

python for data science.

- Pandas
- Matplotlib
- Numpy
- Seaborn.

- **dataframe** comprises of multiple columns.
- **series** comprises of single columns.

preferences > IPython console > Graphics > Backend = automatic.

A Series is used to model one dimensional data, similar to list in python. The series object also have a few more bits of data including an index and a name.

- To Import excel file.
Import pandas as pd
f1 = pd.read_excel("directory")

- To Import csv file
pd.read_csv ("directory")

- difference bet' list, series and array.

① Series

Import pandas as pd

song = pd.Series([12, 11, 10, 9], name = "count")
song.index

O/P: Index count

| | |
|---|----|
| 0 | 12 |
| 1 | 11 |
| 2 | 10 |
| 3 | 9 |

george = pd.Series([10, 7, 1, 22],

index = [1968, 1982, 1945, 1933],
name = 'george songs')

george # display series

george[1968] # To read or select item.

george[1968] = 68 # To update value.

del george[1] # To delete item.

george.mean() # display mean.

① arrays

Import numpy as np.

num_ser = np.array([45, 141, 39, 9])

num_ser[1] # accessing array

num_ser.mean() # display mean.

O/p: George (series)

| Index | George Song |
|-------|-------------|
| 1968 | 68 |
| 1982 | 7 |
| 1945 | 1 |
| 1933 | 22 |

O/p: number (array)

| | |
|---|-----|
| 0 | 0 |
| 0 | 145 |
| 1 | 141 |
| 2 | 38 |
| 3 | 9 |

* **Convert types:**

- string : .astype(str)
- numeric : pd.to_numeric
- Integer : .astype(int)
- datetime : pd.to_datetime

`Songs_66 = pd.Series([3, None, 11, 9],
Index = ["george", "ringo", "John", "Paul"],
name = "counts")`

`Songs_66.dtype`, # display datatypes. float '64'

`Pd.to_numeric(Songs_66.apply(str))`
1st we want to convert in str and then in
numeric It will display error.

To ignore the error use errors = 'coerce'

Pd. To numeric (songs_66.astype(str), errors = 'coerce')

Songs_66.dtypes

To deal with none use .fillna.

songs_66.fillna(-1)

Nan values will be dropped using
.dropna (not good practice)

Songs_66.dropna()

O/P.: float64

* error

- (displays series) dtype : float64

- dtype ('float64')

- george 3.6

- ringo -1.0

- John 11.0

- Paul 9.6

Name: count, dtype : float64

- (drops ringo & -1.0)

* Combining 2 series.

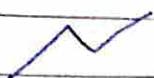
Songs = songs_66.append(songs_69)

O/P: After prints series 1 it will print 2nd
Series vertically.

* Plot

Import matplotlib.pyplot as plt

```
fig = plt.figure() # open canvas
songs_69.plot() OR plot(kind="bar", color='g')
plt.legend()
```



* Pandas Data Frame

- It is two dimensional datastructure.
- It is Immutable, heterogeneous tabular
- data structure with labeled axes rows, and columns.

To upgra upgrade pandas

conda install --upgrade pandas

conda install -c anaconda pandas

Conda Install pandas == 2.0.3

update

To check pandas version (on the spyder)

import pandas as pd

pd.__version__

* Creating data frame using list of list.

Import pandas as pd

technologies = [{"Sports": 2000, "30days"}, {"pandas": 2000, "40days"}]

df = pd.DataFrame(technologies)

O/P

| Index | 0 | 1 | 2 |
|-------|--------|------|--------|
| 0 | Sparks | 2000 | 30days |
| 1 | pandas | 2000 | 40days |

Since we have not given labels to columns and ~~rows~~ indexes, DataFrame by default assigns incremental sequence numbers as labels to both rows and columns, these are called Index.

columns_name = ["course", "fee", "duration"]

row_label = ["a", "b"]

df = pd.DataFrame(technologies, columns=columns_name, index=row_label).

print df

O/P:

| Index | courses | fee | duration |
|-------|---------|------|----------|
| a | Sparks | 2000 | 30 days |
| b | pandas | 2000 | 40 " |

* Creating DataFrame using dictionary.

```
technologies = {
    "course": ["spark", "pyspark", "hadoop", "python"],
    "fee": [2000, 2500, 3000, 4200],
    "duration": ["30 days", "40 days", "25 days", "50 days"]
}
```

```
df = pd.DataFrame(technologies)
print(df.dtypes)
```

O/P : (DataFrame)

courses object

Fee int64

duration object.

* Convert object into string.

```
df2 = df.convert_dtypes()
```

```
print(df.dtypes)
```

O/P : course string

fee int64

duration string.

* Convert all column to object.

```
df = df.astype(str)
```

```
print(df.dtypes)
```

O/P : course object

fee string

duration string

* change specific column:

```
df = df.astype({'fee': int, 'discount': float})
print(df.dtypes)
```

O/P: S

* convert all datatypes for all column in list.

```
pf = pd.DataFrame(technologies)
cols = ["fee", "discount"]
df[cols] = df[cols].astype("float")
df.dtypes.
```

O/P: living leaving object all will be float.

* ignore error

when we will try to convert object into int.

datatype will not change but error will be ignored.

```
df = df.astype({'courses': int}, errors="ignore")
df.dtypes.
```

* generate error

```
df = df.astype({'course': int}, errors="raise")
df.dtypes.
```

O/P: error.

* Convert dataframe to csv.

import pandas as pd.

```
technologies = {
    "Courses": ["sparks", "pysparks", "hadoop",
                "python"],
    "fee": [2000, 2500, 4000, 2500],
    "duration": ["30 days", "40 days", "35 days", "20 days"]
}
```

df = pd.DataFrame(technologies)

df # displays table in kernel.

df.to_csv('data_file.csv')

convert data frame to csv.

Exploratory Data Analysis (EDA).

df.shape

(8,4) (rows, columns)

df.size

32 (rows x columns)

df.columns / df.columns.values

Index(['courses', 'fee', 'duration'], dtype='object')

df df.index

Index(['r0', 'r1', 'r2', 'r3', 'r4', 'r5', 'r6', 'r7'], dtype = 'object').

* accessing one column

`df['Fee']`

* accessing two columns or more.

`df[['Fee', 'discount']]`

* accessing specific rows, column.

`df.iloc [rows, column]`

$df = df.iloc[1:3, 0:2]$

$df = df.iloc[:, 0:2]$ all rows & 0,1 columns

* accessing one column contents

~~`df['fee']`~~

* accessing individual cell.

`df['duration'][3]`

* Subtracting specific number from column

$df['Fee'] = df['Fee'] - 500$

`df['Fee']`

* Shows 5 no. summary `Describe()`

`df.describe()`

* changing column names.

`df.columns = ["A", "B", "C", "D"]`
`df`

rows = `axis=0`

column = `axis=1`

* `df2 = df.rename({‘A’: ‘C1’, ‘B’: ‘C2’}, axis=1)`

`df2 = df.rename({‘C’: ‘C3’, ‘D’: ‘C4’}, axis=1)` # columns

`df2 = df.rename(columns={‘A’: ‘C1’, ‘B’: ‘C2’})`

* To convert series into list.

Import pandas as pd.

`ds = pd.Series([1, 2, 3])`

`print(ds)` # Series

`print(ds.tolist())` # Converted into list

`print(type(ds.tolist()))` # <class ‘list’>

* Tuples, list & dict can be converted into series.

* Numpy array can be converted into series.

`import numpy as np`

`arry = np.array([1, 4, 9, 8])`

`ds2 = pd.Series(arry)`

`print(ds2)`

A) Write pandas program to convert one column as series.

```
d = {"col1": [1, 2, 3],  
     "col2": [4, 5, 6],  
     "col3": [7, 8, 9]}  
     }
```

```
df = pd.DataFrame(data=d)  
print(df) # display dataframe.
```

```
ds1 = df.iloc[:, 0] # all rows 0th column.  
print(ds1) # displays series.
```

① stack()

It is used to stack the prescribed level(s) from column to index.

It converts multilevel Index into single.

eg. a = [[1, 2], [3, 4]]
 .stack() = [1, 2, 3, 4].

eg:

```
s = pd.Series([  
    [1, 2],  
    [3, 4, 5],  
    [6, 7]  
])
```

print("original Series", s)

print("single Series")

```
s = s.apply(pd.Series).stack().reset_index()  
drop = True)
```

o/p:

| | |
|---|-----------|
| 0 | [1, 2] |
| 1 | [3, 4, 5] |
| 2 | [6, 7] |

resetIndex will reset
Index upto expected
value,
drop=True means drop
the original index

| | |
|---|---|
| 0 | 1 |
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 5 |
| 5 | 6 |
| 6 | 7 |

If $S = S \cdot \text{apply}(\text{pd.Series})$

$S = S \cdot \text{stack}()$

$S = \text{resetIndex}(\text{drop=True})$

* Add existing series to the series

$s1 = \text{pd.Series}([2, 3, 4])$

$\text{print}(s1)$

$\text{newS1} = \text{pd.concat}([s1, \text{pd.Series}(["Hi", 5])], \text{ignoreIndex}=\text{True})$

print(newS1)

o/p:

| | |
|---|---|
| 0 | 2 |
| 1 | 3 |
| 2 | 4 |

| | |
|---|----|
| 0 | 2 |
| 1 | 3 |
| 2 | 4 |
| 3 | Hi |
| 4 | 5 |

* Sort the series

```
s1 = pd.Series(["1", "2", "Python", "0"])
s2 = pd.Series(s1).sort_values()
print(s2)
```

O/P: 0 0
1 1
2 2
3 Python.

* ReIndexing Series

```
s1 = pd.Series(data=[1, 2, 3], index=["A", "B", "C"])
```

```
s1 = s1.reindex(index=["B", "C", "A"])
print(s1)
```

O/P: B 2
C 3
A 1.

* dropping rows using row labels.

```
df = df.drop("x1").
```

x1 row will be deleted.

`df = df.drop(df.index[1])` # 1st row deleted
`(df.index[2:5])` # 2nd, 5th row deleted
`(df.index[4:6])` # from 4th deleted.
`[2:5]` # 2-4 deleted.
`df.drop(0)` # 0th Index deleted
`df.drop([1,4])` # 1, 4 deleted
`df.drop(range(2,5))` # 2-4

Dropping Column

Must Mandatory to write axis=1

`df = df.drop(['courses'], axis=1)`
`(labels = ['courses'], axis=1)`
`(columns = ['fee'], axis=1)`

When we are making changes in same DataFrame it is mandatory to use `Inplace=True`

`eg: df.drop(df.columns[1], axis=1, inplace=True)`

- To delete multiple columns.

`df = df.drop(df.columns[[1,2]], axis=1)`
 # It will delete columns Index 1,2.

`l1col = ['courses', 'fee']`
`df = df.drop(l1col, axis=1)`

`df.iloc[[1,3],[1,2]]` # rows, column

* **iLoc (iloc)**

`df.iloc [startrow : endrow, startcolumn : endcolumn]`

`df = df.iloc[1:3,0:2]` # 1st row, 0,1 column

`df = df.iloc[0:2,:]` # 0st row, all column

`df = df.iloc[0:2]` # 2nd Index rows (series)

`df = df.iloc[[2,3,6]]` # It will select 2nd, 3rd, 6th Index

`df = df.iloc[1:5]` # It will select 1st to 4th row

`df = df.iloc[:2]` # 0th, 1st row

`df = df.iloc[-2:]` # -1th & -2th row

`df = df.iloc[: -3]` # all rows till -2th index

`df = df.iloc[:, :3]` # all rows in step of 3.

