

SDS → Unit 4

- Chi-Square test →
- Probabilities of  $k$  items are denoted by:  $p_1, p_2, p_3, \dots, p_k$ .

→ Null hypothesis is of the form

$$H_0: p_1 = p_{01}, p_2 = p_{02}, \dots, p_k = p_{0k}$$

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

→ When expected values are all suff. large  
good approx. possible

$(k-1)$  degrees of freedom denoted by:

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

When each outcome has same probability  
each expected value is  $\frac{\text{Observed value}}{\text{Total no. of values}}$

→ For row & column data:

$$E_{ij} = R_{ij} \times C_{ij} \rightarrow \begin{matrix} \text{Column} \\ \text{sum} \end{matrix}$$

$$\begin{matrix} \downarrow \\ \text{Row} \\ \text{sum} \end{matrix} \quad \begin{matrix} \downarrow \\ \text{Row sum + Column} \\ \text{sum} \end{matrix}$$

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

$$\text{Degree of freedom } f = (k-1) \times (l-1)$$

## Non-parametric tests

- Distribution-free tests →
- Samples not reqd. to come from any specific dist.
- Need assumptions but they are not as restrictive as t-test.
- Wilcoxon Signed-Rank test →
  - Can be used even with outliers

slide 3f:

9.3, 9.0, 0.9, 21.7, 11.5, 13.9

$H_0: M \geq 12, H_1: M < 12.$

X	X - 12	Rank	p-value:
9.3	-2.7	-3	
9.0	-3	-4	$S^+ = 7$
0.9	-11.1	-6	
21.7	9.7	5	
11.5	-0.5	-1	$P(S^+ < 7)$
13.9	1.9	2	

$$P(S^+ < 4) = 0.1094, P(S^+ \geq 7) > 0.1094$$

can't rej  $H_0$

when  $n > 20 \rightarrow$  normal disty

$$Z = S^+ - \frac{n(n+1)}{4} \sqrt{\frac{(n+1)(2n+1)}{24}}$$

→ Ties → replace with avg. of their pos.

$$3, 4, 4, 5, 7 \rightarrow 1, 2.5, 2.5, 4, 5$$

→ 0 → dropped from sample & n reduced by 1

# Mann-Whitney Test

- Wilcoxon Rank-Sum Test →
- 2 assumptions →
  - First → pop, continuous
  - Second → prob, density func must be qd & shape & size
- 2 samples from 2 diff pop, same shape

Slide 45

Value	Pop	Rank	H <sub>0</sub> : M <sub>X</sub> ≥ M <sub>Y</sub>	H <sub>1</sub> : M <sub>X</sub> < M <sub>Y</sub>
20	X	1		
28	X	2		
29	X	3		
34	Y	4		
35	Y	5		P = 0.0260
36	X	6		0.0
38	X	7		0.0
41	Y	8		reject H <sub>0</sub>
46	Y	9		
47	Y	10		
49	Y	11		

→ If m, n ≥ 8 normal dist

$$z = \frac{w - \frac{m(m+n+1)}{2}}{\sqrt{\frac{mn(m+n+1)}{12}}}$$

→ Fixed-level testing →

→ measures plausibility of null hypothesis by prod. P-value.

↓ P-value → ↓ plausible  $H_0$ .

Choose  $\alpha$  such that  $0 < \alpha < 1$ .

If  $P \leq \alpha \rightarrow$  reject, else do not reject  $H_0$ .

Slide 59

5% level,  $\alpha = 0.05 \rightarrow z = -1.645$

Inv Normal

$$Z = \frac{\bar{X} - M}{S/\sqrt{n}}$$

$$\bar{X} = \frac{zs + M}{\sqrt{n}}$$

$$\bar{X} = -1.645 \times 0.6 + 2$$

$$\sqrt{80}$$

$$\bar{X} = 1.89$$

$H_0$  rej if  $\bar{X} \leq 1.89 \rightarrow$  rej. region.

10% level,  $\alpha = 0.1, z = -1.28$

$$\bar{X} = -1.28 \times 0.6 + 2$$

$$\bar{X} = 1.9141$$

$$1.85 < 1.9141$$

$H_0$  rej at 10% level

→ Errors →

→ Type I error: Reject  $H_0$  when true

→ Type II error: Fail to reject  $H_0$  when false

Prob. of type I error never greater than  $\alpha$ .

→ Power of test → prob. of rejecting  $H_0$  when  $q_f$  is false.

$$\rightarrow \text{Power} = 1 - P(\text{Type II error}) = 1 - \beta$$

→ Factors affecting power →

→ Sample size ( $n$ ),  $n \propto P$

Power of → Significance level ( $\alpha$ ),  $\alpha \propto P$

1 tailed test → Type of stat. test

2 tailed test → Effect size (True mean - Hypothesized mean)

greater effect size → Standard deviation ( $\sigma$ ),  $\sigma \propto \frac{1}{P}$

greater power  $\mu$  close to  $H_0$  → ↑ power  
 $\mu$  far from  $H_0$  → ↑ power

Reject  $H_0$   $H_0$  true  $H_0$  false ↑  
Fail to reject  $H_0$  Type I error Correct decision  
Correct decision Type II error

→ In general, power  $> 0.80$  or  $> 0.90$  considered acceptable

→ Computing power →

→ Compute rej, reg,

→ Compute prob. test stat falls in rej reg if null hypothesis true. Page No.

