

Univariate charts

Number: [distribution of number column]

list Plot, distplot, hist, box plot

category: [distribution of category column]

count-plot, bar chart, frequency Plot

Bivariate charts

Number v/s Number: [relative distribution in 2 column]

Scatter Plot, reg Plot

Number v/s category

Box Plot, Strip, Swarm, violin

category v/s category [Cross Tab.]

Grouped bar chart

Bivariate charts with third category (hue)

number v/s number v/s category

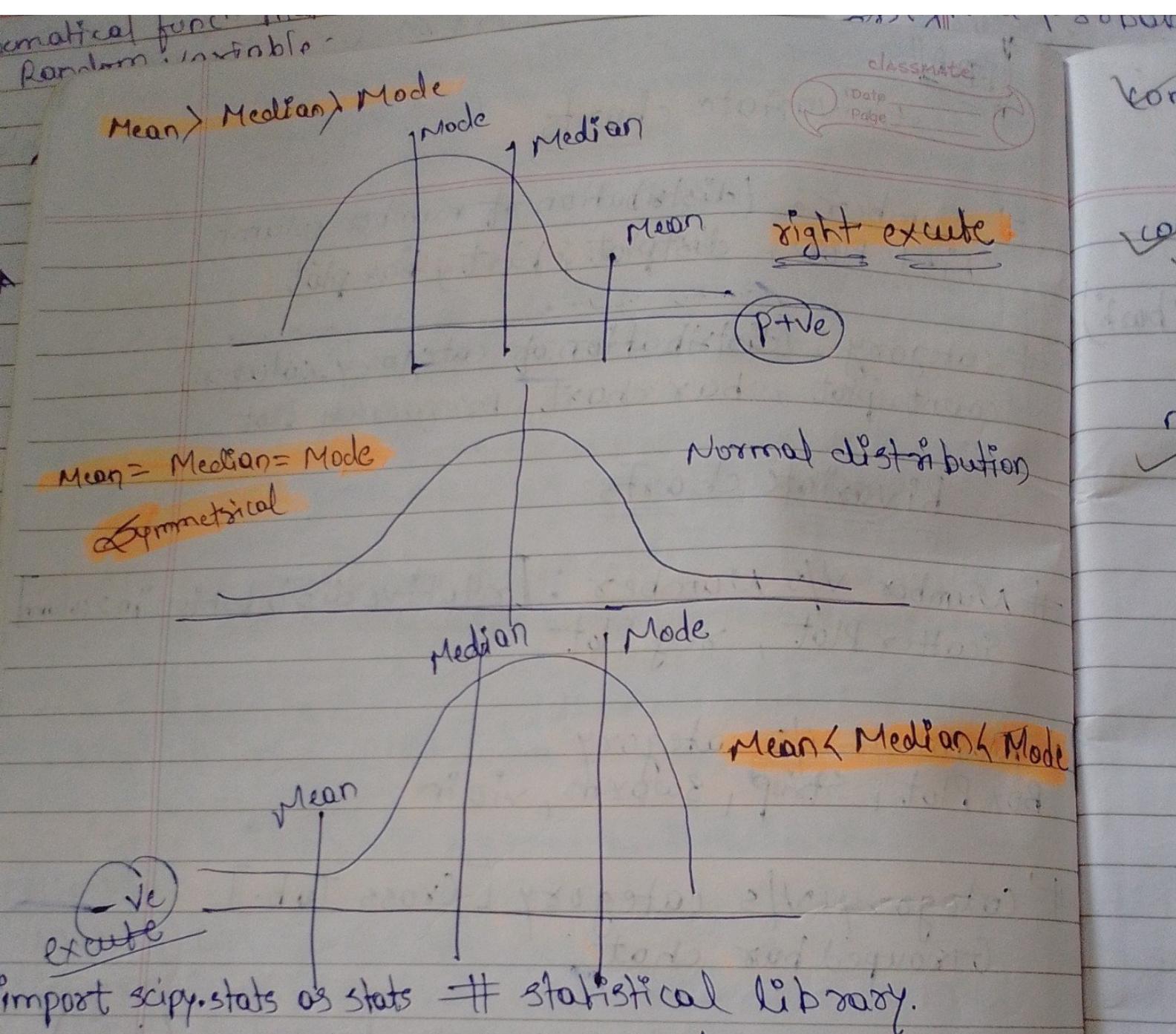
Scatter plot, reg plot, pair plot

category
number v/s number v/s category

Box plot, strip, swarm, violin.

Series Data

◎ Shot on AWESOME A05s



```
import scipy.stats as stats # statistical library.
```

```
import pandas as pd
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns.
```

```
test = pd.read_csv('test-AbST221.csv')
```

```
train = pd.read_csv('train-vgrgXOR.csv')
```

```
sklearn.preprocessing import StandardScaler, MinMaxScaler
```

Shot on AWESOME A05s

a probability to eat

combined = pd.concat([train, test], ignore_index=True)

→ दो data को combine कर दिया

combined.select_dtypes(include=np.number).columns
→ find all numerical variables

Univariate

Variables

num_cols = ['Item-Weight', 'Item-Visibility', 'Item-MRP',
'Item-Outlet-Sales']

nrows = 2 # no. of rows in the plot

ncols = 2 # no. of columns

iterator = 1 # Plot iterator

plt.figure(figsize=(10, 6))

for i in num_cols:

plt.subplot(nrows, ncols, iterator)

sns.distplot(combined.loc[:, i])

iterator = iterator + 1

plt.tight_layout()

plt.show()

Category

combined.select_dtypes(include='object').columns

cat_cols = ['Item-Fat-Content', 'Item-Type',
'Outlet-Identifier', 'Outlet-Size',
'Outlet-Location-Type', 'Outlet-Type']

Shot on AWESOME A05s

natural funcn that des
Rand

... funcn
nrows = 3
ncols = 2
repeater = 1

plt.figure(figsize=(10,8))
for i in cat_cols:
 plt.subplot(nrows, ncols, repeater)
 combined.loc[:, i].value_counts().plot(kind='bar')
 repeater += 1
plt.tight_layout()
plt.show()

Bi- Variates

Number V/s Number

nrows = 2

ncols = 2

iterator = 1

plt.figure(figsize=(12,8))

for i in num_cols:

plt.subplot(nrows, ncols, iterator)

plt.scatter(combined.loc[:, i],
 combined.Item_Outlet_Sales)

plt.title(i)

iterator += 1

plt.tight_layout()

plt.show()

