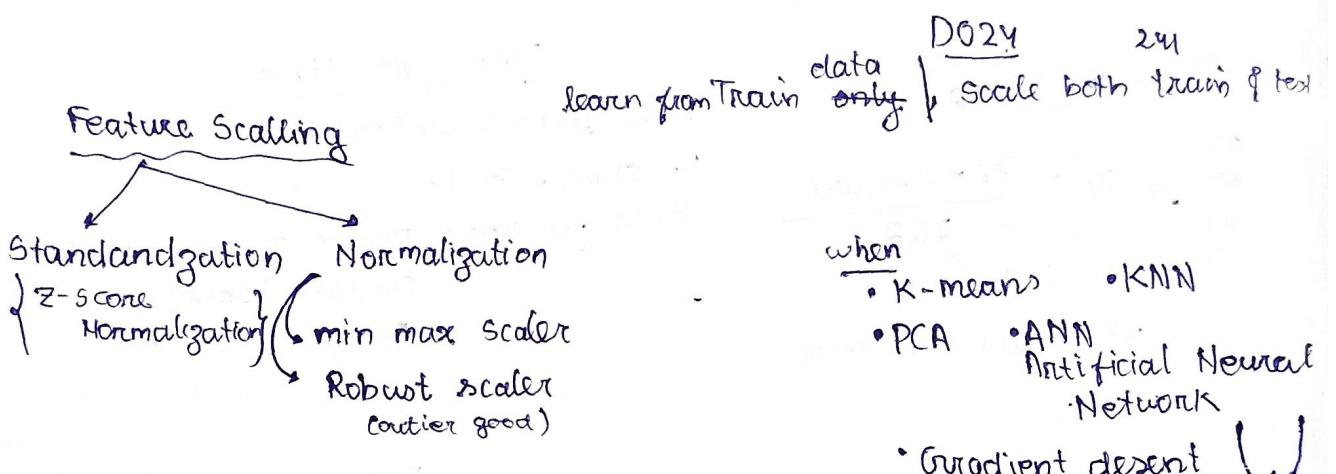
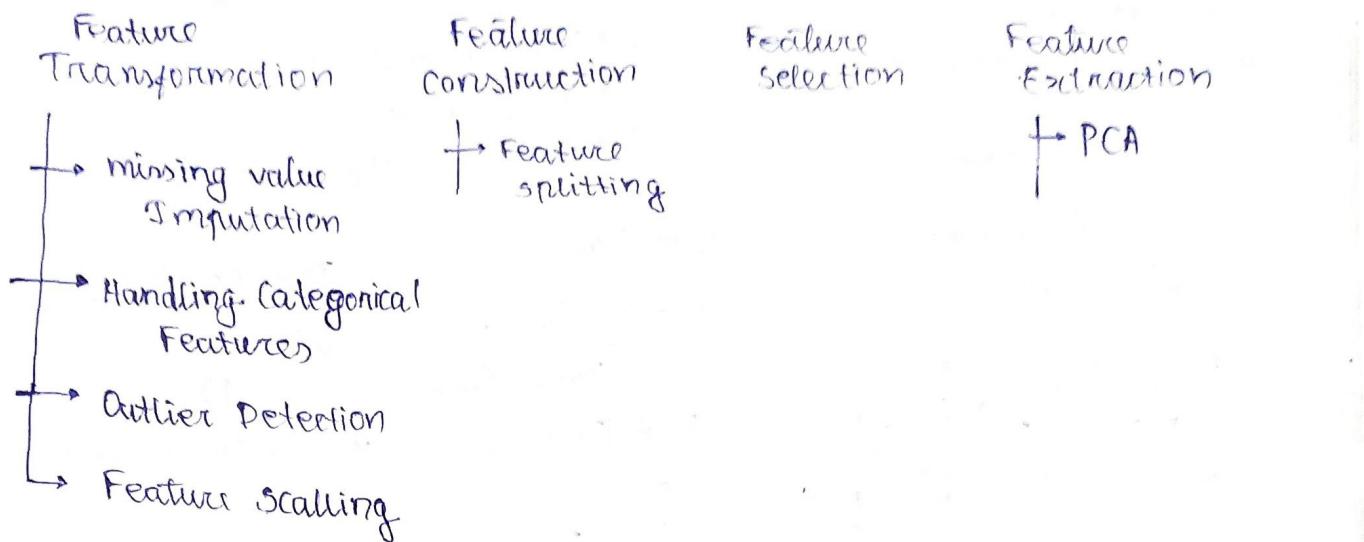


