

Welcome Everyone

Agenda :-

1. ANOVA
2. EDA & Feature Engineering
3. CRISP - DM
4. DS libraries → numpy, pandas, matplotlib
5. Practical of Statistics }
6. Few projects and EDA parts

Hypotheses testing

0. Write everything given. Finalize the test
1. H_0 & H_1
2. Decision Boundary, CI, α
3. Test statistics, p-value or direct, rule
4. Decision & Conclusion

A NOVA \Rightarrow analysis of variance

Anova is a statistical test used to compare mean of 2 or more groups. & see if atleast one of them is significantly different from the other.

2 main important things to understand

1. Factor

G
M
F
M
C

2. Level

G with 2 levels

Gender

M
F
O
O
M
F

Medicine

Dolo \leftarrow factor
500 } \Rightarrow levels
650
900

Say you are given all of these dosages as a part of test group & you are asked to rate out of 10, how much was each dosage effective?

Medicine	factor	levels
500	650	900
7	9	5
8	3	6
7	9	8
8	3	5
3	1	9

{Gender} Factor / Variable
 {M, F, O} Levels / sub-categories

Assumption in Anova

1. Normality of means of sample distribution

⇒ means of sample dist. should be normally distributed

Outliers

2. Absence of Outliers ⇒ needs to be removed

3. Homogeneity of variance of samples

⇒ sample should have same variance

4. Population variance in diff levels of

each independent factor/variable are equal.

5. Samples are random & independent.

Types of ANOVA

1. One way Anova: When we have
One factor with atleast 3 levels
& all the levels are independent.

Gender \Rightarrow M F Others

Dolo \Rightarrow 650 500 900

2. Repeated Measure Anova \div We have one
factor with atleast 3 levels but
levels are dependent.

Sleep

Day 1	Day 2	Day 3
18	6	6

There will be some relationship b/w levels.

3. Factorial Anova - 2 or more factors &
each of which have at least 2 levels.

Levels can be dep. or indep.

<u>Medicine</u>			
Gravell's M	500	650	900
	3	5	9
	3	6	8
	8	7	1
F	2	6	8
	2	6	9
	3	7	10

↳ random variable

One way Anova [F-test]

- ⇒ used for inferential statistics
- ⇒ comparing mean of more than 2 groups

Hypotheses in anova

$$H_0 \Rightarrow \mu_1 = \mu_2 = \mu_3 = \dots = \mu_n$$

$H_1 \Rightarrow$ at least one is not equal

Test statistic

$$F = \frac{\text{Variance between sample}}{\text{Variance within sample}}$$

Q : Researchers want to test a new medicine . They split participants into 3 levels (dosage) & asked them to rate the effect of medication from 1-10

0	50	100	100 } before sample
9	7	4	
8	6	3	
7	6	2	
8	7	3	
8	8	4	
9	7	3	
8	6	2	

Are there any diff. b/w 3 doses
using $\alpha = 0.05$?

$$\textcircled{1} \quad H_0 \Rightarrow \mu_0 = \mu_{50} = \mu_{100}$$

$H_1 \Rightarrow$ mean are not equal

$$\textcircled{2} \quad \alpha = 0.05 \quad \text{C.E. } 95\%$$

\textcircled{3} In f test we calculate d.o.f.

$$n = 7 \quad \{\text{total data points in a sample in a level}\}$$

$$N = 21 \quad \{\text{all data points collected}\}$$

$$a = 3 \quad \{\text{no. of categories}\}$$

In ANOVA, we use sample data to make inference about population parameters.