✓ nRows = 3

nCols = 2

rep = 1

plt.figure(figsize=(12,10))

for i in cat_cols:

plt.subplot(nRows, nCols, rep)

sns.boxplot(data=combined, x=i, y='Item_Outlet_Sales')

plt.xticks(rotation=90)

rep+=1

plt.tight_layout()

plt.show()

{Category CLASSMATE
Date Page ✓/5 ✓
Number}

{Category vs Category ✓/5}

fix the item fat content
combined.Item_Fat_Content.replace(to_replace = ['low fat', 'reg'], value = ['Low fat', 'Low Fat', 'Regular'], inplace = True)

change the name, names change huge

average sales

print("Mean sales:", mean)

what is the 95% values of sales.

print("95% sales values:", combined.Item_Outlet_Sales.quantile([0.95]))

Identify the product that sells most.

combined.Item_MRP.median()[0]

Mode

Find the middle most element observation for the MRP

combined.Item_MRP.median()

Median

Compare the Mean and Trimmed Mean of Sales

print('mean:', combined.Item_Outlet_Sales.mean())

Mean

trimmed_mean = stats.trim_mean(train.Item_Outlet_Sales, proportion剔除 = 0.01)

print("Trimmed Mean", trimmed_mean)

mean: 2181.28891357032

trimmed Mean: 2138.3762

mean
Trimmed mean

what % Diary Items are in the data

print(combined.loc[combined.Item_Type == 'Diary'].shape[0] / combined.shape[0])

0.079977471.

which product has the highest variation in sales

Standard deviation

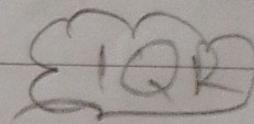
combined.groupby("Item_Type")["Item_Outlet_Sales"].std()

what is the middle 50% of Item NRP

q3 = combined.Item_NRP.quantile(0.75)

q1 = combined.Item_NRP.quantile(0.25)

print('IQR Range:', q3-q1)



find the % of variation in item weight for each product

combined.groupby("Item_Type")["Item_Weight"]

for i in combined.Item_Type.unique():

mean = combined.loc[combined.Item_Type == i, "Item_Weight"].mean()

std = combined.loc[combined.Item_Type == i, "Item_Weight"].~~mean~~.std()

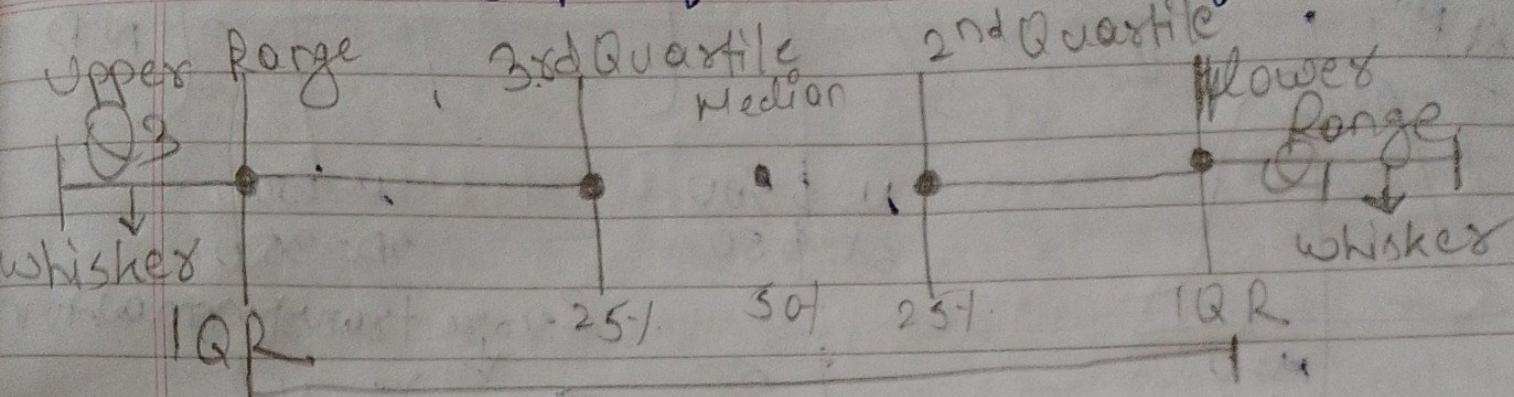
print("(%) for", i, std/mean)

print
sum product has edge.



Shot on AWESOME A05s

$IQR = Q_3 - Q_1$ (Upper quartile (Q_3) - Lower quartile (Q_1))