Feature Engineering

D023

216pm
Tue, Apr 15



$$x_i \quad x'_i = \frac{x_i - \bar{x}}{\sigma}$$

mean SD

$x_1 = 15$

$x_2 = 27$

$x_3 = 10$

$x'_1 = \frac{15 - 30}{10} = -1.5$

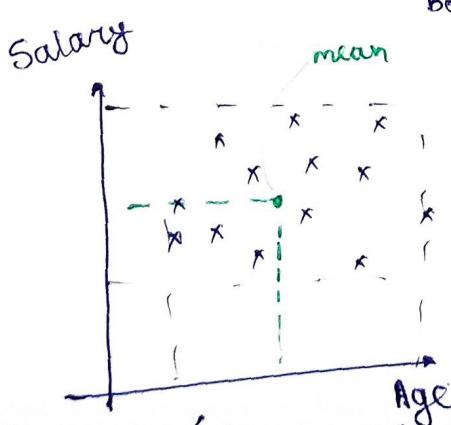
$x'_2 = \frac{27 - 30}{10} = -0.3$

$x'_3 = \frac{10 - 30}{10} = -2.0$

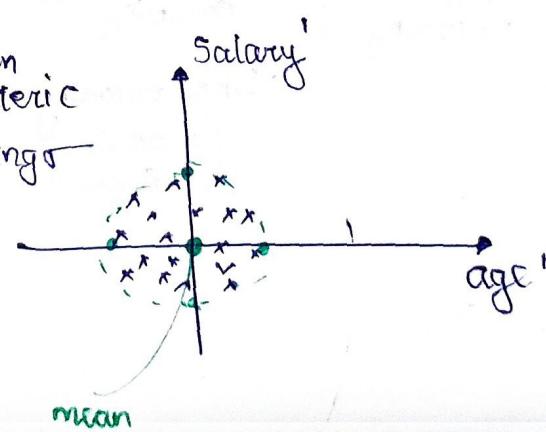
New Series

$$\mu = 0 \quad \sigma = 1$$

mean SD



- mean centeric
- Salating



Normalization

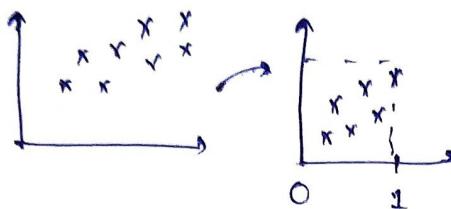
D025

314

Eliminate the unit.
make into common scale.

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

$$x'_i \in [0, 1]$$



mean Norm.

$$x'_i = \frac{x_i - \bar{x}_{\text{mean}}}{x_{\max} - x_{\min}}$$

$$x'_i \in [-1 \text{ to } 1]$$

no clus
wcl in
centeric data

Max Abs Scalling

$$x'_i = \frac{x_i}{|x_{\max}|}$$

wed \rightarrow Spce data
when zero in more

Train data, but Transform
both train & test

Robust scaling

wl
200 $x'_i = \frac{x_i - \text{median}}{\text{IQR}}$
300
100
! { 75th - 25th }

when outlier in more.

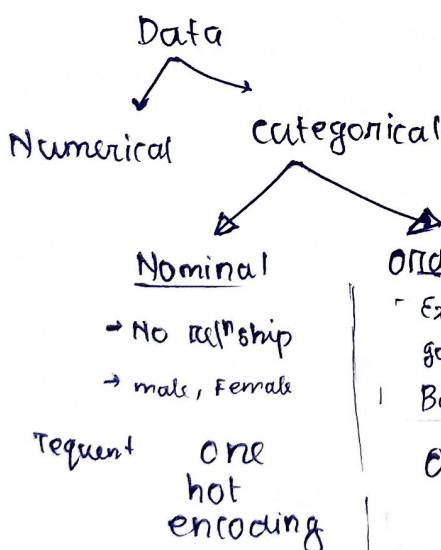
Norm V/s Stand.

