

10.11 Goodness of Fit test (Karl Pearson 1900)

Suppose $X_1 \sim B(m, p_1)$. With $\begin{cases} mp_1 \geq 5 \\ m(1-p_1) \geq 5 \end{cases}$

$$Z = \frac{X_1 - mp_1}{\sqrt{mp_1(1-p_1)}} \sim N(0, 1) \quad (\text{Central Limit Theorem})$$

$$\text{and } Z^2 = \frac{(X_1 - mp_1)^2}{mp_1(1-p_1)} \approx \chi^2(1)$$

$$\text{Now } \frac{a}{mp_1} + \frac{a}{m(1-p_1)} = \frac{a(1-p_1) + ap_1}{mp_1(1-p_1)} = \frac{a}{mp_1(1-p_1)}$$

Therefore letting $a = (X_1 - mp_1)^2$ we get.

$$Q = Z^2 = \frac{(X_1 - mp_1)^2}{mp_1(1-p_1)} = \frac{(X_1 - mp_1)^2}{mp_1} + \frac{(X_1 - mp_1)^2}{m(1-p_1)}$$

Let $X_2 = m - X_1$ and $p_2 = 1 - p_1$.

$$\text{We have } (X_1 - mp_1)^2 = (m - \tilde{X}_1 - m(1-p_1))^2 = (X_2 - mp_2)^2$$

which gives,

$$Q = Z^2 = \frac{(X_1 - mp_1)^2}{mp_1} + \frac{(X_2 - mp_2)^2}{mp_2} \approx \chi^2(1)$$

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

In a multinomial distribution the sample space consists of k (instead of two) mutually exclusive outcomes, and it can be shown that $Q = \sum_{i=1}^k \frac{(X_i - mp_i)^2}{mp_i} \approx \chi^2(k-1)$

$$\text{where } mp_i \text{ must not be too small } (mp_i \geq 5)$$

Example. Tossing a die. Did X be the outcome

Want to find out. Is $P(X=x) = \frac{1}{6}$, $x=1, 2, \dots, 6$

For 120 tosses we observed

$x:$	1	2	3	4	5	6
Observed	20	22	17	18	19	24
Expected	20	20	20	20	20	20

The test statistic for a goodness-of-fit test is

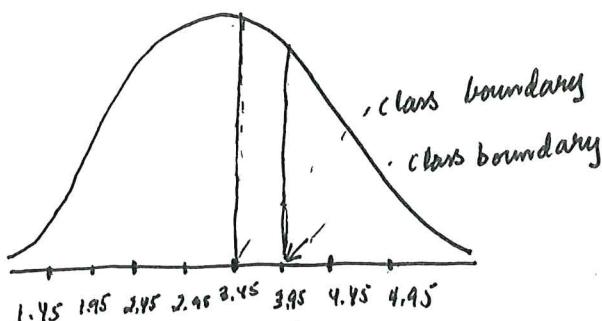
$$Q = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \approx \chi^2_{(k-1)}$$

$$Q_{\text{obs}} = \frac{(20-20)^2 + (22-20)^2 + (17-20)^2 + (18-20)^2 + (19-20)^2 + (24-20)^2}{20} \approx 1.7$$

$$\chi^2_{0.05-1} = 11.07 \quad \text{i.e. no reason to reject } H_0: P(X=x) = \frac{1}{6}, x=1, 2, \dots, 6.$$

Continuous distributions

Lifetime of batteries, $T \sim N(3.5, 0.7^2)$?



<u>Class Boundaries</u>	O_i	e_i	
1.45 - 1.95	2	0.5	
1.95 - 2.45	1	2.1	0.2125
2.45 - 2.95	7	5.9	
2.95 - 3.45	15	10.3	0.2773
3.45 - 3.95	10	10.7	0.2678
3.95 - 4.45	5	7.0	0.2625
4.45 - 4.95	3	3.5	

Expected numbers for the class (cell) 2.95 - 3.45

Probability to fall in the class

$$P(2.95 \leq T \leq 3.45) = P\left(\frac{2.95 - 3.5}{0.7} \leq \frac{T - 3.5}{0.7} \leq \frac{3.45 - 3.5}{0.7}\right)$$

$$= \Phi(-0.07) - \Phi(-0.79) = 0.4721 - 0.2148 = 0.2573$$

Expected number is $40 \cdot 0.2573 = \underline{10.3}$

Cells with expected number less than 5 should be joined with neighbours. We are left with 4 cells.