* data wo column kabhi include nahi kya jata where SD standard deviation = 0 *

$$CV = \frac{\sigma}{\mu}$$

coefficient of variation

जब तक कि सभी जरूरी

$$SD = 0 \text{ हो देता है}$$

ex - 2 players का

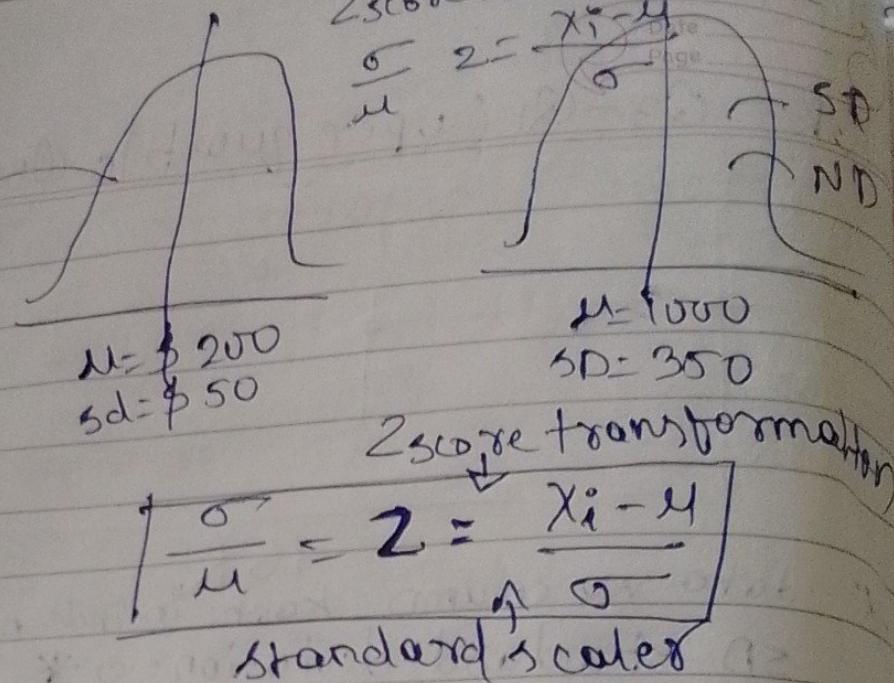
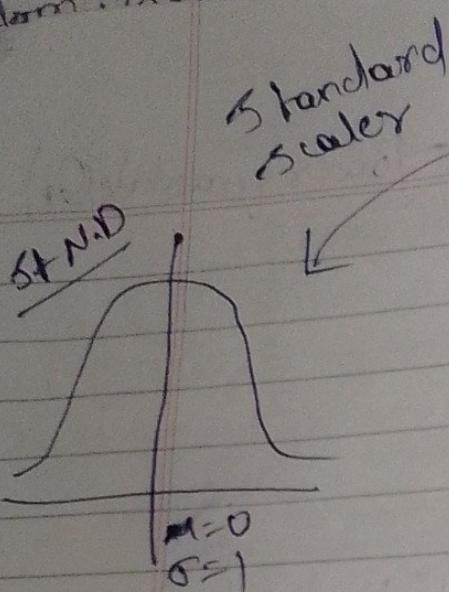
उसका CV कम होगा। यह को
अद्यता player में से उत्तम।
क्योंकि वह stable jada

1	$\frac{1-1}{5-1} = 0$
2	$\frac{2-1}{5-1} = \frac{1}{4} = 0.25$
3	$\frac{3-1}{5-1} = \frac{2}{4} = 0.5$
4	$\frac{4-1}{5-1} = \frac{3}{4}$
5	$\frac{5-1}{5-1} = 1$ $\min = 1$ $\max = 5$

Min Scale =

$$\frac{x_i - \min(x_i)}{\max(x_i) - \min(x_i)}$$

mathematical function that
Random variables



Two Scaling techniques

① Standard Scaler

② Min/Max Scaler

$$\text{Min/Max Scaler} = \frac{x_i - \text{Min}(x_i)}{\text{Max}(x_i) - \text{Min}(x_i)}$$

$$1 \quad \frac{1-1}{5-1} = 0$$

$$2 \quad \frac{2-1}{5-1} = \frac{1}{4} = 0.25$$

$$3 \quad \frac{3-1}{5-1} = \frac{2}{4} = 0.50$$

5

Min/Max Scaler को apply करना उत्तम
Outliers present हैं।

SCALING

classmate

Date

Page

- ⑥ Standard Scaler is one of the popular techniques used for scaling in the Data Science Feature. It involves calculating the mean and std of all the numerical columns and then using Z score to convert the data into standard Normal Distribution where the Mean of all the columns is equal to 0 and SD is 1. Standard Deviation is 1.

⑦ Normalization

- ⑦ This another technique used to scale the data in such a way where the min of each column is equal to 0 and the max of each column is equal to 1.
- ⑧ Since, it takes the range in denominator thus, it is very much influenced by the presence of the outliers in the data. Therefore Min Max Scaler is not a very effective technique in scaling the data.

Note ① Standard Scaler is quite versatile in nature as it can handle the presence of Outliers well

Note ② However, in the presence of outliers in data one must apply ROBUST SCALER to scale the data.

Shot on **AWESOME A05s**

stistical function that describes the random variable.

```
① from sklearn.preprocessing import StandardScaler, MinMaxScaler  
② sc = StandardScaler() # Machine Instance
```

Scale the Data
$$\# \text{sc} = \frac{(x_i - \text{mean})}{\text{std}}$$

.fit() will only learn the mean & std.

.fit_transform() will not only learn the mean & std but also convert/transform

scaled = sc.fit_transform(pd.DataFrame(CombinedItem))
(Error: it is 1D convert it into 2D (Dataframe))

Note: sc takes data in 2D format and hence,
I need to supply the Data Frame.

Apply the standard scaler on the whole Data

```
print(pd.DataFrame(scaled, columns=["ZScore"]).mean)  
print(pd.DataFrame(scaled, columns=["ZScore"]).std)  
ZScore -1.391374e-16  
dtype: float64  
ZScore 1.000035  
dtype: float64
```

Apply Standard Scaler on the whole Data

sc = StandardScaler()

scaled = sc.fit_transform(combined.loc[:, num_cols])

pd.DataFrame(scaled, columns = num_cols).describe()

Min Max Scaler()

mmax = MinMaxScaler()

pd.DataFrame(mmax.fit_transform(combined.loc[:,
output num_cols]), columns = num_cols).min()

Item_Weight 0.0

Item_Visibility 0.0

Item_MRP 0.0

Item_Outlet_Sales 0.0

dtype: float64

def minmax(df, xi):

return ((df[xi] - min(df[xi])) / (max(df[xi]) -
min(df[xi])))

minmax(combined, "Item_MRP").describe()

output count 14204.00000

mean 0.465686

std 0.268529

min 0.00000

25% 0.266244

50% 0.470958

75% 0.856055

1.00000

critical function that describes the probability distribution of the random variable.

