

31

* { descriptive value for population → parameter
 " " sample → statistics
Monday

January 2011

Week 5 • 31-334

#

Day-7.

Different aspects of python packages -

- Numpy
- Pandas
- matplotlib
- scikit-learn
- Scipy.

Day-8.

Introduction to World of statistics.

10

11 Statistics → It is a collection, organization, analysis &
 12 interpretation of data.

- Design of experiments
- Sampling.
- Descriptive statistics
- Inferential statistics
- Probability theory.

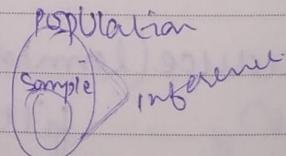
numerical way
 of analyzing
 data
 ↓
 Stats.

Why we want to learn stats.?

- Draw conclusion abt large sets.
- We can make forecasting about business activity.
- Improve business process.

describe/
 presenting
 data
 clearly,

Types of stats.



Mean,
 median,
 mode

Descriptive
 (summarization of
 data)

Inferential.

(using hypothesis, concluding
 result of whole popn.)

take small data → conclude on
 large data.

make prediction

	Mo	Tu	We	Th	Fr	Sa	Su
Wk	1	2	3	4	5	6	
5	7	8	9	10	11	12	13
6	14	15	16	17	18	19	20
7	21	22	23	24	25	26	27
8	28						

Tuesday 1

February 2011

1

Week 5 • 32-333

• Types of Variable.

Variables are predefined

↳ categorical (Predefined categories): Ex. Gender, Season, continuous. (numerical variables) → measured in numbers.

↓
Discrete
(Finite)

↓
continuous
(Infinite range)

Ex. NO. of children ↳ sales, ... Real Numbers.

Method !

Measure of central tendency.

① mean, ② median, ③ mode.

Sample mean: $\bar{x} = \frac{\sum x}{n}$

Population mean: $\bar{y} = \frac{\sum x}{N}$

obtained by deleting
a few small & large values from
a dataset & then
computing the
mean of remaining
values.

Trimmed mean
Removing extreme value and
finding mean. (outlier removal)

Geometric Mean.

When want to check growth rate in financial data

calculated by finding
nth root the product
of n values.

Weighted mean.
Find avg. wrt category wise.

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$

wi similar to
 $\left(\frac{E_f i}{E_f} \right)$

2

Wednesday

February 2011

Week 5 • 33-332

Wk	Mo	Tu	We	Th	Fr	Sa	Su
5	1	2	3	4	5	6	
6	7	8	9	10	11	12	13
7	14	15	16	17	18	19	20
8	21	22	23	24	25	26	27
9	28						

Note

arrange values in ascending order

(2) Median

To find middle values.

- Observation → even.

take two values & find mid pt.

- Observation → odd

take the middle value.

(3) Mode

most frequent (repeated) value.

10 Method 2:11 # Dispersion measures.12 It describes the data spread ~~to~~ ^{or} how far measurement are from center.

- 3 • Range → max.value - min.value.

- 4 • Mean Deviation → How much value is deviated from mean.

- 5 • Percentile → relative performance.

Ex. $i = (P/100) n$

Find 80th percentile → $(80/100) \times 70 = 56$
even no. (56) →

Avg. of 56 and 57 data value:

$$\text{80th percentile} = (535 + 549)/2 \\ = 542$$

6 Quartiles -

First Quartile → 25 percentile

Second Quartile → 50 percentile (Median).

Third Quartile → 75 percentile.

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Thursday

February 2011

Week 5 • 34

Ex. Third Quartile. \rightarrow 75 percentile.

$$i = (\frac{P}{100})n = (75/100) \cdot 40 = 30 \approx 31$$

Third Quartile: $= 525$ \rightarrow value at 33 posn of data.
(Q3)

* Inter Quartile range.

dist. b/w 1st & 3rd quartile.

Q1	Q2	Q3
25%	25%	25%

$\brace{ }_{\text{Q3} - \text{Q1}}$ (Interquartile range)

$$(3rd \text{Quar.}) Q_3 = 525$$

$$(1st \text{Quar.}) Q_1 = 445$$

$$(\text{Inter Quar.}) Q_3 - Q_1 = 525 - 445 = 80$$

* Variance.

difference b/w the value of each observation (x_i) & the mean (\bar{x} for a sample, μ for a population).

Used
to find
variation
in
data

* Standard Deviation.

Use to measure the consistency.

$$\sqrt{\frac{\sum (x - \bar{x})^2}{(n-1)}}$$

Small S.D \rightarrow More consistent

Large S.D \rightarrow Less consistent

4

Friday

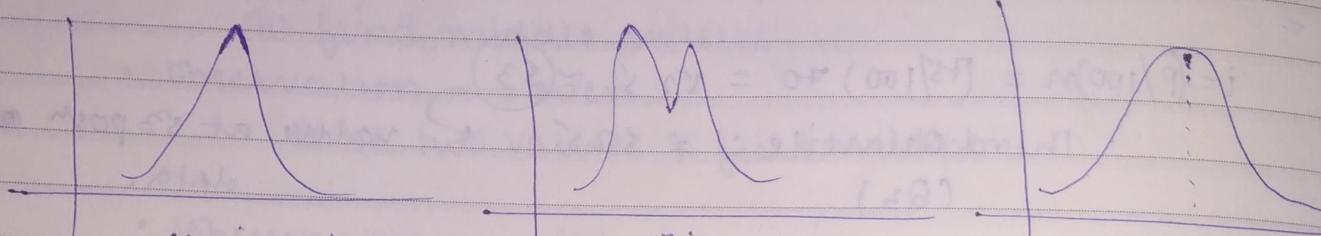
February 2011

February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
5		1	2	3	4	5
6	6	7	8	9	10	11
7		14	15	16	17	18
8		21	22	23	24	25
9		28				26 27

Week 5 • 35-330

Method 3.

