

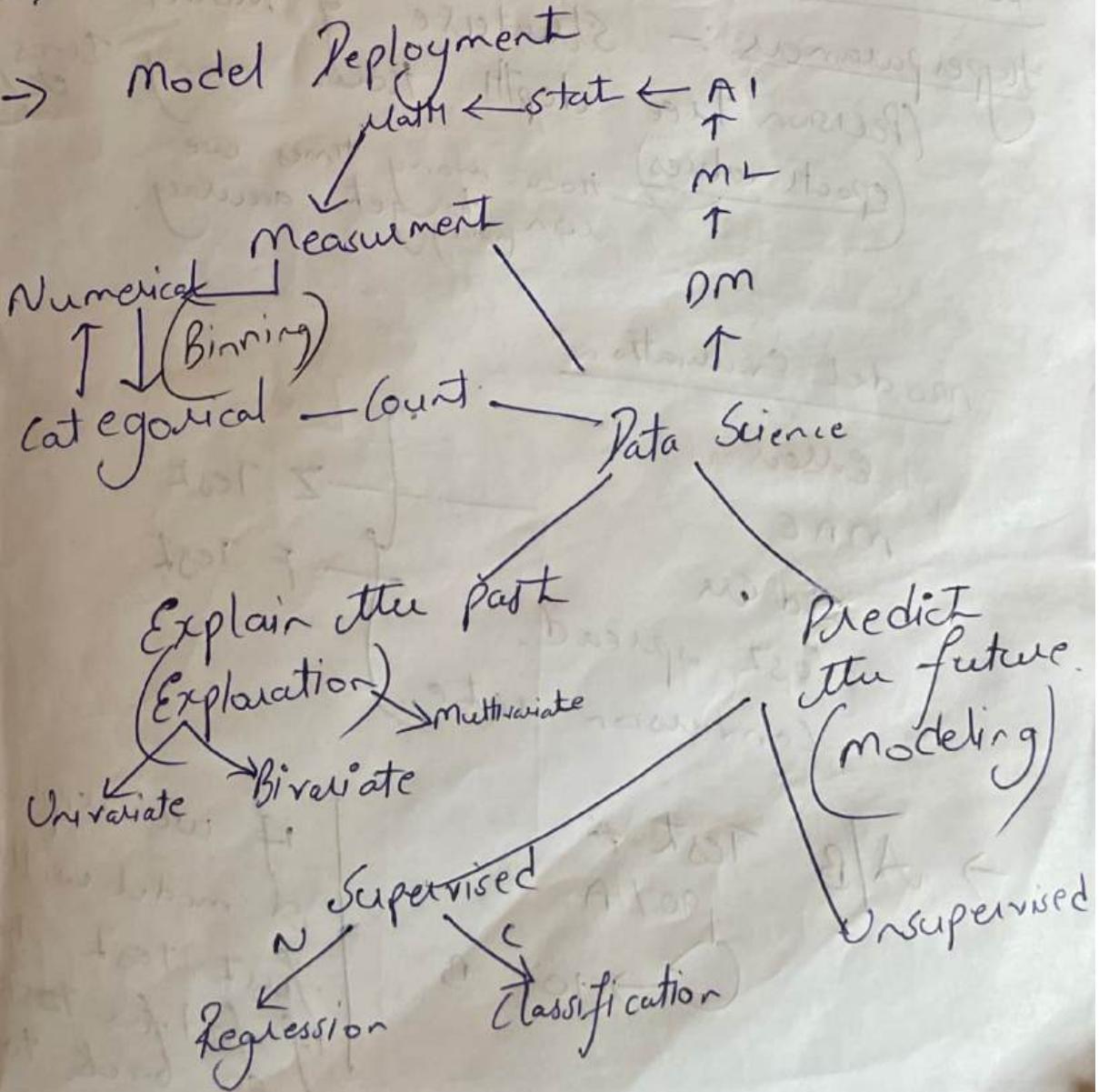
Day-1

[F1L116]

Tell story in eight way with the data

6 steps ★★★

- Problem definition + criteria of success
- ~~Criteria~~ how to get data (Data preparation)
- Data Exploration (with statistics visualization)
- modeling (Regression / Classification)
- Model Evaluation
- Model Deployment



Discovery

→ (junction details was completely empty)

Data present in 3 files &
it had missing values in it

merge - 3 datasets to
(2005-2007) (2009-11) (2012-14)

Structuring

Removed the null values
The column names had special
characters which were replaced

Cleaning

checked all the columns datatype,
Only date column was a character
column, changed it to date

Enriching

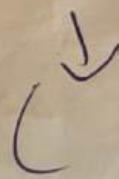
Answering 3rd question (Risk factor
we needed columns like
No of fatalities, accident severity,
frequency of fatalities

followed repeated Discovery, structuring &
cleaning for this dataset file
to our data set.

Validation

We then validate our data &

Publish



Data Normalization (Structuring / Transforming)

Reduces data redundancy & improves?

Data Integrity → [Data stored in logical sense]

Anomaly

Insertion

Update

Deletion

INF

Remove repeating

- Create separate

- Identify primary

Example [Employee & department tables]

- Single valued - each cell
- entries of column are of same type
- Row is uniquely identified.

groups - finding entities

table for each entity

key for each table

Emp Id	Name	Phone	Salary	2 value
		a1087 a1126	30,000 40K	→

→ Atomity → each column should have single value

→ [all attributes depend on key]

2 NF → In INF should not contain partial dependency

Employee id	Dept id	Office location
		→ Non prime attribute

office location only depends on dept id.

→ Split Table into 2

Emp id	Dept id

Dept id	location

3NF → In 2NF
 There should be no transitive dependencies
 for non-prime attributes
 All non prime attributes should only
 depend on prime attributes
 $A \rightarrow C \rightarrow D$ then $A \rightarrow D$

Eg	student id	Student Name	Subject id	Subject	Address

$\text{student id} \rightarrow \text{subject id} \rightarrow \text{subject}$
 $\therefore \text{student id} \rightarrow \text{subject}$

Transform to

stud id	subject id	stud. Name	address	Subject id	Subject

BCNF (or) 3.5 NF [Boyce Codd normal form]
 Every functional dependency $A \rightarrow B$
 then A has to be super key for
 that table

Eg	student id	subject	Professor

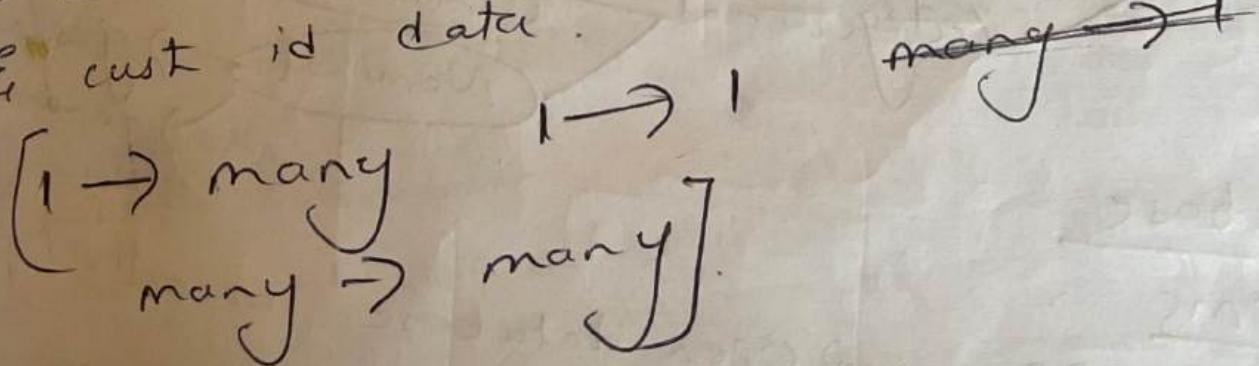
a professor can teach multiple subjects

<u>student id</u>		<u>Professor id</u>	
<u>prof id</u>	<u>subject</u>		<u>Professor</u>
↓ (super key)			

Remove
non prime
attributes
functional
dependency

<u>HNF</u>			
→ No multivalued dependency			
<u>cust id</u>	<u>Cust Name</u>	<u>Address</u>	<u>News letter</u>
20	Alan	10 paul Rob	NYC times
20	Alan	10 paul Rob	TOL
21	smith	AIC philips dr.	TOL

~~if it is one~~
one can subscribe to multiple newsletters
or one can have multiple shipping address
So we have multiple rows of cust name
& cust id data.