TRANSFORMATION

classmate

Date _____

Page _____

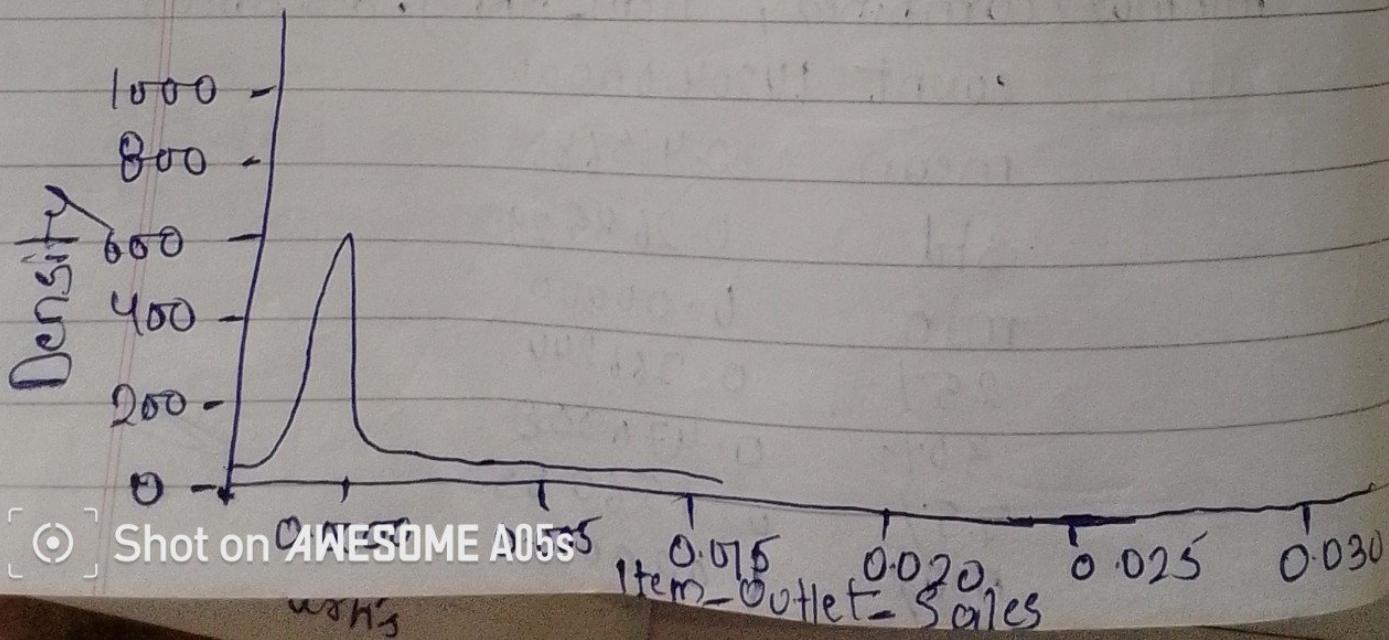
Statistical 1

why Transform?

- ① The Purpose of Transformation is to reduce the SKEWNESS in the data
- ② Most of the statistical models have an assumption about the data that it should be Normally Distributed / Gaussian Curve
- ③ We can apply some of the transformation techniques to reduce the skewness in the data.
- ④ Those Techniques are as follows.
 - ① Log Transformation
 - ② SQRT Transformation
 - ③ Cube Root
 - ④ Reciprocal
 - ⑤ Box - Cox
 - ⑥ Yeo - Johnson.

APPLY RECIPROCAL TRANSFORMATION

```
sns.distplot(np.reciprocal(train.Item_Outlet_Sales))  
print(np.reciprocal(train.Item_Outlet_Sales).skew())
```



Shot on AWESOME A05S

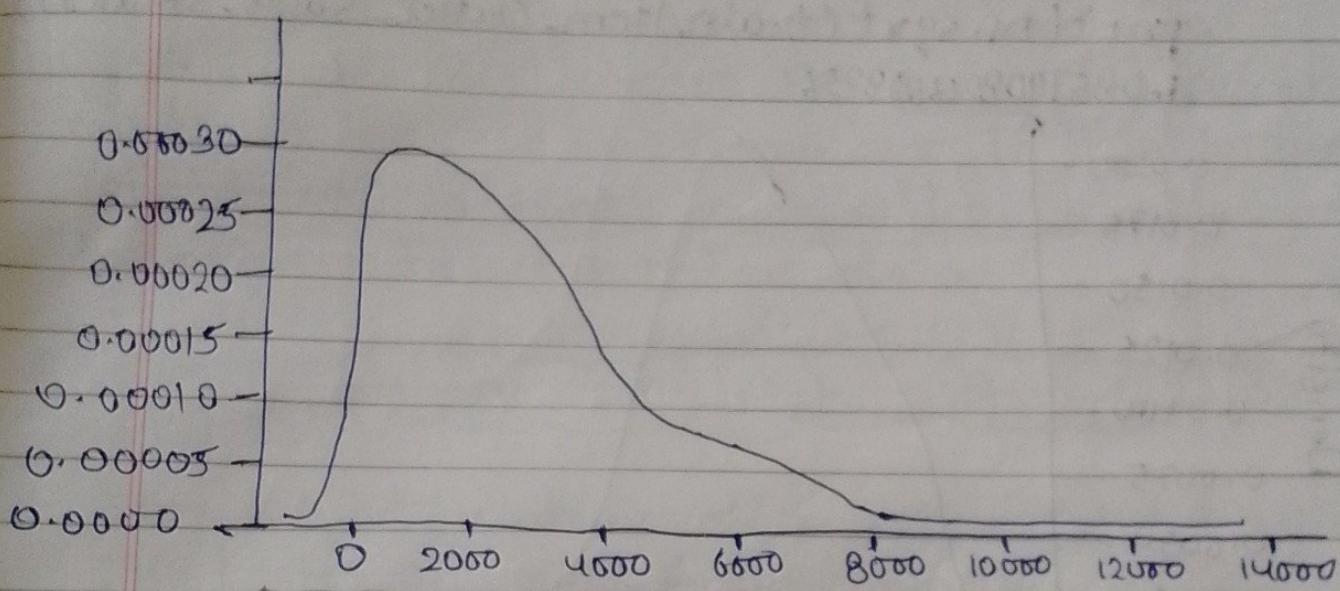
work

mt

(WITHOUT TRANSFORMATION) describes the pd of ⁹ n Dnt

sns.distplot(train.Item_Outlet_Sales)
print(np.log(train.Item_Outlet_Sales).skew())