* **loc** ~~df.loc~~

`df = df.loc[1]` # display 1st row Series

`df = df.loc[['x1', 'x2']]` # display "x1", "x2" row

`df = df.loc[['x1', :, 'x7': 2]]` # display x1 to x7 in step of 2

`df = df.loc['x1': 'x4']` # display from x1 to x4 including

`df = df.loc['x0': 'x5', ['course', 'fee']]`

`df = df.loc['x0', 'fee': 'discount']` # Series

`df = df.loc[:, "duration":]` # all rows & all columns from duration.

* To Add column.

`tutors = ['ram', 'sham', 'ghansham', 'radhesham']`
`df2 = df.assign(Tutor_assigned=tutors)`
 column name values

① Add multiple columns

`mnccompany = ['tata', 'Infosys', 'amazon', 'google', 'TCS']`
`df2 = df.assign(Company=mnccompany)`.

② Derive new columns from existing columns.

`df2 = df.assign(per = lambda x: x.fee * x.disc / 100)`

③ Add column to existing column.

`df['Company'] = mnccompany`.

④ Add new column to specific position

`df.insert(0, 'tutor', tutors)`

* To Add row

`df.loc["f"] = ["avanti", 20, 3, "no"]`.
 or
 Index

* To find no. of column & rows in DF

`rows_count = len(df.index)` or `df.shape[0]`
`= len(df.axes[0])`

`column_count = len(df.axes[1])` or `df.shape[1]`

Assignment

- * no. of attempts must be less than 2.

```
print(df[(df['attempt'] > 2)])
```

```
df2 = df.loc[df['attempt'] > 2]
```

- * To find null value.

```
print(df[df["score"].isnull()])
```

```
print(df.loc[df['score'].isnull()])
```

- * Select Score between 15 to 20

```
df2 = df[df['score'].between(15, 20)]
```

```
df2 = df.loc[df['score'].between(15, 20)]
```

- * no. of attempts is < 2 & score > 15

```
print(df[(df['attempt'] < 2) & (df['score'] > 15)])
```

```
print(df.loc[(df['attempt'] < 2) & (df['score'] > 15)])
```

- * To change specific value in Dataframe.

```
df.loc["d", "score"] = 11.5
```

- * sum(), mean()

```
print(df["attempts"].sum())
```

```
print(df["attempts"].mean())
```

* Sort in ascending, descending

`df = df.sort_values(by=["column name"], ascending=[False])`

sorted in descending order

`df = df.sort_values(by=["column"], ascending=[True])`

* Replace values in DF map()

`df["column name"] = df["column name"].map({ "Hi": "Hello", "By": "Bye" })`

`df["column name"] = df["column name"].replace("Hi", "Hello")`

* Iterate over rows

```
for index, row in df.iterrows():
    print(row["name"], row["score"])
```

* While renaming the values in column, rename all the values else it will be replaced by nan

apply()

```
data = {"A": [1, 2],
        "B": [3, 4]}  
}
```

```
df = pd.DataFrame(data)
```

```
def add_3(x):
    return x + 3
```

$\text{df}^2 = \text{df}.apply(\text{add_3})$ # all column

$\text{df}^2 = ((\text{df}["A"])).apply(\text{add_3})$ # single column

$\text{df}["B"] = \text{df}["B"].apply(\text{add_3})$ # single column

To apply on multiple column list of column
is mandatory to pass

$\text{df}[["A", "B"]] = \text{df}[["A", "B"]].apply(\text{add_3})$
multiple column

Using lambda function. on all. column

$\text{df}^2 = \text{df}.apply(\lambda x: x + 10)$
 df^2

Using lambda fun' on multiple column

$\text{df}[["A", "B"]] = \text{df}[["A", "B"]].apply(\lambda x: x + 10)$
 df

* **transform()** instead of apply

```
df = df.transform(lambda x: x + 5)
```

* **map()** same as apply & transform.

* **Using numpy function**

```
import numpy as np
df["A"] = df["A"].apply(np.square)
print(df)
```

* **groupby()**

Single column.

```
df2 = df.groupby(["courses"]).sum()
print(df2) # Not have Index
```

```
df2 = df.groupby(["courses"]).sum().reset_index()
```

On multiple columns

```
df2 = df.groupby(["courses", "duration"]).sum() # with Index.
```

* **Shuffle**

shuffle DF rows.

```
df2 = df.sample(frac=1) # 100% shuffle.
```

```
df1 = df.sample(frac=0.5) # 50% shuffle
```

drop = True

If we set drop = True, reset_index will delete the index instead of inserting it back in column.

df1 = df.sample(frac=1).reset_index(drop=True).

* Joining DF

tech2 = {

"courses": ["cpp", "Java", "Python"]

"fee": [100, 200, 300]

"duration": ["30 days", "10 days", "20 days"]

}

Index1 = ["r1", "r2", "r3", "r4"]

df1 = pd.DataFrame(tech2, index1)

tech3 = {

"courses": ["C++", "Java", "C"]

"Discount": [300, 400, 500]

}

Index2 = ["r1", "r5", "r6"]

df2 = pd.DataFrame(tech3, Index2)

Left Join

df3 = df1.join(df2, lsuffix="left", rsuffix="right")
print(df3)

y_1 CPP 100 30days 300
 y_2 JAVA 200 10days nan
 y_3 python 300 20days nan

InnerJoin

```
df3 = df1.join(df2, lsuffix="left", rsuffix="right", how="inner")
```

O/p:

y_1 CPP 100 30days 300

Right Join.

how = "right"

O/p:

y_1 CPP 300 100 30days,
 ~~y_2~~ JAVA nan nan nan
 y_3 C nan nan nan

Joining according to column Index don't matter

```
df3 = df1.set_index("courses").join(df2.set_index("courses"), how = "inner")
```

O/p. ~~y_1~~ CPP 100 30days 300
 JAVA 200 10days 400.

* **Merge()** = InnerJoin

$df3 = pd.merge(df2, df3)$

$df3 = df1.merge(df2)$

* **Concat()**

Vertically joins all dataFrame

$data = (df1, df2)$

$df3 = pd.concat(data)$

$df3.$

Numpy

* Create an array.

Import numpy as np

```
arr = np.array([10, 20, 30]) # 1 dimension
print(arr)
```

* Create multidimensional array.

```
arr = np.array([[10, 20, 30], [40, 50, 60]]) # 2dimension.
```

* ndmin - to specify \geq how many dimensions you need.

```
arr = np.array([10, 20, 30], ndmin=3)
```

```
O/P: [[[10 20 30]]]
```

* To change datatype.

```
arr = np.array([10, 20, 30], dtype=complex)
```

```
O/P: [10.+0.j 20.+0.j 30.+0.j]
```

* To get dimension of array.

```
arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
```

```
print(arr.ndim)
```

```
O/P: 2
```

* to get array type, size, shape.

```
arr.dtype # int32
```

```
arr.size # 6
```

```
([10, 20, 30],
 [4, 5, 6])
```

```
arr.shape # (3, 3) (row, column)
```

* create sequence using range fn?

```
arr = np.arange(0, 20, 3)  
print(arr)
```

* accessing array using Index.

```
arr = np.arange(11)
```

```
print(arr)
```

```
print(arr[2])
```

2

```
print(arr[-2])
```

9.

* accessing multidimensional array

```
arr = np.array([[10, 20, 30], [40, 50, 60]])
```

```
print(arr[1, 2]) # 60
```

```
print(arr[0, 2]) # 30
```

```
print(arr[0, -1]) # 30
```

```
arr = [1, 2, 3, 4, 5, 6, 7, 8, 9, 0]
```

```
print(arr[1:8:2]) # start:stop:step.  
# [2, 4, 6, 8]
```

```
print[-2:3:-1] # -2 to 3 in step 1
```

[9, 8, 7, 6, 5]

```
print(arr[-2:1])
```

[9, 0]

`arr = ([[1, 2, 3], [4, 5, 6], [7, 8, 9]])`

`arr[1, 2] = 6`

`arr[1, :]` # one row all column

`arr[:, 2]` # 2nd column represented as array.

`arr[:, :, 2]` # all rows till 3 of alternate columns.

* Integer array indexing, allows selection of arbitrary items

`arr = np.arange(35).reshape(5, 7) (rows, columns)`

* Boolean array Indexing,

`rows = np.array([False, True, True]) # 0=False
1, 2=True`

`wanted_rows = arr[rows, :]` # 1st, 2nd row of all columns

* Convert array into list

`array = np.array([1, 2, 3])`

`lst = array.tolist()`

Multidimensional array → Nested list.

1st, 2nd, 3rd row → 1st, 2nd, 3rd index of list.

* Convert list into array.

lst = [1, 2, 3]

array = np.array(lst) or np.asarray(lst)
print(array)

* Resize array.

array = np.array([[1, 2, 3], [4, 5, 6]])
print(array.shape(3, 2))

P array.shape = (3, 2)
print(array.shape)

* Find if any number is zero

① all()

Import numpy as np.

x = np.array([1, 2, 3, 4])
print(np.all(x)) # True

y = np.array([0, 2, 3])

print(np.all(y)) # False.

If array contains zero then o/p is false.

② any()

z = np.array([1, 0, 0])

print(np.any(z))

O/P: True

```
y = np.array([0,0,0])
print(np.any(y))
o/p : False
```

True if any non zero number in array

④ isfinite()

check element by element if is any not infinite or nan value.

```
x = np.array([-1,0,np.nan,np.float])
print(np.isfinite(x))
# o/p : [True True False False]
```

⑤ isnan()

```
print(np.isnan(x))
op: [False, False, True, False]
```

⑥ Comparision element by element

```
x = np.array([2,3])
y = np.array([1,4])
print(np.greater(x,y))
(np.greater_equal(x,y))
(np.less(x,y))
(np.less_equal(x,y))
```

o/p : [True False]
 [True False]
 [False True]
 [False True.]

* create 3×3 Identity matrix $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

array = np.Identity(3)

* generate random number betn 0 to 1.

syntax: center of location, scale, size

rand_num = np.random.normal(0,1)
print(rand_num).

* create 3×4 array & iterate over it.

arr = np.arange(10,22).reshape(3,4)

print(arr) # display array of 3x4

for x in np.nditer(arr):

 print(x, end=" ")

 print()

$$\begin{bmatrix} 10 & 11 & 12 & 13 \\ 14 & 15 & 16 & 17 \\ 18 & 19 & 20 & 21 \end{bmatrix}$$

* Vector creation of length 5 betn 10 to 50

v = np.linspace(10,49,5)

start, stop, length.

* reverse an array,

x = np.arange(8,10) # 0....9

x = x[::-1]

print(x) # 9....0

* dot product of 2 matrix

```
x = [[1, 2], [3, 2]]
```

```
y = [[3, 2], [3, 5]]
```

```
result = np.dot(x, y)
```

```
print(result)
```

O/p : 9 12
15 16

* cross product.

```
result1 = np.cross(x, y)
```

```
result2 = np.cross(y, x)
```

| | |
|-----------|----------|
| # result1 | result2 |
| [-4 9] | [4 -9] |

* To find determinant of given square array.

```
from numpy import linalg as la
```

```
a = np.array([[2, 2], [3, 1]])
```

```
print(np.linalg.det(a))
```

o/p: -5.0000...1

* Find eigenvalues & eigen vectors.

```
m = np.mat("3 -2 ; 1 0")
```

```
print(m)
```

```
w, v = np.linalg.eig(m),
```

```
print("values, vectors:", w, v)
```

o/p: [[3 -2]

[1 0]]

values = [2. 1.]

vector = [[0... 0...]
0... 0...]]

* Inverse of Matrix

```
result = np.linalg.inv(a)
print(result)
```

* Data science processes

- 1) Setting the research goal
- 2) Retrieving data
- 3) Data preparation
- 4) Data exploration
- 5) Data modeling or model building
- 6) Preparation presentation & automation

Matplotlib.