Problem definition

Understand domain perspective & convert this to data science problem

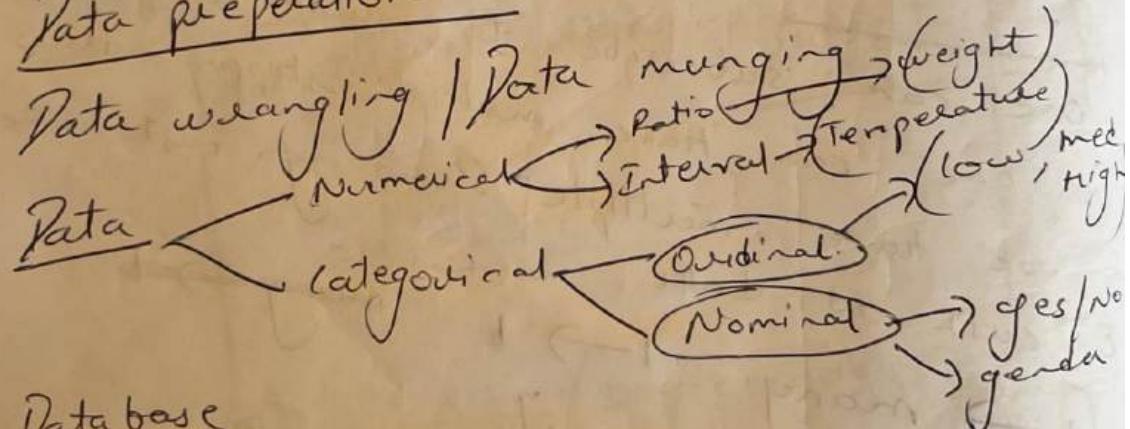
Example of customer hiring a server for my restaurant

Data \Rightarrow Bill amount, tip amount

Linear regression

$X \rightarrow$ Bill amount
 $Y \rightarrow$ Predict tip amount
 on avg per month 200 tables
 \rightarrow total bill (X) \Rightarrow use model to predict
 & hire them in job off

Data preparation



Database

DBMS

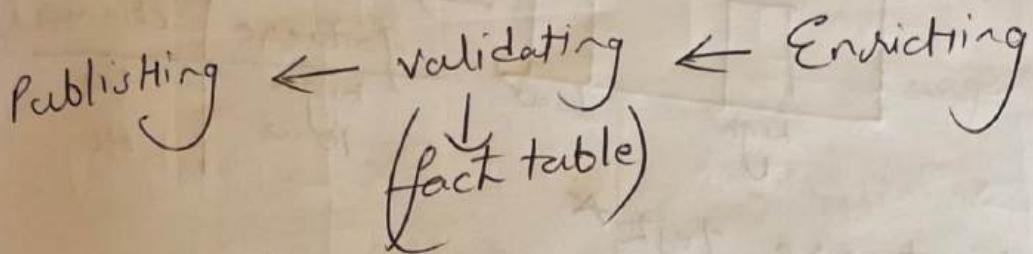
\hookrightarrow ODBC \rightarrow Open database connectivity

JDBC \rightarrow Java database connectivity

<u>SQL</u> (Data definition language)	\rightarrow Create, Alter, Drop Table Create index, drop index	<u>DMQL</u> \rightarrow (Data manipulation language)
		select insert into update delete

6 steps of Data Preparation

Piscovery → Structuring → Cleaning ↴



Data Exploration

Explaining the past

→ Univariate

→ Bivariate

Categorical

→ frequency counts, bar & pie charts

Numerical

count, Mode, Variance, sd, Skewness, → Histogram
Min, max, range, mean, median, quantile → Box plot

Categorical ↗ (binning)
Numerical ↗ (encoding)

binning
1P - 30 }
31 - 50 }

	up	flat	down
up	1	0	0
up	1	0	0
down	0	1	0
flat	0	0	1

Bivariate analysis

→ Numerical - Numerical
Scatter plot



Linear correlation

$$\rho = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \cdot \text{Var}(y)}}$$

If ρ is -1 negative correlation
if ρ is $+1$ positive correlation

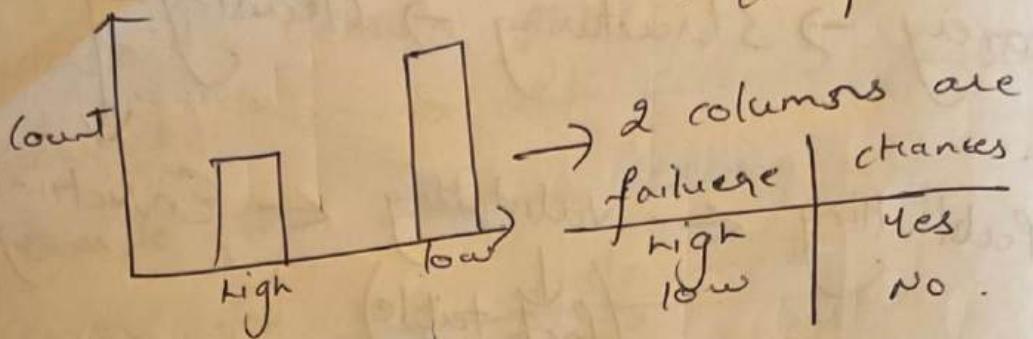
$$\text{Cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

$$\text{Var}(x) = \frac{\sum (x - \bar{x})^2}{n}$$

$$\text{Var}(y) = \frac{\sum (y - \bar{y})^2}{n}$$

Categorical - Categorical

→ Stacked column
combination charts
Chi square test



* Chi square test *

Difference b/w expected frequency (e)

observed frequency (n)

→ Returns a probability for computed Chi square & e $\frac{\text{degree of freedom}}{\text{Total values} - 1}$

$P = 1$ Two categorical variables are completely independent

$P = 0$ Completely dependent

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - e_{ij})^2}{e_{ij}}$$

$$df = (r-1)(c-1)$$

$$e_{ij} = \frac{n_i n_j}{n}$$

$$P = \sqrt{\frac{\chi^2}{n \sqrt{(c-1)(r-1)}}}$$

$r, c \Rightarrow$ rows, columns

Categorical - Numerical

Visuals

line charts

combination charts

Z test, t test, ANOVA.

Z test & T test

decide if averages of 2 groups
are statistically different from
each other.

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}$$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_2 \left(\frac{1}{N_1} - \frac{1}{N_2} \right)}}$$

(We use t test if N_1 or N_2
are less than 30)

If probability of Z it is small, difference
b/w 2 averages is more significant.

ANOVA - (analysis of variance)

Averages of more than 2 groups
are different statistically from each other

	df	ss	ms	F	P value
Regression	k	SSR	$ms_R = \frac{SSR}{df}$	ms_R / ms_E	$p(F)$
Residual	$n - k - 1$	SSE	$ms_E = \frac{SSE}{df}$		
Total		SST			

$SSR = \sum (y_i - \hat{y})^2$ $SSE = \sum (y - \hat{y})^2$ $SST = \sum (y - \bar{y})^2$
 Coefficient of determination $R^2 = \frac{SSR}{SST}$

Missing values - Rectification

Numerical

Average

Mean

Median

Categorical

Missing value as a separate category
Majority value

Numerical

Binning

- Bins of equal width

Temp

Temp-bin

Play off

n

n

y

y

85

80°-90

80

80-90

70

70-80

68

60-70

Numerical
Assumption

variable

that it has a gaussian distribution

(pdf) - mean & sd.

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

probability

distribution

Simple / multiple Regression

SSE - Sum of Squared Errors.

Goal is to create a linear model that minimizes the sum of squares of residuals (SSE)

Reduce SSE

[Comparison with only dependent variable]

Linear Regression Assumptions

- Linear Relationship
- No Multicollinearity
- Multivariate normality - (Residuals are normally distributed)
- No autocorrelation
- Homoscedasticity - (Same variance)

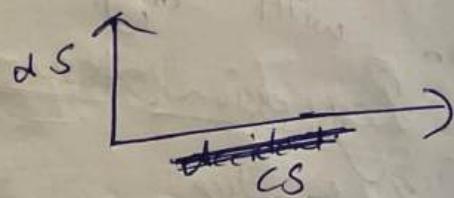
x_1 Number of vehicles	x_2 number of casualties	x_3 Casualty Severity	y Accident Severity
--------------------------------	----------------------------------	-------------------------------	-----------------------------

if we don't have any data & only accident severity then we predict the next accident severity by taking a mean of accident severity

SSE - sum of square of residuals
 \rightarrow Any best fit line would reduce this

Step 1 - Scatter plot

b/w Accident Severity & Casualty Severity



Step 2

find correlation coefficient

(It shows if there is any strong correlation b/w 2 columns)

Centroid $(\text{mean}(x), \text{mean}(y))$
Best fit will pass through centroid

$$\hat{y} = 60 + 6x_i \rightarrow \text{slope.}$$

$$SSR = \boxed{SST}_{\text{with mean}} - \boxed{SSE}_{\text{with best fit line}}$$

$$r^2 = \frac{SSR}{SST}$$

coefficient of
determination

Multiple Regression to multiple regression
Adding more IV's does not mean
the regression will get better.

Multi collinearity

more IV creates more relationships b/w them. If IV's are related to each other then

Ideal \rightarrow IV to be correlated to independent variable but not with each other

this leads to wider confidence intervals & they lower the statistical significance of regression coefficients

Do correlation, scatter plot & simple regression for each IV.
To avoid overfitting & multicollinearity.

multiple
regression
model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \text{error}$$

Multiple Regression Eqn: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

Estimated Multiple Regression Eqn: $\hat{y} = b_0 + b_1 x_1 + \dots + b_p x_p$

b_0, b_1, b_2, \dots estimates of $\beta_0, \beta_1, \dots, \beta_p$

\hat{y} = predicted value of dependent variable
Variance inflation factor (VIF)

$$C_{ii} = \frac{1}{1 - R_i^2}$$

> 5 (Multicollinearity is a problem)

Example

Q3. If the total sum of squares can be explained by using the estimated regression equation

① $SSR \Rightarrow$ Variance explained by the predictors

$SSET \Rightarrow$ Variance that has to be explained by the predictors (Independent variables)

$SSE -$ Variance not explained by the predictors.