Distribution of shape

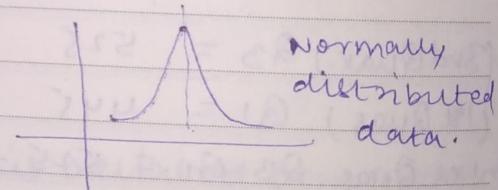
Analyze collection of Data.

(Used to measure → symmetry.)

- 9 o Skewness
- 10 o Boxplot
- 11 o Kurtosis
- 12 o Scatter Plot
- 1 o Coeff. of variance.

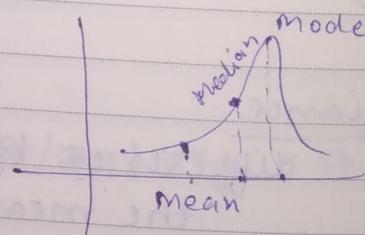
→ If skewness = 0,

Mean = median = mode



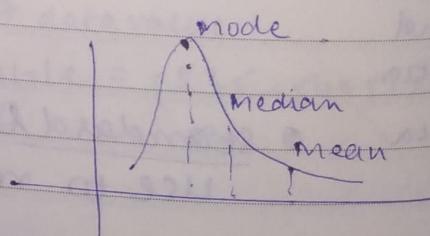
Skewness = -ve

Mean < Median



Skewness = +ve

Mean > Median



$$\text{Skewness}(S) = \frac{3(\text{Mean} - \text{Median})}{\text{Standard Deviation}}$$

all graphs
are
unimodal.

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9	1	2	3	4	5	6
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Saturday

5

February 2011

mean < Median

mean = Median

Week 5 • 36-329

mean > Median

Σx

$$\mu = 23 \quad \{ \quad \uparrow$$

$$Md = 26 \quad \}$$

$$\sigma = 12.3$$

$$S = \frac{3(\mu - Md)}{\sigma}$$

$$= -0.73$$

$$\mu = 26 \quad \{ \quad \uparrow$$

$$Md = 26 \quad \}$$

$$\sigma = 12.3$$

$$S = \frac{3(\mu - Md)}{\sigma}$$

$$= 0$$

$$\mu = 29 \quad \{ \quad \uparrow$$

$$Md = 26 \quad \}$$

$$\sigma = 12.3$$

$$S = \frac{3(\mu - Md)}{\sigma}$$

$$= +0.73$$

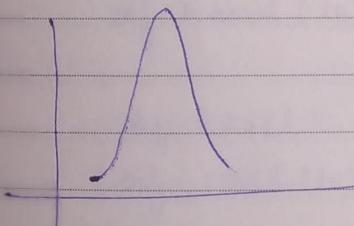
-ve skewed.

Neutral

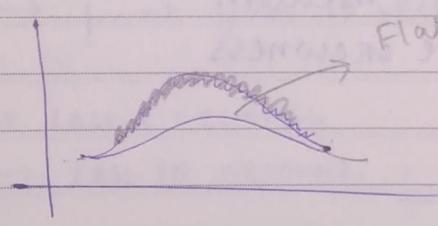
+ve skewed.

② KURTOSIS

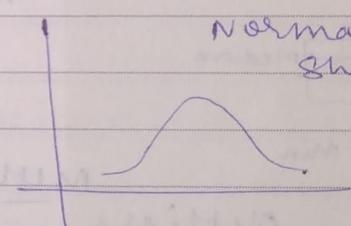
Used to measure peakness of data.



Leptokurtic
peak (High & thin)

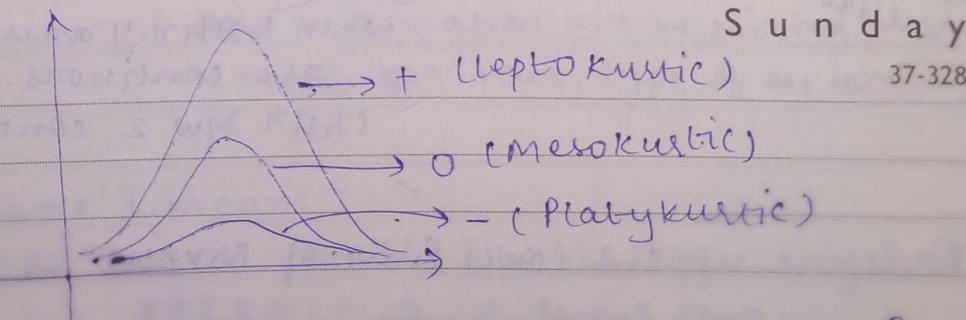


Platykurtic
(Flat & spread out)



Mesokurtic
(Bell shaped)

Measure distribution is flat or peaked.



