

## Mid 2

### Chi-sq test of Goodness of fit

→ Observed counts are different from expected counts

→ How well the <sup>obs</sup> data fit the expected dist

$H_0$ : There is no inconsistency b/w obs and expect counts

$H_A$ : There is inconsistency b/w obs and expect counts

General formula of test statistic is

$\frac{\text{point estimate} - \text{null value}}{\text{SE of point estimate}}$

Chi-square ( $\chi^2$ ) statistic

$$\chi^2 = \sum_{i=1}^{k} \frac{(O-E)^2}{E} \quad \text{where } k = \text{total no. of cells}$$

Why

→ Any standardised diff that is squared will now be positive

→ Diff that already looked unusual will become much larger after being squared

→ To determine  $\chi^2$  statistic is considered high or not we must describe its distribution

→ Chi-sq has one param called **degree of freedom (df)**, which influences the shape, center and spread of the dist

$$df = k-1 \quad : k = \text{no. of cells}$$

p-value < 0.05  $\rightarrow$  reject  $H_0$

p-value > 0.05  $\rightarrow$  ~~accept  $H_0$~~  do not reject  $H_0$

### Condition for chi-sq

#### → Independence:

Each case that contributes a count to the table must be independent of all other cases in the table.

#### → Sample size:

Each particular scenario must have at least 5 expected cases

#### → df > 1

Degree of freedom must be greater than 1

# Example

Candidate Obs Reported

Ahmedinajad	338	63.29%
Mesari	136	34.10%
Minor Candidate	30	2.61%
Total	504	100%

H<sub>0</sub>: The obs counts from the poll follow the same dist as reported vote

H<sub>A</sub>: The obs counts from the poll do not follow the same dist as reported vote

Candidate Obs Exp

i) Ahmedinajad	338	(0.6329 × 504) = 319
ii) Mesari	136	(0.3410 × 504) = 172
iii) Minor candidate	30	(0.261 × 504) = 13

Total 504 504

$$O - E \quad (O - E)^2 / E$$

338	319	31	1.13
136	172	1296	7.83
30	13	299	22.23

$$\sum (O - E)^2 / E = 30.89$$

$$\chi^2_{df=3} = 2$$

$$P\text{-value} = P(\chi^2_{df=2} > 30.89)$$

$$P\text{-value} = 0$$

$< 0.05$  therefore reject  $H_0$ .

## Chi-sq test of Independence

Class	Grades	Popular	Sports
4	63	31	25
5	88	55	33
6	96	55	32

$H_0$ : Goals do not vary by grade

$H_A$ : Goals vary by grade

$$\chi^2_{df} = \sum_{i=1}^k (O - E)^2 / E$$

$$\therefore df = (R-1)(C-1)$$

Class	Grades	Popular	Sports	Total
4	63 61	31 35	25 23	119
5	88 91	55 52	33 33	176
6	96 95	55 54	32 34	183
Total	247	141	90	52 478



$$\text{Expect Count} = \frac{(\text{Row total})(\text{col. Total})}{\text{table total}}$$

Obs	Exp	$(O-E)^2$	$(O-E)^2 / E$
-----	-----	-----------	---------------

63	61	4	0.07
31	35	16	0.48
25	23	9	0.17
88	91	9	0.10
55	52	9	0.17
33	33	0	0
96	95	4	0.04
55	54	1	0.02
32	34	4	0.12

$$\sum = 1.1688$$

$$df = (3-1)(3-1) \\ = 2 \times 2 \\ = 4$$

$$P\text{-value} = \chi^2_{df=4} > 1.16$$

$$P\text{-value} = 0.89 > 0.05$$

therefore we cannot reject  $P\text{-value}$

## One sample mean with $t$ -distribution

### Condition

→ **Independence**

rows are independent

→ **Sample size / Skew**

sample distribution

must be skewed

⇒ we do not know  $\sigma$  and  $n$  is too small  
to assume  $s$  is reliable estimate for  $\sigma$

⇒ As long as dist. are independent and  
the population dist is not extremely  
skewed then large sample could  
ensure