Adjusted n - number of items in the set

k - means no of cell means or groups (No of columns)

$$F = \frac{\text{Explained variance (or) variance}}{\text{Unexplained variance. variance}} = \frac{\text{b/w group variability}}{\text{within group variability}}$$

T-test

To check if 2 sets of data are significantly different from one another. A null hypothesis is used to test for significant difference.

confidence interval for regression coefficients is an interval that contains β_i (true) in 95% of all samples.

VIF estimates how much the variance of regression coefficient is inflated due to multicollinearity in the model.

T test

$$T \text{ value} = \frac{\text{mean}(1) - \text{mean}(2)}{\sqrt{s_e^2 (1)}}$$

If $|t|$ is greater than 1, there is more signal i.e. there is more noise.

Null Hypothesis - No difference between the samples

Critical value if $t < \text{critical}$ don't reject
 $t > \text{reject}$

if t_i is in CI 95%.

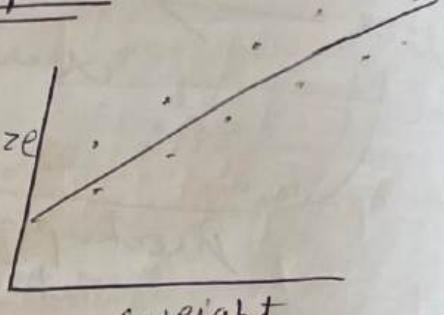
We calculate the P value.

E.g. this would say if we have to reject the null hypothesis.

Ridge & Lasso Regression

Regularization Techniques

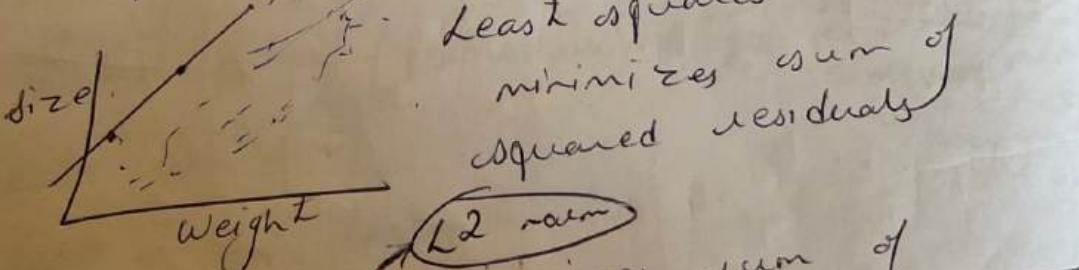
Least squares line size weight
 (what if we only have 2 points)



Training data - least square will be 0
 But testing data has high error (overfitting)

so we introduce a bias to the new line fits the data
 By slightly worse fit, ridge regression can provide better long term prediction.

Least squares method minimizes sum of squared residuals
 Ridge regression minimizes sum of squared residuals + $\lambda * \text{the slope}^2$



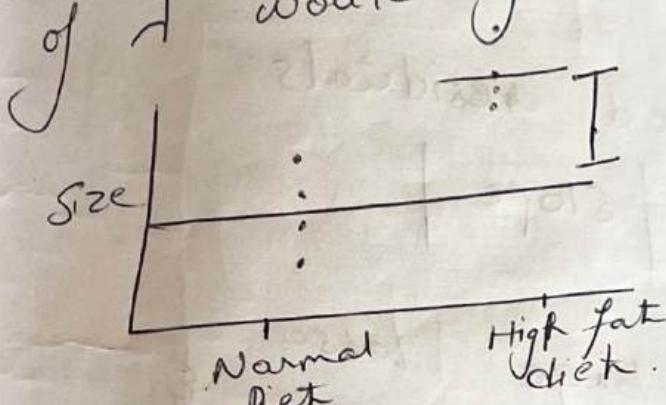
Lambda determines how severe the penalty is
 it has a small amount of Bias but the variance is low

The slope of ridge regression is small, hence less sensitive to weight than least squares line

$$d \Rightarrow 0 \text{ to } \infty \quad \left| \begin{array}{l} \text{If } d \text{ is large} \\ \text{model becomes} \\ \text{less sensitive to} \\ \text{weight} \end{array} \right.$$

↓
sum of
squares
(least squares)

10 fold cross validation - a high value of d would give lowest variance.



$$\text{Size} = 1.5 + 0.7 \times \text{High fat}$$

When applied to logistic regression ridge regression optimizes the sum of the likelihoods instead of squared residuals

$$\text{Sum of likelihood} + d * (\text{slope})^2$$

→ Ridge regression - Every parameter except for the y intercept is scaled by the measurements

$$d(\text{slope}^2 + \text{diet difference}^2 + \dots)$$

$$T = y + \text{slope} \times \text{weight} \rightarrow \text{minimum 2 data point}$$

If we have 100 parameters, we need 100 data points at least to predict the parameters

least squares can't find an optimal solution when we have less data points than parameters (gene data example)

But Ridge regression can find a solution with cross-validation & the regression penalty that favours smaller parameter values.

Lasso Regression

Sum of squared residuals
+ $\lambda \times |\text{slope}|$

Differences

when λ is large Lasso can shrink the slope to 0 & Ridge can only shrink asymptotically close to 0.

$$\lambda(|\text{slope}| + |\text{diet diff}| + |\text{heart rate} + \text{location}|)$$

\swarrow shrinks a little \searrow becomes 0
removes the silly stuff

Lasso can exclude useless variables from equations, so it's better than Ridge regression.

Ridge is better when all variables are useful.

Logistic Regression

[Predicts True/false instead of predicting something continuous like size]

Used for classification

if > 0.5 obese.

< 0.5 not obese

Probability.

Complicated models have more x values.

We check if other variables test to see if they are significantly different from 0.

(Wald's test)

gives probabilities & classifies new samples using continuous & discrete measurements.

Big difference b/w linear & logistic

How line fits to the data

Maximum likelihood - Pick a probability, scaled by weight of observing an obese mouse - just like the curve above & use that to calculate likelihood of non-obese mouse



curve with maximum likelihood is selected.

Scale of coefficient is log(odds)

Logistic Regression & Perception Algo

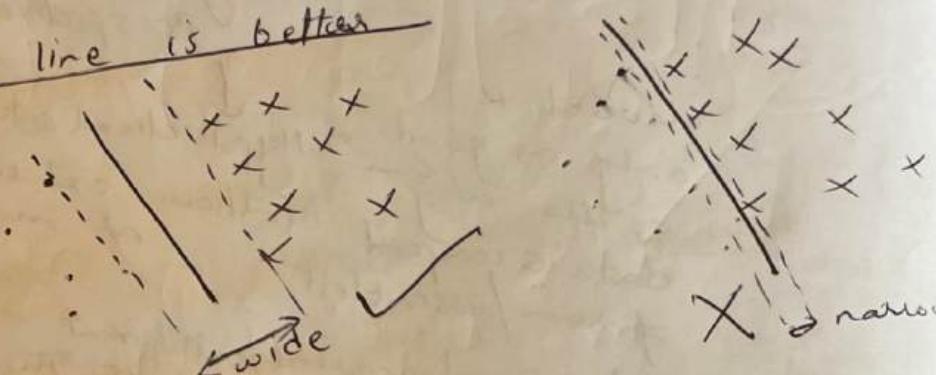
- Start with random line, with blue & red sides.
- Pick a large no - 1000 (No of repetitions / epoch)
- Repeat 1000 time
 - Pick a random point.
 - if classified correctly
continue.
 - else.
move towards the point.

Result - we get a line that separates the data.