Sunday

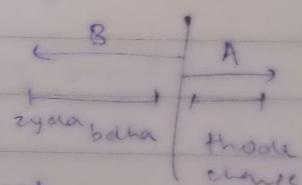
37-328

6

③ Coefficient of Variation.

When we change one variable, how much variation is there in other variable.

$$n = \left(\frac{s}{\bar{x}} \right) \cdot 100\%$$



7

Monday

February 2011

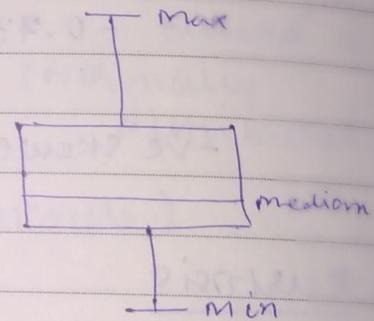
February							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su
5		1	2	3	4	5	6
6	7	8	9	10	11	12	13
7	14	15	16	17	18	19	20
8	21	22	23	24	25	26	27
9	28						

Week 6 • 38-327

① Box plot

visual
representn
of
data

- Minimum
- Q_1
- Median
- Q_3
- Maximum

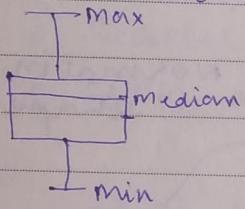


line is below centre \rightarrow +ve skewness
mean > median

line is above centre :

mean < median

-ve skewness

Outliers.

extreme
values.

→ find by Boxplot

→ or by Scatterplot

when both x, y axis
are continuous.

(reln b/w 2 continuous variable)

Redefine upper & lower limit of Boxplot as :

$$\text{Lower limit} = Q_1 - 1.5 \times IQR$$

$$\text{Upper limit} = Q_3 + 1.5 \times IQR$$

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9	1	2	3	4	5	6
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Tuesday

8

February 2011

Week 6 • 39-326

DAY-9.

File Name for today's case study.

Descriptive statistics.xlsx

import pandas as pd

df = pd.read_excel ("Descriptive statistics.xlsx", sheetname=0)

You can also
give sheet
name.

refers 186
sheet

if you have
data on diff. sheet
change no.
acc to
another sheet

Sheet_name = 'Descriptive statistics'

9

df.head

df.head() → print first 5 records

df.head(10) → print first 10 records.

10

df.tail() → last 5 records

df.tail(10) → last 10 records.

11

12

df.columns → print all columns name.

1

df.info() → all information will be printed.

(no. of coln, rows, type of all coln... etc)

2

3

4

5

6

df['colnname'].mean()

df[['coln1', 'coln2', 'coln3']].mean()

print mean
of 3 diff.
coln.

for more than 1 var. use double sq. bracket.

Similarly,

9

Wednesday

February 2011

Week 6 • 40-325

- .drop duplicates
- .isna().sum()
- .dropna
- .duplicated()

Wk	Mo	Tu	We	Th	Fr	Sa	Su
	1	2	3	4	5	6	7
5							
6	7	8	9	10	11	12	13
7		14	15	16	17	18	19
8		21	22	23	24	25	26
9		28					

df[['coln1', 'coln2', 'coln3']].median()

- mode()
- var()
- std()
- describe()

in one shot you get all values

mean, std, var, min, 25 percentile,
(Q1), Q2, Q3, max...

median

skew()

kurt()

compare two/multiple columns →

pd.crosstab(df['Gender'], df['Jobcat'])

coln names.

plt.boxplot(df.coln)

plt.scatter(df.coln1, df.coln2)

plt.hist(df.coln)

Inferential

check by
using
skew,
kurt,
etc

Non-parametric
(data is not normally
distributed)

Parametric
(data is normally
distributed)

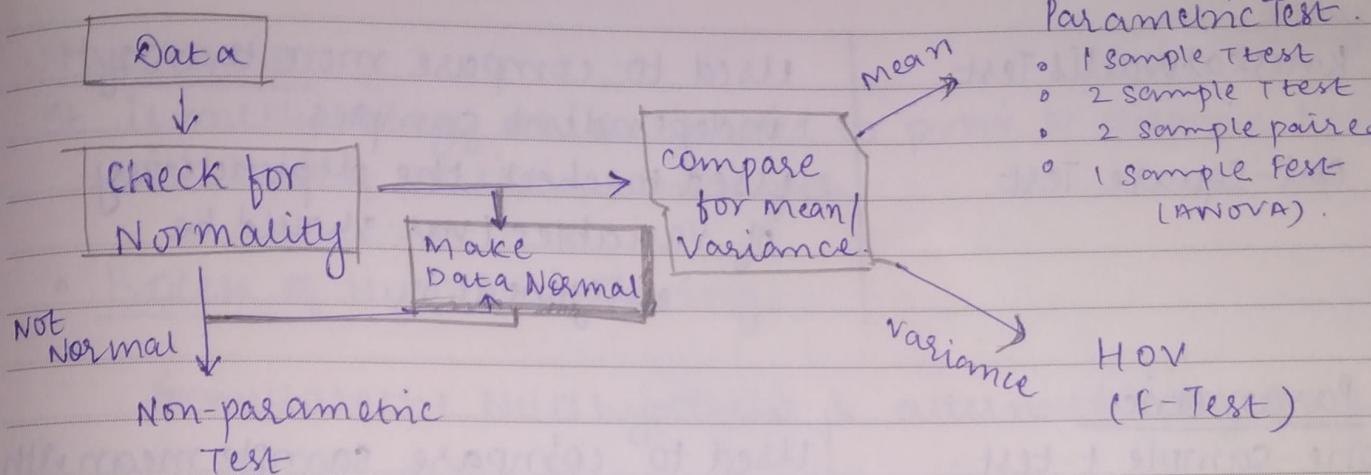
March 2011							
Wk	Mo	Tu	We	Th	Fr	Sa Su	
9	1	2	3	4	5	6	
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Thursday

