

DS203

<https://tinyurl.com/8-2024...>

Books: learning data science
veridical data science
Stats - Business Statistics for contemporary decision making
ml - An intro to Stat. Learning

Test - reasoning and
Prog - internet access
test - assignments

ds203.2023@gmail.com
Prog. four data science
0-2-2-6 , 8-10 assign.
(Assignments - 8-10 - No fraud)
10% 2 sup quizz, 30% midsem,
30% project, 30% endsem.

31st July - (Note: understanding problem/domain knowledge is v. imp)

- Data science from data: ① Broad field
 ② extracting insights
 ③ knowledge
 ④ collecting, cleaning, organizing, analysing, and int data.
 ⑤ to uncover patterns, trends and meaningful information.
- traceaway: Data science is a jumble of many things. we need skills in most of this.

Data science utilizes various techniques, methodology, and tools to extract valuable insight from data. ↳ communication is also an important part

Case Study: The problem: Optimization of processes in a chemical plant.

- goals: ① operating parameters opti.
 ② throughput prediction
 ③ breakdown prediction
 ④ predictive maintenance if eqq. in this most optimised way?
 How to predict if they are sick?

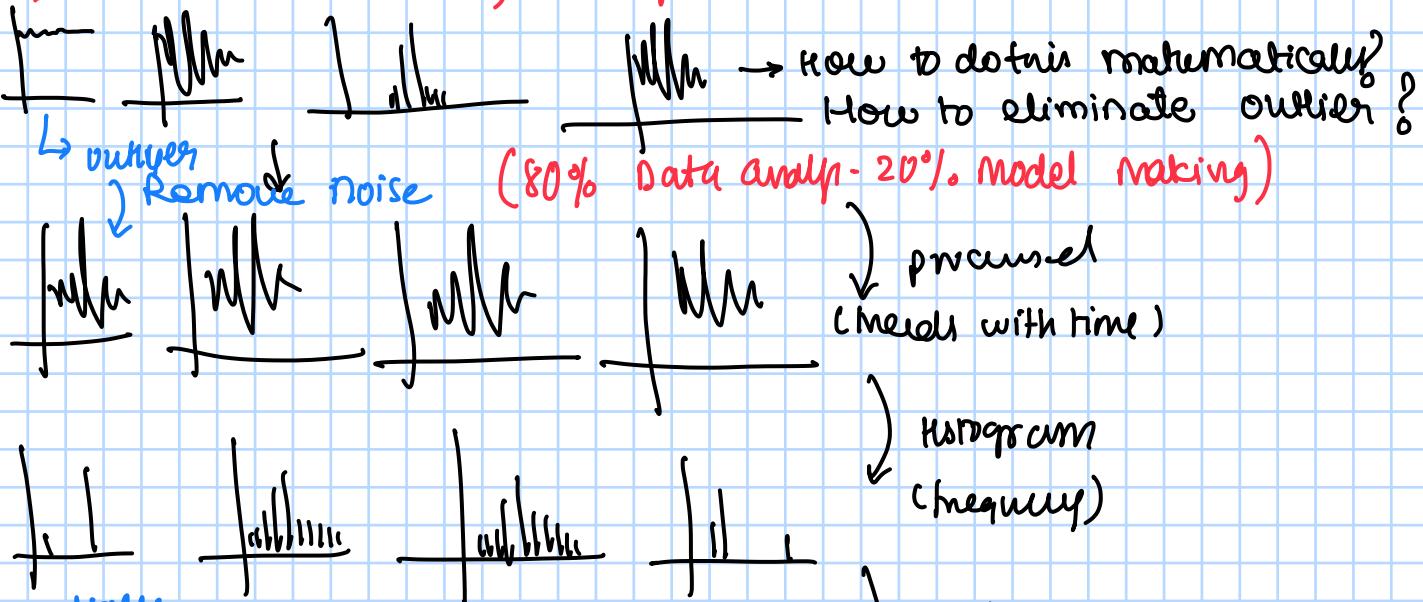
Let's start:

① data given → 132 columns
 (second hand) 7 years of daily average
 ↓
 what tools to use? Excel? deep learning?
 ↓ difference

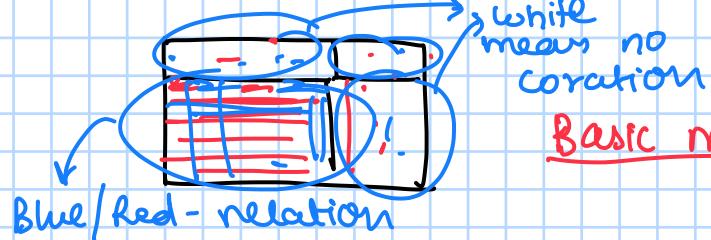
We do preliminary analysis of data: gender?
 ↳ what is data size?
 ↳ most app tools?
 ↳ type of data items present?
 ↳ missing obs? their perc?
 ↳ how to handle missing obs?
 ↳ Are there any outliers?
 ↳ How to handle it?

maturing learning
 $y = f(x)$
 ↳ input variables
 ↳ output variables
 ↳ unknown
 what is y what is x ?
 cause-effect?

* First: understand data, make plots:



essentially two points



(Heatmap)

Basic model: linear regression?



evaluation:
using math and stats, R value, F value ...

data → analysis → prediction / model

2nd Aug: • mastering DS

- ↳ Statistical foundations ↳ machine learning fundamentals (DS 203)
- ↳ knowledge of tools, prog. (DS 203) ↳ Domain knowledge
- ↳ comm. skills

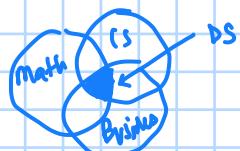
course - intro

↓
Data

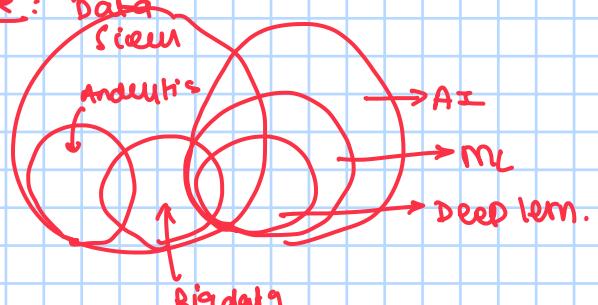
↓
cloud comp

example:

- good ppt vs bad ppt
- ↳ less text ↳ No int. got



Note:



Stages: ① Problem statement

② Prog. for data science

Analytics: ① Systematic comp. analysis of data/stats
 ② Information resulting from analysis of data/stats
 ③ Discovery and comm. of pattern in data.
 ④ raw data → draw conclusion.

AI: AI is ability of machine to perform functions similar to that of a human mind like perceiving, learning, problem solving, etc.

ML: AI technique to handle large and complex data,

supervised learning

unsupervised learning

Reinforcement learning

- linear regression
- logistic regression
- clustering

See Slides

Deep learning - neural network



7th Aug : HW : Do notebooks (moodle)

$$y = f(x)$$

↑ unknown

what is ML?

what do we have? - Data

- Structured data

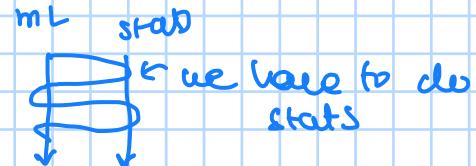
"Table" / "Excel sheet"

eg.	y	x ₁	x ₂	x ₃	x ₄
	:	=	=	=	=
	Something	:	=	=	=

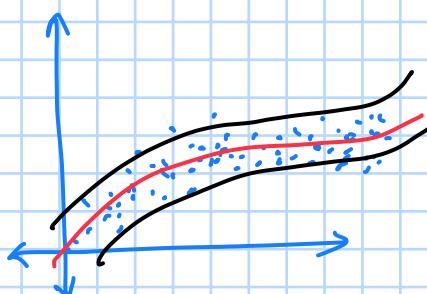
wrong! only 1 variable

we want to find patterns, trends, etc.

we want to predict x is giving us to predict y.

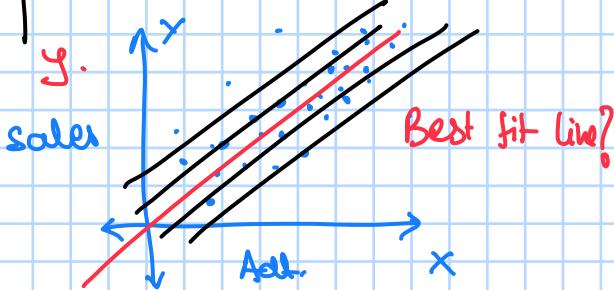


$$y = f(x)$$



$$\text{so } y = f(x)$$

↳ model
↳ line
↳ polynomial
↳ something else



levels of measurement:

see
prob
notes

- NOMINAL . . .

DISCRETE

NO ORDER

DIFFERENTIATE

- ORDINAL . . .

} Gender = M
F
Red, green, Blue
Classification

--- INTERVAL

--- Nominal +

ordering possible

- grades

- tall, mid, short

- DISCRETE x₄ | tall 3

- ORDER 2 | mid 2

3 | short 1



x ₁	x ₂	x ₃
0	0	0
2	2	2
6	6	6

Nominal

continuous
problem?