1.45 - 1.95, 2.95 - 3.45, 3.45 - 3.95, 3.95 - 4.45

$$Q = \frac{(7 - 8.5)^2}{8.5} + \frac{(15 - 10.3)^2}{10.3} + \frac{(10 - 10.7)^2}{10.7} + \frac{(8 - 10.5)^2}{10.5} = 3.05$$

$$\chi^2_{0.05}(3) = 7.815 \Rightarrow \text{No reason to reject H}_0: T \sim N(3.5, 0.7^2)$$

10.12.

Test for Independence for Categorical data.

Example: In order to find out if there are dependence between eye colour and hair colour for humans in a given population, 80 persons are randomly drawn from the population.

Eye colour is classified as blue (B_1) or brown (B_2) and hair colour as light blond (H_1), dark blond (H_2) or dark (H_3) frequencies

The observed data are presented in a contingency table.

Eye-colour	Hair-colour			Total
	H_1	H_2	H_3	
B_1	29	12	7	48
B_2	10	12	10	32
Total	39	24	17	80

This is a 2×3 contingency Table

Let us define $P(B_I) = p_i$, $I=1, 2$ and $P(H_J) = p_j$, $J=1, 2, 3$

Further let $p_{ij} = P(B_i \cap H_j)$, $i=1, 2$, $j=1, 2, 3$

Independence between eye colour and hair colour can be defined as $p_{ij} = p_i \cdot p_j$, \hat{p}_{ij}

$$H_0: p_{ij} = p_i \cdot p_j \quad \hat{p}_{ij}, \quad i=1, 2, \quad j=1, 2, 3.$$

The ^{estimated} expected cell numbers are $80 \cdot \hat{p}_{ij}$, $i=1, 2$, $j=1, 2, 3$

Under H_0 this becomes $80 \cdot \hat{p}_i \cdot \hat{p}_j = \frac{80 \cdot \frac{m_i}{80} \cdot \frac{c_j}{80}}{80} = \frac{\cancel{80} \cdot m_i \cdot c_j}{\cancel{80}^4} = \frac{m_i \cdot c_j}{4}$

$$\text{or } \frac{(\text{row total}) \times (\text{column total})}{\text{grand total}}$$

The 2×3 contingency Table with observed and expected numbers

	H1	H2	H3	Total
B1	29 (23.4)	12 (14.4)	7 (10.2)	48
B2	10 (15.6)	12 (9.6)	10 (6.8)	32
Total	39	24	17	80

$$\text{Ehs. } 23.4 = \frac{48 \cdot 39}{80}, \quad 15.6 = \frac{32 \cdot 39}{80}$$

$$Q = \frac{(29 - 23.4)^2}{23.4} + \frac{(12 - 14.4)^2}{14.4} + \frac{(7 - 10.2)^2}{10.2} + \frac{(10 - 15.6)^2}{15.6}$$

$$+ \frac{(12 - 9.6)^2}{9.6} + \frac{(10 - 6.8)^2}{6.8} = 6.86.$$

number of rows number of columns

The degrees of freedom $v = (n-1)(c-1)$

$$\chi^2_{0.05}(2) = 5.9915$$

i.e. ~~H0~~ H0 (independence) is rejected on a 5% level.

How to find the degrees of freedom. We have nc cells

and we are estimating $n-1 + c-1$ parameters

$$\text{Thereby } DF = nc - 1 - (n-1) - (c-1) = nc - n - (c-1)$$

$$= n(c-1) - (c-1) = \underline{(n-1)(c-1)}$$