1. Is feature scaling reqd?
2. Stand vs beat
3. if you know min, max used MinMax
Outlier Robust
Spce MaxAbs

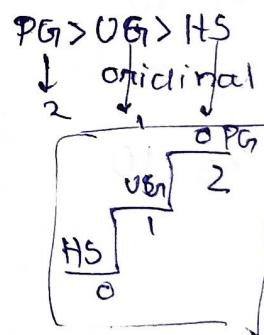
D026

334

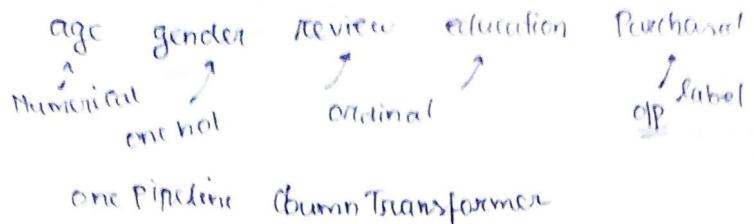
Encoding Categorical variables



Education col	
HS	0
UG	1
PG	2
PG	2
HS	0
UG	1



- Ordinal (if o/p in numerical)
- encoding
- label encoding (if o/p in categorical)



One Hot Encoding

Color	Target	Y	B	R	Target
Yellow	0	1	0	0	0
Y	1	1	0	0	1
Blue	1	0	1	0	0
Red	0	0	0	1	1
B	1	0	0	0	0
R	1	0	1	0	1
		0	0	1	1

D027

3:59 pm

if we do like thi

YBR ML think
012 R>B>Y

But we convert into vector

1,0,0 → Y
0,1,0 → B
—0,0,1 → R

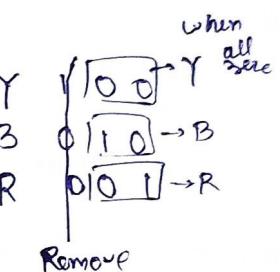
Dummy Variable Trap

* zip column must be independent for o/p

Multicollinearity

• no mathematical Relationship

n categories
(n-1) categories



Column Transformer

D028

4:39

from sklearn.compose import ColumnTransformer

transformer = ColumnTransformer(transformers=[

name object transformer

('tnf1', SimpleImputer(), ['fever']),

('tnf2', OrdinalEncoder(categories=[[['mild', 'strong']], [['cough']]]),

('tnf3', OneHotEncoder(sparse=False, drop='first'), ['gender', 'city']).

), remainder = "passthrough")

→ other col as it is

transformer.fit_transform(X_train)

transformer.transform(X_test)

* if you are not train the model

then used fit_transformer.

* If train the model used fit()

ML Pipeline

D029

```
from sklearn import set_config  
set_config(display='diagram')
```

Mathematical Transformations		DO30
Skew Transf / Power	Quantile	Normal Distribution
log Transf.	Box-Cox	pd.skew()
Reciprocal	Yeo-Johnson	sns.kdeplot()
power(67/89)		Q-Q plot

Log Transform

- Right skewed distribution



- don't use in -ve

$$\begin{aligned} \log &= \log(x) \\ x &= x+1 \\ &\uparrow \\ &\text{add each value w/ } 1 \end{aligned}$$

Reciprocal ($\frac{1}{x}$) ; $\text{sq}(x^2)$; $\text{sqrt}(\sqrt{x})$

left skewed

913

lambda

$$x_i^{(\lambda)} = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(x_i) & \text{if } \lambda = 0 \end{cases} \quad n > 0$$

$$\lambda \in [-5, 5]$$

maximum likelihood
Bayesian stat

Box-Cox Transform

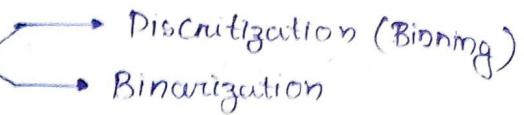
DO31

Yeo-Johnson Transform

$$x_i^{(\lambda)} = \begin{cases} [(x_i + 1)^\lambda - 1] \lambda & \lambda \neq 0, x_i \geq 0 \\ \ln(x_i) + 1 & \lambda = 0, x_i \geq 0 \\ - [(-x_i + 1)^{2-\lambda} - 1] / (2-\lambda) & \lambda \neq 2, x_i < 0 \\ -\ln(-x_i + 1) & \lambda = 2, x_i < 0 \end{cases}$$