Hy hypothesis test in ANOVA compares the population mean, but by using sample mean

$$\text{dof}_{B/W} = a-1 = 2$$

$$\text{dof}_{\text{within}} = N - a = 18$$

$$\text{dof}_{\text{total}} = N - 1 = 20$$

$$\alpha = 0.05$$

f-value $\{\text{dof}_{B/W}, \text{dof}_{\text{within}}\}$

$[2, 18]$

Degrees of Freedom	1	2	3	4	5	6
1	161.448	199.500	215.707	224.583	230.162	233.986
2	18.513	19.000	19.164	19.247	19.296	19.330
3	10.128	9.552	9.277	9.117	9.013	8.941
4	7.709	6.944	6.591	6.388	6.256	6.163
5	6.608	5.786	5.409	5.192	5.050	4.950
6	5.987	5.173	4.757	4.534	4.364	4.284
7	5.591	4.731	4.347	4.120	3.972	3.866
8	5.318	4.499	4.066	3.838	3.687	3.581
9	5.117	4.259	3.863	3.653	3.524	3.374
10	4.965	4.103	3.708	3.474	3.326	3.217
11	4.844	3.982	3.587	3.357	3.204	3.095
12	4.747	3.885	3.490	3.259	3.106	2.996
13	4.667	3.806	3.411	3.179	3.025	2.915
14	4.600	3.739	3.344	3.112	2.958	2.848
15	4.543	3.682	3.287	3.056	2.901	2.790
16	4.494	3.634	3.239	3.007	2.852	2.741
17	4.451	3.592	3.197	2.965	2.810	2.699
18	4.414	3.552	3.160	2.928	2.773	2.661
19	4.381	3.522	3.127	2.895	2.740	2.628
20	4.351	3.493	3.098	2.866	2.711	2.599
21	4.325	3.467	3.072	2.840	2.685	2.573
22	4.301	3.443	3.049	2.817	2.661	2.549
23	4.279	3.422	3.028	2.796	2.640	2.528
24	4.260	3.403	3.009	2.776	2.621	2.508
25	4.242	3.385	2.991	2.759	2.603	2.490
26	4.225	3.369	2.975	2.743	2.587	2.474
27	4.210	3.354	2.960	2.728	2.572	2.459
28	4.196	3.340	2.947	2.714	2.558	2.445
29	4.183	3.328	2.934	2.701	2.545	2.432
30	4.171	3.316	2.922	2.690	2.534	2.421

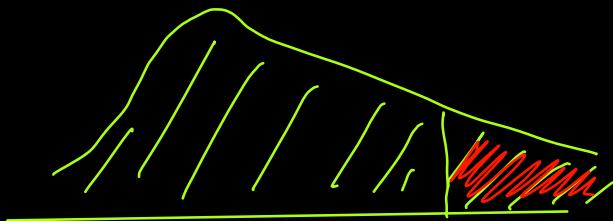
2 dof B/w

dof w/t |

18

$f\text{-value}$
= 3.555

If f-test statistic is greater than f-value
we reject null value



H0 H1 ?

Why -f-test graph
is rightly showed?

Calculate F - statistics

	S.S	D.F	Mean Square S.S/D.F	F	
Between	98.67	2	49.33①		
Within	10.29	18	0.54②		
Total	108.95	20			

57 47 21

$$SS_{\text{Between}} = \frac{\sum (\bar{x}_i)^2}{n} - \frac{T^2}{N}$$

\bar{x}_i = all data having key level

T = total summation of all values

$$\frac{57^2 + 47^2 + 21^2}{7} \Rightarrow 842.71$$

$n = 7$

T = total summation of all values

$$= (57 + 47 + 21)^2 \quad N = 21$$

$$= 125 \quad \Rightarrow 744.04$$

$$SS_{\text{Between}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{125^2}{21}$$

$$\Rightarrow 98.67$$

$$S^2_{\text{within}} = \frac{\sum y^2 - \frac{\sum (\bar{y})^2}{n}}{n}$$

$\sum y^2 = \text{all data points squared}$

$$= 853 - 842.71$$

$$= 10.29$$

$$F = \frac{\text{variance between sample}}{\text{variance within sample}} \Rightarrow \frac{\textcircled{1}}{\textcircled{2}}$$

