

# ONE WAY ANOVA :

One X (categorical) two (more independent levels) - Y numeric

Ex. X 60 individuals - diets A, B, C, D

	sample mean	sample S.D
A	9.18	2.29
B	8.91	2.78
C	12.11	1.79
D	10.54	2.23

$$H_0: \mu_A = \mu_B = \mu_C = \mu_D$$

$$H_1: \mu_p \neq \mu_q \text{ for } p \neq q$$

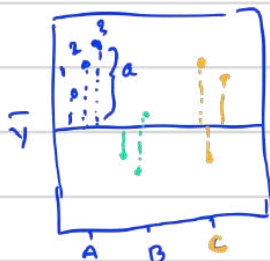
at least one differs, not which one or one or more differ

Assumes - normality

Other possibilities:

① non-parametric alternative: KRUSKAL WALLIS  
(doesn't require normality)

② Bootstrap / Resample



$$a = (\text{individual} - \bar{Y})$$

Sum of squares total

$$SS_{\text{total}} = \sum_{\text{all}} (\text{individual} - \bar{Y})^2 \text{ - total variability in } Y$$

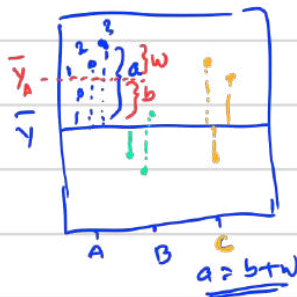
$$s^2_{\text{total}} = \frac{SS_{\text{total}}}{n-1} \text{ | sum of squares}$$

Why d's are different?

- ① Different diet - Explained by diet - called between (b)
- ② People are different - Unexplained by X (diet) (w)  
- Also called within

$$SS_{\text{total}} = SS_{\text{Explained}} + SS_{\text{unexplained}}$$

$$= SS_b + SS_w$$



$$SS_b \text{ or } SS_{\text{Explained}} = \sum_{\text{all}} (\text{group } \bar{Y} - \text{Overall } \bar{Y})^2$$

$$s^2_b = \text{Mean Square bet' groups (MS}_b) = \frac{SS_b}{k-1} \text{ k: number of groups}$$

Signal

$$SS_w \text{ or } SS_{\text{unexplained}} \text{ - also called SS error or SS residual}$$

$$= \sum_{\text{all}} (\text{individual} - \text{Group } \bar{Y})^2$$

Noise

$$s^2_w = MS_w = \frac{SS_w}{n-k}$$

On to test statistic for ANOVA:

	sample mean	sample s.d
A	9.18	2.29
B	8.91	2.78
C	12.11	1.79
D	10.54	2.23

$i = \text{groups (A, B, ... K)}$

$\bar{y}_i = \text{mean of group } i$

$s_i = \text{SD of group } i$

$n_i = \text{sample size of } i$

$\bar{y} = \text{overall mean}$

$j = \text{individual within group}$

$y_{ij} = \text{individual observation in group } i \text{ obs. } j$

$y_{13} = \text{person 3 in 1st group (i.e., A)}$

$$\text{Explained } S_B^2 (MS_B) = \frac{SS_B}{df_B} = \frac{\sum_{\text{group}} n_i (\bar{y}_i - \bar{y})^2}{k-1} = \frac{97.3}{3} = 32.4$$

$$\text{Unexpl.} \rightarrow S_W^2 (MS_W) = \frac{SS_W}{df_W} = \frac{\sum_{\text{all obs}} (y_{ij} - \bar{y}_i)^2}{n-k}$$

$$= \frac{\sum_{\text{groups}} (n_i - 1) s_i^2}{n-k}$$

like pool variance in t-test

$$= \frac{297}{56} = 5.3$$

If  $H_1$  true, we expect  $MS_B > MS_W$

$$F_{\text{statistic}} = \frac{MS_B}{MS_W} > 1$$

$$F_{\text{stat}} = \frac{32.4}{5.3} = 6.1$$

If  $H_0$  true  $\frac{MS_B}{MS_W} \approx 1$

[F dist. with 2 df.  $\rightarrow df_n = k-1$   
 $\rightarrow df_d = n-k$ ]

from: F dist.  $\rightarrow$  p-value = prob. ( $F_{\text{stat}} \geq 6.1$  if should  $\approx 1$ )  
 $= 0.0011$

$\rightarrow$  Reject  $H_0$

$H_1$  just says that diet has a difference

"but not which one"  $\rightarrow$  need multiple comparison

ANOVA multiple comparison: Multiple Testing Correction

$H_0: \mu_A = \mu_B = \mu_C = \mu_D$   $H_1$ : At least one  $\mu$  differs

$$F_{\text{STAT}} = \frac{S_B^2}{S_W^2} = \frac{32.4}{5.3} = 6.1 \rightarrow p = 0.0011 \Rightarrow \text{Reject } H_0$$