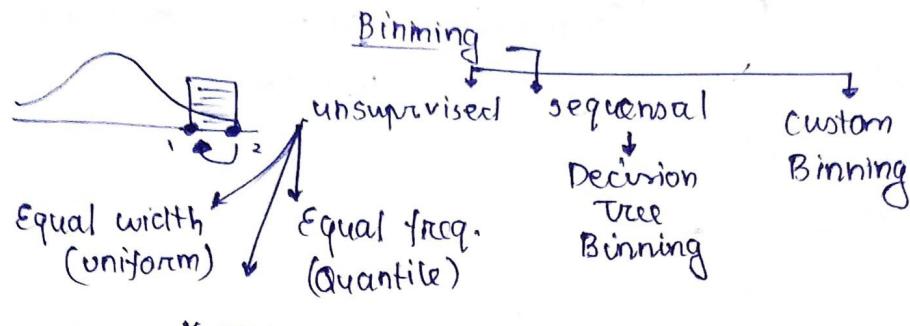
Encoding Numerical Features

DO.3.2

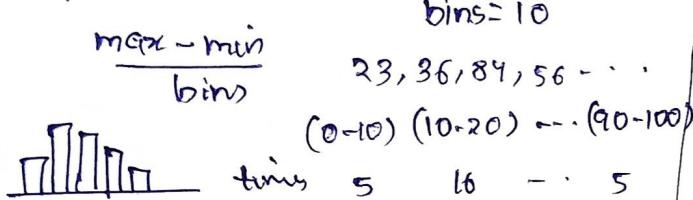


Binning

- To handle Outliers
- To Improve the value spread



Equal width/uniform Binning



Equal freq/Quantile Binning

uniform distribution interval = 10
 Each interval contains 10% of obs.

0-10, 10-20, 20-30, ...

no fixed width 10 4 2

K-means Binning

centroids \rightarrow apply when numerical Discrete

1) Randomly take 5 centroids

2) Find each point dist. w/ each centroid then make into cluster.

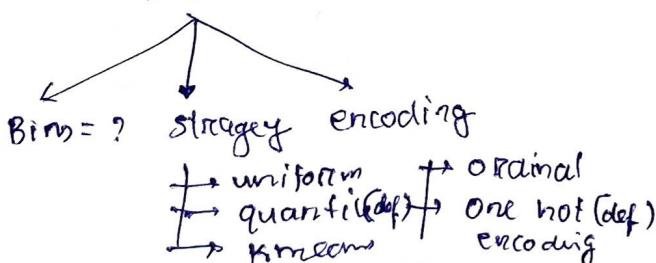
3. or we can bisect the blw two cluster centroid.

4. shift all centroid into mean

5. Repeat all until same

App

KBins Discretize()



Binarization

use during colour image \rightarrow Black white

How EDA helps to decide how many bin in great for accuracy?

Custom Domain
 Use Business Knowledge of Pandas

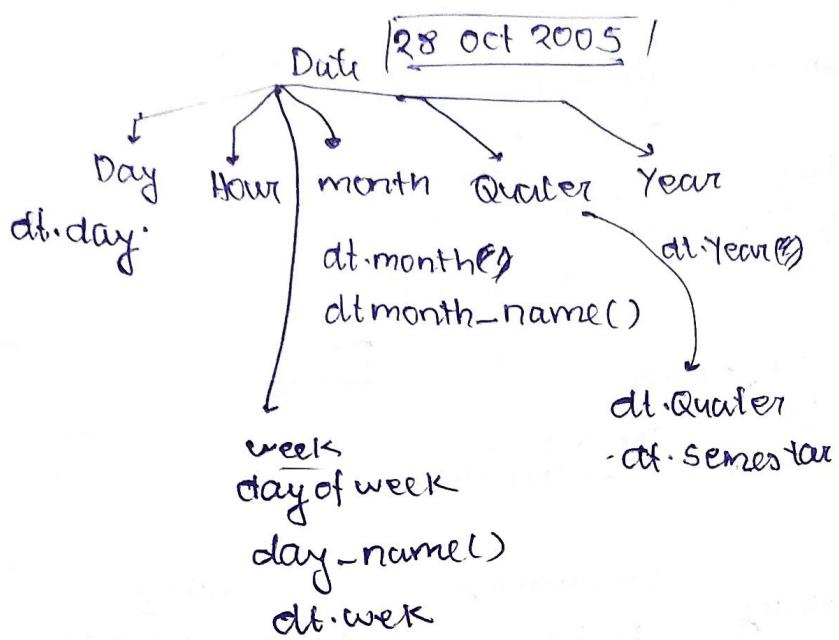
Mixed Data (two types)