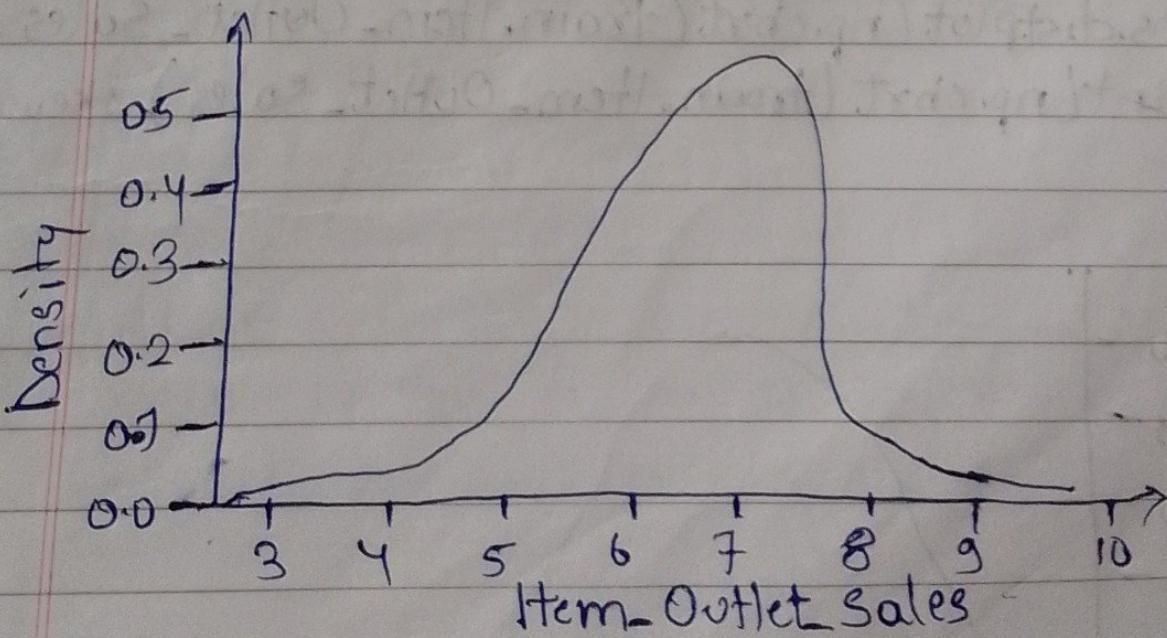
1.775306028542796



APPLY LOG TRANSFORMATION

sns.distplot(np.log(train.Item_Outlet_Sales))
print(np.log(train.Item_Outlet_Sales).skew())

-0.887753343204305



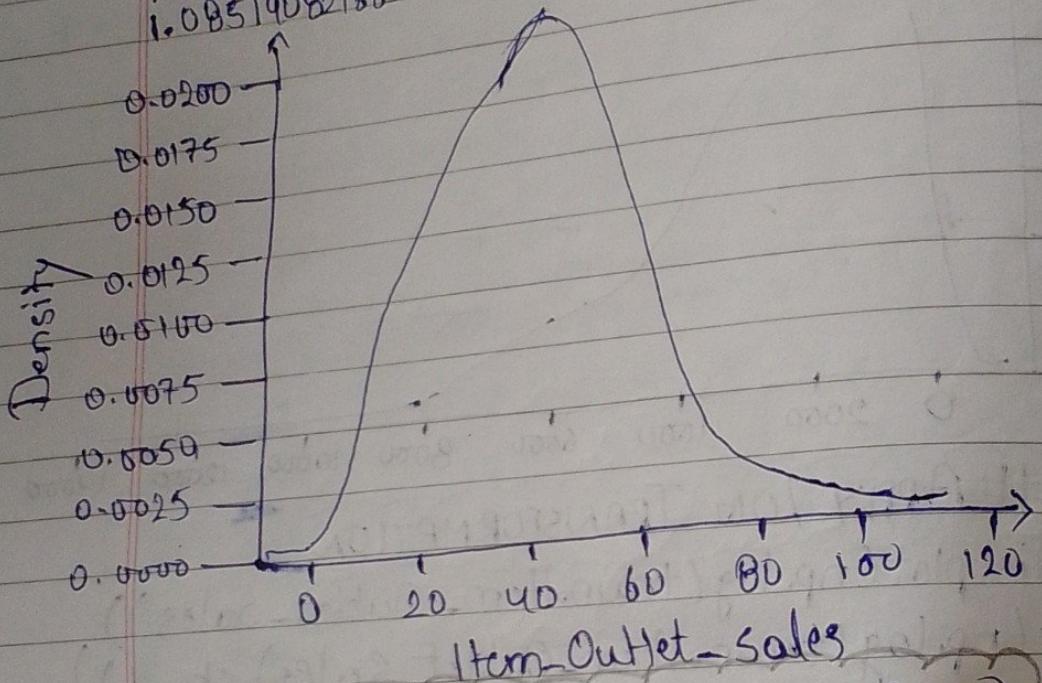
Shot on AWESOME A05s

P-Demo

Date: CLASSMATE
Page:

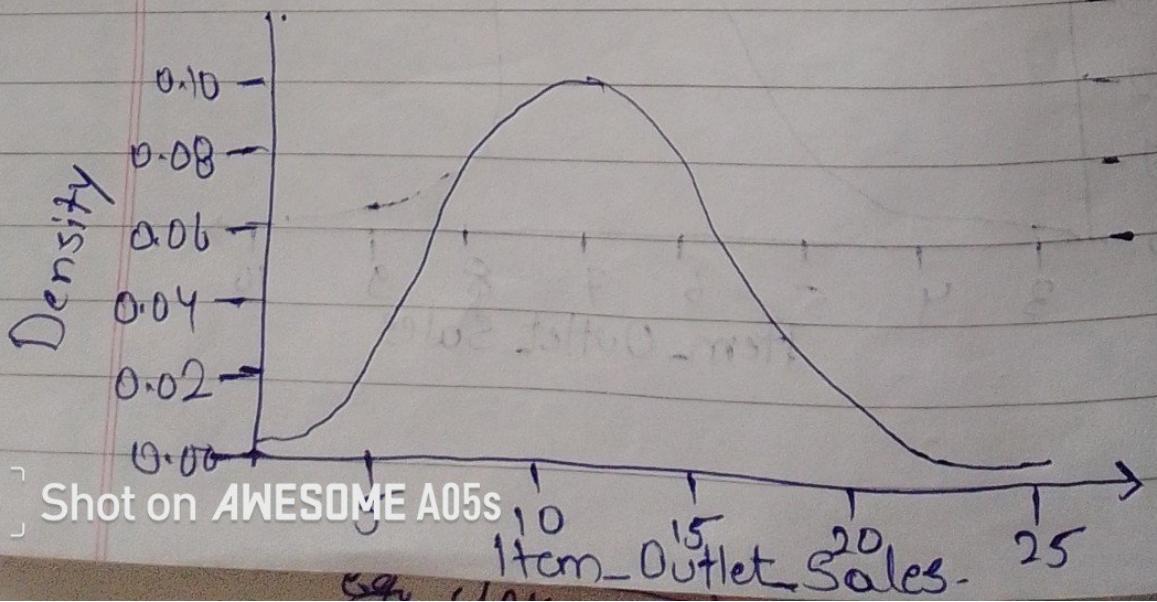
mathematical function but I... # APPLY SQRT TRANSFORMATION

```
sns.distplot(np.sqrt(train.Item_Outlet_Sales))  
print(np.sqrt(train.Item_Outlet_Sales).skew())  
1.0851408216698326
```