- Types of data
  - Univariate
  - Bivariate
  - Multivariate
- Univariate → analysis using mean, median, mode, spread of data, histogram, piechart, etc.
- Bivariate → analysis using correlation coeff., regression analysis, scatter plot.

~~so~~ find relation b/w two sets of values. ( $x, y$ )  
 find  $\downarrow$  dep/vars  
 ordered pair

→ Scatter plot → math diagram that plots pairs of data on  $x-y$  graph to reveal relationship b/w data sets. Shows strength of relationship, direction of relation b/w vars & whether outliers exist.

~~corr coeff~~ → only used when relationship b/w vars → linear

Correlation → strength of relation b/w two linear vars,  $r \in [-1, 1]$

→ Pearson Correlation coeff.

$$r = -0.5 \text{ to } 0.5 \quad r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

$$r = \pm 0.8 \text{ to } \pm 0.5 \quad r = \frac{n}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$r = \pm 0.8 \text{ to } \pm 1 \quad \text{moderate}$$

↓  
correlations → +ve, -ve, no/poor

Strong

$r = \pm 1$  → perfect positive/negative corr

$0 < r < 1$  → +ve corr

$0 > r > -1$  → -ve corr

$r \geq 0$  → no corr

- Anscombe's Quartet → diff. nos, but  $\bar{x}$ ,  
Summary stats. (EX: mean, variance, etc.)
- Conclusions of Anscombe's Quartet →
  - Presence of outlier impacts evaluated  
value of  $r$ .
  - Basic stat, prop, inadeq, for desc  
realistic datasets.
  - Imp to look at data visually before  
any kind of analysis.
- confounding var, → influences both indy  
& dep. var, causing spurious correlation.  
Cause two major prob →
  - ↑ variance in raw plot → allows for  
introduction bias
  - they act as extra indy var that have  
hidden effect on dep. var.
  - can be cause of correlation.
- Controlled exp, reduce the risk of confounding  
study repr, nos, of times & variety of  
cond, before drawing reliable conclusions.
- Regression analysis → Study of set of data  
to make best guess or some sort of pred.
- In stat, modelling, OLS used for est  
relation b/w dep. var, & 1 more indy var.

Regression Models →

- Simple → Linear → Non-linear
- Multiple → Linear → Non-linear

Page No.

When two var,  $\rightarrow$  linear relationship, scatter plot clustered along straight line  $\rightarrow$  least fit line.

least square line:  $y_i = \beta_0 + \beta_1 x_i$   $\beta_0 \rightarrow y$  intercept

where  $\hat{\beta}_1 = \frac{n}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$   $\hat{\beta}_1 \rightarrow$  slope of least square

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

least square line

Residual  $\rightarrow y_{\text{observed}} - y_{\text{predicted}}$

Least square line  $\rightarrow$  sum of squared residuals  $\rightarrow$  min  $\sum_{i=1}^n e_i^2 \rightarrow$  minimum

Goodness of fit  $\rightarrow$  how well model explains given set of data

$$\text{Regression Sum of Squares} \leftarrow \sum_{i=1}^n (y_i - \bar{y})^2 \leftarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Total sum of squares  $\rightarrow$  error sum of squares

$$r^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

coefficient of determination  $\rightarrow$  stat measure of how close obs. data is to fitted regression line!

$$r^2 = \frac{\text{Regression Sum of Squares}}{\text{Total sum of squares}}$$

Value  $r^2 \in [0, 1]$ .  $\uparrow r^2 \rightarrow$  better prediction

- If the vertical spread doesn't vary with fitted value → homoscedastic else, heteroscedastic
- Outlier that makes significant diff to least square line when removed → influential point.

### Comments

- weighted least squares regression → points with  $\neq$  variance → greater influence on computation of least squares line.

→ outlier groups from all groups based on  $\text{new } \hat{y}_i - \text{old } \hat{y}_i$

→  $\hat{y}_i$  is mean of  $y_i$  in group  $j$  →  $\hat{y}_i = \bar{y}_j$

$$\sum ((y_i - \hat{y}_i)^2)_{\text{group } j} + \sum ((\bar{y}_j - \hat{y}_j)^2)_{\text{group } j} \rightarrow \text{min}$$

→ forward plot → outliers points for  $\hat{y}_i$

→ and regression better at  $\hat{y}_i$  when  $y_i$  was

→ group  $j$  and removed =  $\hat{y}_i$

→ group  $j$  and left =  $\hat{y}_i$

→ third with  $y_i$  =  $\hat{y}_i$