$$\frac{49.33}{0.57} = 86.54$$

Denominator Degrees of Freedom	Critical Values of the F-Distribution $\alpha = 0.05$					
	1	2	3	4	5	6
1	161.468	199.500	215.726	224.583	230.747	236.7
2	18.513	19.000	19.164	19.264	19.298	18.330
3	10.128	9.552	9.277	9.117	9.013	8.941
4	7.709	6.944	6.591	6.388	6.256	6.163
5	6.943	6.395	6.099	5.864	5.649	5.47
6	5.987	5.473	4.737	4.534	4.284	4.15
7	5.591	5.73	4.347	4.120	3.972	3.866
8	5.318	5.57	4.200	3.960	3.811	3.71
9	5.121	5.34	3.983	3.740	3.627	3.51
10	4.965	4.103	3.708	3.455	3.217	3.07
11	4.844	3.982	3.587	3.357	3.204	3.095
12	4.747	3.842	3.489	3.236	3.096	2.959
13	4.667	3.806	3.411	3.179	3.023	2.915
14	4.600	3.739	3.344	3.112	2.958	2.848
15	4.542	3.682	3.297	3.060	2.901	2.77
16	4.494	3.634	3.239	3.007	2.832	2.741
17	4.451	3.584	3.197	2.963	2.810	2.699
18	4.414	3.555	3.160	2.928	2.773	2.661
19	4.382	3.526	3.127	2.892	2.736	2.626
20	4.351	3.493	3.098	2.866	2.711	2.599
21	4.325	3.467	3.072	2.840	2.685	2.573
22	4.301	3.443	3.049	2.817	2.661	2.549
23	4.280	3.420	3.040	2.787	2.625	2.515
24	4.260	3.403	3.009	2.776	2.621	2.508
25	4.242	3.385	2.991	2.759	2.603	2.490
26	4.226	3.369	2.975	2.743	2.587	2.474
27	4.210	3.352	2.959	2.721	2.565	2.450
28	4.196	3.340	2.947	2.714	2.558	2.445
29	4.183	3.328	2.934	2.701	2.545	2.432
30	4.171	3.316	2.922	2.690	2.534	2.421

F value from table
→ 3.55

$$86.54 > 3.55$$

⇒ reject H_0

Conclusion :-

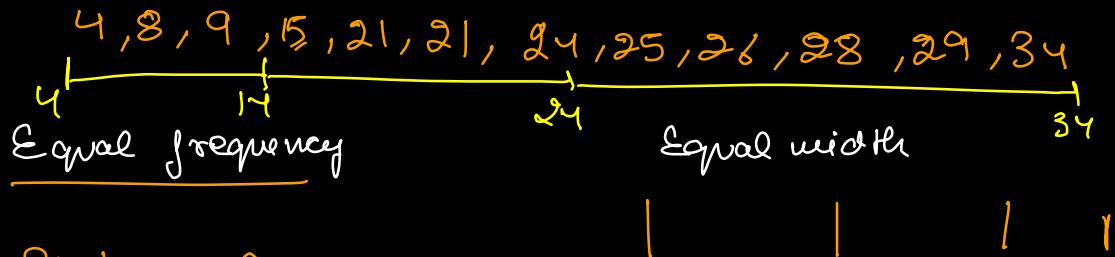
From f test value we reject
null hypothesis, so there are
difference b/w 3 condⁿ
(10, 50, 100)
with 95% interval

E DA :-

Data Scientist Plan :-

1. Load the data into the data mining environment.
2. Understand the data using graphical and non-graphical techniques
3. Document Insights
4. Pre process the data to suit your model
5. Divide into train & test
6. Quickly build the model with small data set and manually verify the error on test set for improving.

Changing numerical to categorical



$$\text{Bin 1} = 4, 8, 9, 15$$

$$\text{Bin 2} = 21, 21, 24, 25$$

$$\text{Bin 3} = 26, 28, 29, 34$$

$$\frac{34 - 4}{3} \quad \frac{(\text{max-min})}{\text{(no. of bins)}}$$

$$\Rightarrow 10 \quad \boxed{4+10} \quad \begin{array}{|c|c|} \hline 4 & 14 \\ \hline 14 & 24 \\ \hline \end{array}$$

Categorical to numerical

$$\begin{array}{ccc} M & \longrightarrow & 1 \\ F & \rightarrow & 0 \end{array}$$

Missing Value

Why missing value :-

- a) human error
- b) Refuse to answer / fill form
- c) Optional box
- d) Technical reasons

95
91
○
98

Should we ignore or impute ?