* Import Matplotlib.pyplot as plt
`plt.plot([1, 3, 2, 4])`
`plt.show()`

* Multiple plot

```
x = range(65)
plt = plot(x, [xi * 1.5 for xi in x])
plt = plot(x, [xi * 3 for xi in x])
plt.show()
```

* Grid in graph.

Import matplotlib.pyplot as plt

Import numpy as np

`x = np.arange(1, 6)`

`plt.plot(x, x * 1.5, x, x * 1.3, x, x / 3.0)`

`plt.grid(True)`

`plt.show()`

* handling axes.

`x = np.arange(1, 6)`

`plt.plot(x, x * 3, x, x * 1.5, x, x / 3.0)`

`plt.axis()`

`plt.axis([0, 5, -1, 13])`

`plt.grid(True)`

`plt.show()`

* adding labels.

```
plt.plot([1,3,2,4])  
plt.xlabel("This is x axis")  
plt.ylabel("This is y axis")  
plt.show()
```

* adding title

```
plt.plot([1,3,2])  
plt.title("simple plot")  
plt.show()
```

* adding legend

```
x = np.arange(1,5)  
plt.plot(x, x**3, label = "Normal")  
plt.plot(x, x**1.5, label = "Fast")  
plt.plot(x, x/3.0, label = "Slow")  
plt.legend()  
plt.show()
```

* color abbreviation

| | |
|---|---------|
| b | blue |
| c | cyan |
| g | green |
| k | black |
| m | magenta |
| r | red |
| w | white |
| y | yellow |

Import matplotlib.pyplot as plt
Import numpy as np

y =

```
y = np.arange(1,3)  
plt.plot(y, "y")  
plt.plot(y+1, "m")  
plt.plot(y+2, "c")  
plt.show()
```

* Specifying style in multiline plot.

```
y = np.arange(1,3)  
plt.plot(y, "--", y+1, "-.", y+2, ":")  
plt.show()
```

* Markers

```
y = np.arange(1, 3, 0.2)
plt.plot(y, "x", y+0.5, "o", y+1, "D", y+1.5, "v", y+2,
plt.show()
```

* histogram charts.

```
import matplotlib.pyplot as plt
import numpy as np
y = np.random.randn(1000)
plt.hist(y)
plt.show()
```

* Bar graph.

syntax: `bar([x cordin], [y cordi])`

```
plt.bar([1, 2, 3], [3, 2, 5])
plt.show()
```

* Scatter plots.

displays values of 2 sets of plot data at x, y axis

```
Import matplotlib.pyplot as plt
import numpy as np
x = np.random.randn(1000)
y = np.random.randn(1000)
plt.scatter(x, y)
plt.show()
```

* change size & color

```
size = 50 * np.random.randn(1000)
color = np.random.randn(1000)
plt.scatter(x,y,s=size,c=color)
plt.show()
```

* adding text

```
x=np.linspace(-4,4,1024)
y=.25*(x+4.)*(x+1.)*(x-2.)
plt.text(-0.5,-0.25,"bracketed minimum")
plt.plot(x,y,c='k')
```

Statistics

| | |
|----------|--|
| Page No. | |
| Date | |

- (i) **Descriptive Statistics**: (inwards)
- (ii) **Inferential Statistics**: (outwards)

- **descriptive** is with processing of data without attempting to draw any **inferences** from it.
- **Inferential statistics** uses mathematical tools to make ~~or~~ **predict** and ~~or~~ **projection** by analyzing the given data.

* Population

The population often consist of large group of specifically defined element.
ex: population of town means that all the people living within that boundaries of town.

A small group of elements from the population is called **sample**.

* Data :

It is collection of numbers gather to give some information.

Types of Data (depending upon source):

- **Primary data**
- **Secondary data**.

D Primary data

The data collected by experimentation.

2) Secondary data.

data collected from someone's experimentatio

* Variable

The quantity entity which can vary from one individual to another is called variable.

* Continuous Variable / data.

A variable which can assume each an every value within given range is called continuous variable / data.

e.g monitoring of temperature.

* Discrete variable or data

A variable which can assume only some specific values within given range is called discrete variable or discrete data.

e.g students marks, population of village.

* Raw data.

The Data obtain from original source but not arranged properly is called raw data.

* Preprocess data / Tabular Data.

The process of placing classified data in tabular form is known as tabular data or processed data.

* Structured / Unstructured data.

Structured data is tabular data, represented in form of numerical values.

* Unstructured all data is in original form and not numerical form.

* Range

The difference between highest & lowest value of observation.

* Frequency of observation

The frequency of particular data value, number of times data value occurs.

* Frequency distribution

It is tabulation of values that one or more variables take in a sample.

DATA

Quantitative

continuous

3.5, 4.6

discrete

38

nominal

rule, beautiful etc.

order

good
bad

slow, fast

Qualitative

* Mean.

Name . monthly Income

| | |
|-------|------|
| Rob | 5000 |
| Rafiq | 6000 |
| Nina | 4000 |
| Sofia | 7500 |
| mohan | 8000 |
| Tao | 7000 |

Average = 6250

| Name | \$ |
|---------|------------------------|
| Rub | 3000 |
| Rafiq | 6000 |
| Nina | 4000 |
| Sofia | 7500 |
| mohan | 8000 |
| Tao | 7000 |
| elomwuk | 10 million |
| | average : 1.43 million |

If ~~outlayer~~ value get's added & mean is wrong.

* Median

| x | y | z | a | b | c | d | e |
|------|------|------|------|------|------|--------|-------|
| 4000 | 5000 | 6000 | 7000 | 8000 | 9000 | 10,000 | 11000 |
| | | | | | 7500 | | |

* Handling missing values

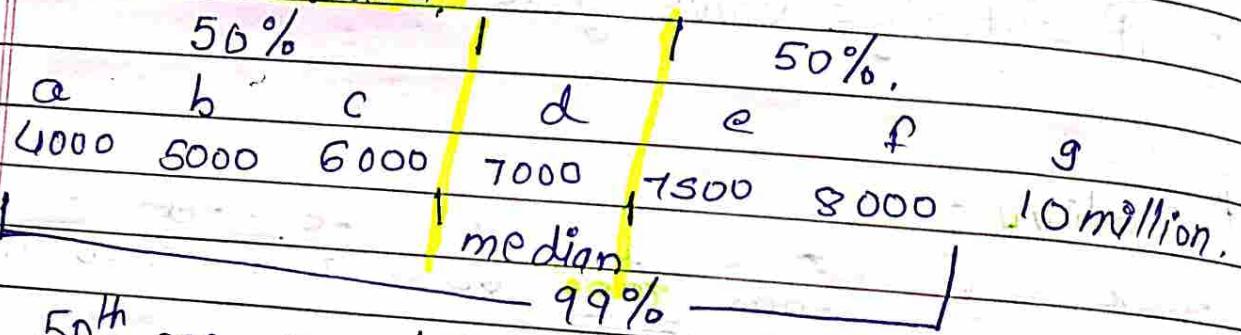
We can replace missing value with mean or median.

* Replacing missing values using mean values

| name | \$ | creditScore | Approv lone? |
|------|-------------|-------------|--------------|
| x | 5000 | 650 | No |
| y | 6000 | 400 | No |
| z | 4000 | 780 | Yes |
| a | 1.6 million | 810 | Yes |
| b | 8000 | 410 | No |
| c | 7000 | 850 | Yes |
| d | 10 million. | 880 | Yes |

- * Replacing missing value with median.
 $6500 \approx 1.6 \text{ million}$
- (* median used in descriptive analysis & Data cleaning (filling NA values))
- * Outlayer is an extreme value of data in given dataset.

* Percentile value,



50^{th} percentile for this dataset is 7000
 0^{th} percentile for this dataset is 4000.

* Mode

most frequently appearing value.

Name Restaurant vote.

a mexican

b mexican

c Italian

d thai

e mexican

f India

Mode = Mexican.

* Dealing with outlier. (Mean/Median absolute deviation) (MAD)

perfect example:
History Test

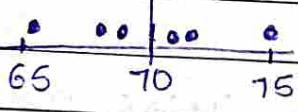
Not perfect:
Math Test

| name | score | Abs |
|------|-------|-----|
| a | 75 | 5 |
| b | 72 | 2 |
| c | 63 | 2 |
| d | 65 | 5 |
| e | 67 | 3 |
| f | 73 | 2 |

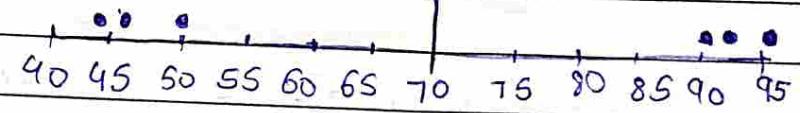
$$\text{average} = 70 \\ \text{mean} = 3.36.$$

| name | score | Abs (Score-avg) |
|------|-------|-----------------|
| a | 91 | 23 |
| b | 96 | 26 |
| c | 73 | 27 |
| d | 47 | 23 |
| e | 53 | 19 |
| f | 90 | 20 |

$$\text{average} = 70, \\ \text{mean} = 23$$



here data(score)
is near to mean
value.



here all the scores are far
from mean value \therefore model
will be poor.

After calculating
Absolute value:
(score-avg)

The mean is
3.36 & deviation
is less.

After calculating Absolute value,
mean is 23 & deviation is less.

The Above was mean Absolute
deviation (MAD)/Var.

If there is still outlier after Applying mean
Absolute deviation then go for median absolute
deviation.

$$G = \sqrt{\frac{\sum (x_i - \bar{X})^2}{N}} \quad \text{Standard deviation} = \sqrt{\text{Var.}}$$

Page No. _____
Date _____

let us consider another table,

History Test.

| Name | Score | Nbs | (score-avg) ² |
|------|-------|-----|--------------------------|
| a | 76 | 5 | 25 |
| b | 72 | 2 | 4 |
| c | 68 | 2 | 4 |
| d | 65 | 5 | 25 |
| e | 67 | 3 | 9 |
| f | 73 | 3 | 9 |

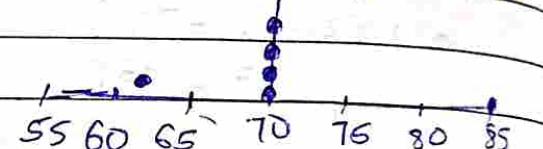
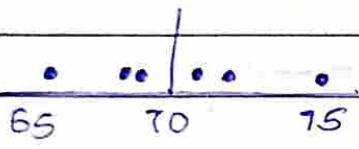
$$\text{MAD} = 3.33 \quad \text{avg}$$

$$\text{avg} = 12.66 \quad \sqrt{\text{avg}} = 3.55$$

| Name | Score | Abs | (score-avg) |
|------|-------|-----|-------------|
| a | 83 | 13 | 16.9 |
| b | 70 | 0 | 0 |
| c | 70 | 0 | 0 |
| d | 63 | 7 | 4.9 |
| e | 70 | 0 | 0 |
| f | 70 | 0 | 0 |

$$\text{MAD} = 3.33 \quad \text{avg}$$

$$\text{avg} = 36.33 \quad \sqrt{\text{avg}} = 6.01$$



here, after MAD the deviation spread is lower, which is good for model.

here, even after mean is the deviation spread is higher, which is not good altho MAD is same for both table but it's standard deviation is more, If mean is't spread is more

so, if standard deviation is more than we can say spread is also more

* L1 norm & L2 norms

L1 norm that is calculated as the sum of absolute values of vector. L2 norm that is calculated as square root of sum of squared vector values.

mean absolute = 61
standard deviation = 20

| | |
|----------|--|
| Page No. | |
| Date | |

* Transforming data on spread and centre.