Cabin	two col	class	num	
B5		B	5	
C23		C	23	
D41		D	41	

class	num
7	N/A
3	N/A
A	A
N/A	A
N/A	A

ID
Tue, AM 17

DO33



Time 08:
Hr min sec

109am
FRI, APR 18

1) convert into
Datetime
type

DO34

Missing Value

Remove
complete case Analysis
(CCA)

Impule

Univariate

Multivariate

Numerical

Categorical

mean/median

mode

Random

mining
indication

Arbitrary

Automatic

End q. distribution

Selection
(Both num/cat)

KNN
impute

Iterative
Impute
MICE

CCA

Assumption complete at
missing, Randomly

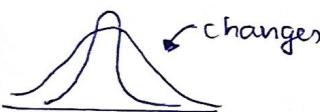
MCAR

when

1. MCAR
2. 5% <

Arbitrary

Create difference
of other data
in same column



when MCAR
not mining
at Random

1. mean / median

Normal
distribution

skewed
distribution

DO36

Benefits

1. simplicity

Disadvantages

1. shape of distribution
2. Outliers
3. covariance
(change the Reln)
w/ other col

when

1. MCAR
2. 5% <

Normality

skewed

$$\mu + 3\sigma$$

$$Q_1 - 1.5 \text{ IQR} \quad 25^{\text{th}}$$

$$\mu - 3\sigma$$

$$Q_3 + 1.5 \text{ IQR} \quad 75^{\text{th}}$$

* most frequently value Imputation

- replace w/ mode
- because more efficient in

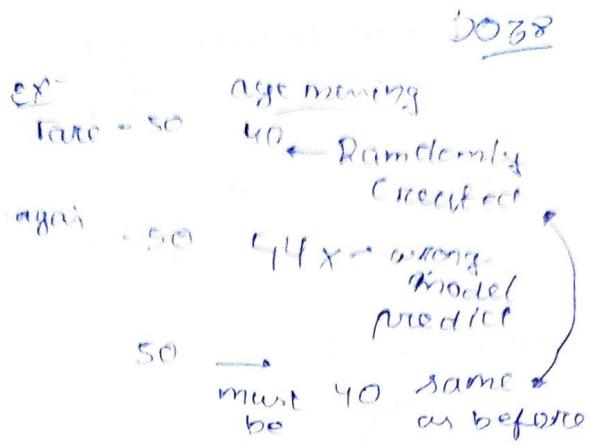
* Missing Category Imputation

Replace null w/ "missing" word

DO37

Random Imputation

- If take the train data on server
- If we use random imputation then the for same IP off course imputation maybe same.

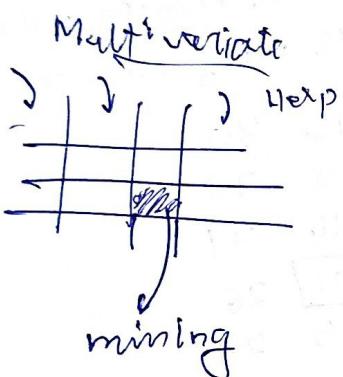


Missing Indicator

<u>Age</u>	<u>Fare</u>	<u>Indicator</u>
27	31	F
41	35	F
NA	31	T
62	89	F

KDD

simple indicator (add indicator = True)
Rather than using missing indicator



kNN

DO39

nan-euclidean

$$\text{dict}(x, y) = \sqrt{\text{wt} * \sum_i \text{dist from present co-ordinate}}$$

weight = $\frac{\text{Total } \# \text{ of co-ordinates}}{\# \text{ of present co-ordinates}}$

F_2, F_3, F_4 for $(2, F_1)$
But present only F_3, F_4 → 3

$$\sqrt{(68-67)^2 + (12-21)^2}$$

<u>sno</u>	<u>F-1</u>	<u>F-2</u>	<u>F-3</u>	<u>F-4</u>
1	33	□	67	21
2	□	45	68	12
3	23	51	71	18
4	40	□	81	□

calculate $(4, F_2)$
present F_1, F_3, F_4 → 3

$$\sqrt{(40-23)^2 + (81-71)^2}$$

KNN Imputer (weight = 'distance' or 'uniform')
 $(1, 2) (2, 3)$

distance $\frac{1}{10} \sqrt{(50-33)^2 + (50-23)^2}$

$$\frac{\left(\frac{1}{30} \times 33\right) + \left(\frac{1}{50} \times 23\right)}{2}$$