February 2011

10

Week 6 • 41-324



wilcoxon - sign test

- wilcoxon - sign test
- Friedman
- mann - whitney
- kruskal - wallis
- chi - square.

Non-parametric:

- wilcoxon - sign test
- Friedman test
- Mann - Whitney test
- Kruskal - wallis test
- chi - square test

parametric

- One - sample T - test
 - Two sample Paired t - test
 - Two sample Separate t - test
 - One sample F test
- # one way (ANOVA)
- two sample F test
- # two way (ANOVA)

Non-parametric

① wilcoxon - sign test

It is used to compare two paired samples.

It is used to compare more than two paired samples.

It is used to compare two independent samples.

② Friedman Test

③ Mann - Whitney Test

Purpose

11

Friday

February 2011

Week 6 • 42-323

February							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su
5		1	2	3	4	5	6
6		7	8	9	10	11	12
7		14	15	16	17	18	19
8		21	22	23	24	25	26
9		28					

④ Kruskal-Wallis Test

Used to compare more than 2 independent samples.

⑤ Chi-Square Test

Used to check the dependency of variables (var. should be categorical).

Parametric

① One-sample t-test

Used to compare sample mean with population mean

② Two-sample paired t-test

Used to compare mean of two paired sample

③ Two-sample separate or t-test (independent)

Used to compare mean of two independent samples.

④ One-sample F-test

Used when dependent var. is continuous & independent var is categorical.

⑤ Two sample F-Test
or (Two-way ANOVA)Used when dependent variable is continuous & independent variable is categorical.
in which there are 2 independent variable.Non-Parametric

- without using population parameters → you are analyzing data
- data → not normally distributed.
- No-use of mean/median/mode.

	March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Saturday

February 2011

12

Week 6 ■ 43-322

Day 10.

Hypothesis.

It is methodology which is used to prove or disprove statement with help of statistical analysis.

Process of Hypothesis Testing.

Formulate the Null hypothesis & alternative hypothesis.



Select appropriate test stats.



choose level of significance, α , & degree of freedom



confidence Interval

(CI)

(DF)

compute calculated test value of test statistics



compute table test value of test statistics



Compare calculated values & table values.



Make statistical decision & state the Sunday managerial conclusion.

44-321

13

I) Hypothesis formulation.

• Null hypothesis (H_0).

→ state claim or assertion to be tested.

What we are assuming.

e.g. Avg. no. of TV sets in US Homes equal to 3

$$H_0: \mu = 3$$

- It is about population parameter, not about sample stats.

14

Monday

February 2011

Wk	Mo	Tu	We	Th	Fr	Sa	Su	2011
5		1	2	3	4	5	6	
6	7	8	9	10	11	12	13	
7	14	15	16	17	18	19	20	
8	21	22	23	24	25	26	27	
9								

Week 7 • 45-320

$$H_0: \mu = 3 \quad \checkmark$$

$$H_0: \bar{X} = 3 \quad \times$$

Always contain: (= or \leq or \geq) sign.

• Alternative Hypothesis (H_1)

→ opposite to Null Hypothesis.

Ex Avg no. of TV sets in US homes is not equal to 3.

$$H_1: \mu \neq 3$$

Never contain: (= or \leq or \geq) sign.

May or may not proven.

2.) Select appropriate Test:

Hypothesis Test
stat. form

when both
sides
rejection
is there

Two tail test

when there
is only
1 rejection

one tail test

is there

known
(Z test)

Unknown
(t test)

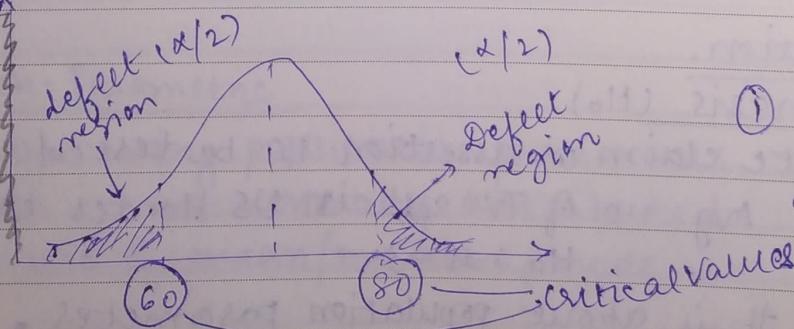
Left Tail Test

Right Tail
Test

① Null hyp. should have
equal sign

$$H_0: \mu = 3$$

$$H_1: \mu \neq 3$$



Both sides has rejection area.