What is the mean weight and what is standard deviation?

$$\bar{x} = 135.6 \quad \sigma = 31.75$$

| | | | | |
|-----|-----|-----|-----|-----|
| a | b | c | d | e |
| 105 | 156 | 145 | 172 | 100 |

D) Transform data with addition.

$$\begin{array}{l} a \quad 105 \quad + \quad 5 \quad = 110 \quad \bar{x}_{\text{new}} = \frac{a+b+c+d+e}{5} \\ b \quad 156 \quad + \quad 5 \quad = 161 \\ c \quad 145 \quad + \quad 5 \quad = 150 \quad \bar{x}_{\text{new}} = 140.6 \\ d \quad 172 \quad + \quad 5 \quad = 177 \quad \sigma_{\text{new}} = 31.75 \\ e \quad 100 \quad + \quad 5 \quad = 105 \end{array}$$

When we are going to add same value in all data point then there is no transformation in given data set.

Measure of centre

Affected by:

(+, -, ×, ÷)

(mean, median,
mode)

Measure of spread

Affected by:

(×, ÷)

(standard deviation)
range

Pie, Bar = categorical data
stamp plot, histogram = continuous (numerical) (quantitative)
Time plot



2) a 105 $\bar{x} = 135.6$ $(\times) 2.5 + 750$

b 156

c 145 $\bar{x}_{new} = (135.6)(2.5) + 750 = 1089$

d 172

e 100 $s = 31.75$ $(\times) 2.5$

$s_{new} = (31.75)(2.5) = 79.38$

25 ml of water is for every pound they

measure of centre

$$\text{Centre new} = (\text{centre old})(x) + B$$

$$\text{Spread new} = (\text{Spread old})(x)$$

* Impact of Outlayer on mean, median, Mode & Range

| a | b | c | d | e | f | g. |
|------|------|------|------|---------------|------|------|
| 26.0 | 15.0 | 20.5 | 31.0 | <u>-350.0</u> | 31.0 | 30.5 |

outlayer.

measure with outlayer without outlayer

Affected mean - 28 25.667

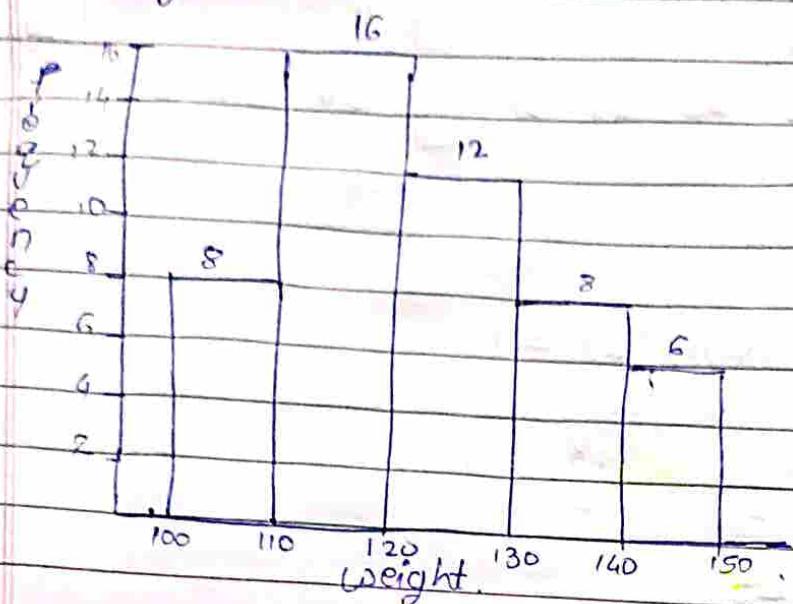
Resistant median 26 28.25

Resistant mode 31 31

Affected range 381 16.

max - min

* Histogram



* frequency distribution Table:

| weight | frequency |
|-----------|-----------|
| 100 - 110 | 8 |
| 110 - 120 | 16 |
| 120 - 130 | 12 |
| 130 - 140 | 8 |
| 140 - 150 | 6 |

Use of df. describe

1) 5 number summary

2) minimum, 1st Quartile, median, 3rd Quartile, maximum.

Quantile value 1st (Q_1)

Quantile value 3rd (Q_3)

4 8
 10 11 12 25 25 27 31 33 34 34 35 36
 = ↓
 Q₁ Q₃

mean.

left side of mean (L) = 7.

$$Q_1 = \frac{L+1}{2} = \frac{8}{2} = 4.$$

$$\therefore Q_1 = 25.$$

Right Side of mean (R) = L

$$Q_3 = \frac{R+1}{2} = \frac{8}{2} = 4.$$

$$Q_3 = 36$$

* Frequency distribution

weight frequency, cal

100-110

8

$8 \div 50$

110-120

16

$16 \div 50$

120-130

12

$12 \div 50$

130-140

8

$8 \div 50$

140-150

6

$6 \div 50$

Sum = 50

Relative frequency distribution

relative frequency, freq (%)

0.16

16%

0.32

32%

0.24

24%

0.16

16%

0.12

12%

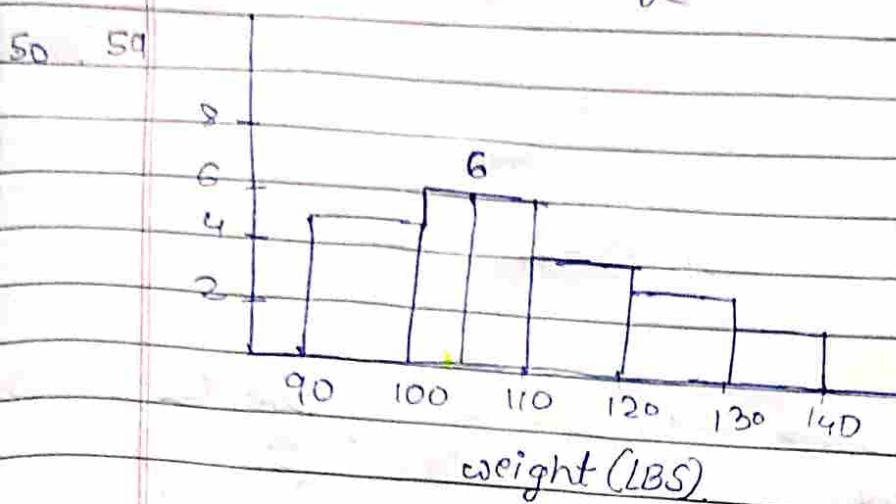
Sum = 1

100%

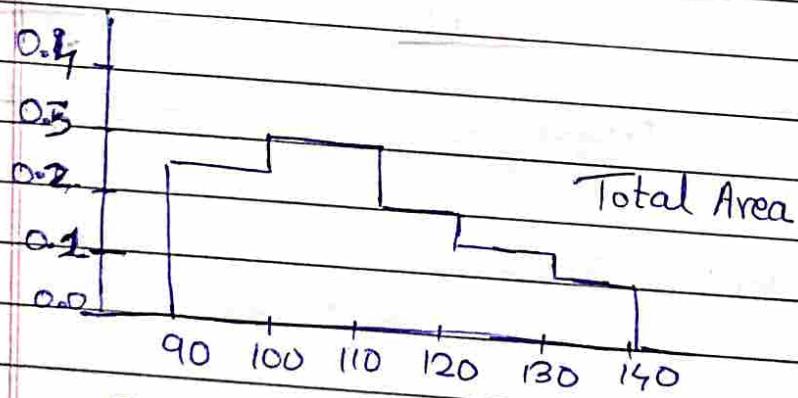
$$\text{relative frequency} = \frac{8}{50} = 16\%$$

Sum must be always 1 or 100%.

* 'Regular' frequency distribution ~~table~~ graph.

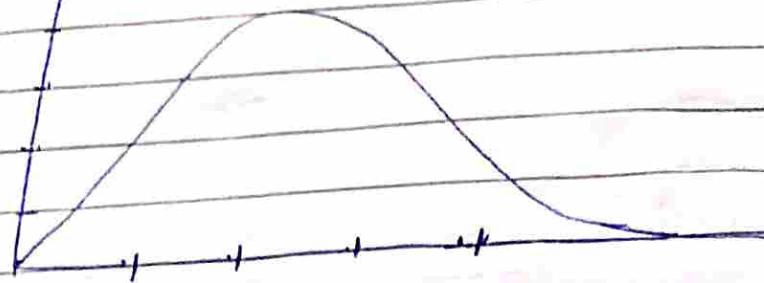


In regular frequency distribution, if we want corresponding values of weights ~~between~~ 100 & 110 you can easily plot that values as shown in Fig. But This is not true with relative frequency distribution.



For example If we want to find corresponding weight in relative frequency distribution betw 100 & 110 then In that case we will have to find out area under the curve.

Relative frequency distribution is transformed in equivalent density curve

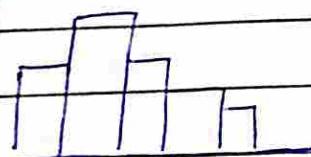
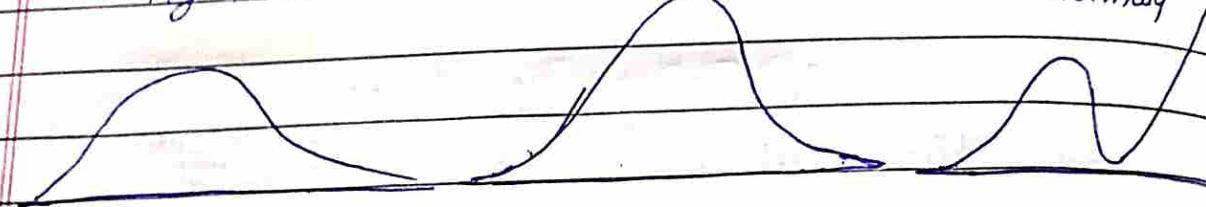


* Various forms of density curves.

Normaly distributed
but right skewed

symetrically
normal distributed

Bimodel
-cally distributed
normally



Total area = 1

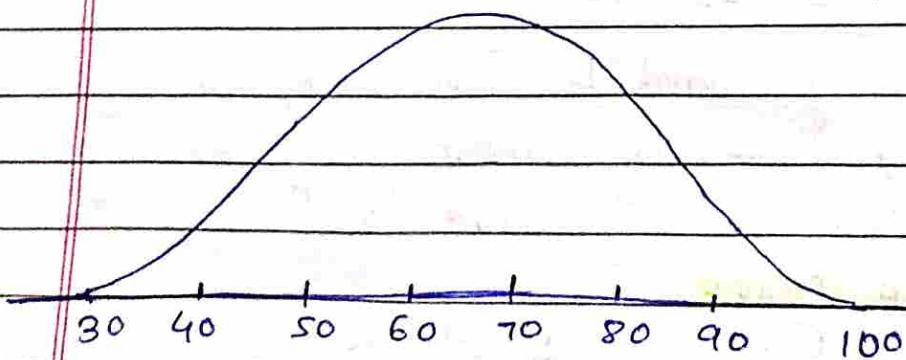
Total area = 1

Total area = 1

* How to read density curves

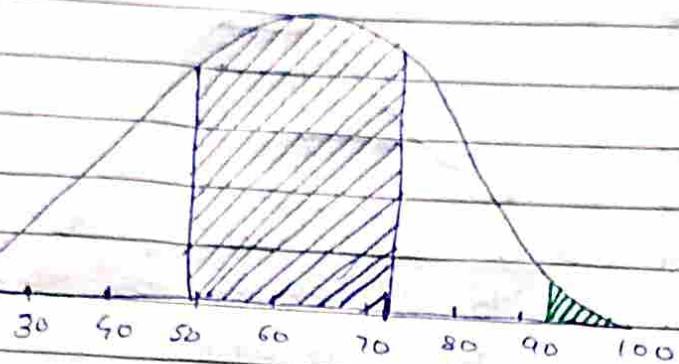
Test scores .

$n = 1000000$



If n (number of entries) are more than curve is defined) but if n is less then there are empty spaces in curves.