$$\frac{33+23}{2}$$

1. more Accurate

Iterative Imputer / MICE

DO 40

Multivariate Imputation by Chained Eqn

MCAR
MAR
MNAR

Assumption / accuracy
MAR memory slow

51. Fill NAN w/ mean of same col.

52. Remove all col missing values

1.	80	15	30	
2		5	20	
3	15	10	41	
4	12	11	26	
5	2	15	29	

Feat

Now
Train
data
predict

15 30
10 u1
11 26
15 29

53. Predict the missing
values of col 1
using other cols.

80	15	30
23	5	20
15	10	u1
12	11	26
2	15	29

54. Remove col 2 missing values

Train	80	15	30	
	23	5	20	
	15	10	41	
	12		26	
Test	2	15	29	

Train data
predict

8 30
23 20
15 41
2 29

55. Predict the missing values
of col 2 using other col

80	15	30
23	5	20
15	10	41
12	11	26
2	15	29

so on ..

Outliers

DOI 1

11:18 am
Sun Apr 20

what are Outliers?

when is outlier dangerous?

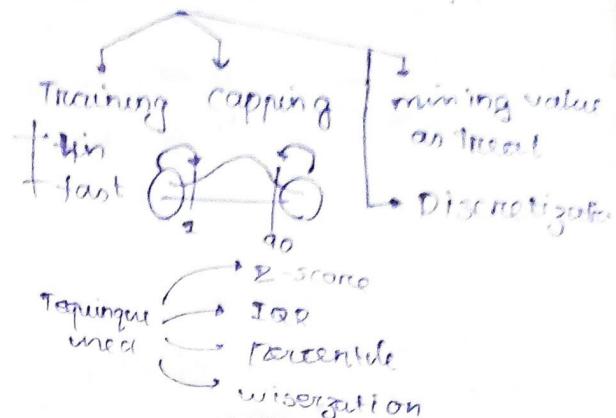
- Anomally: outlier detection.

what to do with outlier?

Effect of outliers

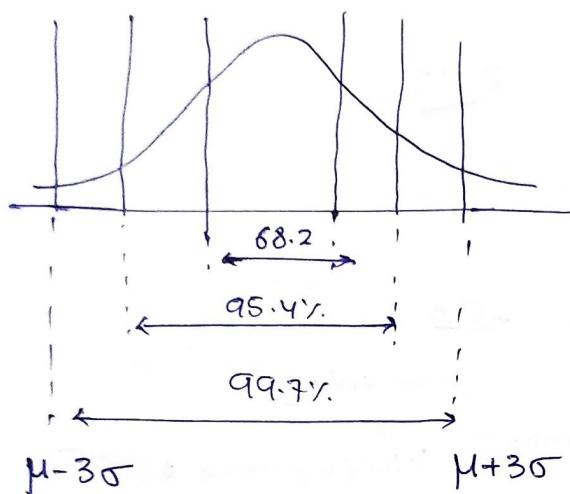
- Does we work w/ weight?

How to treat Outliers?

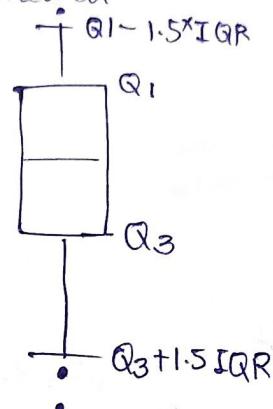


How to detect Outliers

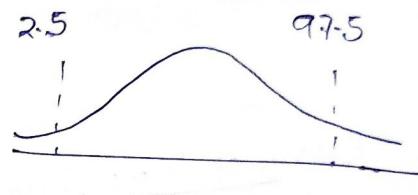
① If data in Normal dist



② Skewed



③ Other distributions
Percentile based



Outlier removal using Z-score

D042

Z-score

$$X_i = \frac{X_i - \mu}{\sigma}$$

$$\mu - \sigma, \mu + \sigma \rightarrow 68.2\%$$

$$\mu - 2\sigma, \mu + 2\sigma \rightarrow 95.4\%$$

$$\mu - 3\sigma, \mu + 3\sigma \rightarrow 99.7\%$$

Trimming \rightarrow Remove all

Capping \rightarrow 85, 3, 90

80, 5, 80 ^(lower)