Support Vector Machine (SVM)

It tries to separate points of 2 classes using a line (Best line, farthest from pts)

which line is better



How to separate line?

$$\text{Eqn} \Rightarrow 2x + 3y - 6 = 0$$

Expanding rate (0.99)

Algorithm

1) Start with a line &
2 parallel lines to it

2) Pick a large no (No of repetitions / epoch)

3) Expanding factor - Pick close to 1

4) Repeat 1000 times:

- Pick random point

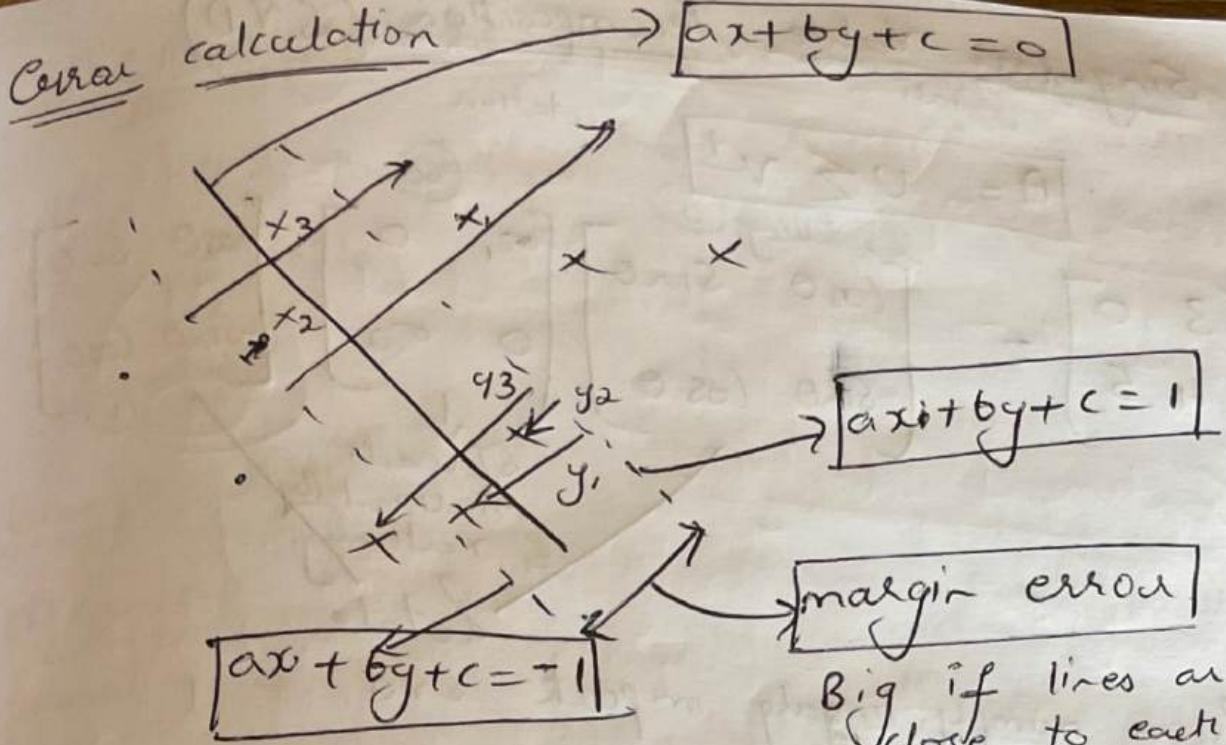
- If point is classified correctly
continue

else

- Move line towards point.
- Separate lines using expanding factor

$$0.99(2x + 3y - 6) = 1$$

$$0.99(2x + 3y - 6) = -1$$



Big if lines are close to each other
Eq vice-versa.

$$\text{Margin} = \frac{2}{\sqrt{a^2 + b^2}}$$

$$\text{margin error} = a^2 + b^2$$

$$\text{Error} = C [\text{classification error}] + \text{margin error}$$

$$\frac{d(\text{Error})}{da} = 2a \quad \frac{d(\text{Error})}{db} = 2b$$

$$a \rightarrow a - \eta 2a = a(1 - 2\eta)$$

$$b \rightarrow b - \eta 2b = b(1 - 2\eta) \quad \begin{array}{l} \eta \Rightarrow \text{learning rate} \\ \text{Expanding factor.} \end{array}$$

C parameter \Rightarrow we use it to choose if classification error is important or margin error is important.
Eg Large C \Rightarrow importance to classification error.

Hyperparameters \Rightarrow C, expanding factor, learning rate

Singular Value Decomposition (SVD)

rotation rotation

$$A = U \Sigma V^+ \quad \text{2 scalings} \quad \text{③}$$

$$\begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \quad \text{①}$$

~~stretch horizontally & vertically
rotate~~

from numpy.linalg import svd.

Rank-1 matrix

$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ -1 & -2 & -3 & -4 \\ 2 & 4 & 6 & 8 \\ 10 & 20 & 30 & 40 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 10 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$

8 numbers

16 numbers

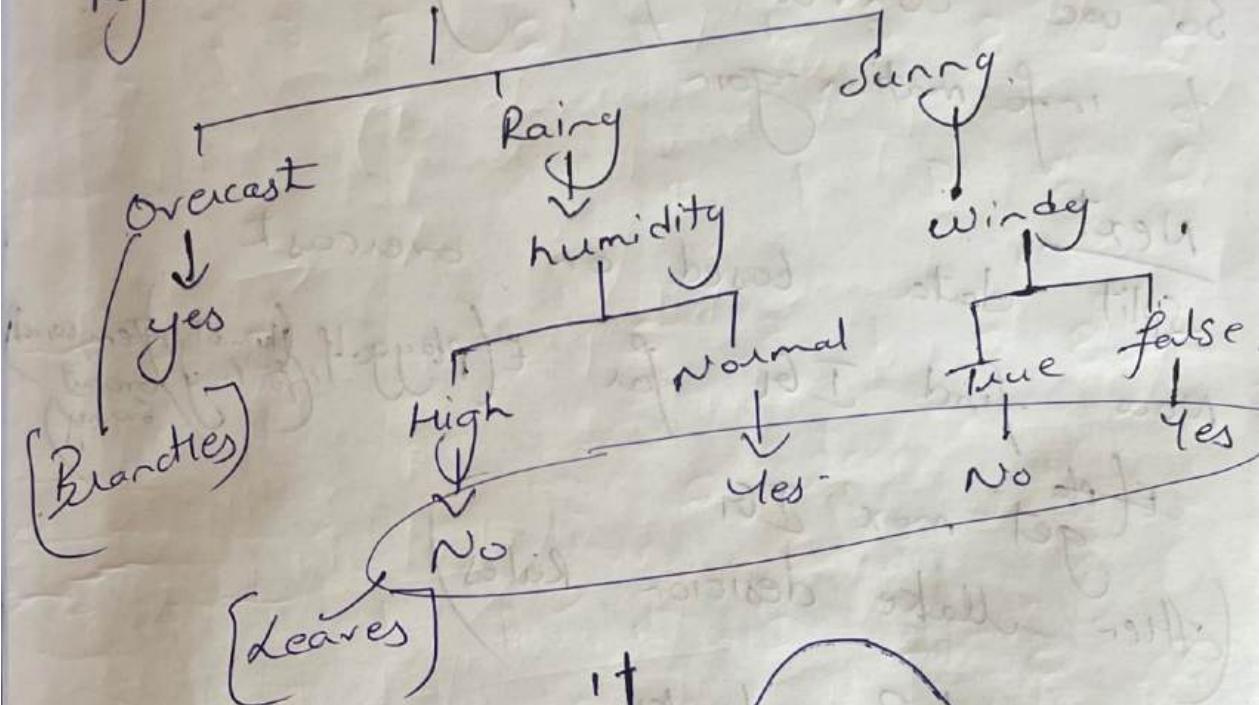
$$\begin{bmatrix} \times & \times \\ \times & \times \\ \times & \times \\ \times & \times \end{bmatrix} \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix} \boxed{\text{Rank 2}}$$

(ignore lower σ values
approximate original matrix)

Svd gives 3 outputs $\Rightarrow U, \Sigma, V^T$

Decision Trees

→ can build both classification & regression models.



Entropy

$$E = -p \log_2 p - q \log_2 q$$

$$= -0.5 \log_2 0.5 - 0.5 \log_2 0.5$$

$$\Rightarrow \text{Entropy}(3,2) = \text{Entropy}\left(\frac{3}{5}, \frac{2}{5}\right)$$

$$= -(0.6 \times \log_2 0.6) - (0.4 \times \log_2 0.4) = 0.971$$

frequency Table

outlook		play golf		E(play golf, outlook)
		Yes	No	
Overcast		4	0	
Sunny		3	2	
Rainy		2	3	

$$= 0.693$$

$$\text{Information gain} \rightarrow \\ \text{gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

So we convert frequency table
to information gain

Next, split data based on overcast!
Now find I_G for $E(\text{play golf}, \{\text{Humidity, Temperature, Windy, Sunny}\})$

If we get max I_G
(then make decision rules)

Issues / Drawbacks

missing values - 2 sol.

Numerical attributes \rightarrow binning

overfitting & pruning

attributes with many values

(super attributes)

$\text{gain}(T, X)$ is tighter for super attribute

& super attributes will easily be selected as root. To avoid this

we calculate gain ratio

$$\text{split}(T, X) = - \sum P(x) \log_2 P(x)$$

gain
split

$$1 - \sum (P_i)^2 = G_{\text{uni}}$$

Gini index - demonstrates degree of inequality in distribution
0 - only 1 class 1 - random dist $| 0.5 - equal dist$

Decision Trees - Regression

~~SDR~~ =
stop if coefficient of variation
is $\leq 10\%$.