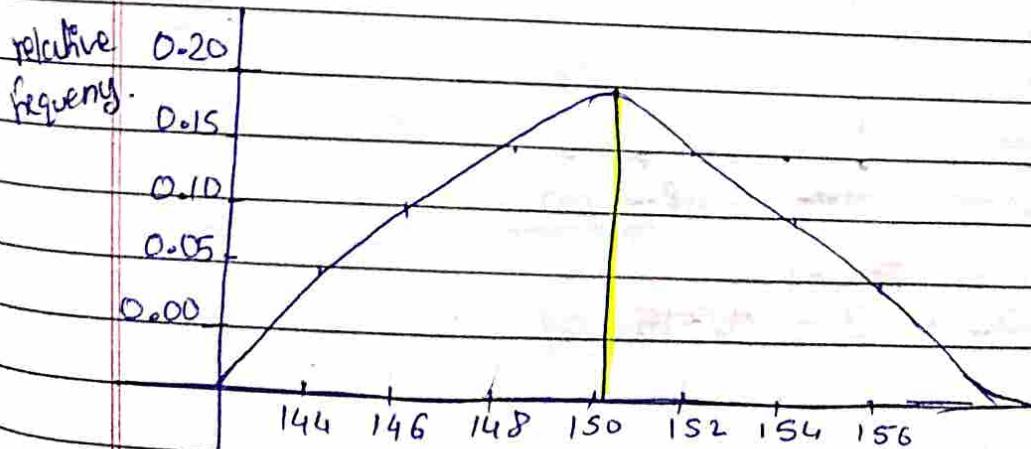
Large majority of students scored score between 50 - 70



Very few student score betⁿ 8 95 to 100

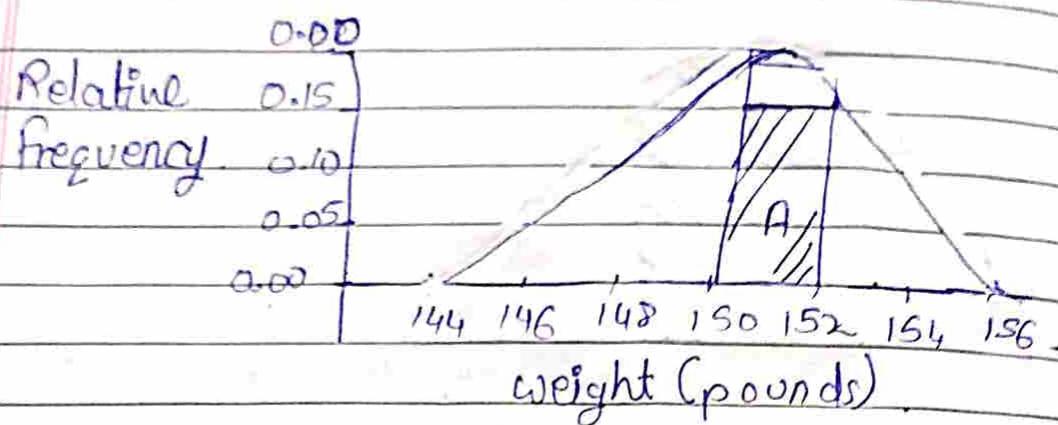
* Properties of density curve.

- 1) A density curve must lie on x axis or ^{above} \wedge horizontal axis.
- 2) The total area under the curve is always 1.
3. For the density curve below, approx what is percentage of people weight exact 150 pounds?

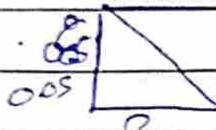


Area of line is 0 \therefore 0 people weight 150 pound.

* What percentage of people weight below 150 pounds?



~~0.05~~



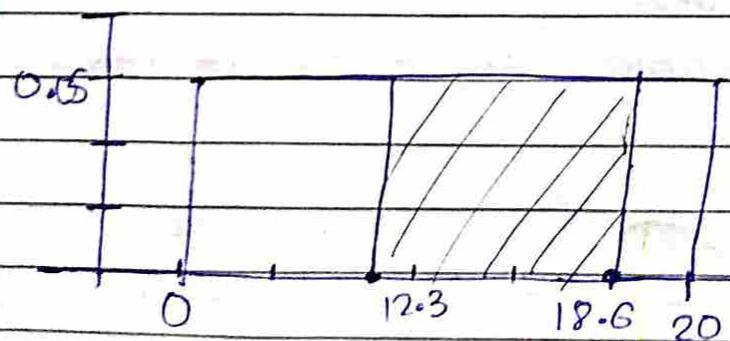
$$A = 1 \quad A/2 = 0.05$$

~~0.05~~

$$A = 2 \times 0.15 = 0.30.$$

$$A_{\text{Total}} = 0.30 + 0.05 =$$

* For the ~~non~~ uniform distribution below, what proportion of values are located between 12.3 & 18.6?



$$18.6 - 12.3 = 6.3.$$

$$A = 0.05 \times 6.3 = 0.315 \text{ or } 31.5\%.$$

* Normal distribution/Gaussian distribution.

- Parameter

A number that describes the data from a population.

- Statistics

A number that describes the data from a sample.

Sample population.

mean \bar{x} μ

standard deviation s σ

Statistic Parameters.

If σ is small then bell shape curve is taller (A)
If spread is higher then curve will be wider.

- The normal distribution is unimodal.
- The normal curve is symmetric about its mean
- The parameters μ and σ completely characterize the normal distribution.

$$X \sim N(\mu, \sigma^2)$$

↓ ↓
Variable mean

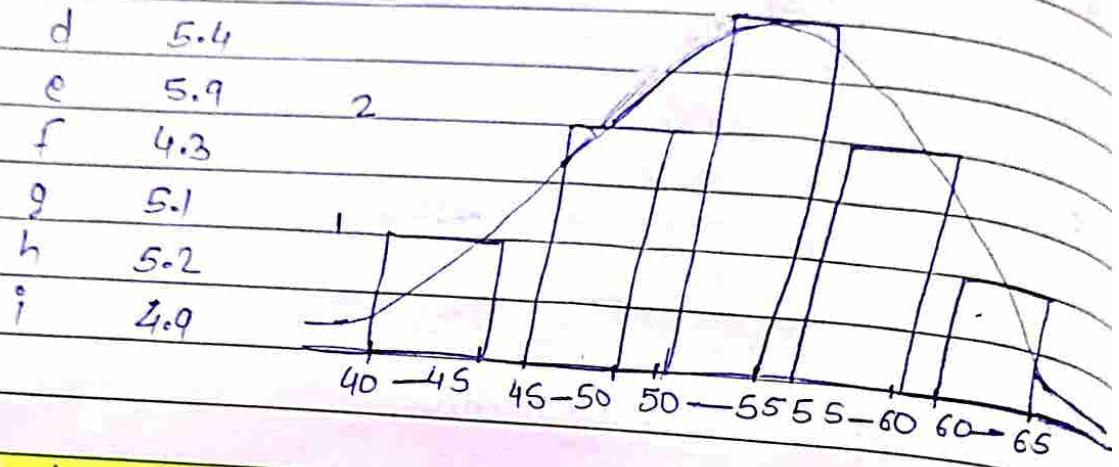
normal
distribution.

when data follows symmetric distribution *normally*
then means median are same.

| | |
|----------|--|
| Page No. | |
| Date | |

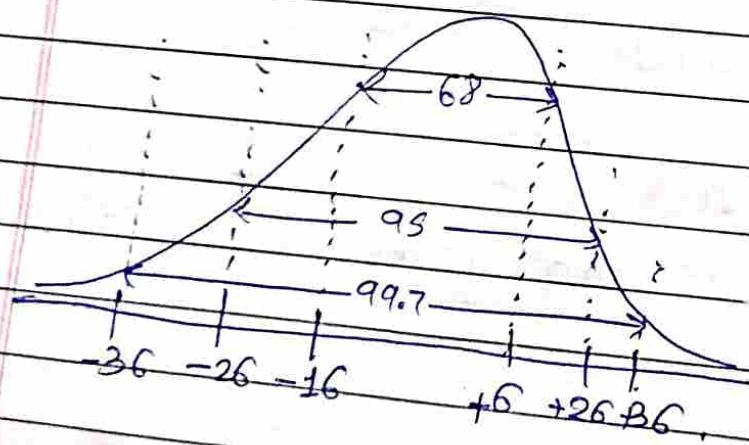
Name height

| | | |
|---|-----|-------|
| a | 6.2 | count |
| b | 5.7 | |
| c | 4.6 | 3 |
| d | 5.4 | |
| e | 5.9 | 2 |
| f | 4.3 | |
| g | 5.1 | 1 |
| h | 5.2 | |
| i | 4.9 | |



* What formula we use to remove outliers with normal distribution.

68-95-99.7 RULE



* normal distribution (file)

54.2631 63.50 66.36 69.17 73.99
min Q1 mean Q3 max

* To eliminate some outliers.

1) Z score.

$$\text{Z score} = \frac{x - \text{mean}(\mu)}{\text{standard deviation}(\sigma)}$$

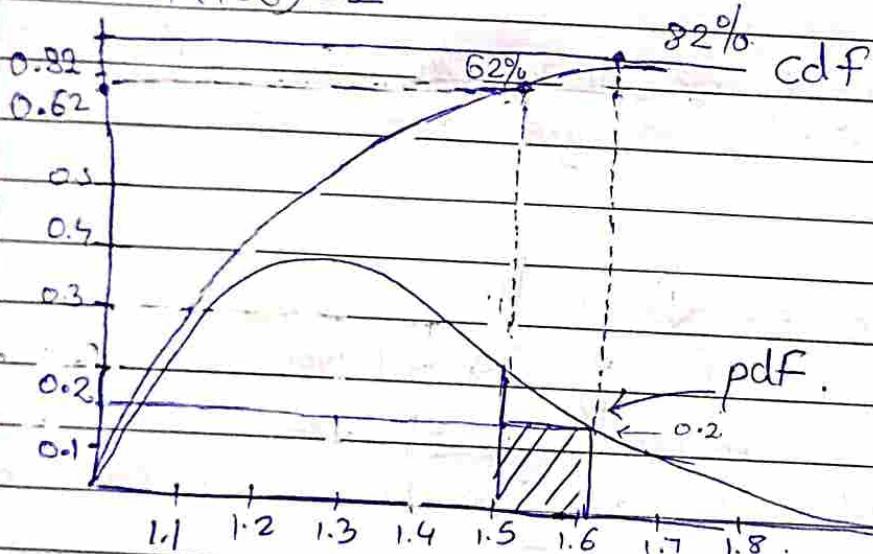
* Cumulative distribution (CDF)

Cumulative distribution function (CDF) for Random variable (x) is denoted by.

$$F(x) = P(X \leq x)$$

$$f(-\infty) = 0$$

$$f(+\infty) = 1$$



20

pdf says. 20% of flower having petal length 1.5 to 1.6.

- 82% of setosa flower having petal length ≤ 1.6
- There are 62% of setosa flower having petal length ≤ 1.5
- 30% of setosa flower having petal length having length ≤ 1.1

cdf calculation.

- 1) There are total 50 setosa flowers & there are total 42 flowers having length = < 1.6 .

$$\frac{42}{50} = 0.82\%$$

- 2) You have histogram of setosa flowers cdf value of Corresponding of 1.6 will be cumulative sum of 1.1, 1.2, 1.3, 1.4, 1.5, 1.6 histogram.

- * If We want to find the count of variable in particular range then we use pdf or cdf

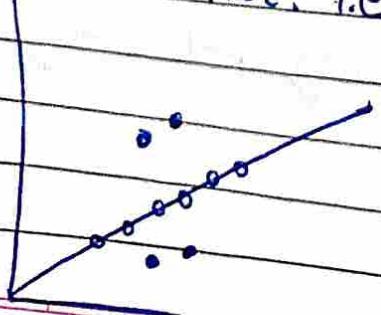
Variable $\sim N(\mu, \sigma)$

\downarrow mean \rightarrow std deviation.