All pairwise comparison  $\binom{4}{2} = 6$  possible  $\rightarrow$  independent 2 sample t-test

$$(\bar{y}_1 - \bar{y}_2) \pm t_{df}^{\text{crit}} \times SE(\bar{y}_1 - \bar{y}_2)$$

$$\text{or } t_{\text{STAT}} = \frac{\bar{y}_1 - \bar{y}_2}{SE(\bar{y}_1 - \bar{y}_2)}$$

When we do more test,  $\rightarrow$  probability of Type I error increases  
 One test  $\rightarrow \alpha = 0.5\%$ , another test another  $0.5\%$  error

Eg. AB, AC, AD, BC, BD, CD - assume each comparison independent

each test: we  $\alpha = 0.05 \rightarrow$  prob. of making Type I error

$$P(\text{Type I}) = 0.05,$$

$$\Rightarrow P(\text{no Type I}) = 0.95$$

$$\text{Overall } P(\text{At least 1 Type I}) = 1 - P(\text{making no Type I error})$$

(also called

Family wise Error Rate)

$$= 1 - (0.95)^6$$

$$= \underline{\underline{0.265}}$$

Bonferroni approach:

$$\text{use adjusted } \alpha^* = \frac{0.05}{\# \text{ comparison}} = \frac{0.05}{6} = 0.00833$$

$$\text{for C.S.} = 95\% \Rightarrow 99.167\%$$

$$\text{Now, Overall } P(\text{at least 1 Type I error}) = 1 - (0.99167)^6$$

$$= 0.049$$

Pearson's' Chi-square test of Independence:

independence test X & Y - both categorical, 2 / more levels.

Non parametric, but relies on chi-square dist<sup>2</sup>

like other parametric

X - vaccinated for MMR

Y - diagnosed for Autism

Any relationship?  $\chi^2$  - just tells if / not there  
 is an assoc., but not how STRONG?

	Y		
	Yes	No	
X	Yes	440,034	440,655
	No	96,51	96,648
	238	536,565	539,303

Call this observed table "O"

other way to  
 do  $\rightarrow$  Two proportion  
 test

$$H_0: X, Y \text{ no association / Independent} \Rightarrow P_1 = P_2 \Rightarrow P_1 - P_2 = 0$$

$$H_1: X, Y \text{ associated} \Rightarrow P_1 \neq P_2 \Rightarrow P_1 - P_2 \neq 0$$

$$\hat{P}_1 = P(\text{AUT} | \text{Vacc}) = \frac{621}{440,655} = 0.00141$$

$$\hat{P}_2 = P(\text{AUT} | \text{no Vacc.}) = \frac{117}{96,648} = 0.00121$$

Expected table "E"  
fill the values by taking the values in blue from "O" table

	Y	y	N
X	Y	605.25	440,655
	N	96,648	537,303
		738	536,565

These are not free  $\therefore$  1 d.f  
(only 605.25 is free)

Now we know one value, rest can be filled

$$\begin{aligned} O &= 238 - 605.25 \\ P &= 440,655 - 605.25 \\ \star &= 96,648 - O \end{aligned}$$

Assumption of  $\chi^2$  test:

- Groups, Obs<sup>2</sup> independent
- All cells  $\geq 1$
- All cells in "E"  $\geq 5$

% if not met then

Actually this is simply  
row total  $\times$  column total  
overall total

one can

→ Fisher Exact Test

→ Bootstrapping

Test - statistic: compare "O" to "E" table

$$\chi^2_{STAT} = \sum_{all\ cells} \frac{(O - E)^2}{E} = \frac{(621 - 605.25)^2}{605.25} + \dots + \frac{(96531 - 96515.25)^2}{96515.25}$$

$$= 2.28$$

$\sim \chi^2$  dist<sup>n</sup> (if  $H_0$ : TRUE) d.f = (# rows - 1)(# cols - 1)

So, here p-value =  $P(\chi^2_{STAT} \geq 2.28 \text{ if } H_0 \text{ TRUE}) = 0.1309$

Accept  $H_0$

NOTE:

- For large sample sizes, small differences may show up to be statistically significant
- For small sample sizes, big " " " NOT " " " " " "

CONTINUITY CORRECTION: - Above example, had numbers of persons  
- discrete case

- In a discrete case  $P(X \geq 10)$  is same as  $P(X \geq 11)$

- In a continuous case  $10 \neq 10.5 \neq 10.3$ , suggestion is use a midpoint

$$P(X \geq 10) = P(X \geq 10.5)$$

\* To know how strong is assoc.<sup>n</sup>  $\left\{ \begin{array}{l} \text{Risk RATIO} \\ \text{Risk DIFFERENCE} \\ \text{ODDS RATIO} \end{array} \right.$