx} = \sum x_i^2 - \frac{1}{n} \left(\sum x_i \right)^2, \quad S_{yy} = \sum y_i^2 - \frac{1}{n} \left(\sum y_i \right)^2$$

$$S_{xy} = \sum x_i y_i - \frac{1}{n} \left(\sum x_i \right) \left(\sum y_i \right)$$

- ② Obtain b_1 and b_0 by

$$b_1 = \frac{S_{xy}}{S_{xx}}, \quad b_0 = \bar{y} - b_1 \bar{x}.$$

- ③ Calculate other quantities as required, e.g.,

$$SS_T = S_{yy}, \quad SS_E = S_{yy} - \frac{S_{xy}^2}{S_{xx}}, \quad R^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}}.$$

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Distribution of Estimator for Variance

An unbiased estimator for variance σ^2 is given by

$$S^2 = \frac{SS_E}{n-2} = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{\mu}_{Y|x_i})^2$$

The statistic

$$\chi_{n-2}^2 = \frac{(n-2)S^2}{\sigma^2} = \frac{SS_E}{\sigma^2}$$

follows a *chi-squared distribution* with $n-2$ degrees of freedom.

Distribution of B_1 with Estimated Variance

The least squares estimator B_1 for β_1 follows a normal distribution with

$$E[B_1] = \beta_1, \quad \text{Var}[B_1] = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}.$$

The statistics

$$T_{n-2} = \frac{B_1 - \beta_1}{S/\sqrt{S_{xx}}}$$

follows *T-distributions* with $n - 2$ degrees of freedom.

The $100(1 - \alpha)\%$ *confidence interval* of β_1 is given by

$$B_1 \pm t_{\alpha/2, n-2} \frac{S}{\sqrt{S_{xx}}}.$$

Distribution of B_0 with Estimated Variance

The least squares estimator B_0 for β_0 follows a normal distribution with

$$E[B_0] = \beta_0, \quad \text{Var}[B_0] = \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

The statistics

$$T_{n-2} = \frac{B_0 - \beta_0}{S \sqrt{\sum x_i^2} / \sqrt{n S_{xx}}}$$

follows *T-distributions* with $n - 2$ degrees of freedom.

The $100(1 - \alpha)\%$ *confidence interval* of β_0 is given by

$$B_0 \pm t_{\alpha/2, n-2} \frac{S \sqrt{\sum x_i^2}}{\sqrt{n S_{xx}}}.$$

Distribution of Estimated Mean

The estimated mean $\hat{\mu}_{Y|x}$ follows a normal distribution with mean and variance

$$E[\hat{\mu}_{Y|x}] = \mu_{Y|x}, \quad \text{Var}[\hat{\mu}_{Y|x}] = \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}} \right) \sigma^2.$$

Therefore, the statistic

$$T_{n-2} = \frac{\hat{\mu}_{Y|x} - \mu_{Y|x}}{S \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}}$$

follows a *T-distribution* with $n - 2$ degrees of freedom. A $100(1 - \alpha)\%$ *confidence interval* for $\mu_{Y|x}$ is given by

$$\hat{\mu}_{Y_x} \pm t_{\alpha/2, n-2} S \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}.$$

Distribution and CI for Predictor

The statistic $Y | x - \widehat{Y} | x$ follows a normal distribution with mean and variance

$$E[Y | x - \widehat{Y} | x] = 0, \quad \text{Var}[Y | x - \widehat{Y} | x] = \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}\right) \sigma^2.$$

Therefore, the statistic

$$T_{n-2} = \frac{Y | x - \widehat{Y} | x}{S \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}}$$

follows a *T-distribution* with $n - 2$ degrees of freedom. A $100(1 - \alpha)\%$ ***prediction** interval* for $Y | x$ is given by

$$\widehat{Y} | x \pm t_{\alpha/2, n-2} S \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