If Variable is Gaussian distributed then we can find CDF & PDF.

* Summary:-

To understand any dataset we should have Q-Q plot of that dataset. from Q-Q plot we can check whether the data is normally distributed or not. i.e



- Normally distributed - main points lie on the s line only few of them lie outside line.

Find cdf & pdf from Q-Q plot

* Z score :- (Standard Normal Distribution)

$$Z = \frac{x - \mu}{\sigma}$$

How many standard deviation away a datapoint is from mean.

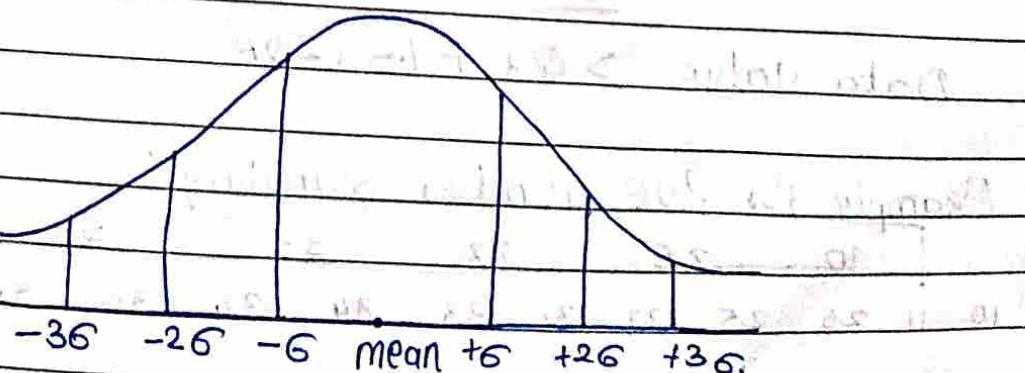
| Name | Height. | $Z \text{ score} = (x - \mu) / \sigma$ |
|------|---------|--|
| A | 6.2 | $(6.2 - 5.25) / 0.61 = 1.53$ |
| B | 5.7 | $(5.7 - 5.25) / 0.61 = 0.72$ |
| C | 4.6 | $(4.6 - 5.25) / 0.61 = 1.06$ |
| D | 5.4 | $(5.4 - 5.25) / 0.61 = 0.23$ |
| E | 5.9 | $(5.9 - 5.25) / 0.61 = 1.04$ |
| F | 4.3 | $(4.3 - 5.25) / 0.61 = 1.55$ |
| G | 5.1 | $(5.1 - 5.25) / 0.61 = 0.25$ |
| H | 5.2 | $(5.2 - 5.25) / 0.61 = 0.09$ |
| I | 4.9 | $(4.9 - 5.25) / 0.61 = 0.58$ |

Average = 5.25

Standard deviation = 0.61

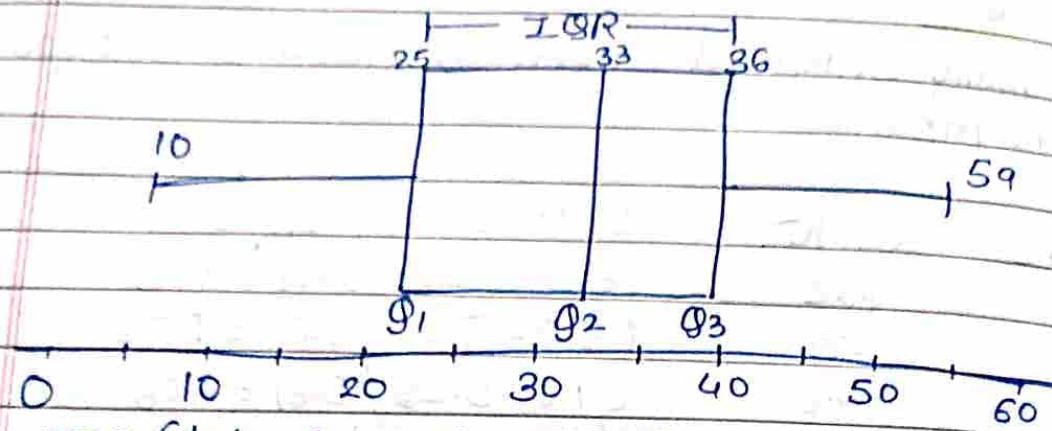
$$Z = \frac{x - \mu}{\sigma}$$

- Z score is also called as Standard Normal Distribution.



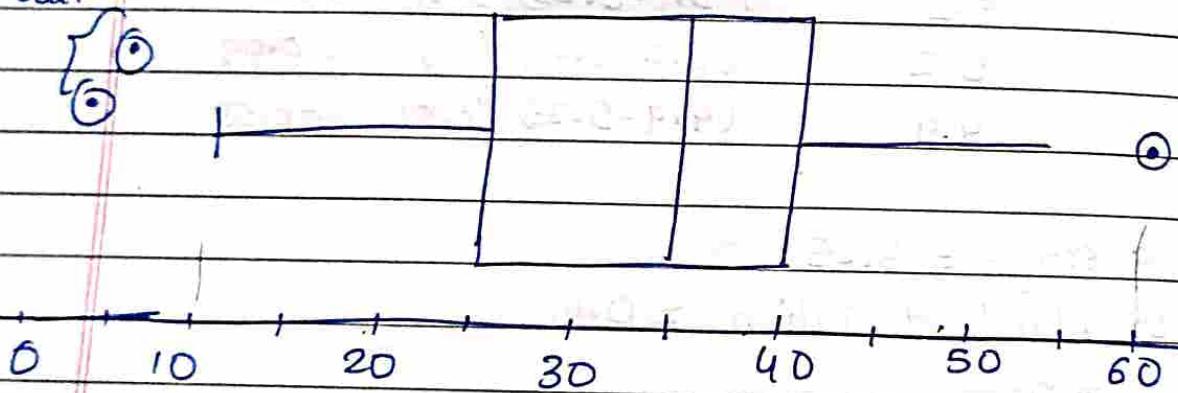
* Boxplot :

Gives us a visual representation of the number summary.



$$\text{IQR (Interquartile Range)} = Q_3 - Q_1$$

* Modified Boxplot
outliers.



- A data value is considered to be an outlier if
Data Value $< Q_1 - 1.5 \times \text{IQR}$
OR

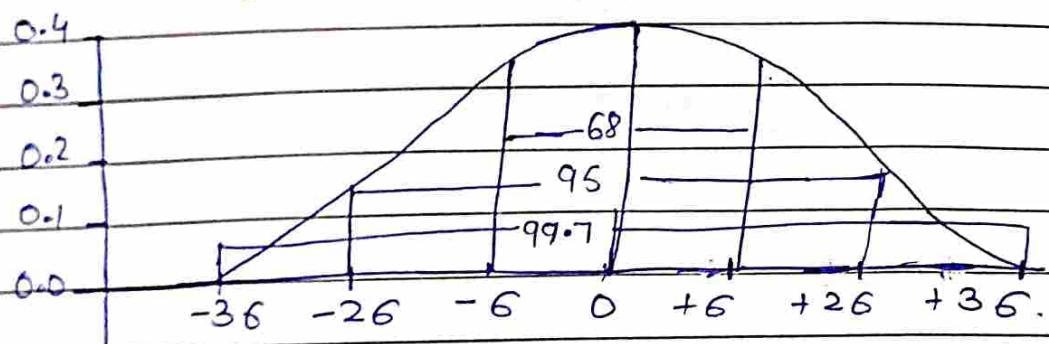
$$\text{Data Value} > Q_3 + 1.5 \times \text{IQR}$$

* Example for Five number summary :

| | 10 | 25 | 33 | 36 | 59 |
|----|----|----|----|----|----|
| 10 | 11 | 25 | 27 | 31 | 33 |

* Chebychev's Inequality (Contd). Ex?

We know that if data is normally distributed we have 68, 95, 99.7 Rule i.e. 68% of data points lie in the Range -1σ to $+1\sigma$, 95% of data points lie in the Range -2σ to $+2\sigma$ & Similarly 99.7% data points lie in the Range -3σ to $+3\sigma$ as shown in the Fig.



Q Let us assume height of students data follows normal distribution.

here it follows normal distribution given that $\mu = 6$

given that σ is 10cm & σ is 10cm. Then probability of

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 95\%$$

$$\therefore P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 95\%.$$

[Ans]

This is possible bcz we know that data is normally distributed. What if we don't know the distribution of data but we know μ & σ where μ is finite & σ is nonzero & finite. can we find % of data betw the range

If we don't know distribution, like normal distribution we did?

In this case chebyshev's inequality equation will help to solve this problem.

Chebyshev inequality equation is given by

at

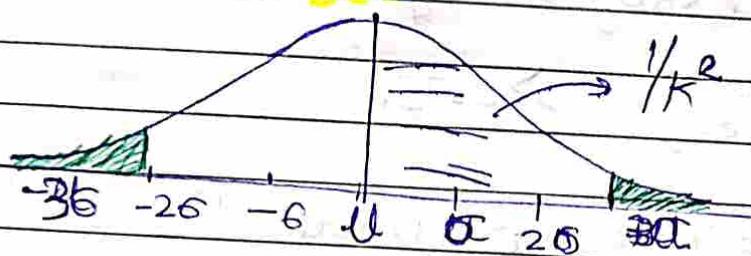
$$P(|X - \mu| \geq k\sigma) \leq 1/k^2$$

where $k = -1, -2, -3, 1, 2, 3$

This eqn can be further simplified

$$X \geq \mu + k\sigma$$

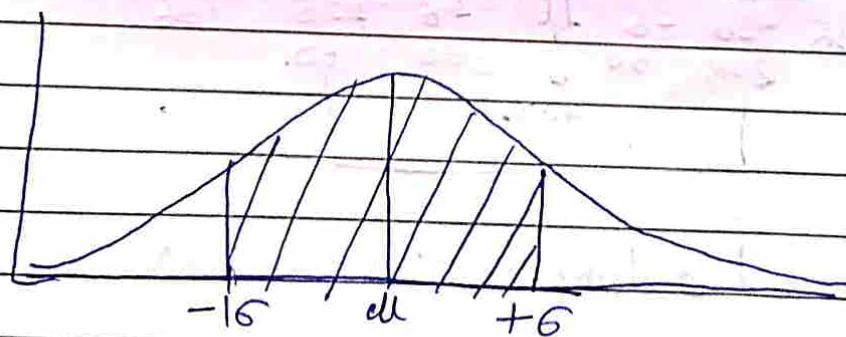
$$X \leq \mu - k\sigma$$



If $k=2$ the data points will be as shown in fig

2nd eqn of chebyshev eqn :-

$$P(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - 1/k^2$$



Don't Know distribution

Q. $\mu = 40, \sigma = 10k$ \rightarrow salary betⁿ 20k, 60k?

$$20k = \mu - 2\sigma$$

$$40k \text{ or}$$

$$60k = \mu + 2\sigma$$

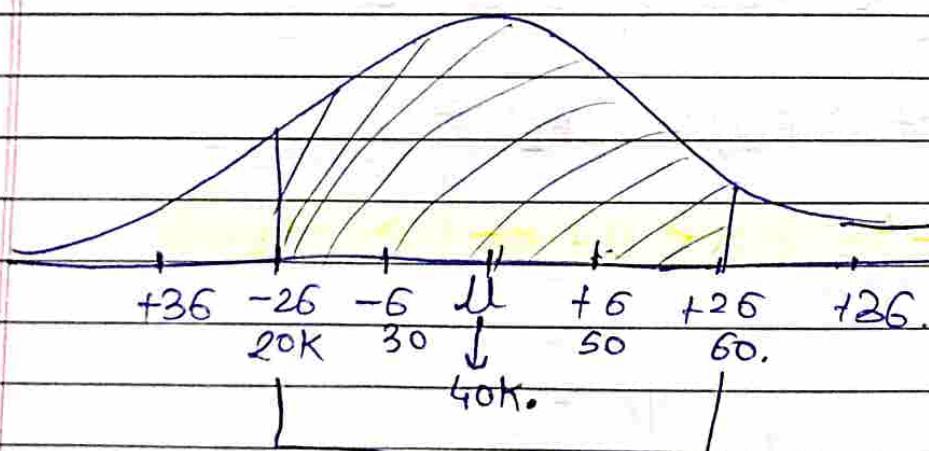
$$P(40-2\sigma < X < 40+2\sigma) > 1 - \frac{1}{K^2} \rightarrow \text{2nd eq chp}$$