Average is assigned to the leaves

split is done by calculating

$$\text{Average} = \bar{x}$$

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

Coefficient of variation (C_v) = $\frac{s}{\bar{x}}$

SD b/w 2 attributes

$$s(T, x) = \sum P(c) s(c)$$

standard deviation Reduction

$$SDR(T, x) = s(T) - s(T, x)$$

[Split on max SDR]

Random forest

Binary classification problem (Eg)

Decision trees are highly sensitive to training data & thus could result in high variance

- Hence random forest - which is much less sensitive to training data
- We Bootstrap datasets \rightarrow making multiple datasets at random from the dataset
- Train decision tree on each of the bootstrap dataset independently
(with feature subsets)
(random) \leftarrow this reduces correlation b/w trees
- Pass new data through each tree & note down predictions
- We combine all the predictions & take majority predicted value as the answer

Bootstrapping + Aggregation

(Bagging)

Ideal size of feature size = $\log/\sqrt{\text{size of features}}$

for regression problem - while combining the predictions, just take the avg

Bagging

DT 1

DT 2

DT 3

DT 4

DT 5

Bootstrap + Aggregation

Split data - n Tree

majority
prediction
are taken

Random forest

Randomly

feature bagging - Select subset of p features dim at each split. This is to deliberately avoid the features that lead to similar splits in DT.

Therefore it reduces correlation b/w each DT.
 $\Leftrightarrow \uparrow$ predictive accuracy on average.

Use \sqrt{p} features (rule of thumb)

The proportion of out of bag samples that were incorrectly classified is the out of bag error.

is

$P_{\text{incorrect}} = \frac{\text{Number of misclassified samples}}{\text{Total number of samples}}$

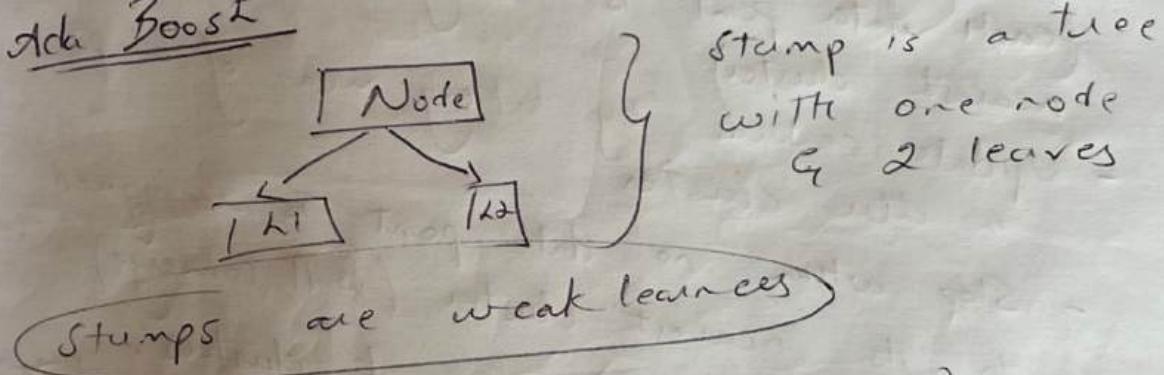
\Leftrightarrow round off at each split

best for different bunches of setting to get correct P value for

Boosting

Idea is to iteratively learn weak machine learning models on a continual updated training data set & add them to produce final result
 Differs from bagging (simply take avg of entire models)

Ada Boost



In fastest of stumps (ada boost), some stumps have higher say than others
 → order is important.

Calculate gini index.

$$\text{Total error} \Rightarrow \frac{\text{Error}}{n} = \frac{1}{2} \log \left(\frac{1 - \text{Total error}}{\text{Total error}} \right)$$

Amount of say for each stump

(If stump does a good job dos is large +ve value)
 New sample weight = $\begin{cases} \text{sample weight} \times \text{amount of say} \\ (\text{to increase}) \end{cases}$

$$(\text{to decrease}) = \text{s.w} \times e^{-\text{amount of say}}$$

→ normalize it \rightarrow random number \downarrow
 take subset

Batch Normalization

Attention

Gradient Boosting

- Boosting \Rightarrow converting weak learners to strong learners

\Rightarrow Easy Explanation is data boost

3 - 32 Leaves (Trees are larger)

Step-1 - Calculate the average of target label (680)

Step-2 calculate residuals for each sample

Residual = Actual - predicted value.

Step-3 \Rightarrow construct D.T
Decision trees predict target variables,
in GB it tries to predict the residuals

If some residuals end up in same leaf
we take the average.

Step-4

Predict target labels using all the trees within ensemble.

Predict value = avg price + $\frac{\text{Learning rate} \times \text{Residual}}{650}$ predicted by DT

Step 5

Compute new residual

$$\frac{\text{Actual value} - \text{Predicted value}}{650}$$

GB performs better on unbalanced data (dist.)
Random forest builds each tree
 \rightarrow multi class obj detection

Regression

Step 6
 Repeat 3-5 until for a given hyperparameter - No. of trees does not significantly reduce the size of residuals (i.e. no of estimators)

Step 7
 Once trained, we use all the trees in ensemble to make a prediction as to the value of target variable

$$\text{Avg} + \frac{0.1 \times \text{Residual}_{gt1}}{T_1} + \frac{0.1 \times \text{Residual}_{gt2}}{T_2} + \dots$$

Classification

$$\text{Prediction} = \frac{\log(\text{odds})}{\text{Residual} = \text{Actual} - \text{Predicted}}$$

Eg
 1 - Yes
 2 - No

$$\log\left(\frac{4}{2}\right) = 0.7$$

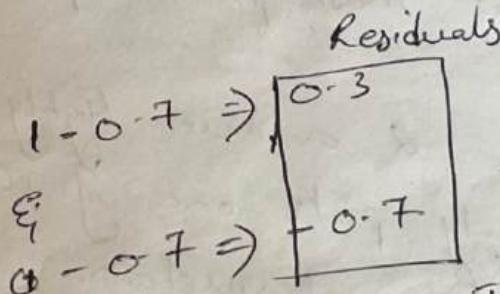
Initial prediction

$$P(Y) = 1$$

$$P(N) = 0$$

$$\text{Residual} \Rightarrow$$

(Yes)



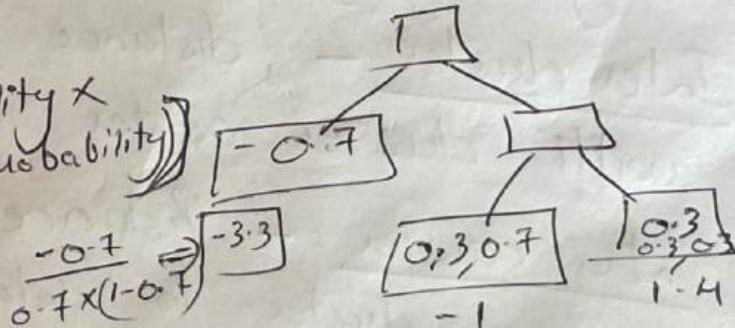
Tree

$$P = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$$

$$P = 0.7$$

$\sum \text{Residual}$

$$\sum \left[\text{Previous Probability} \times \frac{\text{Residual}}{(1 - \text{Previous Probability})} \right]$$



$$\Rightarrow 0.7 + 0.8 \times \frac{(-3.3)}{0.7 \times (1 - 0.7)} = 0.7 + 0.8 \times 1.4 = 1.8$$

Calculate residual

$$\text{Probability} = \frac{e^{1.8}}{1 + e^{1.8}} = 0.9 \Rightarrow \boxed{\text{Yes}}$$

XGBoost

Extreme Gradient Boost

Speed of processing, performance
Advanced - gradient boost

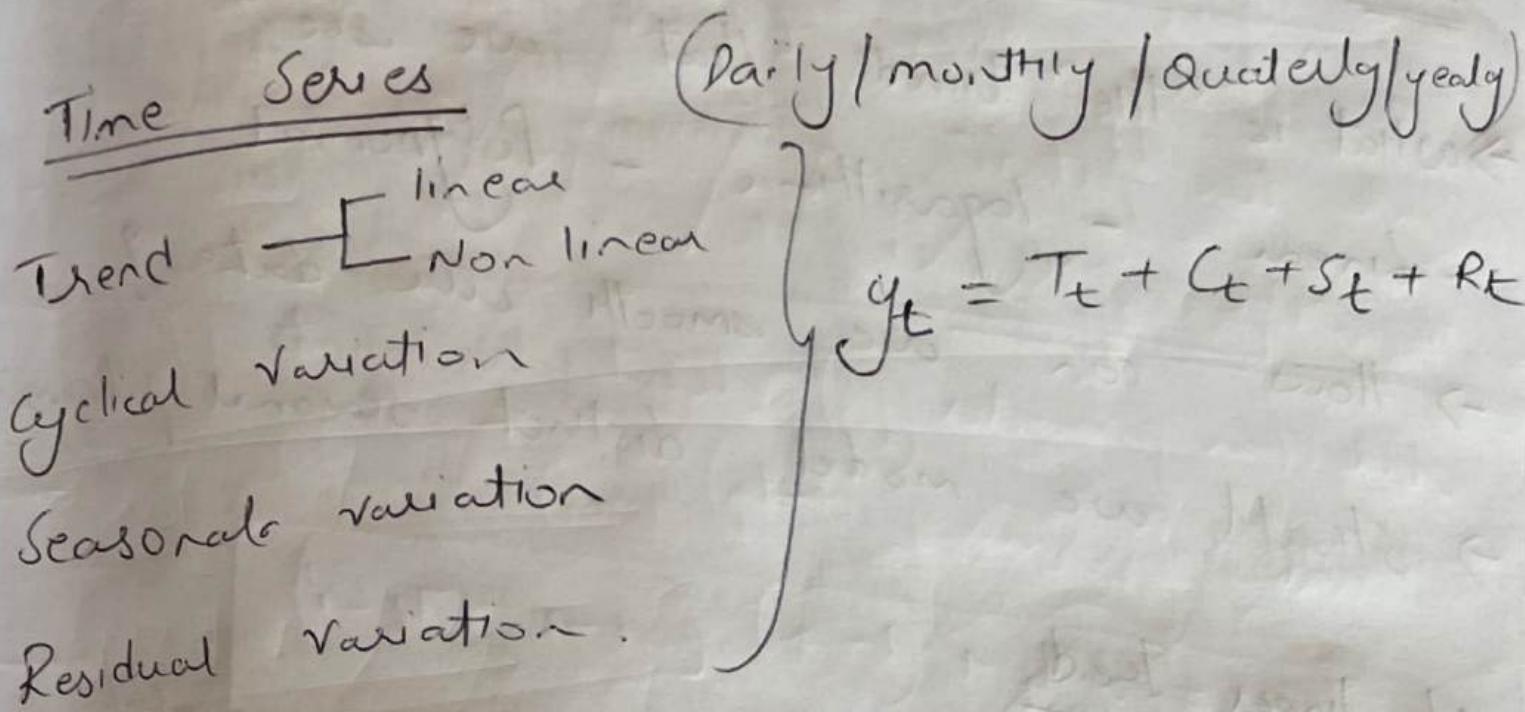
parallelization
CPU's
cache optimization
Out of memory computation