20°C v/s 10°C

- unit issues

Height 4 feet, 2 feet

- addition/subtraction

multiplication and
division

captures
the input
correctly

{ x₃ into
3 columns }

x₃₁ x₃₂ x₃₃

0 0 1

1 0 0

0 1 0

label / target /
response depends

$$y = f(x)$$

↳ prediction/
features

$$y = f(x)$$

↳ nominal / ordinal / intt / ratio

DISCRETE
"Classification"

CONT
"Regression"

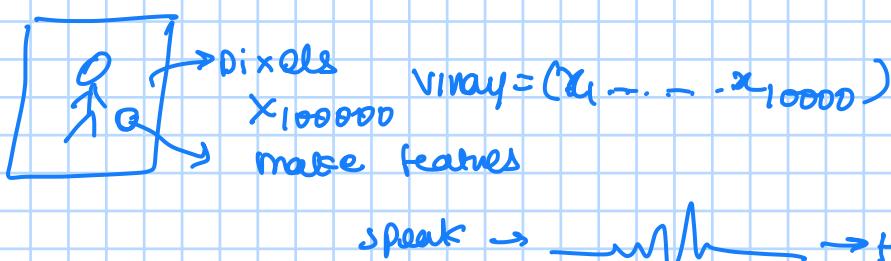
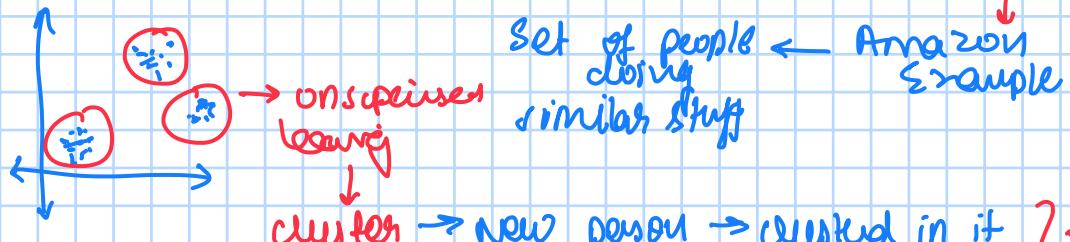
How to handle these cases?

PURCHASING-POWER = $f(\text{AGE}, \text{GENDER}, \dots)$
"DISTANCE" → Imp. role in ML.

"DISTANCE" → Imp role in ML.

Based on particular distances, can we group stuff together?

Note: If y is known → "Supervised Learning"
If y is not there → "Unsupervised Learning" → clustering



These features are very important
But Big ML models can make features for us.

speak → → floats → x₁, x₂, ..., x_n

Population → Entire possible set of data points

Sample → we can always get - As data is sample space

What to do with sample?
→ use it to understand

- ① Understand the population,
- ② Predict the behavior of the pop

9th Aug:

0 on a scale of 10 test - does not mean no sig

time 1 am vs 2 am issue (similar to temp - int vs ratio)

} Basically 0 can have different meaning for different int vs ratio.

$$y = f(x)$$

refers
point?

Classification v/s Regression

↳ discrete

↳ cont values

→ entire data (very abstract)

Population v/s Sample

subset of population

Number of Bank accounts - entire data!

Population?

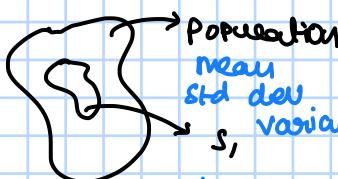
∴ Any data is 'always' a sample

use as much data as we can

Note 80/20 is not comp

Training

Testing

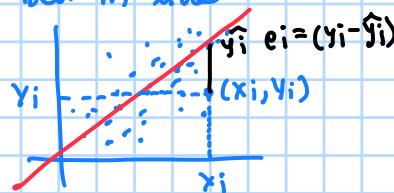


Parameters

Statistics

Using statistics we want to ESTIMATE parameters.

Best fit line



THINGS DON'T END WITH $f(x)$

$\sum e_i^2 = 0$ (Best fit line)

do

$\sum e_i^2 = \sum (y_i - \hat{y}_i)^2 \rightarrow$ why e_i^2 and not $|e_i|$?
↳ sum of square of errors

$$\downarrow y = \beta_0 + \beta_1 x$$

we have to find β_0 and β_1
s.t. we minimise β_0 and β_1 .

$$\sum e_i^2 = \sum (y_i - (\alpha x_i + b))^2$$

$$\frac{\partial}{\partial a} (\text{SSE}) = 0 \quad \frac{\partial}{\partial b} (\text{SSE}) = 0$$

$$\sum 2(y_i - (\alpha x_i + b))(-x_i) = 0$$

$$\sum (-y_i x_i + \alpha x_i^2 + b x_i) = 0$$

$$\sum y_i x_i = \alpha \sum x_i^2 + b \sum x_i \quad \text{--- ①}$$

$$\frac{\partial}{\partial b} (\text{SSE}) = 2 \sum (y_i - (\alpha x_i + b))(-1) = 0$$

$$\sum y_i = \sum a x_i + \sum b \quad \Rightarrow \quad b = \bar{y} - a \bar{x} \quad \text{--- ②}$$

mean of y

mean of x

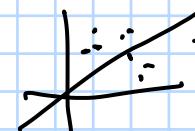
∴ reg. line passes

$$\text{as } y = \alpha x + (\bar{y} - a \bar{x})$$

$$\bar{y} \bar{x}_i = a \bar{x}^2 + (\bar{y} - a \bar{x}) \bar{x}_i$$

$$\bar{y} \bar{y} = a \bar{x}^2 + (\bar{y} - a \bar{x}) \bar{y}$$

Supervised vs unsupervised

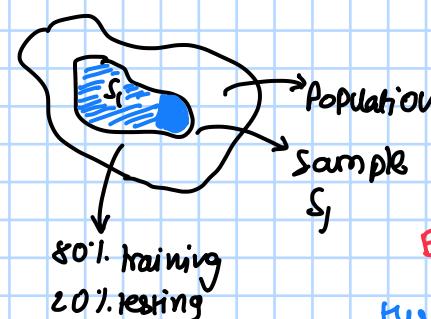


Steps: ① from data train f
② from f predict y.

$$y = \beta_0 + \beta_1 x$$

* Neural network

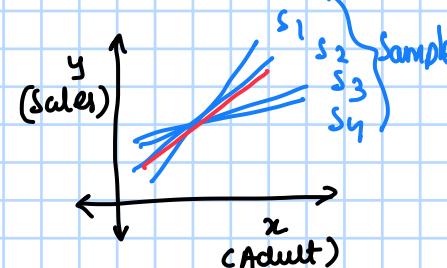
- * Regression
- * Classification
- * Clustering



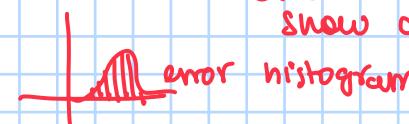
Expected value:

true will always be errors of estimation & it depends on the sample size 'n'!

if $n_1 < n_2$
 $\text{error}(n_1) > \text{error}(n_2)$

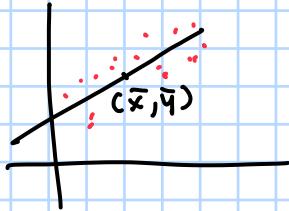


Note: If model is correct - random otherwise errors will show a trend.

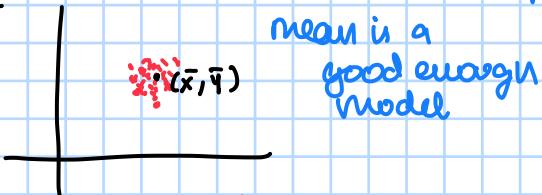


$$\text{as } y = \alpha x + (\bar{y} - a \bar{x})$$

Mean as a model? → Σ



Can we use it?
i.e. $f(x) = x$, y is
solution.
 $f(x) = y$, x .



↑ finally
it comes
down to
error

so see the 'smile' of
data, and then
use / make
functions.

If we are able to visualize data / small data - use basic models.
If huge / big data - use big / heavy models.

14th Aug : see Beta 0 vs Beta 1 - 10

$$a = \frac{\bar{xy} - \bar{x}\bar{y}}{(\bar{x}^2 - \bar{x}^2)}$$

$$b = \frac{\bar{y}(\bar{x}^2 - \bar{x}\bar{x})}{\bar{x}^2 - \bar{x}^2}$$

$y = ax + b$ dependent

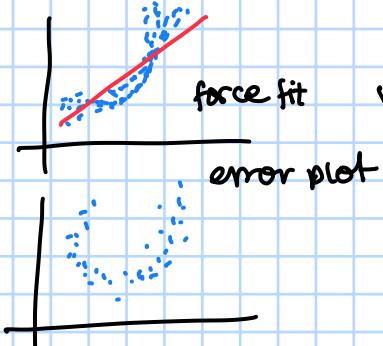
$$= \beta_1 x + \beta_0$$

independent variable

S. Linear Regression

Why arrows are imp?

- if model correct
→ Residual errors
have to be random.

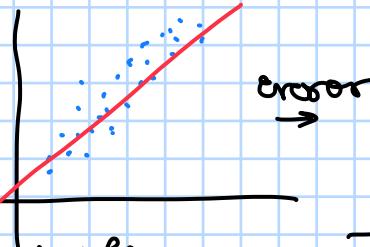


gradient descent

$a, b \rightarrow$ algebraically
but for M.L.R we
use numerical method to find β_i

gradient decent

if random :

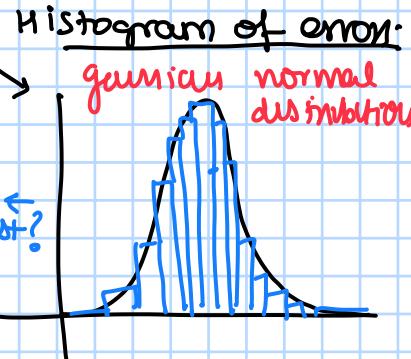


How to use error metric

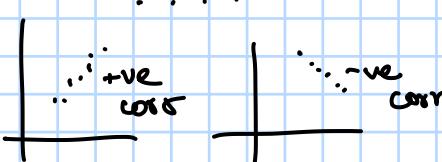
Absolute

MAE
MSE
RMSE

Relative
if RMSE = 634
or 2.35
then $2.35 \downarrow 6.34$
for error



$$\sum (x_i - \bar{x})(y_i - \bar{y}) \sim \text{correlation}$$



R square -

SST = measure of total var

= $\sum (y_i - \bar{y})^2$

= $\sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$

= $\sum (y_i - \hat{y}_i)^2$

+ $\sum (\hat{y}_i - \bar{y})^2$

$$SST = SSE + SSR + \cancel{+ 2C(x)} \rightarrow \text{why zero?}$$

Orthogonality (Sel)

$$I = \frac{SSE}{SST} + \frac{SSR}{SST}$$

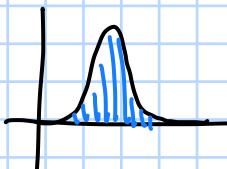
$$1 - \frac{SSE}{SST} = \frac{SSR}{SST} = R-\text{squared}$$

Ques : - correlation
- correlation coeff

$R^2 = \text{sq of correlation coeff}$
w/o x and y

- See all stat usage.