$$P(40-2\sigma < X < 40+2\sigma) > 1 - \frac{1}{22}$$

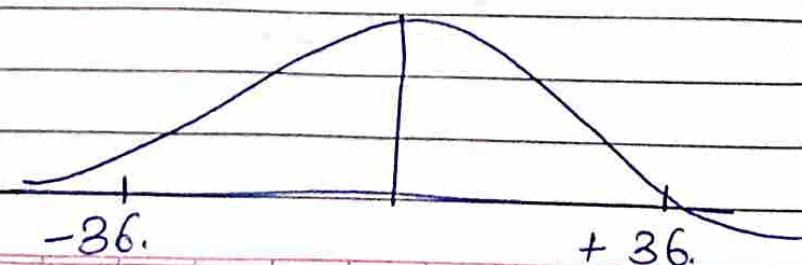
$$P(20 < X < 60) > 1 - \frac{1}{4}$$

$$P(20 < X < 60) > 0.75.$$

i.e. ~~75%~~ 75% data point lies betⁿ 20k-60k



Q. If we want salary betⁿ 10k to 70k.



$$P(10 - 3\sigma < X < 10 + 3\sigma) \geq 1 - \frac{1}{3^2}$$

$$P(40 - 3\sigma < X < 40 + 3\sigma) \geq 1 - \frac{1}{9}$$

$$P(10 < X < 70) \geq 0.89$$

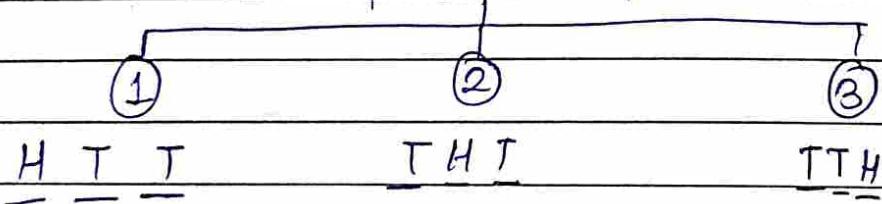
89%

BINOMIAL PROBABILITY DISTRIBUTION

- The number of trials (n), must be fixed.
- There are only two possible outcomes for each trial:
 - success
 - Failure
- The probability (P) of success must be constant for every trial
- Each trial must be Independent

Q If you flip a regular coin three times. What is probability of getting exactly one head and this a binomial experiment?

3 possible ways

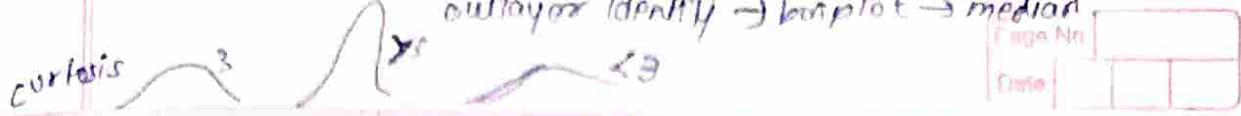


$$P(HTT) = 0.125 \quad P(THT) = 0.125 \quad P(TTH) = 0.125$$

$$\begin{aligned}P(\text{exactly one head}) &= 0.125 \times 0.125 \times 0.125 \\&= 0.375\end{aligned}$$

- 1) The number of trials (n) must be fixed.
- 2) There are 2 possible outcomes for each trial
 - Success = Heads
 - Failure = Tails.

- 3) The probability (p) of success must be constant for every trial
 $P(E) = 0.5$
- 4) Each trial must be independent
- 5) Binomial Experiment



* 1st movement business decision.

- ① calculate mean value, Identifying outliers
It measures ~~se~~ central tendency of data!

* 2nd movement business decision

- ② finding Variance & standard deviation.
spread of data.

* 3rd movement business decision.

- ③ finding skewness. ~~skew~~ (measure of asymetric in the distribution).

* 4th movement business decision.

Finding curtosis (measuring flatness of distribution)

Curtosis of normal distribution should be 3.

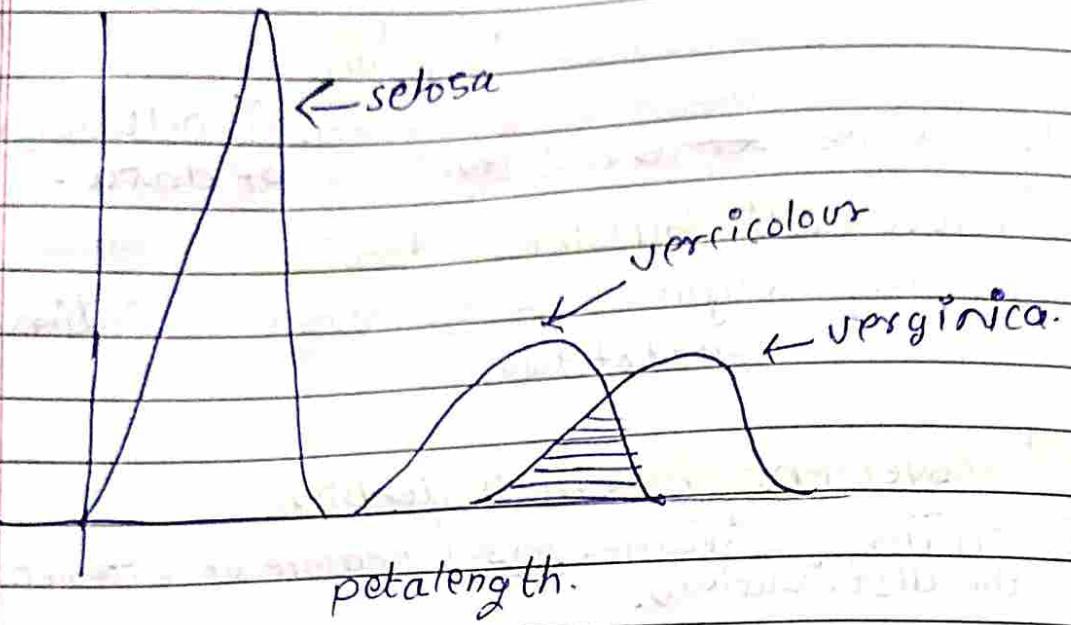
If given distribution of curtosis ≤ 3 then it is said platykurtic which means it reproduces less outliers

If given distribution as curtosis > 3 it is called leptokurtic. It means that it produces more outliers.

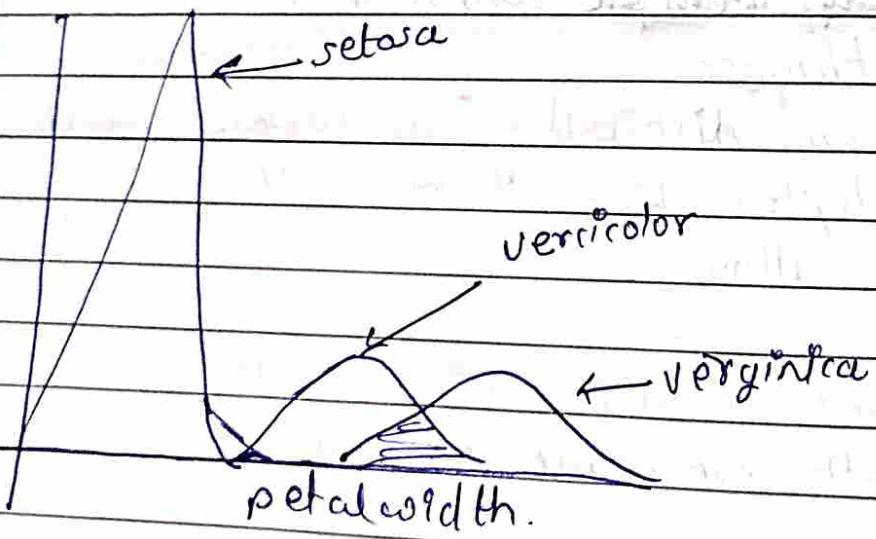
* Univariate analysis using pdf

- ④ Why we want to study univariate analysis?

In univariate analysis using pdf one can understand which feature quantitatively or which features are linearly separable or which features are not linearly separable.



In above diagram we can ~~conclude~~ conclude that setosa flower are linearly ~~separable~~ separable where as vericolor & vergina are not ~~be~~ linearly separable.

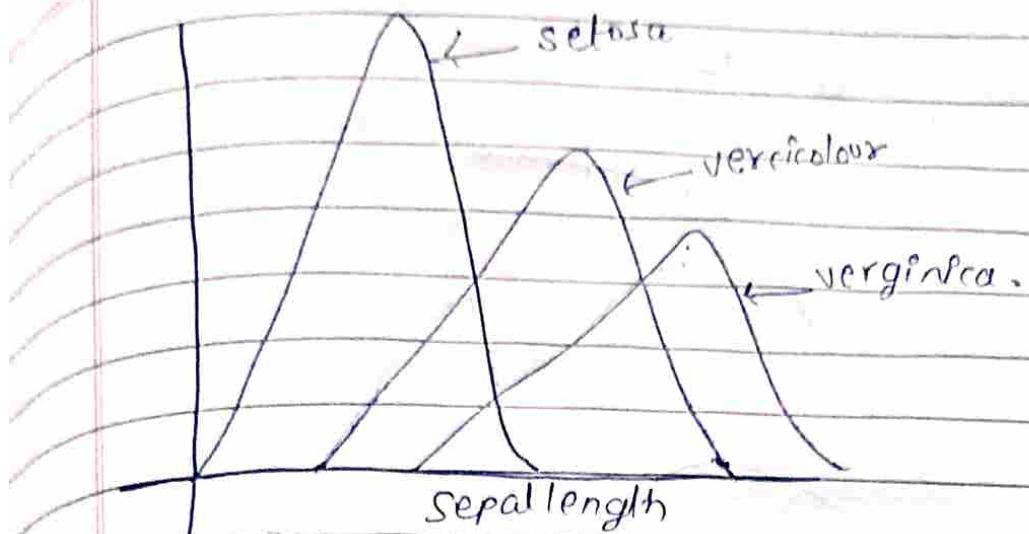


In above diagram we can see that setosa is linearly separable were as ~~vercolor~~ vericolor

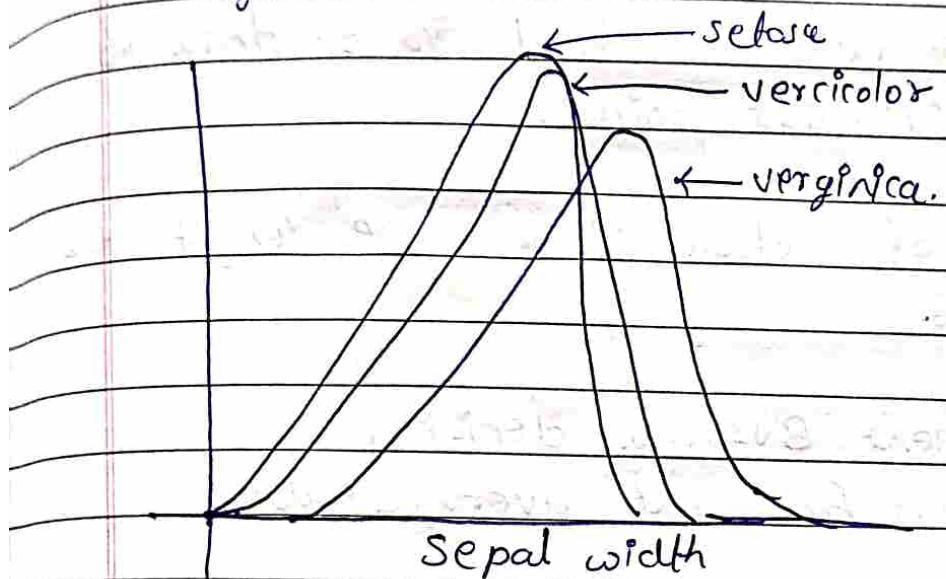
Compare to petal length, petal width is less linearly separable, however

Cdf gives % of that point
pdf gives point on that point.