- Each tree boosts attributes that lead to misclassifications of previous trees
- (Regularized boosting to avoid overfitting)

[Hyperparameters - booster.
(gbtree, gblinear, dart)]

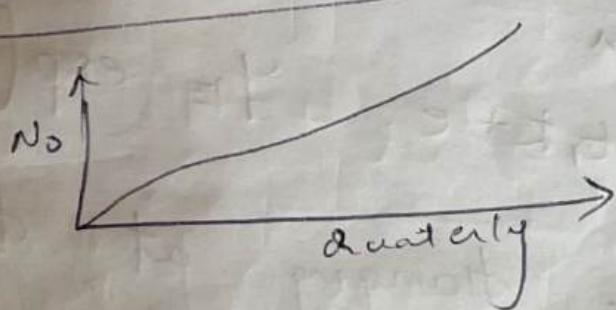
Verbosity - 0, 1, 2, 3 (quality of using more words than needed)

nthread - max number of threads to use XGBoost.



Linear Trend

$$y = b_0 + b_1 t + e$$



Increase by a constant amount at each successive time period - linear trend

Increase by increasing amount - curvilinear trend

Increase by equal percentages - can be made linear by applying logarithmic transformation

$$\log(y) = b_0 + b_1 t + e$$

Time Series model Building

→ What is the trend that we see?

- Linear - logarithmic - Polynomial

→ How can we smooth the data?

→ Should we model distinct seasonal patterns?

Non linear Trend

- Time series changing at increasing rate over time, then logarithmic model in Y works better

$$\log(Y) = b_0 + b_1 t + e \quad Y = \exp(b_0 + b_1 t + e)$$

- Time series changing at decreasing rate over time, then logarithmic model in t works better

$$Y = b_0 + b_1 \log(t) + e$$

Curvilinear Trends

Model with polynomial

Linear (first order) $\Rightarrow Y = b_0 + b_1 t + e$

Quadratic (second order) $\Rightarrow Y = a_0 + b_1 t + b_2 t^2 + e$

Cubic (third order) $\Rightarrow Y = a + b_1 t + b_2 t^2 + b_3 t^3 + e$

* moving Averages *

Compute average of last m consecutive observations.

They merely smooth the fluctuations in the data. It works well when data has

- fairly linear trend
- definite rhythmic patterns of fluctuations

$$Y_{ma}(4) = \frac{Y_t + Y_{t-1} + Y_{t-2} + Y_{t-3}}{4}$$

Exponential smoothing

It is a weighted average that assigns positive weights to the current values \propto to past values of the time series. It also assigns weights to recent values \propto

- Greater weights decrease exponentially as we go further back in time.

$$S_1 = Y_1$$

$$S_t = \alpha Y_t + (1 - \alpha) S_{t-1}$$

* smaller the value of α smoother the plot.

$$\text{Eg } \alpha = 0.5$$

$$S_1 = Y_1$$

$$S_2 = 0.5 Y_2 + 0.5 S_1$$

$$S_3 = 0.5 Y_3 + 0.25 S_2 + 0.25 Y_2$$

$$S_4 = 0.5 Y_4 + 0.25 S_3 + 0.125 Y_3 + 0.125 Y_2 + 0.125 Y_1$$

Two parameters Exponential smoothing
 (α, β) - (Holt's method)

One Parameter (α)

forecasted values beyond the data, into the future remain the same.

Two parameters

Adds a parameter (β) that accounts for trend in the data.

$$\boxed{\begin{aligned} S_t &= \alpha Y_t + (1-\alpha)(S_{t-1} + T_{t-1}) \\ T_t &= \beta(S_t - S_{t-1}) + (1-\beta)T_{t-1} \end{aligned}}$$

forecasted value

$$Y_{t+1} = S_t + T_t$$

S_t = weighted average of current observation & previous forecast value

T_t = weighted average of change in S_t & previous estimate of the trend parameter.

Holt-winters Seasonality
 (α, β, γ)

$$I_p = \gamma \frac{Y_t}{S_t} + (1-\gamma) I_{t-p}$$

The value I_t represents seasonal index at a point p in the season

$$\text{Accuracy of forecast}$$

$$MSE = \frac{\sum_{t=1}^n (y_t - \hat{y}_t)^2}{n}$$

mean absolute Deviation

$$MAD = \frac{\sum_{t=1}^n |y_t - \hat{y}_t|}{n}$$

We choose weight w to minimize MSE
or MAD

Akaike information criteria (Estimator of prediction error)

$$AIC = 2k - 2\ln(\hat{L})$$

\hat{L} - maximum likelihood fn of the model
k - No of estimated parameters

→ It tells us about the relative quality of statistical models for given data set. (Estimates the information criteria representing the best models) (Schwarz IC)

Bayesian Information

$$BIC = k \log(n) - 2 \log(\hat{L})$$

| n - no of data points in observed data

Models with low AIC & BIC are selected.

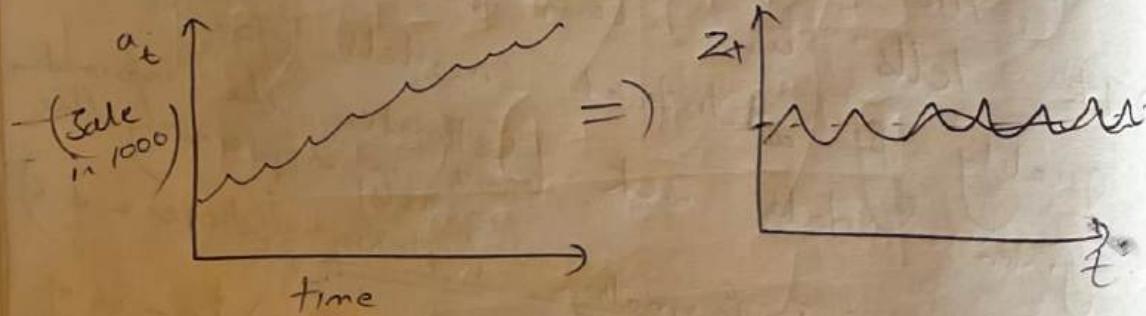
$$x - \text{observed data} \quad \hat{L} = P(x|\theta)^n$$

$\hat{\theta}$ is parameter values that maximize likely hood.

When fitting models, it is possible to increase the likelihood by adding parameters, & but by doing so may result in overfitting. Both BIC & DIC attempt to resolve this problem by adding introducing a penalty term (No of parameters) in the model. Penalty term is large in BIC than DIC.