Histogram of means:



16th Aug:

= RAND() → variable by definition is random 'e.g. x_i '

statistical numbers:

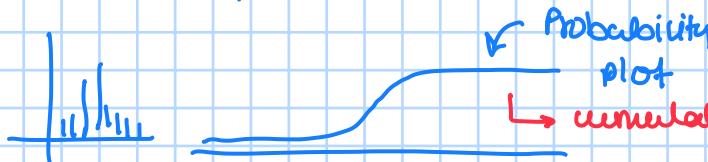
$R^2 = \text{ratio of } \frac{SSR}{SST} \rightarrow \text{do: minimum value}$

$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \sim \text{correlation coeff} \sim x \text{ varies w.r.t its mean}$

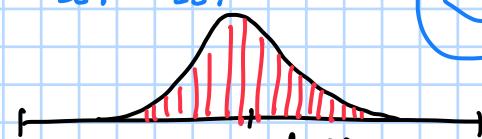
r^2 is sq of correlation coeff



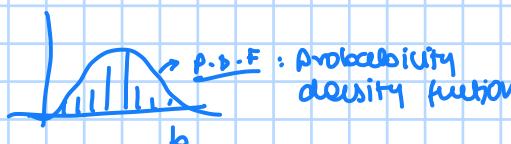
Given - 1000 data points - distribution



$$1 - \frac{SSE}{SST} = \frac{SSR}{SST} = R^2$$



will form the kind of distribution



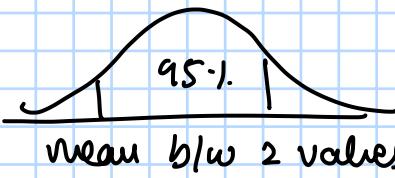
$$\int_{a_1}^{a_2} \text{PDF} dx = \text{area} = \text{prob } b/w a \text{ and } b$$

$$y = \beta_0 + \beta_1 x \rightarrow \text{Statistics} \rightarrow \text{used to estimate parameters}$$

parameter - associated with pop
statistics - associated with sample

confidence interval:

95% confidence interval:

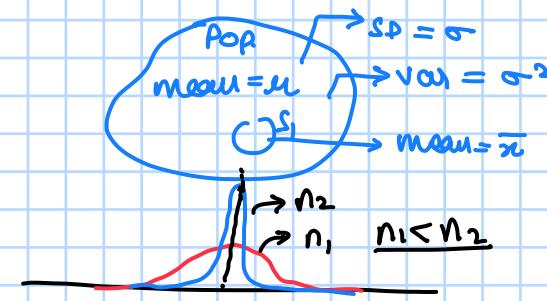


→ 95.1.

mean b/w 2 values

95.1.

central limit theorem:
for a suff. large sample



"dist of sample means → normal dist as n → large"

$$\mu_{\bar{x}} = \mu$$

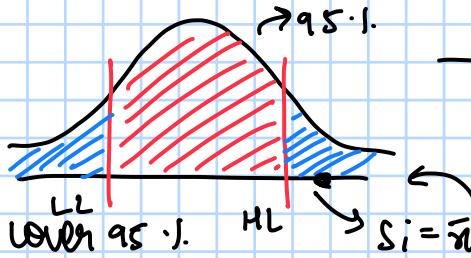
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

if $n \rightarrow \infty$ then straight line

$$\beta_1 = 0 \text{ or } 5$$

if S true is it different from 0?

Hypothesis statistics

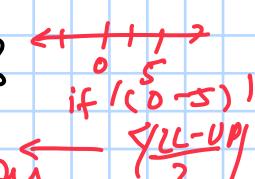


upper tail statistically significant

95.1.

LL

UL



No neg version

$\frac{Y_{LL} - Y_{UL}}{2}$

21st Aug:

$S_1 \leftarrow$ Random representative sample

$S_i^o \leftarrow$ we will get different a_i, b_i values everytime.

standard deviation / standard errors

using a_1, b_1 :

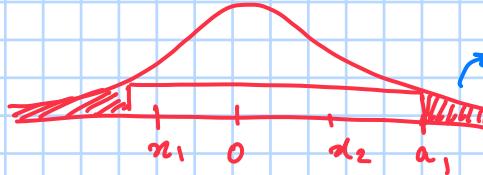
$$SE(a) = \sqrt{\frac{\sigma^2}{\sum(x_i - \bar{x})^2}} \quad SE(b) = \sqrt{\frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)}$$

use $\hat{a} = a_1^o$
as true no regression

σ from S_1 (close)

null hypothesis - Is a^* statistically different from ZERO

P-value:



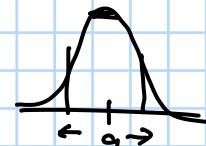
red area is P-value

Two-tailed

vs

one-tailed

< 0.05 is good



F value : multiple linear regression

$$F = \frac{SSR}{n-k-1} \quad n = \text{number of observations}$$

$F_{\text{static}} / \frac{SSE}{n-1}$

move F, better

significance F - p-value associated with F-static ($p_{\text{good}} < 0.05$)

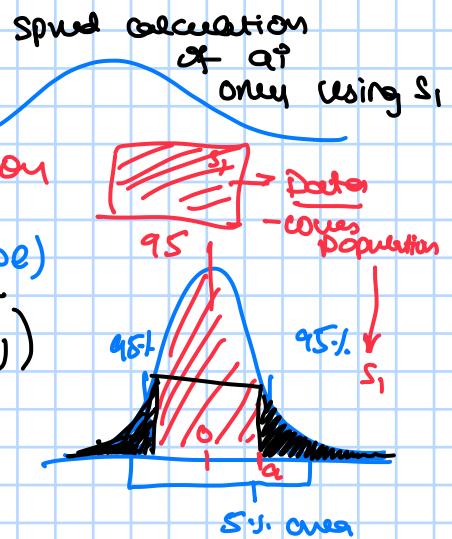
AIC (Akaike Information Criteria) and BIC (Bayes Information)

$$\text{① } AIC - BIC = 2 \text{ good}$$

\downarrow is good
for AIC and BIC

Durbin-Watson test - DLS

- check autocorrelation



Condition numbers: - OLS

- sensitivity of my model

↓ - coeff. near minimal < 10

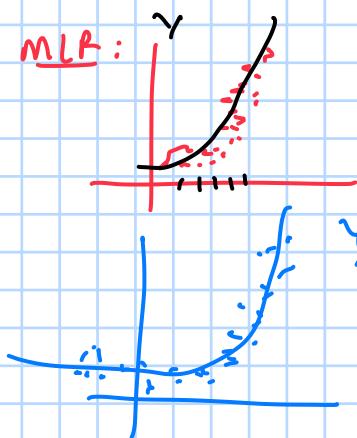
↑ - coeff. near maximal > 30

> 100 very bad

Omnibus: combination of all test (lowvalue)

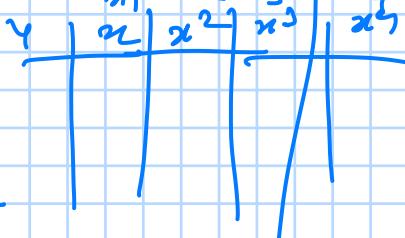
- P > 0.05

MLR:



$y | x_1, x_2, x_3$

$$y = a x_1 + b x_2 + c x_3$$

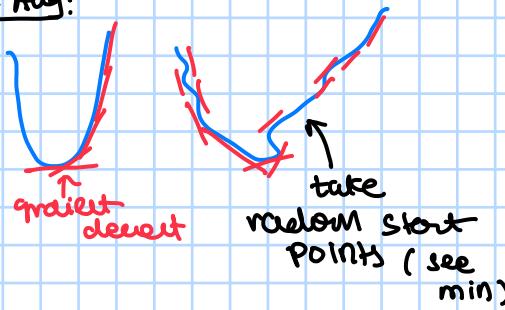


forward feature elim.

vs

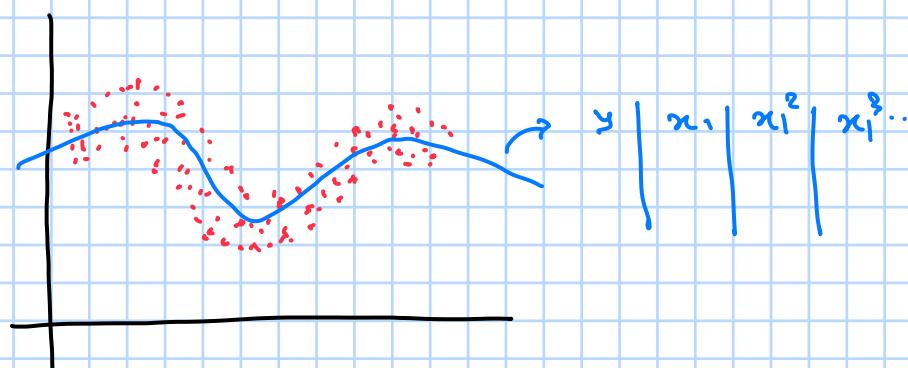
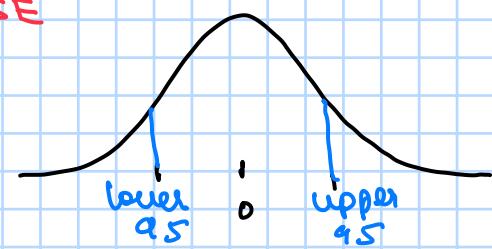
backward feature elim.