APPLY CUBEROOT TRANSFORMATION

```
sns.distplot(np.cbrt(train.Item_Outlet_Sales))  
print(np.cbrt(train.Item_Outlet_Sales).skew())
```



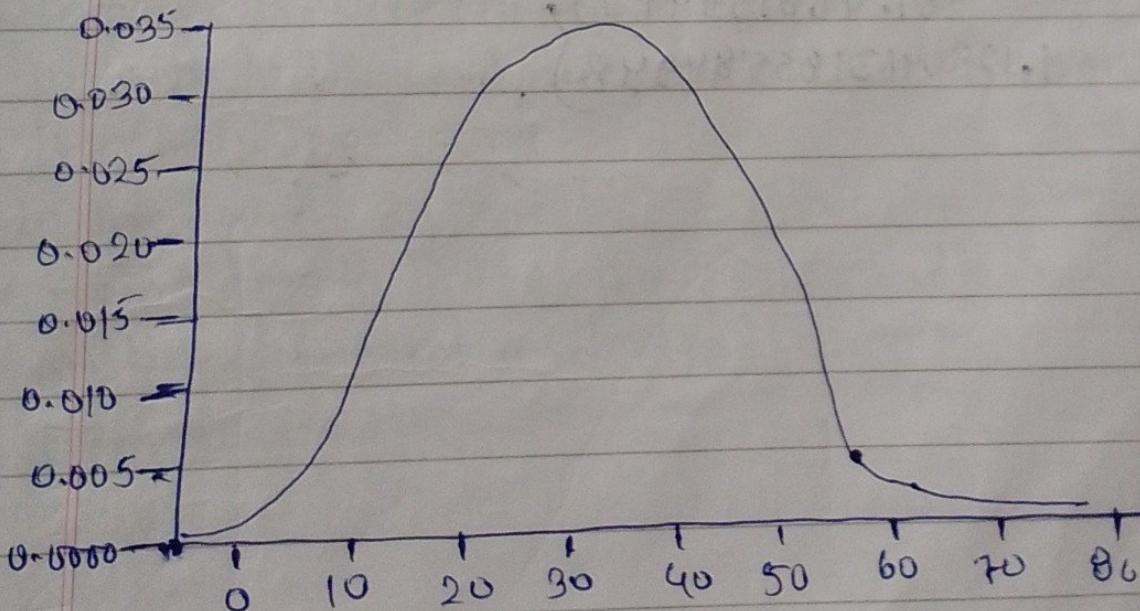
Shot on AWESOME A05s

~~a mathematical function that describes the relationship between Box Cox Transformation & Yeo-Johnson Transformation~~

- ① Box-Cox Transformation works on transforming the Positive Values (greater than 0) only as the function entails usage of Log Transformation as per the formula if the function encounters 0 in the data.
- ② Whereas as the Yeo-Johnson Function can work on any number be it 0, -ve or +ve number. In totality, the Yeo Johnson is more flexible and a versatile transformation bet on the data.

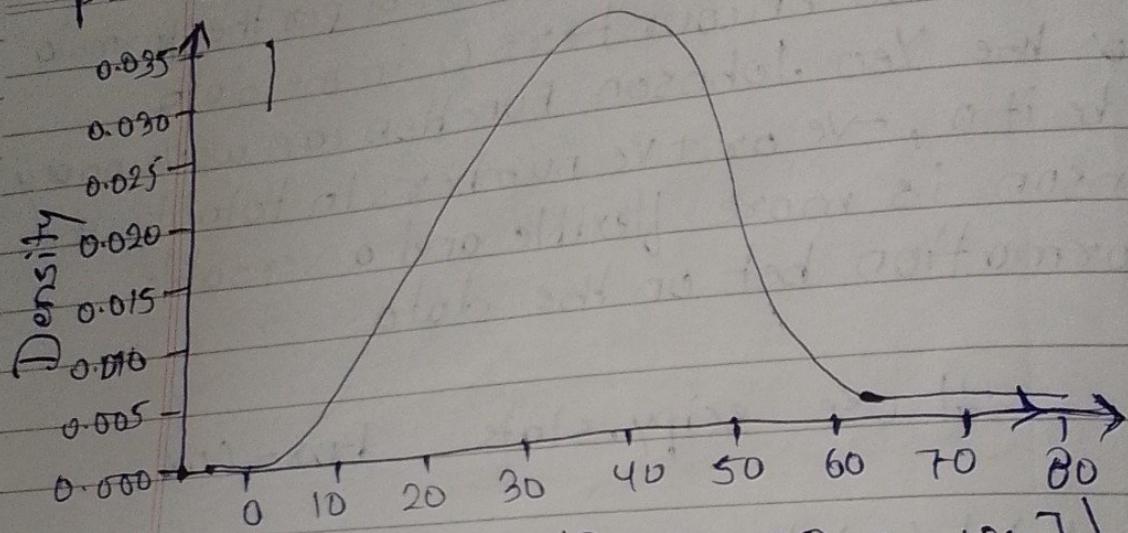
Use `scipy.stats.Box-Cox`

```
sns.distplot(stats.boxcox(train.Item_Outlet_Sales)[0])  
print(stats.skew(stats.boxcox(train.Item_Outlet_Sales)[0]))  
-0.074964223767276
```



a mathematical function now)
Use Yeo-Johnson

sns.distplot(stats.yeojohnson(train.Item_Outlet_Sales)[0])
print(stats.skew(stats.yeojohnson(train.Item_Outlet_Sales)[0]))



stats.yeojohnson([23, 54, 98, 0, -134])

output

(array([31.07413761, 80.56195565, 157.18079993,
-81.46815947]),
1.1280463185584948).

MISSING VALUES

- ⑨ Rule says if 95% or more data is missing in the column, then it is best to drop that column.