Page No. _____
Date _____



here all 3 Flower are not linearly
separable. If you are going to use Sepal
length.



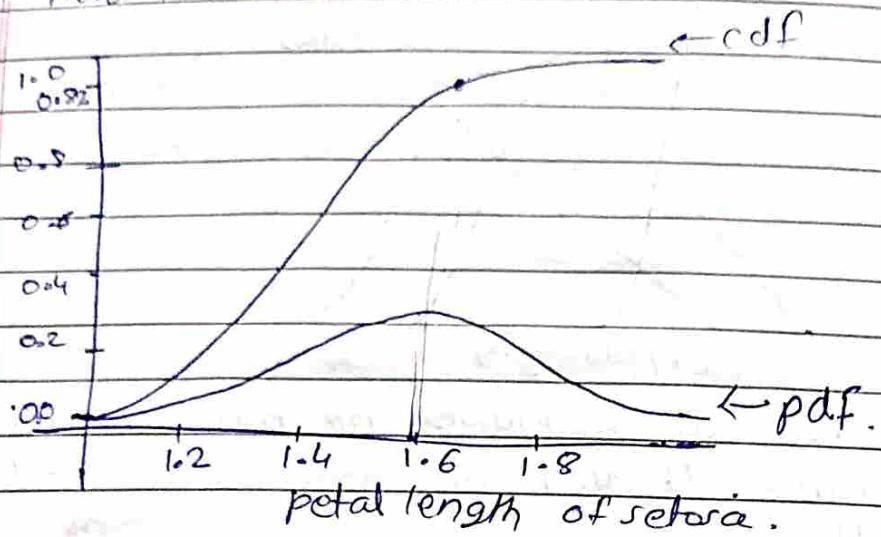
here all 3 Flower are not linearly ~~sep~~
separable. If you are going to use sepal length

* plotting histogram & pdf.

* plotting of cdf

from pdf we can visually see that
what percentage of setosa/virginica/vicolor
Flower have petal length less than some
value

* how to read cdf



We can find % of data at any particular point.

82% of setosa flower have length less than 1.6.

* 1st movement Business decision

- we can find out average value.

* 2nd movement BD,

Instead of variation we will used standard deviation bcoz $\text{Var} = \frac{\sum (x - \bar{x})^2}{N}$ where

$(x - \bar{x})^2$ gives the higher value \therefore we use

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

* equation of Variance for population (complete dataset) & Sample statistic,

$$\text{Var} \sigma^2 = \frac{\sum (x - \mu)^2}{N} \quad (\text{population})$$

parameters.

~~s~~

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

$$\text{Var } s^2 = \frac{\sum (x - \bar{x})^2}{n-1} \quad \text{Sample statistic.}$$

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

* Bar plots : (Univariate analysis),

- data less - Bar plot

- data more - histogram, - (frequency distribution table)

* 5 Number summary,

minimum, 1st Quartile, median, 3rd Quartile, max
(F. describes).

* Box plot

- 1) Understanding distribution

- 2) Finding outliers

Data Preparation
sweetviz, autoviz, dftale, dataprep.

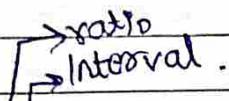
| | |
|----------|--|
| Page No. | |
| Date | |

Data Analytics, analyse past data to make conclusion.

stages of Analytics.

- prescriptive analysis
- predictive analysis
- diagnostic analysis
- descriptive analysis

* Data types



- 1) Continuous
- 2) discrete
- 3) Qualitative & Quantitative
- 4) Structure, semi, non
- 5) Big Data, not big data
- 6) Cross sectional vs Time series vs longitudinal / panel data
- 7) Balanced / Imbalanced Data
- 8) offline / live data.

* Business problem

- " object
- " success criteria
- ML "
- economics "

* Data preparation.

Data cleaning/
preparation/man
-aging/organizing

1. Type casting
2. Handling Duplicates
3. Outlier treatment
4. Zero & near zero variable
Features
5. Discretization / Grouping
6. Dummy variable creation
7. Missing value
8. Transformation
9. Feature scaling/
feature shrinking
10. string manipulation
(unstructure Textual Data)

$$ZQR = Q3 - Q1 \quad \text{then outlier}$$

$$< Q1 - 1.5(ZQR) \text{ or } > Q3 + 1.5(ZQR)$$

| | |
|----------|--|
| Page No. | |
| Date | |

Outliers treatment.

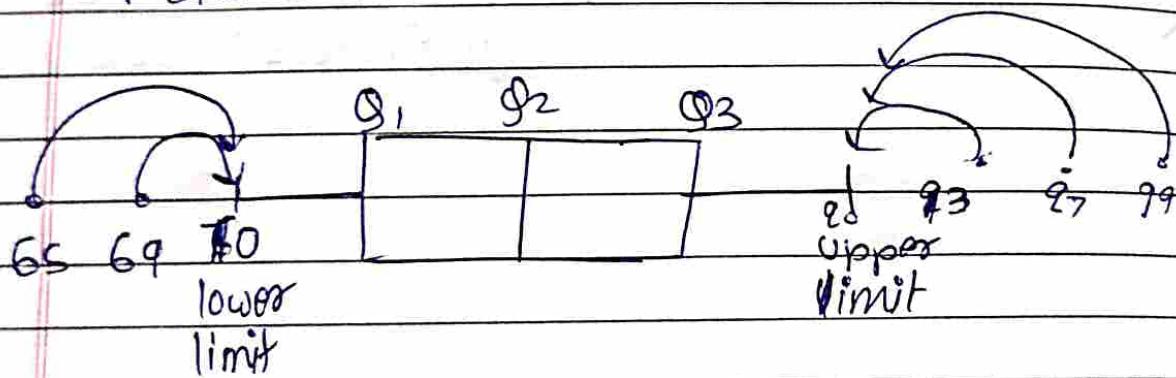
3R Rule

- Replace
- Retain
- Remove
- Trimming: Removal
- Winsorization: Rounding up / masking

Reify

correcting the data / asking data entry operator to check.

Retain



This process is called masking in which we convert outlier into its respective lower or upper limit.

$$\text{lower limit} = Q_1 - 1.5 \times iqr.$$

$$UL = Q_3 + 1.5 \times iqr.$$

pip install feature_engine

In Trimming we remove the outlier, so instead of that we use replace.

* Winsorizer

```
from feature_engine.outliers import Winsorizer  
wWinsor = Winsorizer(capping_method='iqr',  
tail="both",  
fold=1.5,
```

variables ["Salaries"])

```
df_t = wWinsor.fit_transform(df[["Salaries"]])  
sns.boxplot(df_t[["Salaries"]])
```

and min & max of outliers are
i.e. tail value of boxplot is
min & max of boxplot is
which is at 1.5 fold of iqr, i.e.

Final code is

W = Winsorizer(tail="both")

: output

df_t = W.fit_transform(df)

df_t[["Salaries"]]

df_t[["Salaries"]]

df_t[["Salaries"]]

min & max

outliers removed

min & max

outliers removed

min & max

outliers removed

* Discretization (otherwise known as quantization or binning) provides a way to partition continuous features into discrete value.

* necessary to have binning?

Whenever data is fall in small quantity it is very easy to visualize & overall computational complexity is less. however if data quantity is less then it is very difficult to visualize using bar graph. then It is suggested to divide data into quantums/ bins so as to make it visualable and apply to ML model.

* What are bins?

Count 0 1 1 2 5 10 11

Notation:

$$[0,2] \rightarrow 0, 1, 1, 2$$

$$(0,2) \rightarrow 1, 1$$

$$(0,2] \rightarrow 1, 1, 2$$

$$[0,2) \rightarrow 0, 1, 1.$$

* Three Strategies use for binning for discretization

- 1) Uniform strategy
- 2) Quantile Strategy
- 3) Kmean Strategy.

Data Transformation

standardization & normalization

- Why it is necessary.

| x_1 | x_2 | y |
|-------|-------|------|
| 6 | 160 | 21 |
| 6 | 160 | 21 |
| 4 | 108 | 22.8 |
| 6 | 258 | 21.4 |
| 8 | 360 | 18.7 |

When we are applying data to ML model all features should have appropriate mean value & standard deviation if there is drastic variation of x_1 & x_2 across all features then accuracy will be deteriorated

- If data is continuous then distribution is normal and having following properties
 - 1) Bell shape curve will be symmetric
 - 2) Range will be $-\infty$ to $+\infty$
 - 3) Area under curve is 1
 - 4) probability of $x \in (-\infty, -\delta)$ will be equal to zero

Transformation Techniques

1) Normalization] popular

2) Standardization

3) Quantile Transformation

4) log Transformation

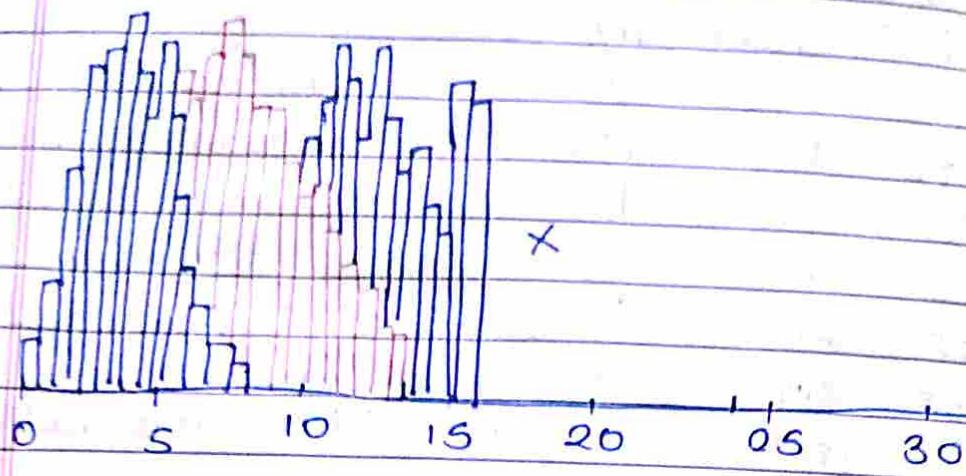
5) binarization

6) Power Transformation

7) Unit Vector scaling

necessary of Transformation

Faster conversion & computing distance appropriately



min-max Normalization $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$

Standardization $x' = \frac{x - \mu}{\sigma}$

- If we apply standard normalization of skewed data Resultant will be skewed
- In standard normal distribution $\mu=0$ and $\sigma=1$.

Winsorizer

from feature_engine.outliers import Winsorizer

winsor = Winsorizer(capping_method='iqr',
tail='both',
fold=1.5)

variables = ["Salaries"]
df_t = winsor.fit_transform(df[["Salaries"]])
sns.boxplot(df_t[["Salaries"]])