ARIMA model

Auto Regressive Integrate Moving Average



$$d=1 \quad [z_t = a_{t+1} - a_t] \quad (\text{Differences from one timestamp from its previous})$$

ARIMA - 3 parameters (p, d, q)

$$d=2 \quad z_t = z_{t+1} - z_t$$

simpliest form $\begin{pmatrix} 1 & 1 & 1 \\ p & d & q \end{pmatrix}$

$$z_t = \phi_1 z_{t-1} + \theta_{t-1} + \varepsilon_t$$

get back at

$$a_k \Rightarrow z_{k-1} + a_{k-1} \Rightarrow z_{k-1} + z_{k-2} + a_{k-2} \dots$$

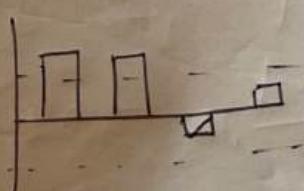
$$a_k = \left[\sum_{i=1}^{k-d} z_{k-i} + a_i \right] \rightarrow \text{last } a \text{ value we know}$$

$$\hat{y}_t = \beta_0 + \beta_1 d_{t-1} + \phi_1 e_{t-1} + e_t$$

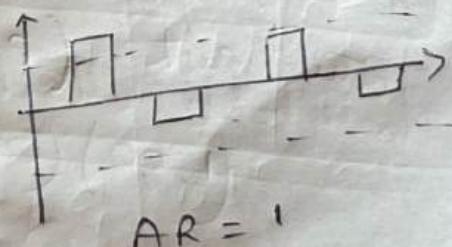
Prediction

$$\hat{L}_t = \beta_0 + \beta_1 d_{t-1} + \phi_1 e_{t-1}$$

MA \Rightarrow ACF, AR \Rightarrow PACF



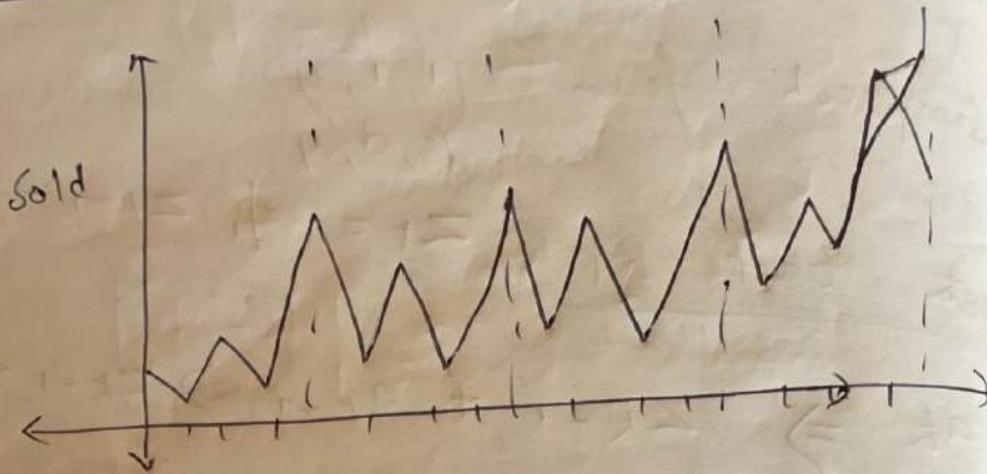
MA $\Rightarrow 2$



AR = 1

SARIMA

Model



ARIMA (P, d, q)

Seasonality - repeating pattern over the year that happens over & over again
 7 parameters - $(P, d, q) (P, D, Q)_m$

$$\frac{(1 - \phi_1 B)^P}{(1 - \theta_1 B)^q} \frac{(1 - \Phi_1 B^4)^P}{(1 - \Theta_1 B^4)^q} \frac{(1 - B)^D}{(1 - B^4)^d} \frac{(1 - B^4)^Q}{(1 - \Theta_4 B^{4m})^Q} y_t = \epsilon_t$$

Back shift time series by 4 coerces in
 other past. $(P \& D)$

$$(1, 90)(0, 1, 1)_4$$

$$z_t = y_t - y_{t-4}$$

$$z_t = \phi_1 z_{t-1} + \theta_1 \epsilon_{t-4} + \epsilon_t$$

Association rules

$$x \Rightarrow y$$

Support Confidence Lift

$$\frac{f_{\text{req}}(x,y)}{N}$$

$$\frac{f_{\text{req}}(x,y)}{f_{\text{req}}(x)}$$

$$\frac{\text{Support}}{\text{supp}(x) \times \text{supp}(y)}$$

first pass - Support of each individual items

second pass

third pass

Apriori Algorithm item set is generated

- Candidate item set of target itemset of previous pass. without considering items in previous pass.

- Large itemset of previous pass is joined with itself. to generate itemsets whose size is higher than

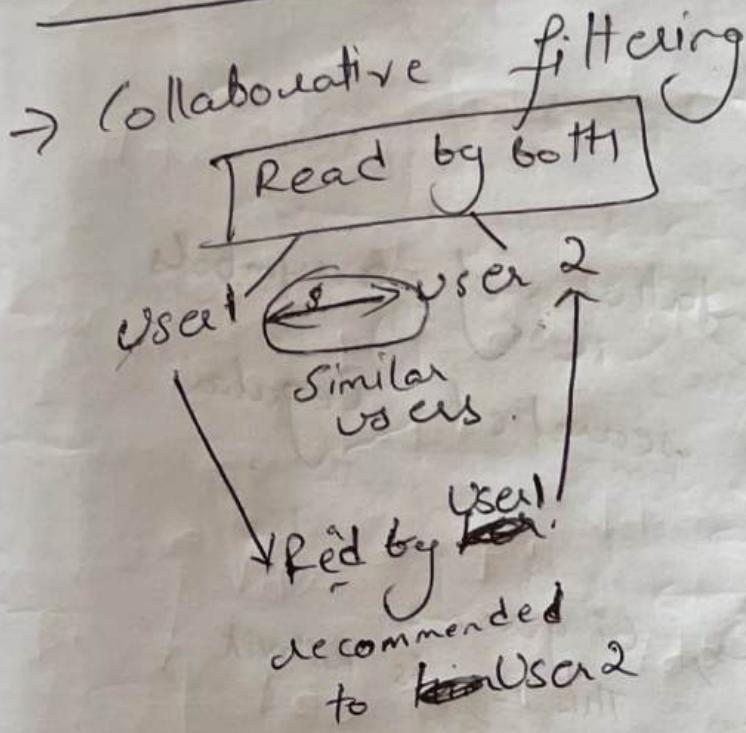
- Each item set that has

- subset with support 1 is deleted

Example
E remaining are taken.

1 2 1	1	2	1 2	2	1 2 3
2 3 4	2	4	1 3	3	1 2 5
2 3 5	3	4	1 5	5	2 3 5
1 2 3 5	1	2	2 3	3	=>
	5		2 5	2	
			3 5	2	
					2 3 => 5

Recommendation Systems



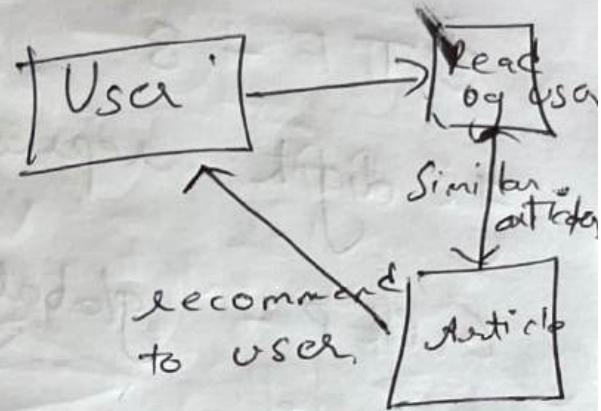
Similar users will like similar things.

- Amazon recommendation
- People who bought ~~Item 1~~ also bought ~~Item 2~~ & ~~Item 3~~
 - Cosine similarity b/w users.

Similarity functions

Cosine, Euclidean, Manhattan, Pearson

→ Content based filtering



Similar content would be liked by user

YouTube recommendation

$$v(x, i) = \cos(x, i) = \frac{x_i}{\|x\| \|i\|}$$

large number \rightarrow small angle
large similarity.

feature extraction is key

3x

2.

general equation of Perceptron

$$\left\{ \begin{array}{l} w_1x_1 + w_2x_2 + \dots + w_nx_n + b \geq 0 \\ y = \text{else} - 0 \end{array} \right.$$

- correct answer - don't change weights
- wrong answer - adjust weights till it classifies it correctly.

$$error = +y - (w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$

$$\epsilon = 1 \quad \left\{ \begin{array}{l} y = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \\ t = 1 \end{array} \right.$$

$$\epsilon = 1 \quad \left\{ \begin{array}{l} y = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \geq 0 \\ y = 1 \end{array} \right.$$

Learning rate (η)

Rate at which learning rate converges to correct answer.

$\eta = 1$ (typically)

Too small - slow convergence

Too large - η can cause oscillation in the process.

$$\Delta w = \eta \times \epsilon \times x \quad | (\epsilon = y - g)$$

XOR will fail with linear model (perception)

Multi layer perception

Input layer Hidden layer Output layer

(Transfer function = Activation function)

$$\text{output layer} - \epsilon_0 = g \times (1-g)$$

$$\text{Hidden layer} - \epsilon_i = g_i \times (1-g_i) \times (w_i \times \epsilon_0)$$

$$\text{weight adjustment} - \Delta w = \eta \times \epsilon \times x$$

Steps

Summation

$$s = \sum w \cdot x$$

To solve local minima
they use relu.

Transfer/activation

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

Sigmoid function
Returns $(0-1) \rightarrow$ Predicts probability of op

Softmax

Back propagation

Propagates the error backward by approximating them to

each unit according to the amount of this error (the unit is responsible for).

3 Neural Networks

(3 Blue 1 Brown)
YouTube

image - 3

Recognize handwritten digits

Plain vanilla concept \Rightarrow Multilayer perception

② \rightarrow neuron

$$28 + 28 \text{ pixels} = 784$$

each pixel
holds grey scale
value
0 - Black
1 - White.