Notes This is a theoretical rule. The point is to find the pattern in the data in such a way that you can handle most of the missing values in the data.

- ⑩ Now, basis discussion, we see that prima facie we can consider the mode of the categorical variable & Median in Numerical variable to impute the
- ⑪ However, this approach might ^{not} work because we have not explore the pattern from the rest of the columns' feature.
- ⑫ Therefore, we need to do the pattern exploration on priority to find out the closest possible pattern for the missing value and then impute it accordingly.
- ⑬ In the event of the pattern not being available or unable to mine the pattern, it's always good to build the model and deal with the missing values accordingly.

|| FEATURE ENGINEERING ||

- ① Bin the Item Identifier and Item Type

Note: The purpose of creating new features and binning is to make sure that the machine is able to capture the pattern well. This concept is known as **generalization**.

	Item-Identifier	Item-Weight	Item-Fat-Content	Item-Visibility
0	FDA15	9.30	Low-Fat	0.016047
1	DRC01	5.92	Regular	0.1919278

Item-Type	Item-MRP	Outlet-Identifier	Outlet-Establishment-Year
Diary	249.8092	OUT049	1999
Soft	48.2692	OUT018	2009
Drinks			
Outlet-Size	Outlet-Location-Type	Outlet-Type	Item-Outlet-Sales
Medium	Tier1	Supermarket	37351380
Medium	Tier3	Supermarket	473.4228

Bins = columns of bins based on food items
 combined-Item-Type-unique = 16 categories
 # Extract the first two letters from Item ID

ids = []

```
for i in combined[Item-Identifier]:
    ids.append(i[:2])
```

```
combined["Item-IDS"] = pd.Series(ids)
```

ids
soybean
dairy

Variables

|| Steps that I Follow ||

- ① Univariate
- ② Bivariate
- ③ Feature Engineering
- ④ Missing Values
- ⑤ Outlier Analysis and Removal = Why?
The step where I split the data in train
and test back again
- ⑥ Scaling & Transformation
- ⑦ Categorical Encoding

Why I split the Data in train and test: Before Outliers ??

- You remove the Outliers from training data only and not from test
- Train : that dataset on which the model will be trained
- Test : is that dataset on which the model will predict.

Outlier Analysis

Here at this stage, we will split the data Train & Test

- * The no of columns / Features in the Test set always $n-1$ where n is the total no. of columns in the Train set.
 - * Train set is the data from where the model will learn the pattern because it has the predictors and the target variable.
 - * Test is the data where we will do the prediction and check the performance of the model.

#Splitting the Data back in train and

train.shape, test.shape

$$((8523, 12), (5681, 12))$$

combined. shape

(14204, 13)

7 original data
split into train

Dough
newtrain = combined.loc[0:train, shape[0]:]
newtest = combined.loc[train, shape[0]:]

newtrain.shape, newtest.shape
(8523, 13) (5681, 13)

(5681, 13)

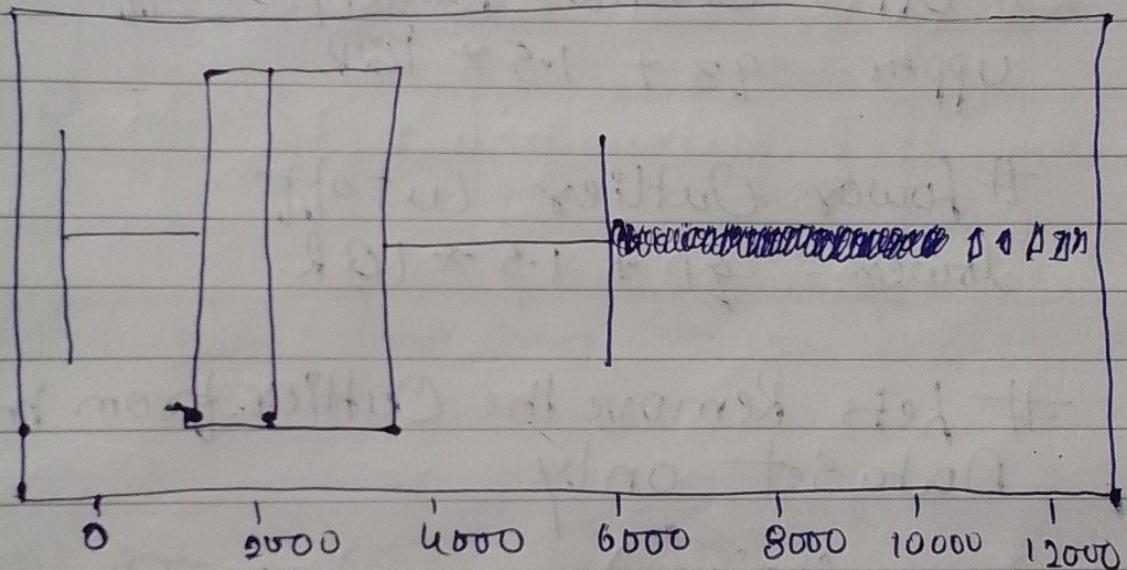
Rule of Identifying the Outlier

Date _____
Page _____

- Any value greater than $(Q_3 + 1.5 * IQR)$ is known as Outlier. This is called Upper Outlier
- Similarly, Any value less than $(Q_1 - 1.5 * IQR)$ is known as Lower Outlier

See the Outlier

sns.boxplot(data=newtrain, x="Item_Outlet_Sales")



$Q_3 = \text{newtrain}.\text{Item_Outlet_Sales}.\text{quantile}(0.75)$

$Q_1 = \text{newtrain}.\text{Item_Outlet_Sales}.\text{quantile}(0.25)$

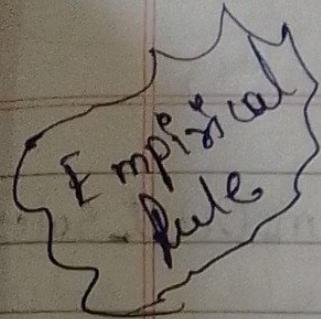
$$IQR = Q_3 - Q_1$$

$$Q_3 + 1.5 * IQR$$

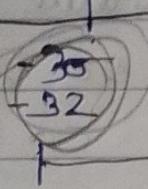
$$6501.8699$$

Any value greater than this is an outlier.