23rd Aug:



$$\frac{\sigma}{\sqrt{n}} = SE$$

β_0
 β_0 , lower 95%
 β_0 , upper 95%



Hypothesis testing
- Null hypothesis

- same?
- we do P test
to find out.

overfitting:



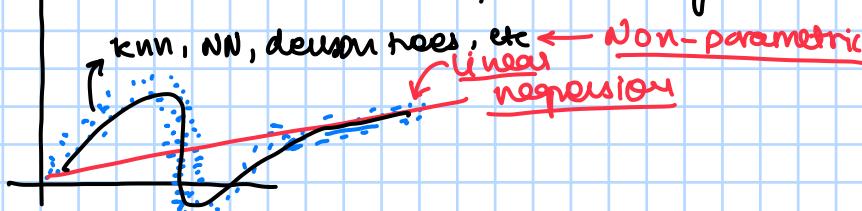
If errors are drastically different from training data, then overfitting is occurring.

$$y | x | \sin(x) | \cos(x) | \tan(x) \leftarrow \text{feature engineering}$$

Backward feature eng.

vs

Forward features eng.



Decision tree /

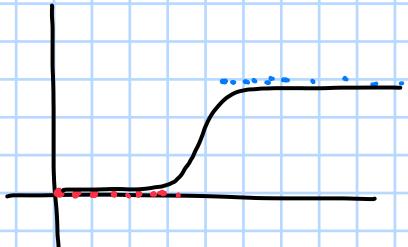
ensemble method

- combine trees together

parametric - what if questions
non-parametric - predict value but
don't give sensitivity analysis directly.

SKLearn - many models available, see models and make them.
- we would have to decide a good model.
- 80% training, 20% testing.
- we are interested in both!

Logistic regression :



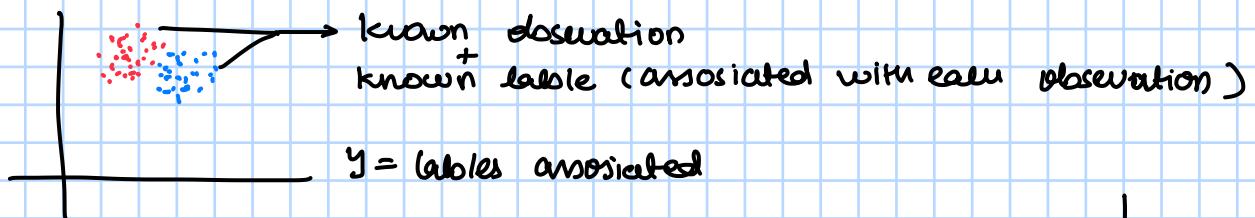
If I only choose 10 from 1000 sample

Each line from different sample

a_i b_i

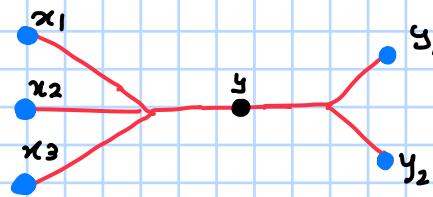
→ 10 data points
→ 1000 data points

- Nominal } 'y' classification → Example, whose face is this?
- Ordinal }
- Interval } 'y'
Ratio regression → we saw this



Problem: given observations, we want to establish a boundary between both.

$y = f(x)$ given x 's, we want to know if it should be classified 0 or 1.



Class boundary
b/w '0' and '1'
for y

or y in general.

0 state and 1 state

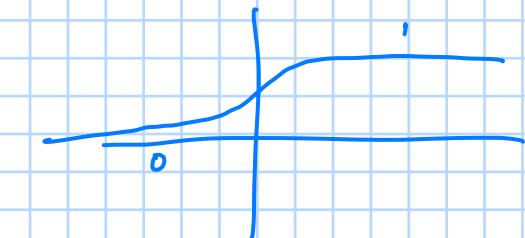
$$a = w_1x_1 + w_2x_2 + w_3x_3 + b$$

0 or 1

sigmoid function

$$f(a) = \frac{1}{1+e^{-a}}$$

f bias/ w_0



$$f(a) < 0$$

sigmoid function

$$f(a) = \frac{1}{1+e^{-a}}$$

← we treat this as probability

$$s(a) = \frac{1}{1+e^{-a}}$$

$\left\{ \begin{array}{l} 0; s(a) \in (0, 0.5) \\ 1; s(a) \in (0.5, 1) \end{array} \right.$

so we have to find w_0, w_1, w_2, \dots , pass it to sigmoid function and get result.

$$f(a) = \frac{1}{1+e^{-a}}$$

$$P(y=1|x) = \frac{1}{1+e^{-a}} \quad (0 \leq p \leq 1)$$

x ₁	x ₂	x ₃	t	y _{act.}
.	.	.	0	1
.	.	-	0	0
.	-	.	1	1

goal: if $t=1 \rightarrow$ we want to maximise P
 $t=0 \rightarrow$ we want to maximise $1-P$

t = observed outcome

$$L = \text{likelihood} = \prod_{n=1}^N p_n^{t_n} (1-p_n)^{1-t_n}$$

we want to maximise this \rightarrow we get weights

$$p = \frac{1}{1+e^{-a}} = \frac{1}{1+e^{-(\omega_0x_0 + \omega_1x_1 + \omega_2x_2 + \dots)}}$$

$$L = \prod_{n=1}^N \left(\frac{1}{1+e^{-(\omega_0x_{0n} + \omega_1x_{1n} + \dots)}} \right)^{t_n} (1-p_n)^{1-t_n}$$

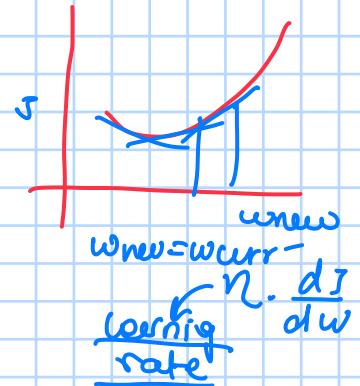
$$\log(L) = \sum_{n=1}^N t_n \log p_n + (1-t_n) \log(1-p_n)$$

\hookrightarrow log likelihood, $\max L \Rightarrow \max \log L$
 $\max L \Leftarrow \max \log L$

$$J = -\log(L)$$

minimising J

$$\begin{aligned} \frac{\partial J}{\partial w} &= x^T(p-T) \\ \text{or} \quad \frac{\partial J}{\partial w} &= (p-T)^T x \end{aligned} \quad \} \text{gradient of objective function}$$



28th Aug :

$$\sigma(a) = \frac{1}{1+e^{-a}} \leftarrow \text{sigmoid function}$$

helps to arrive at weights, till we make boundary

$$\frac{\partial J}{\partial w} = x^T(p - T)$$

$$\frac{\partial T}{\partial w} = (p - T)^T X$$

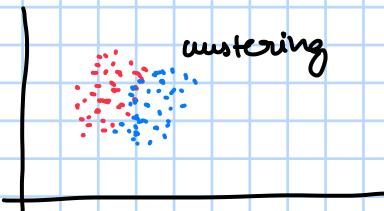
Actual obs(0,1)
predicted prob

$$y = w^T x + w_0 + b$$

η : value

decreasing small value

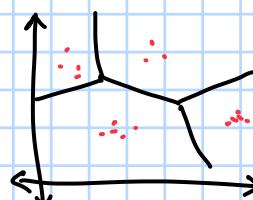
Observations \rightarrow train model \rightarrow user



sk learn linear model \leftarrow your logistic reg. model

Errors in logistic regression: ① we can use .score function

② confusion matrix:



	0	1	2	3	support
0	345	9	6	0	360
1	7	342	8	3	360
2	8	6	332	8	360
3	0	6	8	346	360

this diagonal
is very important

Note: if we had only 25 instead of 360
Accuracy \downarrow This is imbalance

Solution: $360 - 60$
 \downarrow support $360 - 60$

$360 - 360$ $25 - 25$

$360 - 360$ $360 - 60$

$25 - 200$ $360 - 360$

$360 - 360$ $360 - 360$

\uparrow support

CONFUSION MATRIX:

		1	0
Actual	1	TP	FN
	0	Fp	TN

Belong class 1
and class 1

True positive

Accuracy = $\frac{TP + TN}{\text{Total}}$

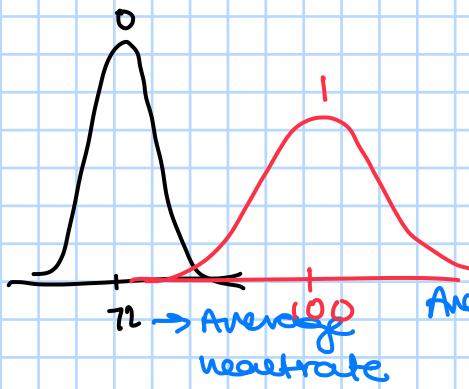
Precision = $\frac{TP}{TP + FP}$ Positive prediction value of the classification, how many are correct

Recall = $\frac{TP}{TP + FN}$ Sensitivity / True positive rate of classification How many are correctly labelled.

$$F_1 \text{ score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

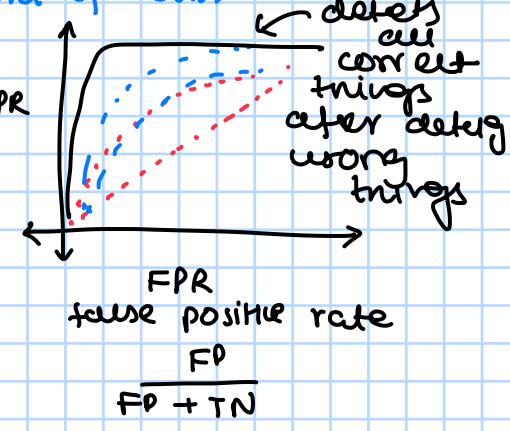
Harmonic mean

③ "ROC": Receptor operator characteristic:
curve created for a classifier (and in the context of CNN)



