

ANOVA

Analysis of Variance

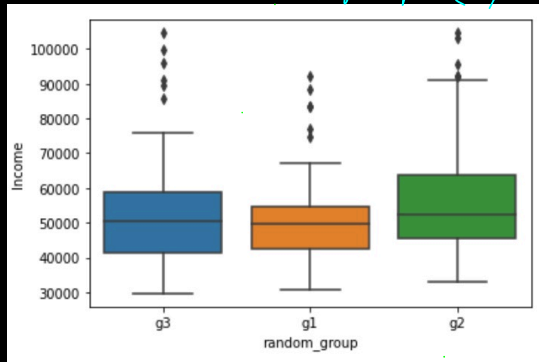
	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
0	KP281	18	Male	14	Single	3	4	29562	112
1	KP281	19	Male	15	Single	2	3	31836	75
2	KP281	19	Female	14	Partnered	4	3	30699	66
3	KP281	19	Male	12	Single	3	3	32973	85
4	KP281	20	Male	13	Partnered	4	2	35247	47

Income Vs Gender → T-test  
(Num) (Categorical - 2 categories)

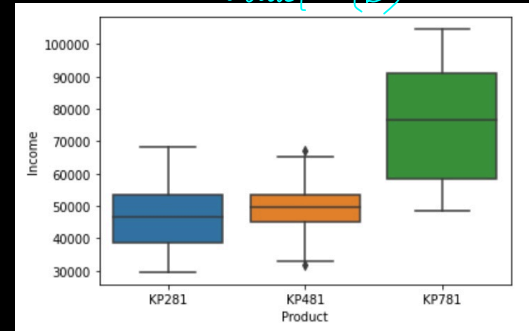
Gender Vs Product → Chi-Squared Test  
(Cat) (Cat)

Income Vs Product → ANOVA  
(Num) (Categorical 2 or more categories)

Random group (A)



Product (B)



$H_0$ : all means are equal  $\rightarrow$  high  $p$ -value

$H_a$ : Atleast one mean is different from others  $\rightarrow$  low  $p$ -value

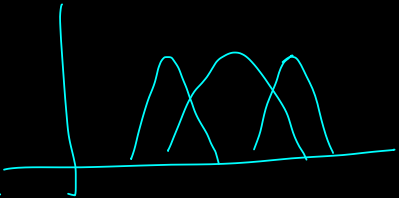
ANOVA : deep dive Height

Setup 1

American Basketball players  
Indonesian college students  
Indian cricket team

① variance within group

② variance b/w group



Setup 2

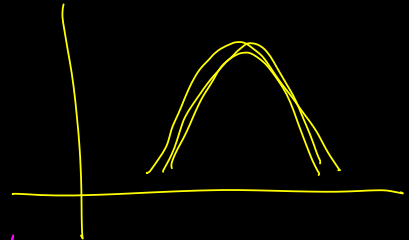
group people alphabetically

A to G

H to N

O to Z

$$f\text{-ratio} = \frac{\text{Variance between groups}}{\text{Variance within groups}}$$



f-ratio will be high in setup 1  $\rightarrow H_a \rightarrow$  right tailed  
will be low in setup 2  $\rightarrow H_o$

## Assumptions of ANOVA:

① Data should be Gaussian  $\nearrow$  QQ plot  
 $\rightarrow$  Shapiro

② Independence

③ Equal variance in each group  $\rightarrow$  Levene

When assumptions of ANOVA don't hold  
we use Kruskal-Wallis

QQ Plot: Test if data is Gaussian

$x_1, x_2, x_3, \dots, x_{10000}$

Empirical rule?  $68/95/99$   
 $1\sigma/2\sigma/3\sigma$

$y_1, y_2, y_3, \dots, y_{10000}$   $\rightarrow$  actually Gaussian

1<sup>st</sup> percentile of  $x \approx$  1<sup>st</sup> percentile of  $y$

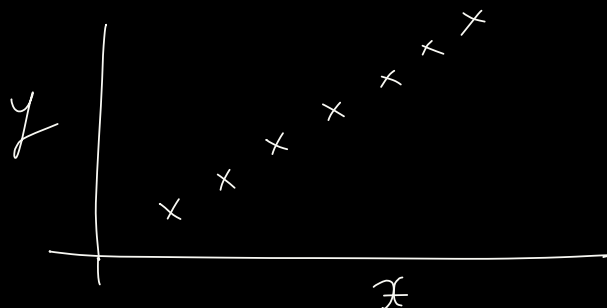
2<sup>nd</sup> percentile of  $x \approx$  2<sup>nd</sup> per of  $y$

$\vdots$

100<sup>th</sup>

$\approx$

100<sup>th</sup> per of  $y$



Shapiro: Test whether data is Gaussian or not

$H_0$ : Data is Gaussian

## Post Read (Optional)

iPhone sales in 3 stores				
A	B	C		
25	30	18		
25	30	30		
27	25	29		
30	24	29		
23	26	24		
20	28	26		
25 $\bar{Y}_1$	26.5 $\bar{Y}_2$	26 $\bar{Y}_3$	25.83 $\bar{Y}$	

$H_0$ : All means are equal       $H_a$ : Means are different  
 Step 1 Compute individual group means     $\bar{Y}_1 = 25$      $\bar{Y}_2 = 26.5$      $\bar{Y}_3 = 26.5$   
 Step 2 Compute mean of these 3 values     $\bar{Y} = \frac{25 + 26.5 + 26}{3} = 25.83$   
 Step 3 Between groups  
 $SSB = 6(25 - 25.83)^2 + 6(26.5 - 25.83)^2 + 6(26 - 25.83)^2 = 6.9$   
 $DF = 3 - 1 = 2$   
 $MSB = \frac{SSB}{DF} = \frac{6.9}{2} = 3.49$   
 Step 4 Within groups  
 $SSW = (25 - 25)^2 + (25 - 25)^2 + (27 - 25)^2 + \dots + (20 - 25)^2$   
 $\quad + (30 - 26.5)^2 + (30 - 26.5)^2 + (25 - 26.5)^2 + \dots + (28 - 26.5)^2$   
 $\quad + (18 - 26)^2 + (30 - 26)^2 + (29 - 26)^2 + \dots + (26 - 26)^2$   
 $\quad = 223$   
 $DF = 18 - 3 = 15$   
 $MSW = \frac{SSW}{DF} = \frac{223}{15} = 14.9$

$F = \frac{3.49}{14.9} = 0.23$   
 $F = \frac{MSB}{MSW}$