- ⇒ Understand why missing value is present
- ⇒ Plot the graphs



Delete the observ. or var.

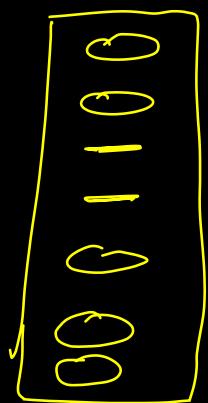
- 1) Drop obs
- 2) Drop var

1) Consider imputing values in var if
missing value < 10%

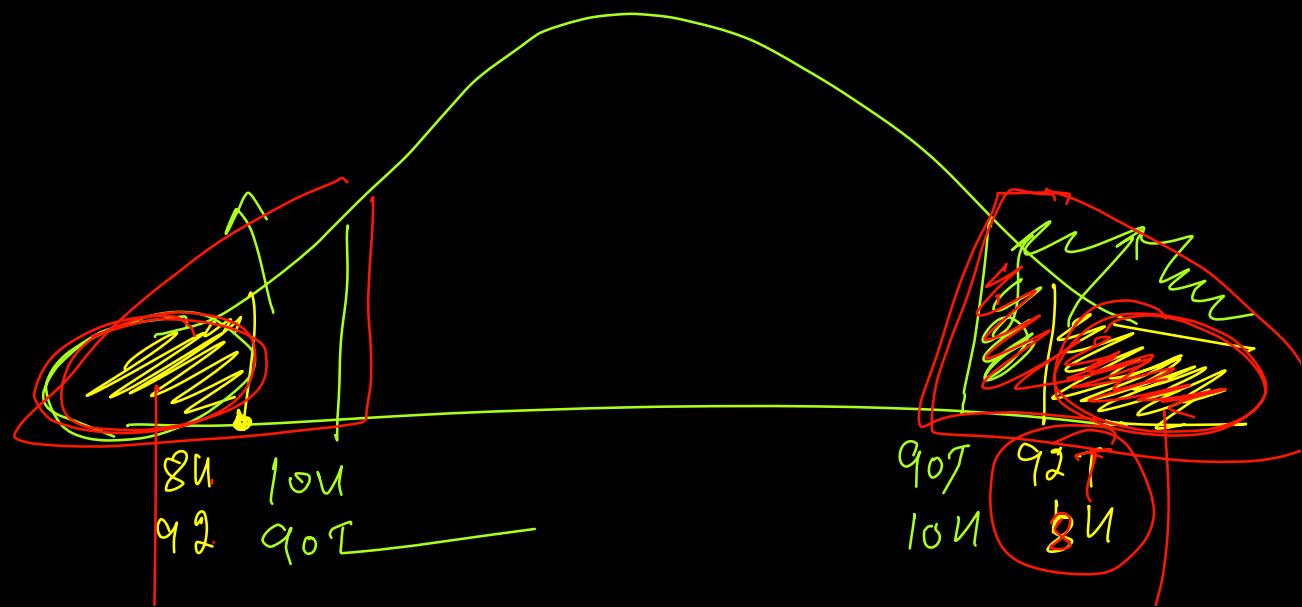
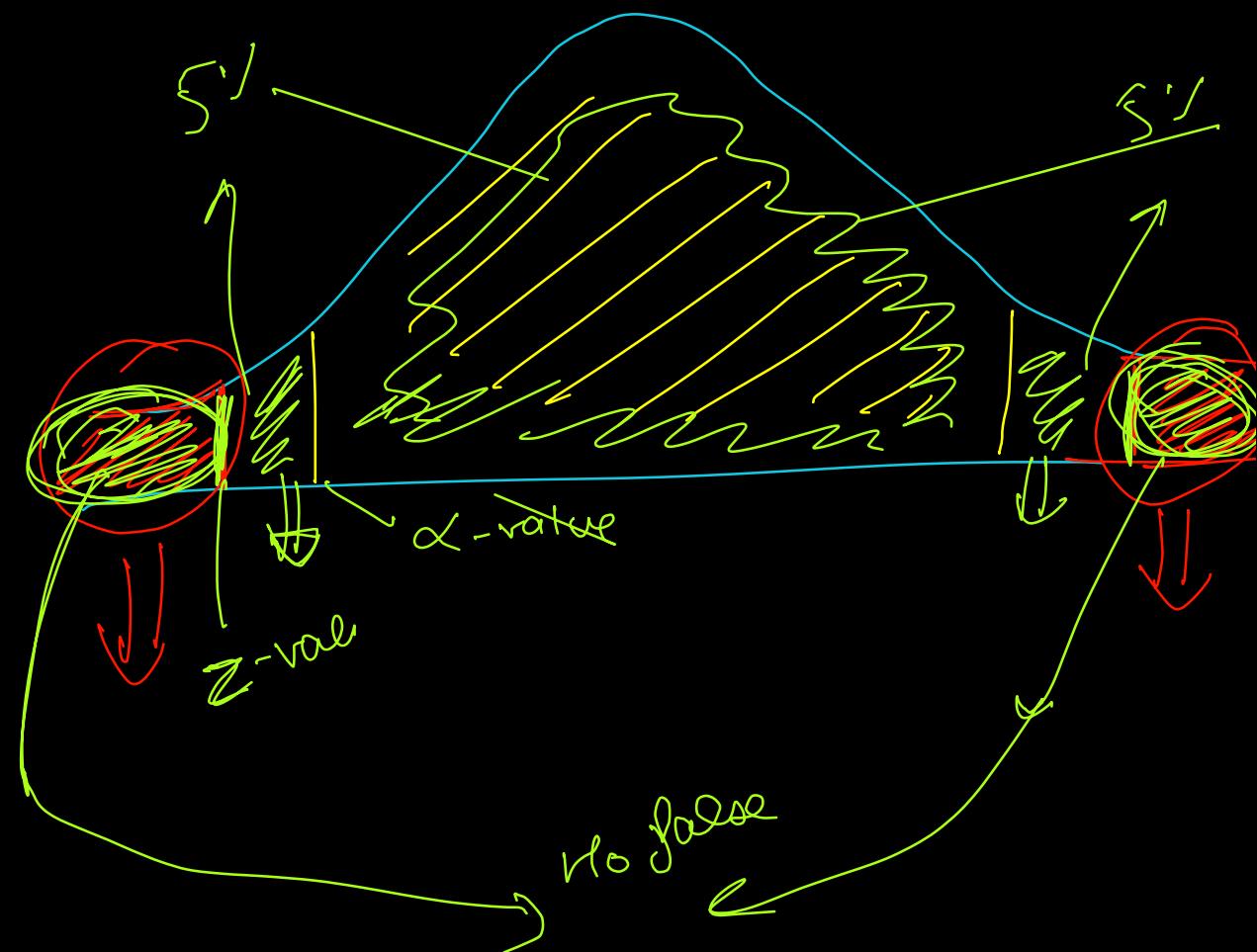
2) If more than 50% , drop var

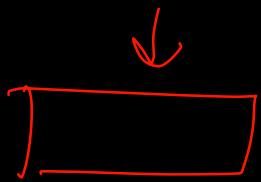
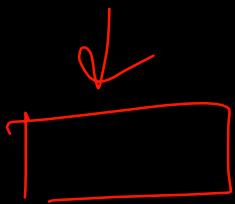
3) If we have huge data & say 10-20%
are missing , drop observation

Marks } \Rightarrow dropping the whole variable



} \Rightarrow dropping observations





41

0

100 52 79

$$\frac{100 - 0}{2} > \boxed{50}$$

$\boxed{(0, 50)}$ 0, 41, 49
 $\boxed{(50, 100)}$ 52, 100