Average neutral
pulsar

True pos. rate $\frac{TP}{TP + FN}$



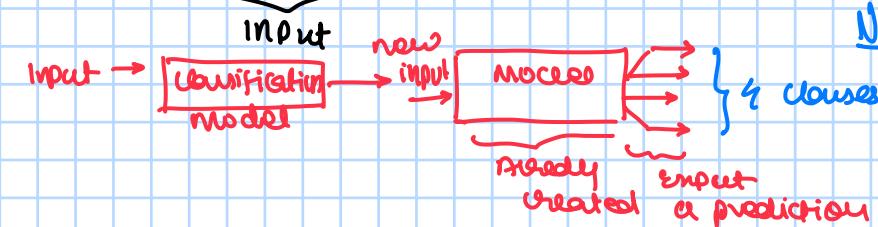
Actually 0 - detected as 1
is Fp \rightarrow not fatal

Fn \rightarrow fatal

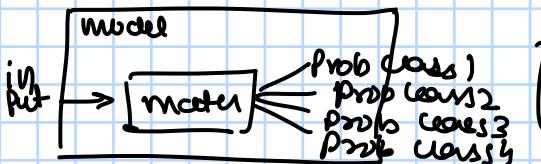
30th Aug:

Classification:

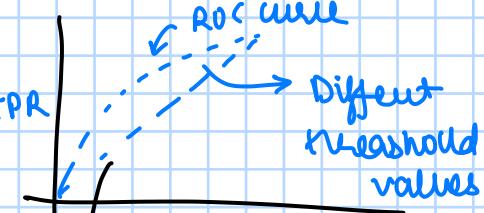
observations (label, imp var-values) $\rightarrow \hat{y} = \text{SOFTMAX}$ (all probabilities)



Note: sigmoid - binary classification problem.



$$\text{SOFTMAX} = \frac{e^{p_i}}{\sum e^{p_i}}$$



When you use threshold value - ROC curve diag.

AREA UNDER CURVE = 1
(ideal case)
Random classifier
(Toss a coin)

Clustering:

unsupervised learning

goal: discovers new aspects of data
mathematically visualise data
subgroups among variables

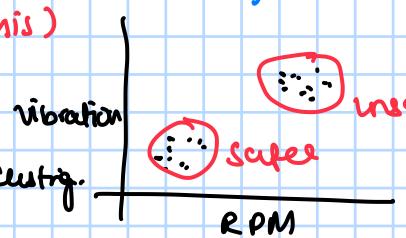
Types:

① PCA (principal component analysis)

② Clustering

Issue: ① subjective process
② metric related to centri.

Note: we can always convert regression problem to classification problem. S, M, L from numbers.



Clustering - part of pre-processing as after classifying, we will make two boxes.

- ① **Kmeans** - Pre-specified no. of clusters — 3
- ② **Hierarchical** - one ad only one cluster

if true randomly 3 is selected



each observation - own center

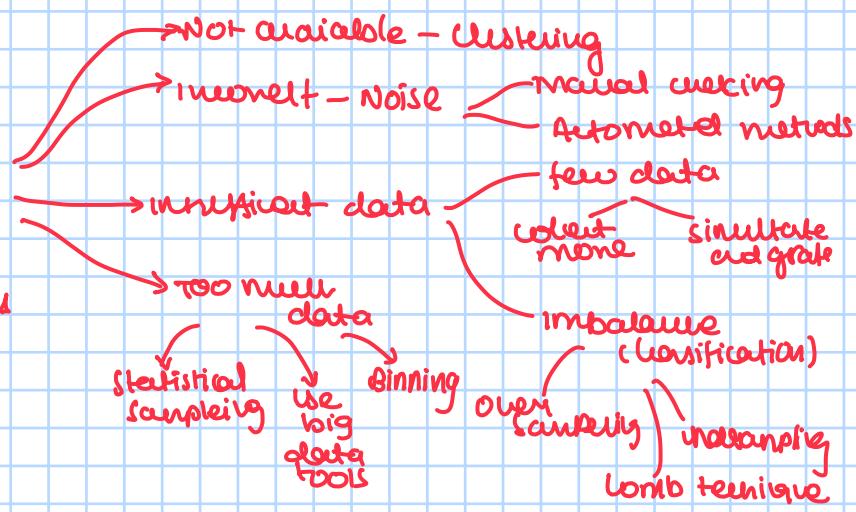
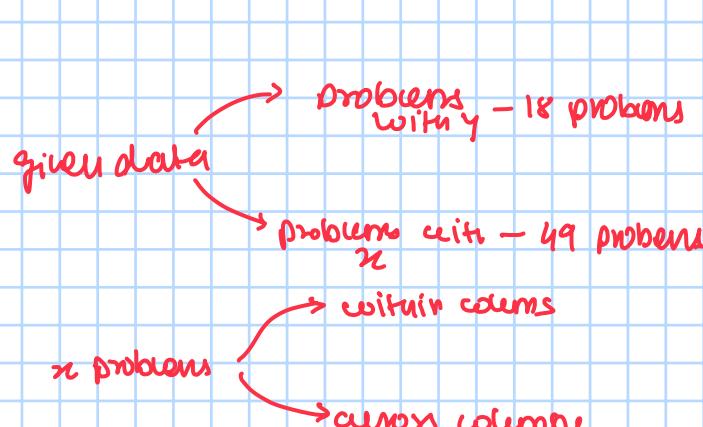
if 1000 data points \rightarrow 1 cluster
1000 cluster

Based on distance

- or **linkage** — complete max dissimilarity $\rightarrow 0$
- single min dissimilarity
- Average
- Centroid



NOT available - clustering



4th sept :

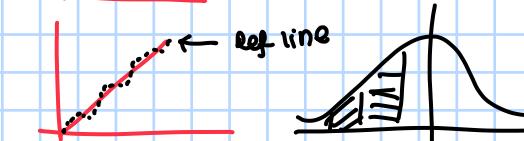
Data analysis

- graphs like ① heatmaps, box graph
② descriptive statistics

EDA

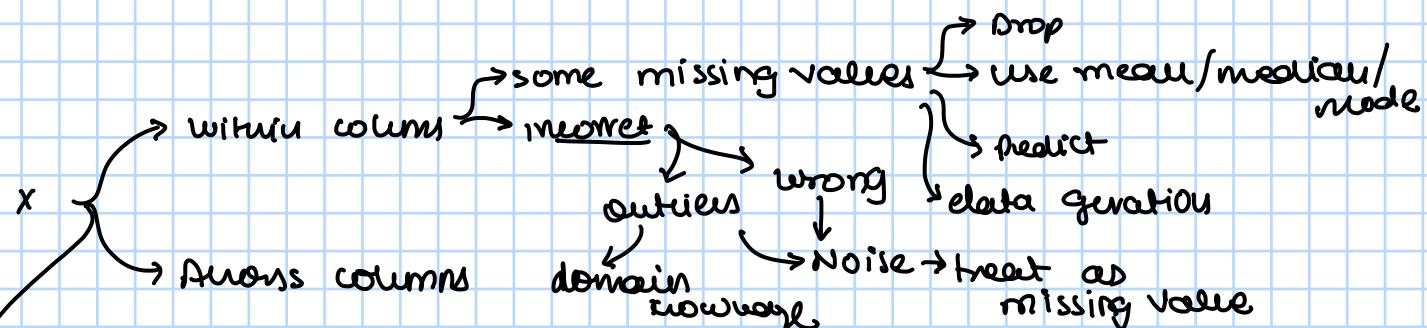
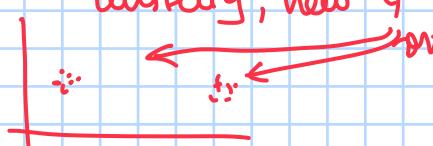
- exploratory data analysis
- gain insights
- spot anomalies
- Q-Q plots

Q-Q plot



given data w.r.t gaussian dist

curving, now $y \leftarrow \text{not } \text{exponential}$

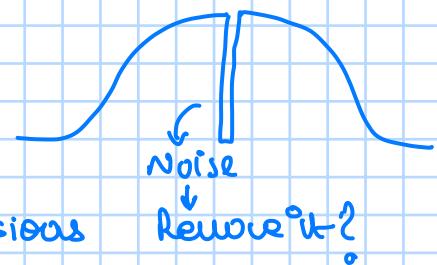


$$f(y_1, y_2, \dots, y_n) = (x)$$
$$f(y_1) \leftarrow (x)$$

incorrect rep
- use encoding methods

using x to produce y_1, y_2, \dots, y_n

negatives
normalization



Action counts

- features not there
- too many features → Feature scaling problems
- Correlation / multi-collinearity

PCA (Principal component analysis)

6th Sept:

example of data cleaning project

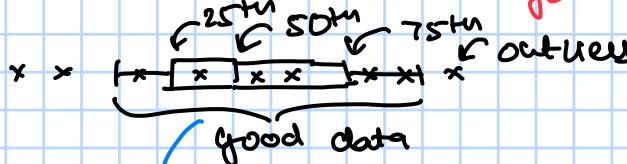
- ① Bird's eye view
- ② finding bad data
- ③ transforming timestamps to require features
- ④ feature eng
- ⑤ predict data for bad data

descriptive statistics:

outliers:

25th percentile, 50th percentile, 75th percentile
median

median close to mean is a good thing



is this wrong/right?