Day 5



Outlier



-22 -15 -12 12 15 22

99.72%

Z score

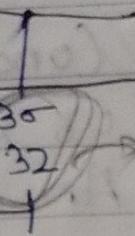
Normal

classmate

$$68.3\% = \pm 1 \text{ S.D}$$

$$95.44\% = \pm 2 \text{ S.D}$$

$$99.7\% = \pm 3 \text{ S.D}$$



Outlier

$$\text{Z score} = \frac{x_i - \mu}{\sigma}$$

st. Normal

Outliers Using Z score

- ④ Generally the Z score is applied on the Target Variable

Empirical Rule

④ 68.3%. Data lies within $\pm 1 \text{ S.D.}$

④ 95.44%. Data lies within $\pm 2 \text{ S.D.}$

④ 99.7%. Data lies within $\pm 3 \text{ S.D.}$

ENCODING

CLASSMATE

Date _____

Page _____

Category → Numeric - binary

- * ONE HOT ENCODING / N-1 ENCODING
- * LABEL ENCODING / ORDINAL ENCODING
- * FREQUENCY ENCODING
- * TARGET ENCODING

Big Mart Sales

House Price
No Broker

Item ID	Outletsize	#B	Area	BSMT	DG	HFX
FD	Small	1	2072	Y	30hn	
NC	Medium	2	1100	N	1m	
DR	High	3	5000	Y	0.5M	
FD		5	3000	N	2mns	
DR		Rn	1080	N	N	

ORDINAL

NO ORDINAL

one
hot
encoding

≈ pd.getdummies

Sex	M	F
M	1	0
F	0	1
M	1	0
F	0	1
M	1	0



Shot on AWESOME A05s

Encoding function / PDL

SEX	M
N	1
F	0
N	1
F	0
M	1
F	0

Class variable
Date
Age

No different
from previous
table

No
Met
SIL
Ry

One Hot Encoding Handles Multicollinearity.

ORDINAL LABEL ENCODING

dict_x = { 'Small': 1, 'Medium': 2, 'Large': 3 }
data['outlet-size'].map(dict_x)

1
2
3

sklearn - LabelEncoder
→ Alphabetical sequence

Nomical	FREQUENCY	ENCODING
Metformin	1500	classmate
Sitagutin	5000	Set 2 item की
Rybelsus	3500	उभकी frequency दोनों and उभकी
		disadvantage of Hash की कमी item की उभकी frequency same Hogalda में fail हो
		04/03/21

TARGET ENCODING

mean_sales = training.groupby(["Outlet_Identifier"])["Item_Outlet_Sales"].
mean().to_dict()

OUTPUT
{
'OUT010': 839.35,
'OUT013': 2298.935,
'OUT017': 2340.675,
'OUT018': 1995.498,
'OUT019': 340.3297}

Target Encoding mein ~~सभी~~ ① outlet-
identifier की encoding ~~करेंगे~~ वह on the
basis of Item_Outlet_Sales and Target
encoding Mein encoding ~~करेंगे~~ Target
variable according.

25

ONE-HOT ENCODING

CLASSMATE

dummytrain = pd.get_dummies(~~wt-obj~~^{age}, drop-first = True)

dummytest = pd.get_dummies(newtest, drop-first = True)

dummytrain.shape, dummytest.shape

S. No.	Date	Title	Page No.	Teacher's Sign / Remarks
		Standard Deviation		

Standard Deviation is a measure which shows how much variation (such as spread, dispersion, spread) from the mean exists. The standard deviation indicates a "typical" deviation from the mean.

Standard Deviation, the most widely used measure of dispersion, is based on all values. Therefore a change in given one value affects the value of Standard Deviation. It is independent of origin but not of scale. It is also useful in certain advanced statistical problems.

Skewness is a measurement of the distortion of symmetrical or asymmetrical distribution in a dataset. Skewness is demonstrated on a bell curve.

Standard Normal Deviation = A standard normal deviate is a normally distributed deviate. It is a realization of a standard normal random variable, defined as a random variable with expected value 0 and variance 1.

S. No.

Date

Title

Page
No.Teacher's
Sign /
Remarks

Standard Deviation is a measure which shows how much variation (such as spread, dispersion, spread) from the mean exists. The standard deviation indicates a "typical" deviation from the mean.

Standard Deviation, the most widely used measure of dispersion, is based on all values. Therefore a change in given one value affects the value of standard deviation. It is independent of origin but not of scale. It is also useful in certain advanced statistical problems.

Skewness is a measurement of the distortion of symmetrical or asymmetrical distribution in a dataset. Skewness is demonstrated on a bell curve.

Standard Normal Deviation = A standard normal deviate is a normally distributed deviate. It is a realization of a standard normal random variable, defined as a random variable with expected value 0 and variance 1.

It's a mathematical function that describes continuous Random...

Pm 5

EDA is used by data scientists to analyze and investigate data sets and summarize their main characteristics, often employing data visualization methods.

Mean = the average of the given set of values. It denotes the equal distribution of values for a given data set.

The mean (average) of a dataset is found by adding all numbers in the data set and then dividing by number of values in the set.

Median = Median is the middle value when a data set is ordered from least to greatest.

Mode = Mode is the number that occurs most often in a data set.

Variance = is the measure of how notably a collection of data is spread out. If all the data values are identical, then it indicates the variance is zero. All non zero variances is zero. All non zero variances are considered to be positive. A little variance represents that the data points are close to the mean and to each other, whereas if the data points are highly spread out from the mean and from one another indicates the high variance. In short, the variance is defined as the average of the squared distance from each point to the mean.

Page No. _____
Title _____
Sign / _____
Date _____

EDDA is used by data scientist to analyze and summarize their main characteristics, often employing data visualization methods.

Mean = the average of the given set of values.
It denotes the equal distribution of values for a given data set.

The mean (average) of a dataset is found by adding all numbers in the data set and then dividing by number of values in the set.

Median = Median is the middle value when a data set is ordered from least to greatest.

Mode = Mode is the number that occurs most often in a data set.

Variance = is the measure of how notably a collection of data is spread out. If all the data values are identical, then it indicates the variance is zero. All non zero variances is zero. All non zero variances are considered to be positive. A little variance represents that the data points are close to the mean or to each other, whereas if the data points are highly spread out from the mean and from one another indicates the high variance. In short, the variance is defined as the average of the squared distance from a point to the mean.