→ Sample dist of the mean is nearly normal

→ The estimate of  $SE$ , as  $s/\sqrt{n}$ , is reliable

## When to use

⇒ When population standard deviation is unknown

⇒ Dist has a bell shaped, but its tails are thicker than the normal model

### Note

Null value is zero because is the null hypothesis we set  $M_{null} = 0$

$$t_{df} = \frac{\text{point estimate} - \text{null value}}{SE}$$

18358

$$\text{point estimate} = \bar{X} = \text{mean}$$

$$SE = \frac{\text{S}_{df}}{\sqrt{n}} - \text{stand deviation}$$

$$df = n-1 : n \text{ is sample size}$$

## Example

	X	X - $\bar{X}$	$(X - \bar{X})^2$
1	694		311875600
2	1104		297700516
3	1037		
4	1889		
5	1911		
6	2416		
7	2761		
8	4382		
9	1839		
10	321		
	$\Sigma x = 18358$		

Ans 2

$$\text{variance} = \text{Ans 1} / 10$$

$$\text{Stand Deviation} = \sqrt{\text{variance}}$$

## Z test statistics

$\Rightarrow$  Sample size of at least 30, one can use the Z-distribution in place of t-dist

# Pair t-test or Dependent sample t-test

## When to use

⇒ When your data values are paired measurements. For example you may have before and after measurements

## Conditions

⇒ Subjects must be independent. Measurement of one subject doesn't affect the other

⇒ Each of the paired measurements must be obtained from the same subject. for example before and after weights of a smasher should be of the same person

⇒ The measured differences are normally distributed

$$\text{or } \frac{\cancel{\sum d^2 - \frac{(d)^2}{n}}}{(n-1)(n)} \quad \text{sample size}$$

$$SE = s_d / t_h$$

$$df = n-1$$

$t = \frac{\text{Average Diff}}{SE}$

$$H_0 = \text{Md}_{\text{diff}} = 0$$

Example  $\rightarrow H_A = \text{Md}_{\text{diff}} \neq 0$

Student	$U_1$	$U_2$	diff $X = U_1 - U_2$	ER
Bob	63	69	6	36
Ning	65	65	0	0
Tim	56	62	6	36
"	100	91	-9	81
"	88	78	-10	100
"	83	87	4	16
"	77	79	2	4
"	92	88	-4	16
"	90	85	-5	25
"	84	92	8	64
"	68	69	1	1
"	74	81	7	49
"	87	84	-3	9
"	64	74	11	121
"	71	84	13	169
"	88	82	-6	36
			$\sum X = 21$	$\sum TQ = 703$

$$\bar{X} = 21 / 16 \\ = 1.31$$

$n=16$  (sample size)

$$SE = \sqrt{\frac{\sum (X_i - \bar{X})^2}{n-1}} = \sqrt{\frac{763 - 21^2}{16}} \\ = \sqrt{\frac{115}{16}} = 1.75$$

$$t = \frac{\bar{x}}{SE} = \frac{1.31}{1.75} = 0.750$$

$$df = n-1 = 16-1 = 15$$

$$p\text{-value} = T_{df=15} > 0.750$$

$$p\text{-value} 2.131$$

$0.750 < 2.131$  we cannot reject  $H_0$

## Difference in two means

$$T_{df} = \frac{\text{point estimate} - \text{Null value}}{\text{SE}}$$

$$\text{point estimate} = \bar{x}_1 - \bar{x}_2$$

$$SE = \sqrt{\frac{\sum x_1^2 - (\sum x_1)^2}{n_1} + \frac{\sum x_2^2 - (\sum x_2)^2}{n_2}}$$

$$\frac{(n_1-1) / (n_1)}{(n_2-1) / (n_2)}$$

OR

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$df = \min(n_1-1, n_2-1)$$

$\Rightarrow$  To compare means of 2 groups  
we use a Z or T test

$\Rightarrow$  To compare means of 3+ groups  
we use ANOVA or F test

## ANOVA

### Conditions

$\Rightarrow$  data are a simple random sample from less than 10% of the population

$\Rightarrow$  Data should be independent

$\Rightarrow$  Observation within each group should be nearly normal

$$F = \frac{\text{variability bet. groups}}{\text{variability within group}}$$

$df_g = k-1$  where  $k$  is the no. of groups  
 $df_T = n-1$ , where  $n$  is the total sample size

$$\text{error} = df_T - df_g$$

$$SSG = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2$$

$\bar{x}_i$  = avg of each group

$\bar{x}$  = overall mean

$$SST = \sum_{i=1}^k (\bar{x}_i - \bar{x})^2$$

$$SSE = SST - SSG$$

$$MSG = \frac{SSG}{df_g}$$

$$MSE = \frac{SSE}{df_e}$$

$$F = \frac{MSG}{MSE}$$