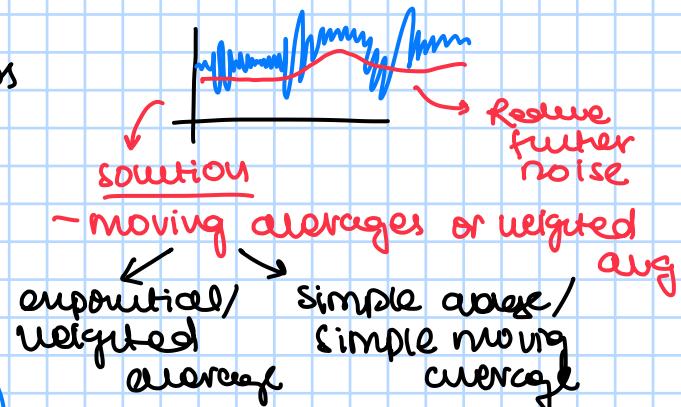
$$\begin{aligned} LB &= Q_1 - 1.5 \times (I.Q.R) \\ UB &= Q_3 + 1.5 \times (I.Q.R) \end{aligned}$$

all data b/w LB and UB is our 'good data' (generally, this 1.5 can also be changed)

Trick: we can recursively apply box plot to filter data

scale before algorithm -

- ① K-nearest neighbors
- ② SVM
- ③ linear regression
- ④ NN



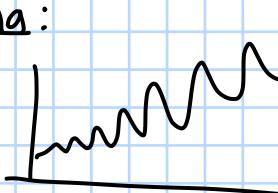
strategies for bad data :

- ① winsorization - replaces extreme value with good value (nearest)
- ② imputation - replace with mean / median
- ③ trimming - removing outliers
- ④ capping - cap the data
- ⑤ other prescriptive methods

CRISP-DM →
Industry Standard for data mining

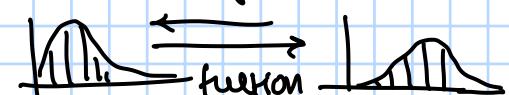
scaling - some algorithms are sensitive to scaling
- some are not

ARIMA:



→ multiplicative → log transformation
→ additive → normal/no transformation

Bonferroni transformation - maximum likelihood estimation



11th Sept:

Rescale data

$\begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$ we do this for better clustering

Homoscedasticity - we need to eliminate this as many digits have assumed this not present

$\log N \rightarrow$ transformation

Data imbalance:

confusion matrix

0	1
0	$253 \ 5^- = 258$
1	$5^+ 12 = 17$

Problem → solution

inc data
dec data

Highly imbalanced data
accuracy = 96.3%
 $= \frac{12+253}{17+12} \approx 70\%$

(less 1)
(less 1 has no value)

SMOTE - Synthetic minority oversampling technique

technique:

- ① randomly choose \oplus
- ② least distance point \ominus
- ③ connect both points and on this line make new \ominus in between lines.

Note: There is an issue to this, as we make points that are intermixing with our data.

ENN - (Noise reduction) - Edited noise-reduction

Algorithm

- ① k nearest points, we find their class
- ② take a vote of majority class, eliminate the instance with minority as major class.

Tomek links (undersampling) -

- different classes
- remove the majority b/w each link

Features / columns - Arises the feature problem

- more the number of features, better model → Not true
- As we inc features, too many variables causing too less data
- overfitting might occur

We need to detect direct correlation and multicollinearity in the system

x_3, x_4 $x_i = f(x_1, x_2, \dots)$

① heat maps
Person w.r.t. factor
heat map

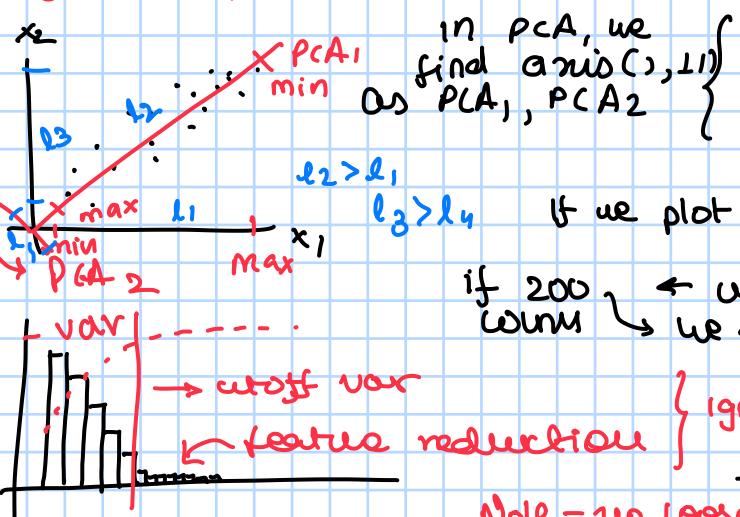
② multicollinearity - VIF - Variance inflation
③ feature exceeding VIF
④ feature eng.
⑤ feature red - PCA

$VIF(x_i) = \frac{1}{1-R_i^2}$

$R_i^2 \ll R^2$ by noisy w.r.t x_i vs other x

① UNDERSTAND STRUCTURE OF DATA
 ② MINIMISE FEATURES

} PCA: principle component analysis



In PCA, we find axis (l_1, l_2)
 as PCA_1, PCA_2

Are these the
 most perfect axis
 for maximally capturing
 variance.

If we plot var by PC_1, PC_2, \dots, PC_n

if 200 \leftarrow we can ignore this \leftarrow too less
 WMM \leftarrow we get first PC_n $n=199, \therefore 1$ dim less

{ eigen value / vector problem }
 \therefore scaled v/s AOs \rightarrow { PC_1, \dots, PC_n }
 200 columns not needed

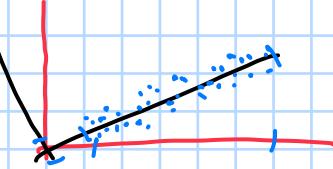
Note - we loose generalisation of model.

$$\text{as } PC_1 = \lambda_1 x_1 + \lambda_2 x_2$$

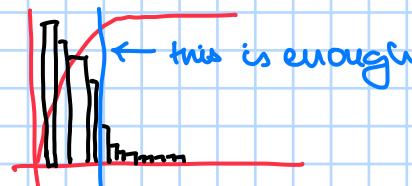
$$PC_2 = \lambda'_1 x_1 + \lambda'_2 x_2$$

13^m Sept:

PCA:



loses features, but we loose explainability

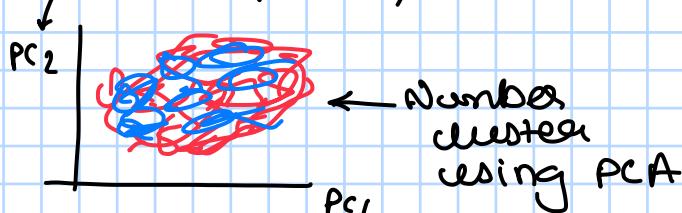


} if this occurs then
data was not
normalised, Normalise
data

multidim data — 2 to 3 visualise data

MNIST Dataset:

28×28 pixels / rows



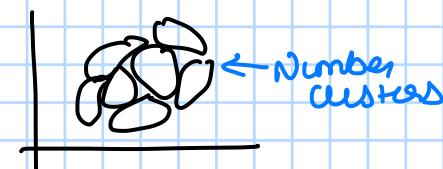
see - Dendrogram

Number
clusters
using PCA

This is
better /
more like
it

New method for 748 dim \rightarrow
2 dim

t-SNE:



← if problem has / uses
distance we cannot use
T-SNE

Feature encoding -

- label encoding
- one-hot encoding
- Binary encoding
- Integer encoding
- frequency encoding
- Target encoding

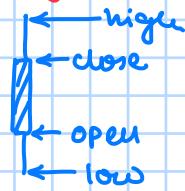
creates
similarity
but does not
have distance

25th Sep:

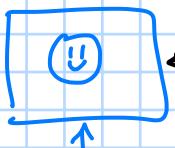
- ① Feature creation methods
- ② Project
- ③ Assignments

midsem how to solve

feature eng: How to convert text, image, Audio, video, etc.



one more example of data

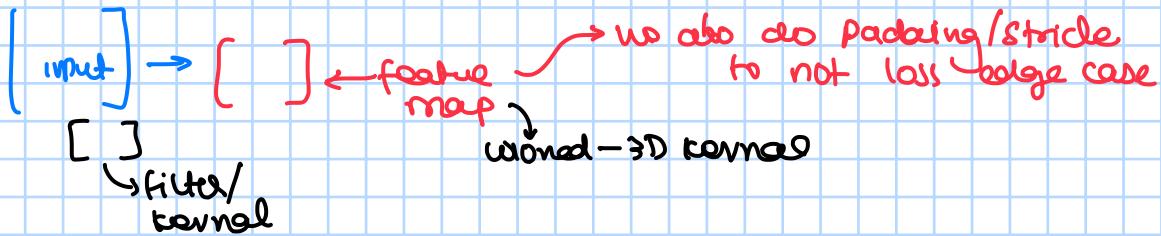
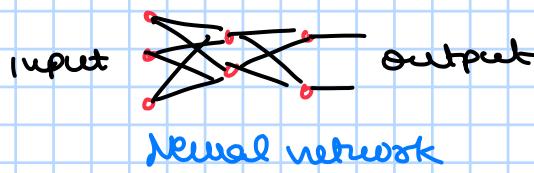


→ model/dataset for pluto

→ Best - convolutional
neural
network

↓
low we
split into
pixels - and the
do it

filter: $\begin{matrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{matrix}$



convolution are further pooled - Pooling

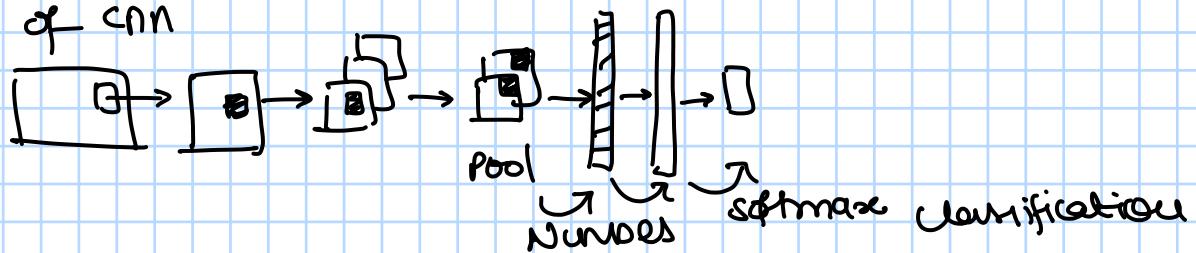


→ [] (we do all this to reduce feature / feature eng)

27th Sept -

feature creation:

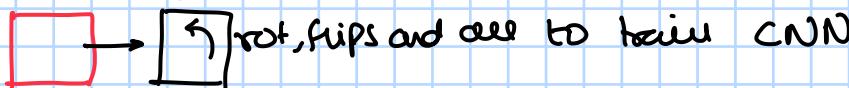
In example of CNN



Example:

- Cars
- ① car 1 [1, 0, 0]
 - ② car 2 [0, 1, 0]
 - ③ car 3 [0, 0, 1]

data training:



Text processing -

Translation, Auto-crop, predict next word, tend from handwritten note

4th Oct:

vectorisation
'embeddings'

Dense — vector representation $(1, 0, 1, 0, \dots, 0)$

Word2Vec — A type of word embedding which converts words to vectors

- distances / similarity search etc can be done

9th Oct:

Audio feature :

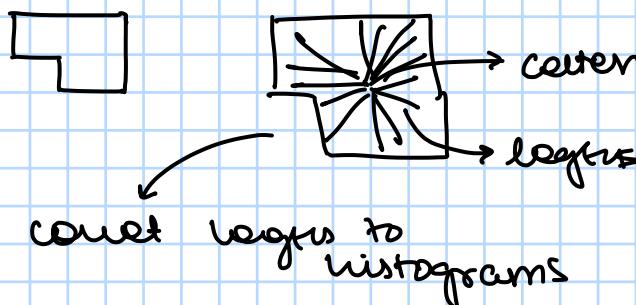
- Audio is finally a signal
- Interested in $y = f(x)$
so we can convert audio into vectors
- Our ear is sensitive to $20\text{Hz} \sim 20\text{kHz}$, we need to have this feature
- Images also get generated 

Note: project : Image processing with audio signals

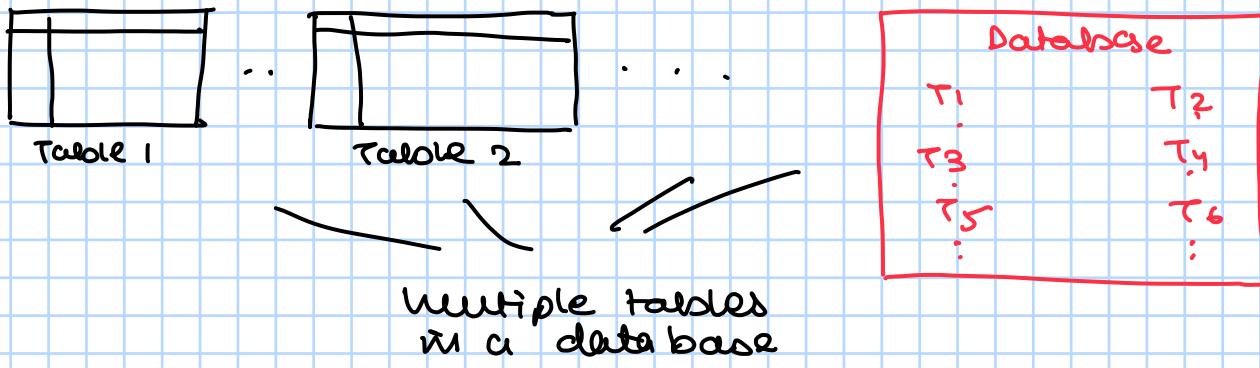
Mel-frequency cepstral coefficients (MFCC)

- features:
- | | | |
|-------------------|---------------------------------------|--|
| Image
Spectrum | ① Mel- <u>FCC</u> | - useful for peculiarity of beats of music
- convert Audio to Image - then do image process |
| | ② Spectrogram - 2D | |
| | ③ Cepstral - Fourier transform of log | |

Shape in Computer:



11th Oct :



Example:

T1: Student

Schema
 {
 first name - Alpha
 name → ref
 DOB
 Address
 :
 :

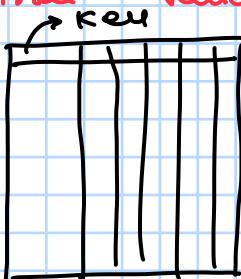
The tables T's
are converted
to
tree in
order to
search better

normalised form

Ist normal → every cell - only one data

IInd normal \rightarrow particular data format

- IIIrd normal \rightarrow validity more than 2nd normal
↑ keen



After 3rd normal

-unique key

- keeps don't repeat

- every cell - one data

The diagram illustrates a many-to-many relationship between two entities: **Students** and **Courses**. It features two separate tables, each with multiple columns. A horizontal double-headed arrow connects the two tables, indicating that many students can be associated with many courses, and vice versa.

given you no, wait a name
and year 2024-Autumn all
courses

```
graph TD; DB[Data Base] --> RDBMS[R.D.B.M.S - system]; DB --> RM[Relational Management]
```

Everything at
Baited links all
of this

Question: How do we define schema for photos? videos? Audio?
(*Apriori/Beliehand*)

SQL - Structured Query language

↳ Query language supported by RDBMS

- 5th gen
- Declarative - lower than nat gpt
but higher than python (3rd gen)
- Describe what to do but there is a syntax

Eg: Select RollNo, Name from db.students where

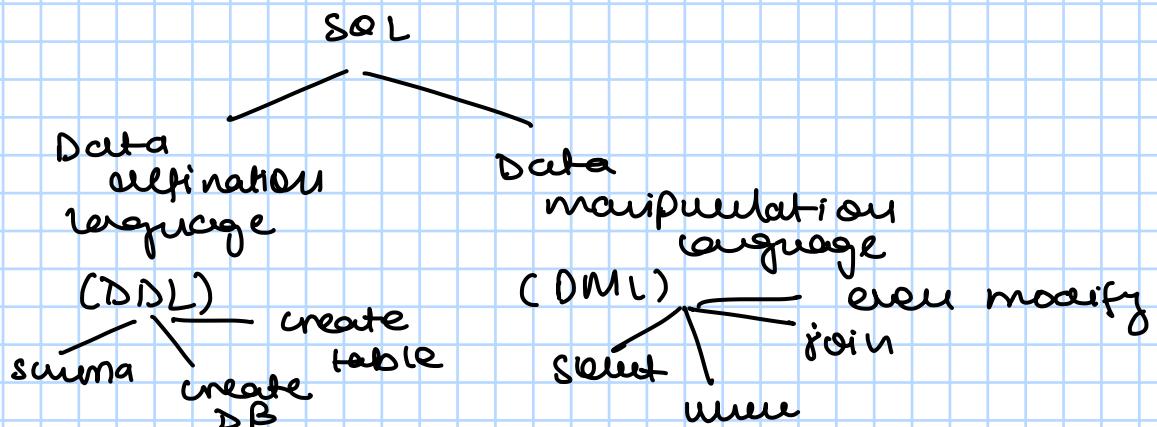
DOB > "2000/01/01"

Syntax: Table is a set

σ select
Π projection
∪ union
∩ intersect
- diff
× Cartesian product
▷ join

select * from table where CONDITION;
select column1, column2 from table;
select * from table1 union select * from table2;
select * from table1 intersect select * from table2;

Example: (Tstock × Tstock) ↗ (use Widow & use TATA)



16th Oct :

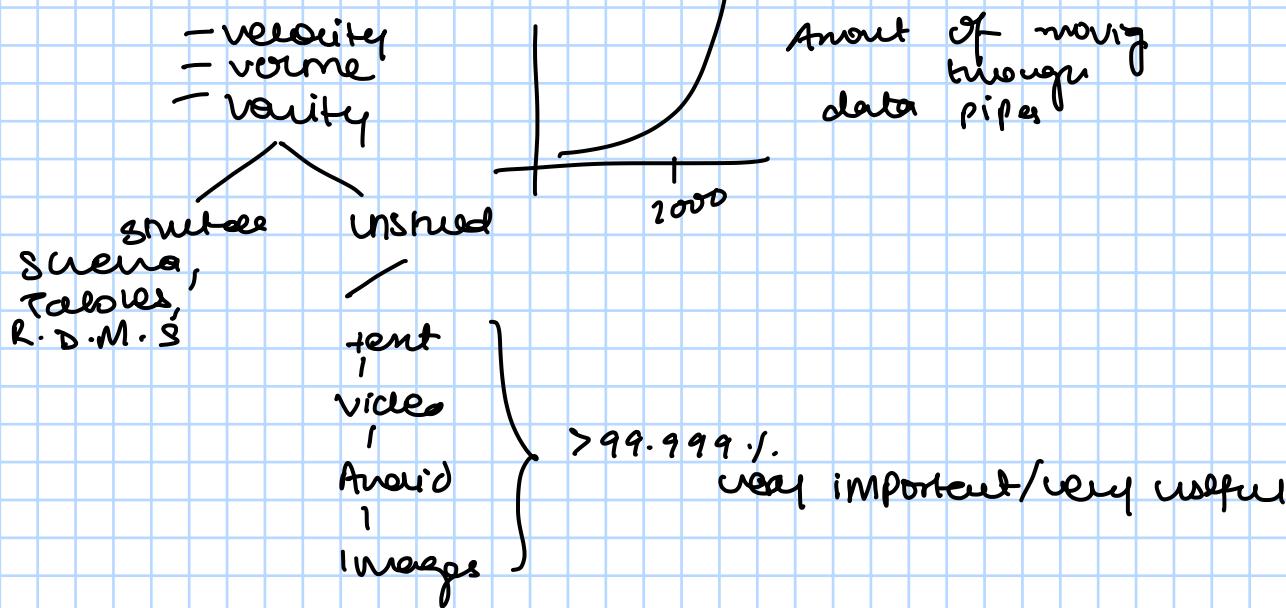
- Project : ① audio files (songs)
② unsupervised data
③ more project - feature engg & feature production

Songs - clusters

Note : What is MFCC - understand that first

Azme : do yourself

Big data : What is big?