percentile method capping
in called mitigation
winsorization

Feature construction

Domain Knowledge v/s Experience
on data

D043/44

D045

Curse of Dimensionality

features

After a certain point the
no of features in the
not improve the model
efficiency

How to find the

optimal Dimensionality

D046

sparsity

1. performance
2. computation

Dimensionality

Reduction

Feature
Selection

Feature
Extraction

PCA

1. Forward Selection LDA
2. Backward Elimination t-sene
3. Bidirectional elimination

PCA

DO47

7:22pm

Sun, Apr 20

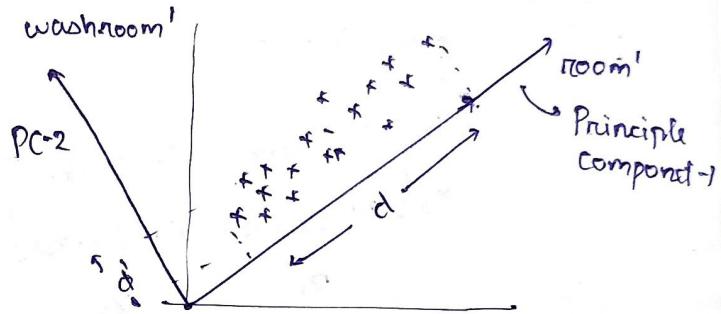
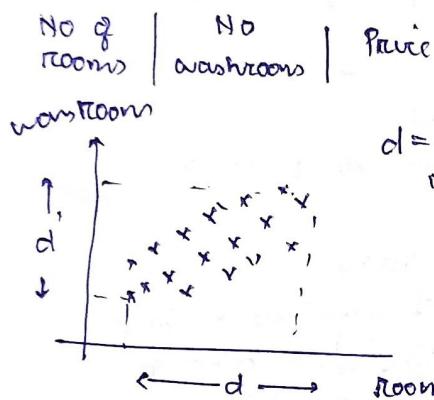
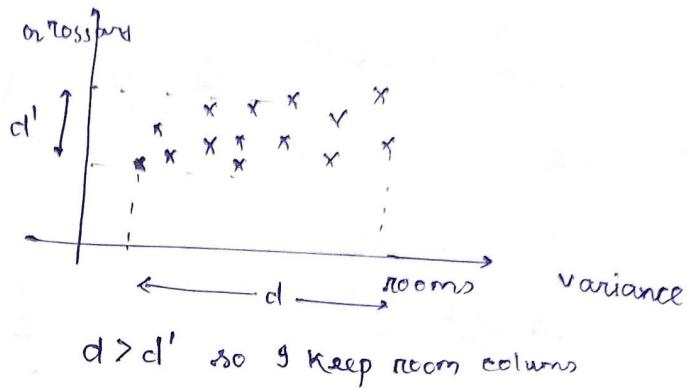
It is technique which can help to reduce higher dimensions to lower dimensions w/o losing essence of data.

Benefits

1. faster execution
2. visualization

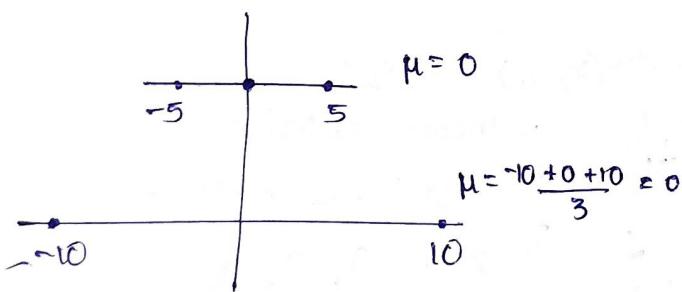
Geometric Intuition

No of rooms	No of grocery shops	price (v)
3	2	60
4	0	130
5	6	170
2	10	90



why variance is imp?