Quantities

- ① *Total sum of squares:*

$$SS_T = S_{yy} = \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

- ② *Error sum of squared:*

$$SS_E = \sum_{i=1}^n (Y_i - (B_0 + B_1 x_i))^2 = S_{yy} - B_1 S_{xy} = S_{yy} - \frac{S_{xy}^2}{S_{xx}}.$$

- ③ *Coefficient of determination:* the proportion of the total variation in Y that is explained by the linear model.

$$R^2 = \frac{SS_T - SS_E}{SS_T} = \frac{S_{xy}^2}{S_{xx} S_{yy}}.$$

Test for Significance with B_1

Let $(x_i, Y \mid x_i), i = 1, \dots, n$ be a random sample from $Y \mid x$. We reject

$$H_0 : \beta_1 = 0$$

at significance level α if the test statistic

$$T_{n-2} = \frac{B_1}{S/\sqrt{S_{xx}}}$$

satisfies $|T_{n-2}| > t_{\alpha/2, n-2}$.

Test for Significance with R^2

Let $(x_i, Y \mid x_i), i = 1, \dots, n$ be a random sample from $Y \mid x$. We reject

$$H_0 : \beta_1 = 0$$

at significance level α if the test statistic

$$T_{n-2} = \frac{B_1}{S/\sqrt{S_{xx}}} = \frac{R\sqrt{n-2}}{\sqrt{1-R^2}}$$

satisfies $|T_{n-2}| > t_{\alpha/2, n-2}$.

Test for Correlation with R^2

Let (X, Y) follow a *bivariate normal distribution* with correlation coefficient $\rho \in (-1, 1)$. Let R be the estimator for ρ . Then we reject

$$H_0 : \rho = 0$$

at significance level α if the test statistic

$$T_{n-2} = \frac{R\sqrt{n-2}}{\sqrt{1-R^2}}$$

satisfies $|T_{n-2}| > t_{\alpha/2, n-2}$.

Testing for Lack of Fit

SS_E is the variance of Y explained by the model.

- Error sum of squares due to pure error.

$$SS_{E,pe} := \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \sum_{i=1}^k \frac{1}{n_i} \left(\sum_{j=1}^{n_i} Y_{ij} \right)^2.$$

The statistic $SS_{E,pe}/\sigma^2$ follows a *chi-squared distribution* with $\sum_{i=1}^k (n_i - 1) = n - k$ degrees of freedom.

- Error sum of squares due to lack of fit:

$$SS_{E,ff} := SS_E - SS_{E,pe}.$$

The statistic $SS_{E,ff}/\sigma^2$ follows a *chi-squared* distribution with $k - 2$ degrees of freedom.

Testing for Lack of Fit

Test for lack of fit. Let x_1, \dots, x_k be regressors and $Y_{i1}, \dots, Y_{in_i}, i = 1, \dots, k$ the measured responses at each of the regressors. Let $SS_{E,pe}$ and $SS_{E,lf}$ be the *pure error* and *lack-of-fit sums of squares* for a linear regression model. Then we reject at significance level α .

H_0 : the linear regression model is appropriate if the test statistic

$$F_{k-2, n-k} = \frac{SS_{E,f}/(k-2)}{SS_{E,pe}/(n-k)}$$

satisfies $F_{k-2, n-k} > f_{\alpha, k-2, n-k}$.

Outline

1 Categorical Data

- Pearson Statistics and Multinomial Distribution
- Goodness-of-Fit Test
- Test for Independence

2 Simple Linear Regression

- Basic Calculation
- Estimations and Predictions
- Model Analysis

3 Supplementary Materials

- Prepare MMA File

It is suggested that you solve problems in assignments using Mathematica. It's the best way to prepare for Final Exam. This notebook file, for your reference, is credited to previous TA Zhang Xingjian and Joy Dong. It would be better to write your own notebook file.

End

Credit to Zhanpeng Zhou (TA of SP21)

Credit to Fan Zhang (TA of SU21)

Credit to Liying Han (TA of SP21)

Credit to Zhenghao Gu (TA of SP20)