March 2011						
Mo	Tu	We	Th	Fr	Sa	Su
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

Tuesday

February 2011

15

Week 7 • 46-319

popul'n SD is given

Z-test. (σ known).

$$Z_{\text{stat}} = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \quad \begin{array}{l} \xrightarrow{\text{Sample mean}} \\ \xrightarrow{\text{Pop. mean}} \\ \xrightarrow{\text{popul'n SD}} \\ \xrightarrow{\sqrt{n}} \text{sample.} \end{array}$$

Based on quality of data $\rightarrow z$ is defined.
or in project client define 'Z'.

t-test (σ unknown).

$$t_{\text{stat}} = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \quad \text{sample SD.}$$

Ex: Two tail test (σ unknown).

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168$$

i.e., don't reject H_0 .
insufficient evidence that true mean
not 168.
note is diff. from $\mu = 168$.

Given:

$$n = 25$$

$$\bar{x} \rightarrow \$172.50$$

$$s \rightarrow \$15.40$$

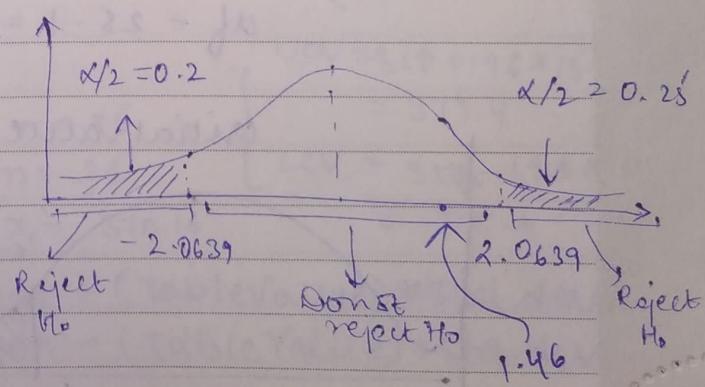
$$\alpha = 0.05 \quad (\text{error rate})$$

$$t \text{ test} \rightarrow \frac{172.50 - 168}{\frac{15.40}{\sqrt{25}}} = 1.46 \quad (\text{fall inside boundary})$$

degree of freedom $\rightarrow df = 25 - 1 = 24$
 $(\text{no. of rows} - \text{no. of coln})$

(from table) critical value:

$$\pm t_{24, 0.025} = \pm 2.0639$$



16

Wednesday

February 2011

February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
5		1	2	3	4	5
6	6	7	8	9	10	11
7	12	13	14	15	16	17
8	18	19	20	21	22	23
9	24	25	26	27	28	

Week 7 • 47-318

One Tail test.

only one side rejection is there.

(either left side or right side)

rejection area is in left side

Left tail test.

$$H_0: \mu \geq 3$$

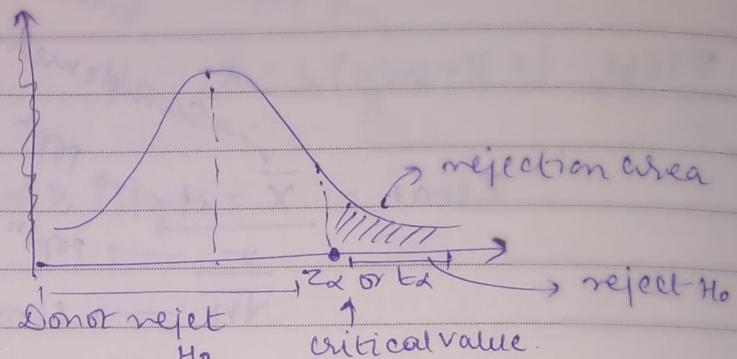
$$H_1: \mu < 3$$

Right tail test

$$H_0: \mu \leq 3$$

$$H_1: \mu > 3$$

rejection area is in right side.



lower-tail test since alternative hypothesis (H_1) is focused on the lower tail below mean of 3

upper-tail test since the alt. hypothesis (H_1) is focussed on upper-tail above mean of 3.

Ex. $H_0: \mu \leq 52$ → Right-tail test
 $H_1: \mu > 52$

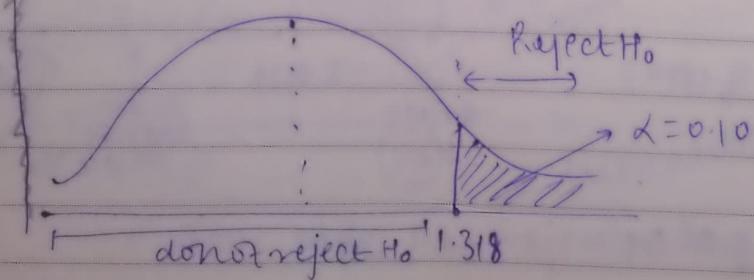
Given.

$$\alpha = 0.10 \quad (10\% \text{ error} \rightarrow 90\% \text{ accurate data}).$$

$$n = 25$$

$$df = 25 - 1 = 24$$

critical value from table: 1.318



Reject H_0 if
 $t_{\text{stat}} > 1.318$

March 2011						
Mo	Tu	We	Th	Fr	Sa	Su
1	2	3	4	5	6	
9						
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

* 95% $\rightarrow \alpha = 0.05$ (~~df~~ \times test table value = 1.96)

Thursday

17

February 2011

Note:

If α is not given ($\alpha = 0.05$ is default) Week 7 • 48-317

3.) Choose level of significance (α)

Level of significance is also called error rate

$\alpha = 0.01$ (1% error, 99% accurate)

$\alpha = 0.05$ (5% error, 95% accurate)

$\alpha = 0.10$ (10% error, 90% accurate)

how much data is accurate.

confidence interval

(limit) \rightarrow Boundary defining (left/right)

Method 1.