Due to not match w/ other point distance we can't see

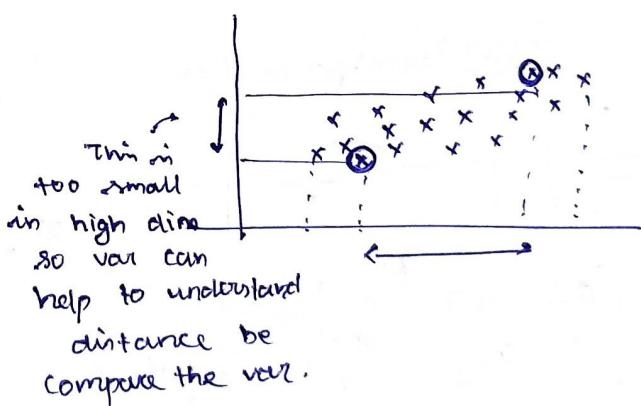


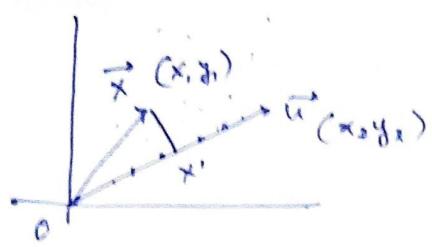
then how to compare

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n} \quad \text{Variance}$$

$$\text{Var}_1 = \frac{25^2 + 25^2}{3} = \frac{50}{3}$$

$$\text{Var}_2 = \frac{100 + 0 + 100}{3} = \frac{200}{3}$$





\vec{x} is a point
 \vec{u} is unit vector

$$\vec{x} \text{ projection} + \frac{\vec{x} \cdot \vec{u}}{\|\vec{u}\|} \cdot \vec{u} = \vec{u}^T \vec{x}$$

$$\therefore \vec{x} \cdot \vec{u} = [x_1, x_2] \cdot [u_1, u_2]$$

$$[u^T x_1] [u^T x_2] \dots [u^T x_n]$$

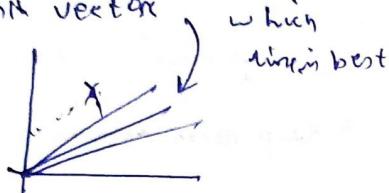
$$[x_1, y_1] \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = x_1 y_2 + y_1 x_2$$

= scalar

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \sum_{i=1}^n \frac{(u^T x_i - \bar{u}^T \bar{x})^2}{n} = \text{variance}$$

Now PCA

Find the Max
unit vector



Covariance matrix

$$x_1 | x_2 \quad \text{cov matrix } [2 \times 2]$$

$$x_1 \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_2, x_1) \\ \text{cov}(x_1, x_2) & \text{cov}(x_2, x_2) \end{bmatrix}$$

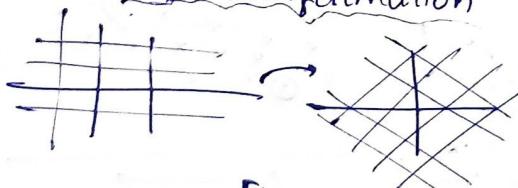
$$= \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_2, x_1) \\ \text{cov}(x_1, x_2) & \text{var}(x_2) \end{bmatrix}$$

symmetric

$$\therefore \text{cov}(x_1, x_2) = \text{cov}(x_2, x_1)$$

- steps theory
1. find covariance matrix
 2. eigen decomposition & find the 2 big values.
 3. the eigenvector is the unit vector

Linear Transformation



$$\begin{bmatrix} 0.5 & -1 \\ -1 & 0.5 \end{bmatrix} \quad \text{Eig vector}$$

dim in same
but magnitude
changes

1. mean centrix
2. covariance matrix
3. find $\lambda_1, \lambda_2, \lambda_3$

$$\text{Eig vector} \quad A\vec{v} = \lambda \vec{v} \rightarrow \text{The largest}$$

Eig value
that data
have highest
eig vector
& spread.

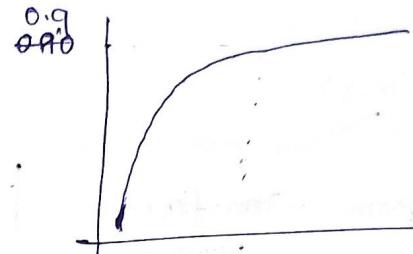
Eig value

magitude
changes.

(x, y, z)
 Eigen value (2, 3)
 784 dim (3, 784)
 if 100 vector (100, 784) $\frac{\text{vector}}{\text{Eigen vector}}$

Finding optimum no. of Principle Component
 $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_{784}$
 Cutoff for original eigen value
 but we found a
 so
 $\left(\frac{\lambda_1}{\sum \lambda}\right)$ $\rightarrow (90\%)$

plt.plot(np.cumsum(pca.explained_variance_))



How much component we need to explain 90%?

when PCA doesn't work



when data spread in same for all axis.

when projection in same

when higher dim. data lossing its originality