How to store/process it :

- How to store?
- How to process?
- What kind of info can be extracted

18th Oct :

Login to VM using
tree

- ① DB inside VM
- ② LSE data table in Stock

Big data:

- size (relative)

Structured v/s Unstructured

Note: semi-structured data in data-scheme

Unstructured - Image
Audio
Video
Text

Picture data - in RDBMS

Storing large data

How to process it
now

23rd Oct:

R.D.B.M.S

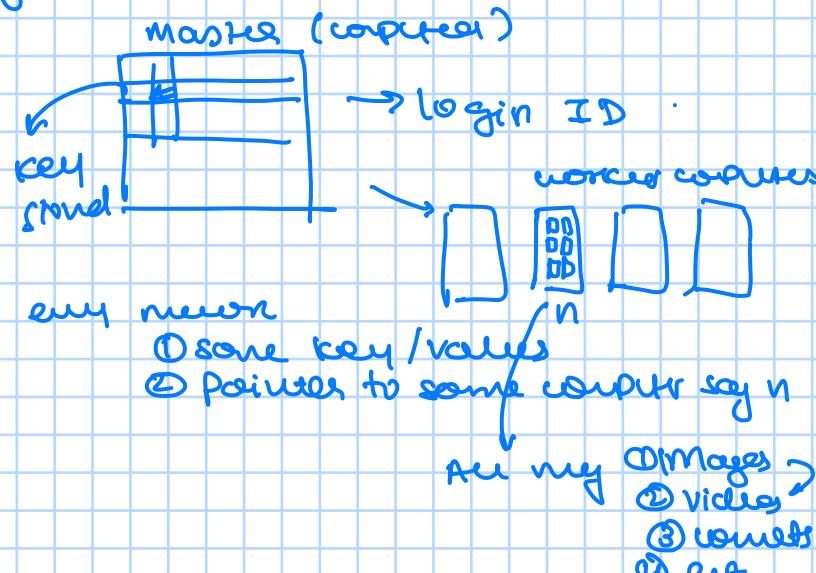
ID	
----	--

key, schema

Database \rightarrow key + value

Database is collection
of key value pairs

Example: facebook \rightarrow $\langle K, V \rangle$
Login key / value place

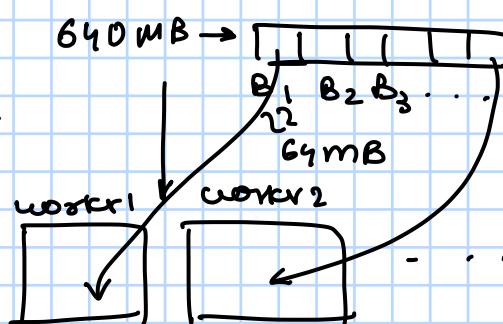


Queries example ① How many videos got uploaded

Note - we don't need database,
we need to use RAM

Then any worker computer gets answer and sends it to master

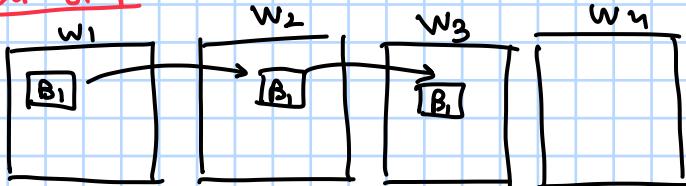
HADOOP - Parallel processing for big data



It stores blocks into worker computer

64 mb - As data is read for 64 mb data chunk

Next Step:

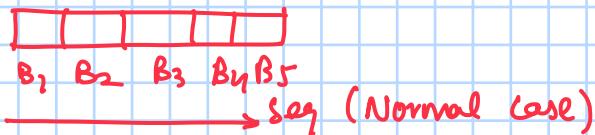


- min 3 blocks (schema)
- Stand for safety
- communication within Racks is much higher than across

\therefore 1 file \Rightarrow file is stored & information available in master.

What if master crashes \rightarrow so there is one more master whose job is to copy master!

Outages: Recovery takes place



Process of pipe: Instead of sequential we can also run them in parallel as every chunk on different computer

HADOOP is a cluster (if m+m+w₁+...)

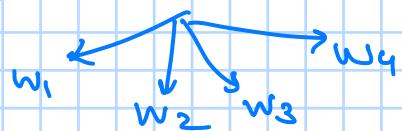
High Availability (HA)

Depends on how much data is needed

Processing:

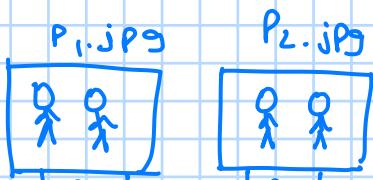


Programme $P_1 \rightarrow$ key/value pair \rightarrow Master



now P_1/P_2 programme gives/rakes Queries and now can be stored to store data or stored in RDBMS.

Ex:



Answer ① How many people
② Who are in the photograph

image process programme (map)



Here P₁ ① identifies no of faces (Image processing)
② identifies who

- Takes you to have building brick by brick
- not conventional way
- we will be able to process almost all kind of data

Project:

Varying columns depending on
length of volume

① Every file has {
20 rows}

② creating right kind of features to solve the problem

③ Note: No of columns not uniform but columns per sec
 $= \frac{44100}{512} = 86$
as $\frac{44100}{512}$ samples per second
 $\frac{44100}{512}$ samples per column

map / Questions ① Identify song

② Identify who sang it

Questions to solve ① How to have same number of columns
↓
solve by clustering (think of features)

② Seems to be 7 groups of songs / clusters

Heatmap → Rearrange it to see clusters together

for first problem need feature set ← some number of features

③ Identify songs sung by P₁, P₂, P₃ ← so now we have categories

for this take your own song → make a testing set /
train data

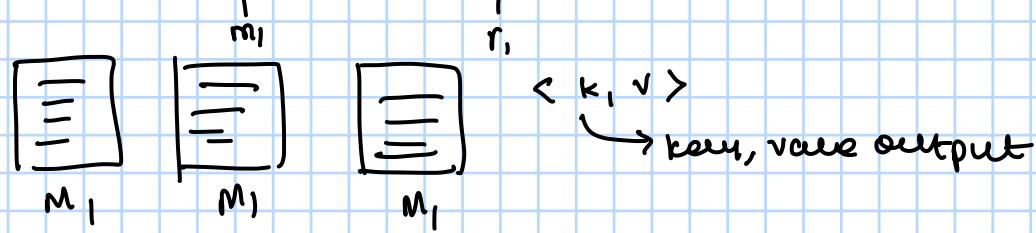
∴ we need training data → make MFCC coefficients

↓
Count into relevant set of features

↓
Use classification
as a tool to identify which cluster does what

25th Oct :

Map Reduce - Input data \rightarrow Map \rightarrow Sort \rightarrow REDUCE

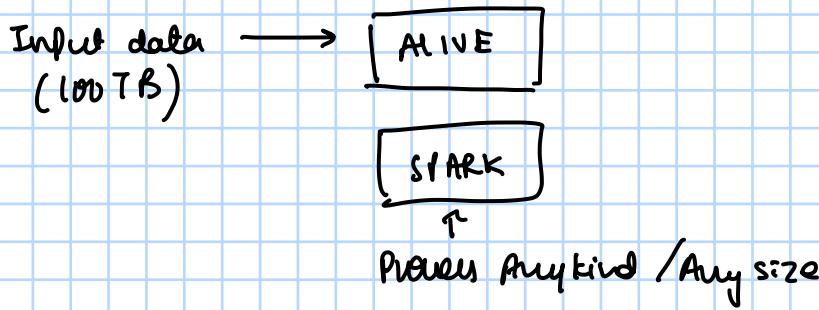


even for SQL queries - we can have key/value using Map/Reduce

- Helps us deal with big data

- Sound/Video/Text/Audio/Images all drawable

\therefore does not matter if data is structure or not



30th Oct :

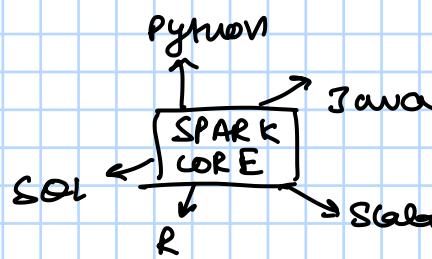
map reduce - strategy - start from output
- input to reduce is what

Spark -
- SQL
- ML
- graph processing

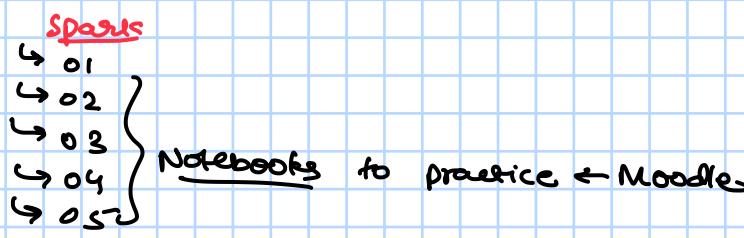
Speed of Spark - in-memory memory

↓
Cloud computing - cluster

lazy execution - code executed only where required
- trace off - errors don't get caught immediately



1st Nov :

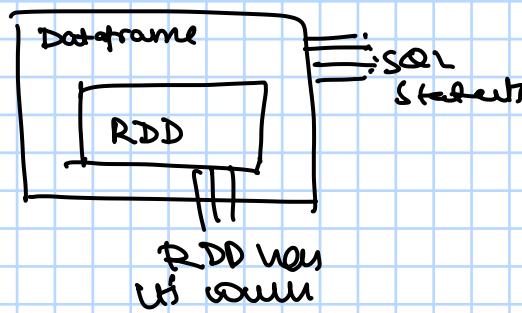


Python list → R D D

↑
Special variable
in spark

6th Nov:

Part -4 - Spark, Load the CSV TO Dataframe



Extract RDD to df:

```
df = df.rdd.map(lambda x:1).reduce(lambda x,y:x+y)
```

We can create a temporary view

```
df = createTempView("stocks")
sq2res = ss.Sq.e ("~~~")
sq2res.show()
```

We can also define the schema using spark

Note - 5th ipynb ← custom schema

6th ipynb: linear reg with spark

```
df1.corr("open","close")
```

Note: we can do linear reg in SPARK

↓
convert to python list
↓
use matplotlib



8th NOV:

Recap

- ① some basics
then done