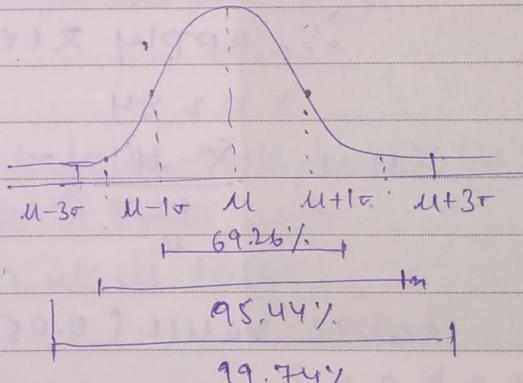
empirical rule tells accuracy of model..

Ex.

mean consumption

of various houses
is 200 units,

SD = 20 units.



• 68.2% consume energy

b/w 180 to 220 unit

• 99% save their energy

consumption b/w 140 to 260 units

Method 2.

Using formula.

$$H_0: \bar{x} = 350$$

$$H_1: \bar{x} \neq 350$$

$$\sigma = 75$$

$$n = 25$$

$$x = 370.16$$

$$UCV = \bar{x} + 1.96 \times \frac{\sigma}{\sqrt{n}}$$

$$= 379.4$$

$$LCV = \bar{x} - 1.96 \times \frac{\sigma}{\sqrt{n}}$$

$$= 320.6$$

- * upper confidence interval = Mean + (Table value) * std.dev.
- * lower confidence interval = Mean - (Table value) * std.dev

18

Friday

February 2011

Week 7 • 49-316

	February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
5		1	2	3	4	5	6
6		7	8	9	10	11	12
7		14	15	16	17	18	19
8		21	22	23	24	25	26
9		28					

Key

Degree of Freedom. $df = \text{no. of row} - \text{no. of coln.}$ 4.) Critical Value approach to Testing.

Ex.

$H_0: \mu = 3$

$H_1: \mu \neq 3$

$\alpha = 0.05$

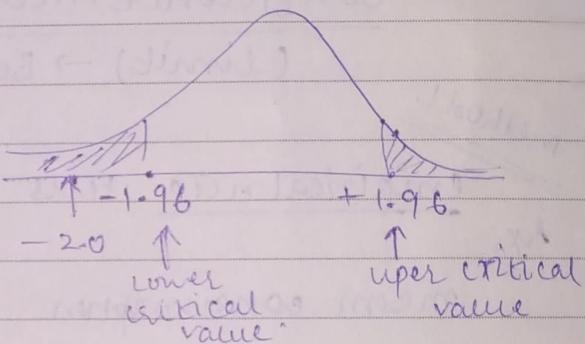
$n = 100$

$\sigma = 0.8$

 \therefore , apply Z-test

$X = 2.84$

$$Z_{\text{start}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{2.84 - 3}{\frac{0.8}{\sqrt{100}}} = -2.0 \quad (\text{fall outside the boundary}).$$

table value ($0.05 \rightarrow \alpha$) $\rightarrow -1.96$ \therefore reject hypothesis② P-Value Approach (Prob. Value approach.)• compare p-value with $\alpha = 0.05$ p-value < 0.05 (reject H_0)p-value ≥ 0.05 (don't reject H_0)

If judgment
is wrong,
Type I or II
error
occurs.

Decision
based on
sample

Truth about population

	H_0 true	H_1 true
Reject H_0	Type I error	correct decision
Accept H_0	Correct decision	Type II error.

2011

	Mo	Tu	We	Th	Fr	Sa	Su
March							
1	2	3	4	5	6		
9	10	11	12	13			
10	14	15	16	17	18	19	20
11	21	22	23	24	25	26	27
12	28	29	30	31			
13							

Saturday

February 2011

19

Week 7 • 50-315

$H_0 \rightarrow$ "Person is innocent"

		Decision	
		Jailed	Set Free
True state	Innocent	Type I correct	Correct
	Guilty	Correct	Type II

Q. How to control Type I & Type II error?

- increase α value,

• Type I error (α)

↳ ★ Reject null hypo. (H_0) when it is True.

considered serious type of error.

The prob. of a Type I error is α
called level of significance of test.

set by researcher in advance

• Type II error (β)

↳ ★ Failure to reject false null hypothesis.

The prob. of Type II error is β .

↳ ★ Accept H_0 when it is False.

Sunday

51-314

20

21

Monday

February 2011

Week 8 • 52-313

February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
5		1	2	3	4	5
6	6	7	8	9	10	11
7		14	15	16	17	18
8		21	22	23	24	25
9		28			26	27

if we want to find
more variation
we go for
Regression.

Day 11

Correlation.

Used to find the relationship b/w only 2 variables.
(not for more than two variables).

how we find relationship b/w 2 variables.

→ By corelⁿ coefficient (γ) or (R)

→ covariance

~~standard dev.~~
when want to measure consistency.

$$\gamma = \frac{\sum_{i=1}^n ((x_i - \bar{x})(y_i - \bar{y}))}{\sqrt{\sum_{i=1}^n ((x_i - \bar{x})^2)(y_i - \bar{y})}}$$

When we are changing one variable how much variation is there in other variable

y ↑

+ve corelⁿ5 (directly proportional)
6 to each other.-ve corelⁿ(inversely proportional)
to each other

} graphically
judge corelⁿ.

Coefficient of Determination (R^2)

$$R^2 = \frac{\text{Explained variance}}{\text{Total variance}}$$

→ When we change one variable how much variation is there in another variable

$$\Sigma \quad r = -0.635 \quad R^2 = .403 \quad (\text{About}$$

40% of variance is there in var.)

March	2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

corrⁿ should always be considered with significance levels (p -values). If corrⁿ is not significant, it is of little use.

Tuesday February 2011

22

Week 8 • 53-312

Correlation Matrix.

↙ * same variable taken in row & colⁿ.
 diag. always = 1.

find corrⁿ
 b/w one var.
 with another
 var. in the
 multiple
 combinations. in one shot.

Below diagonal is filled
 above diag. is empty / get same result.

corrⁿ Classification.

1. Product - Moment corrⁿ
2. Partial correlation.
3. Non-metric correlation

1. Product moment correlation

→ Both var. are continuous.

Ex (Attitude Vs Duration)

$$r = \frac{\text{covariance } xy}{\sqrt{SD_x \cdot SD_y}}$$

covariance formula:

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

standard deviation of x

$$y \cdot n$$

$$\left(\sum_{i=1}^n (y_i - \bar{y})^2 \right)^{\frac{1}{2}}$$

$$\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)$$

pearson

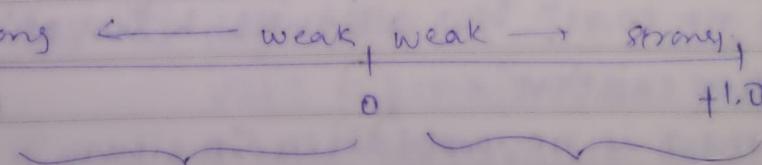
$r \rightarrow$ corrⁿ coefficient or or

corrⁿ coe

if $r=0$, there is
no relationship
b/w variables.

if $r=1$ which means
if one value ↑,
the other ↑ in
same proportion.

if $r=-1$ which
means
if one value +1.0
the other
↓ in
same
proportion.



-ve corrⁿ

+ve corrⁿ

$$(-1) < r < 1$$

perfect
-ve

perfect
colinearity.

23

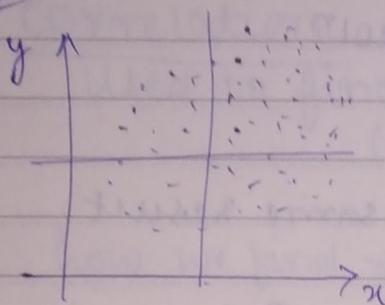
Wednesday

February 2011

Week 8 • 54-311

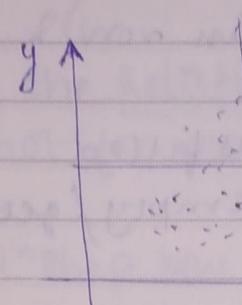
NOTE { +ve coreln \rightarrow graph in upward direction
-ve coreln \rightarrow " " downward "

Sx.05 R.



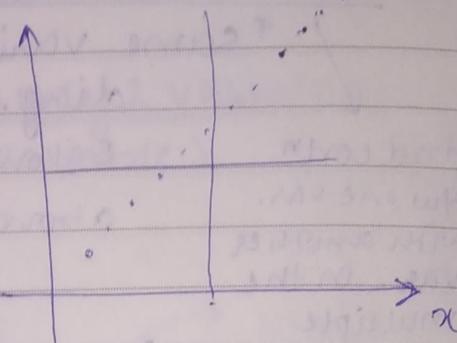
$$\delta = 0.30$$

low +ve coreln



$$\delta = 0.80$$

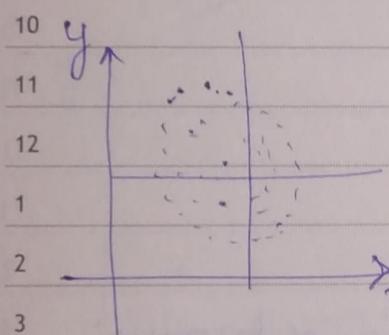
High +ve coreln



$$\delta = 1.0$$

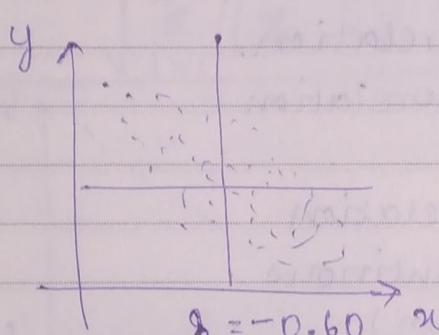
Perfect +ve coreln.

9



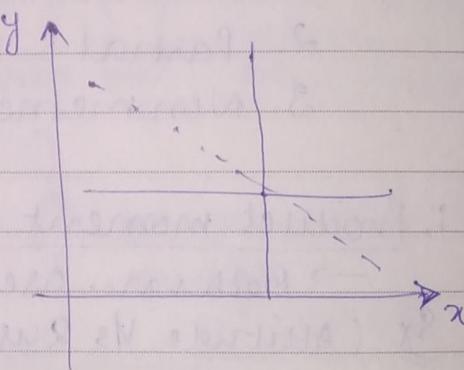
NO coreln

$$\delta = 0$$



Moderately
-ve coreln

$$\delta = -0.60$$



Perfect -ve
coreln

6

* Null hypothesis \rightarrow always be the -ve statement.

H₀: $\rho = 0$

H₁: $\rho \neq 0$

$$\Rightarrow \sigma^2 \neq 0$$

1.) There is no coreln b/w attitude & duration

2.) It is two tail test.

- there is no population SD given \therefore it is t-test.

3.) $\alpha = 0.05$ (5% error, 90% accurate data)

$$df = \text{low-coreln} \rightarrow 12 - 2 = 10$$

	February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
5		1	2	3	4	5	6
6	7	8	9	10	11	12	13
7	14	15	16	17	18	19	20
8	21	22	23	24	25	26	27
9	28						

25

Friday

February 2011

February 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
5		1	2	3	4	5
6	6	7	8	9	10	11
7	14	15	16	17	18	19
8	21	22	23	24	25	26
9	28					

Week 8 • 56-309

Practical Implementation in Python.

~~df = pd.read_excel('correlation.xlsx', sheet_name=1)~~
~~df → view dataset.~~ fn name for corr.

from scipy.stats import pearsonr

stats, p = pearsonr(df.Attitude, df.Duration)

print(stats, p)

0.9360778..., 7.545161167...e-06

≈ 0.9360778
r value

≈ 7.545161167...e-06
p-value

check

p < 0.05

7.5451611...e-06 →

∴, H_0 is rejected

0.0000075

H_1 is accepted

which is less
than 0.05

plt.scatter(df.Attitude, df.Duration)

+ the corelⁿ

df.corr() # Make corelⁿ matrix.

HW → find corelⁿ with Attititon & other variables.

② Statistical Test

(diagram + definations → written in day 9).

Refer day 9 Notes.

	2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10		7	8	9	10	11	12
11		14	15	16	17	18	19
12		21	22	23	24	25	26
13		28	29	30	31		

Check examples
from notes
(provided by LN)

Saturday

February 2011

26

Week 8 • 57-308

Non-Parametric

① Wilcoxon-Sign Test

Ex Same people → 2 different situations.

② Friedman Test → compare more than 2 paired samples

Ex same people → 3 different situations
(more than 2.)

③ Mann Whitney Test → compare 2 independent samples.

↙ 20 → Morning }
20 - Morning 20 → Afternoon } , more than 2 indepen.
20 - Evening 20 → Evening. Kruskal-Wallis

④ Chi-square → applicable only for categorical

variable. Gender : M, F

smoker : Postsmoker, Current smoker,
Non smoker.

Parametric

Sunday

27

① One Sample Test.

compare sample mean with population mean.

② Sample paired t-test.

compare mean of 2 paired samples

Ex same student compared mean of 2 subjects

③ Sample independent t-test.

compare mean of 2 indep'n samples

Ex Marathon Ex.

= Mean of Athlete & Non-Athlete duration.

28

Monday

February 2011

Week 9 • 59-306

Note. Stats value → critical value.

(every test has critical value & Table value).

calculated value → critical value.

February							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su
5		1	2	3	4	5	6
6		7	8	9	10	11	12
7		14	15	16	17	18	19
8		21	22	23	24	25	26
9		28					27

Day 12.

Practical implementation of Non-parametric

Wilcoxon Test.

→ H_0 : There is no significant difference in calcium level of patient initial & after 2 weeks.

```
df = pd.read_excel('wilcoxon.xlsx', sheet_name=0)
df.head()
```

```
from scipy.stats import wilcoxon
stats, p = wilcoxon(df.TOTALCIN, df.TOTALCW2)
print(stats, p)
```

↓ ↓

29.5 0.0025

$p < 0.05 \rightarrow$ reject H_0
accept H_1 ✓

Friedmann Test.

```
from scipy.stats import friedmanchisquare
stats, p = friedmanchisquare(df.TOTALCIN, df.TOTALCW2,
                               df.TOTALCW4)
```

```
print(stats, p)
```

↓ ↗ 27.9277

8.62133 e-07

exponential
power
 e^{-7}

$p < 0.05$

$\therefore H_0$ rejected
 H_1 accepted

Mann-Witney Test

```
df1 = pd.read_excel('mann whitney.xlsx', sheet_name=1)
df1 =
```