84	0	0	0	0	0
	0	0	0	0	1
	0	0	0	0	2
	0	0	0	0	3
	0	0	0	0	4
	0	0	0	0	5
	0	0	0	0	6
	0	0	0	0	7
	0	0	0	0	8
	0	0	0	0	9

Hidden layers

activations in 1 layer determine activations in another layer.

$\sigma(w_1 a_1 + w_2 a_2 + \dots + w_n a_n)$

compute weighted sum ϵ to make it

a value between 0-1

we use Sigmoid $\sigma(x) \Rightarrow \frac{1}{1 + e^{-x}}$

other activation functions are relu, tanh

Let's say we need weighted sum to be active only if weighted sum > 1.

bias to be inactive

$$\sigma(w_{1,0} + w_{2,0} \dots + w_{n,0} - 10)$$

Layer 1 \rightarrow Layer 2 784×16 weights
16 biases

Learning \Rightarrow finding eight weights & bias to all the weight neurons

math

$$\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n,0} & w_{n,1} & \dots & w_{n,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}$$

(matrix multiplications)

ReLU - Rectified linear unit

$$\text{ReLU}(a) = \max(0, a)$$

if a is < 0 then inactive or else its active

Learning Part (Gradient Descent)

Train - Test

initialize w, b randomly

Cost function

$$\text{odd } (x_0 - y_0)^2 + (x_1 - y_1)^2 + \dots + (x_n - y_n)^2$$

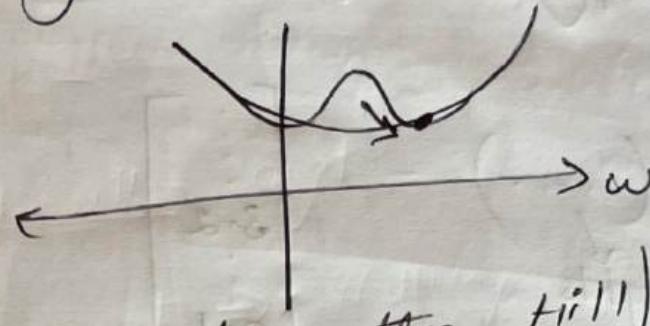
$x_0 \Rightarrow$ last layer neuron activation value

$y_0 \Rightarrow$ Expected value

$$\text{example } (0.43 - 0)^2 + (0.03 - 1)^2 + (0.5 - 0)^2$$

→ compute average cost over many examples

goal is to find min of cost f(\vec{w})



($c(\vec{w})$) = (Ball rolling down the hill)
gradient → the direction of steepest increase $\nabla c(x, y)$

⇒ So we take "-" of gradient gives direction to step.

$$-\nabla c(\vec{w}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.4 \\ 0.82 \end{bmatrix} \xrightarrow{\vec{w}_0(r)} \xrightarrow{\vec{w}_2(t)} \text{How to nudge } \vec{w} \text{ w/ bias.}$$

Gradient descent

Properly labelled data will have faster gradient descent, & they all have many local minima roughly of equal quality.

Since we are using structured data (image + labelled). This is by Human.

Supervised learning.

Back propagation

Concept

$$-\nabla C(\dots) = \begin{bmatrix} \vdots \\ 3.2 \\ \vdots \\ 0.1 \\ \vdots \end{bmatrix}$$

↓
weights

So weights in proportion to a_i will
~~have~~ a Fibonacci theory

[Neurons that fire together
wire together]

Average over all training data
 $w_0, w_1, \dots, w_{13001}$ & change them

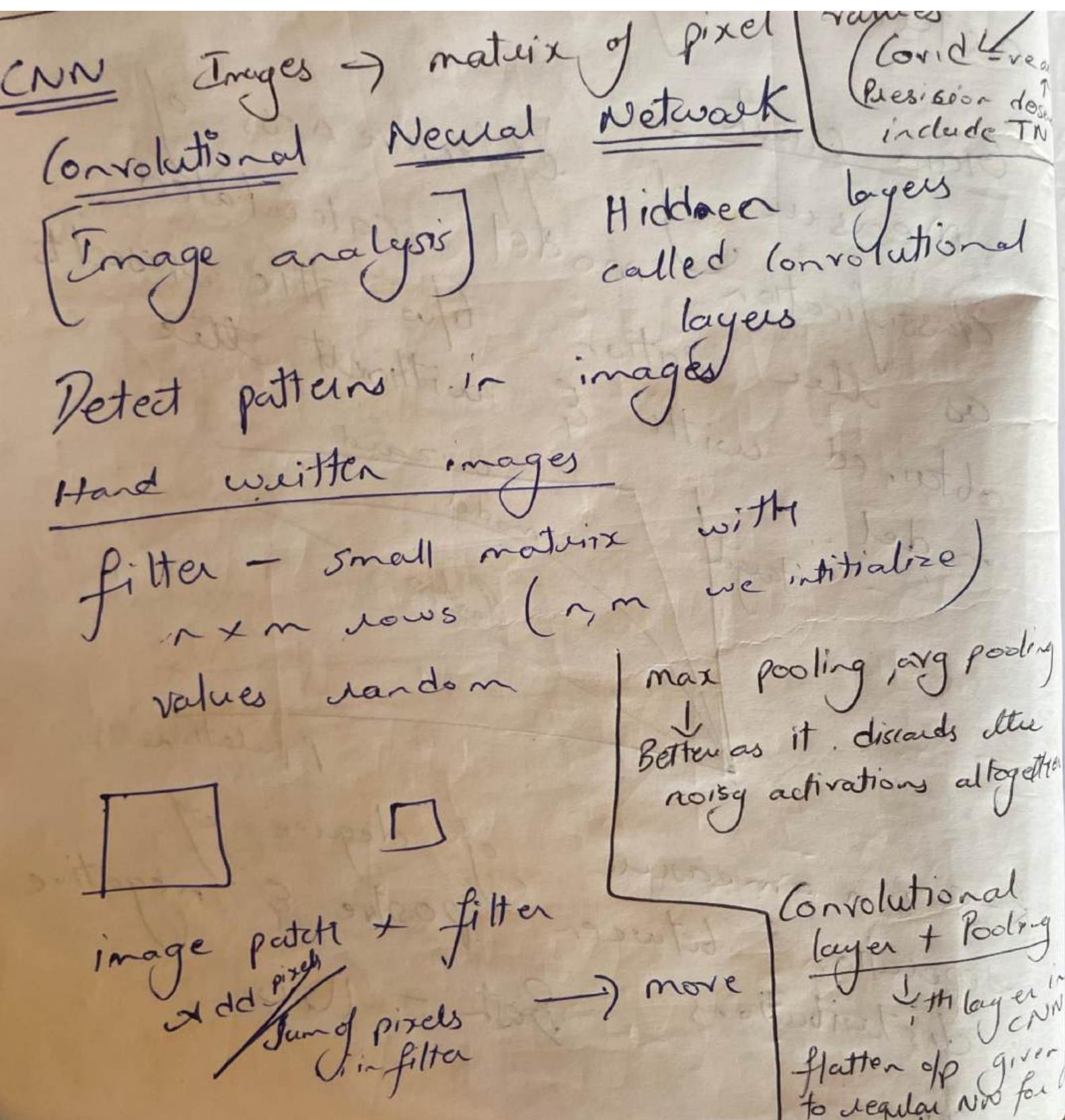
10,000 images

$\rightarrow 100 \times 100 = 100$ mini batches of 100 images each

\rightarrow compute gradient decent step with each
mini batch.

Stochastic gradient decent

Calculation / Calculus of Back propagation



~~Pooling layer~~ \Rightarrow Extract dominant features
stick of images become smaller

- Pick window size
- Pick stride
- walk window across filtered images
- from each window pick the max value (to decrease the computational power required to process data)

* Relu \rightarrow No negative values.
Layers can be repeated several times

Convolution Relu Convolution Relu pooling
(Convolution Relu)
Convolution Relu

Back propagation

$$\text{Error} = \text{(right answer)} - \text{(wrong answer)}$$

Minimize error.

Adjust weights up & down, & see how error changes

Padding → Amount of pixels added to an image when it is being processed by kernel
Extends the area of which CNN of CNN processes an image.

Kernel is a NN filter which moves across the image, scanning each pixel & converting data into a smaller/larger format

In order to assist the kernel to process the image padding is done & it leads to more accurate analysis of images.

Batch Normalization

General technique to the layer (makes NN more stable during training)
before or after activation function

$$\text{log odds} \quad \log \left(P(A) / P(\bar{A}) \right) \quad (\text{Easier for calculations for new data})$$
$$P = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$$

Hiref

$$\log \left(\frac{5}{100} / \frac{95}{100} \right) \Rightarrow -2.94$$

Breaking

$$\log \left(\frac{1}{2} / \frac{1}{n} \right) \Rightarrow \frac{0.69}{-2.25}$$
$$P = \frac{e^{-2.25}}{1 + e^{-2.25}} \Rightarrow 0.1053$$
$$\Rightarrow \frac{0.1053}{1 + 0.1053} \Rightarrow 0.095$$
$$\Rightarrow 9.5\%$$

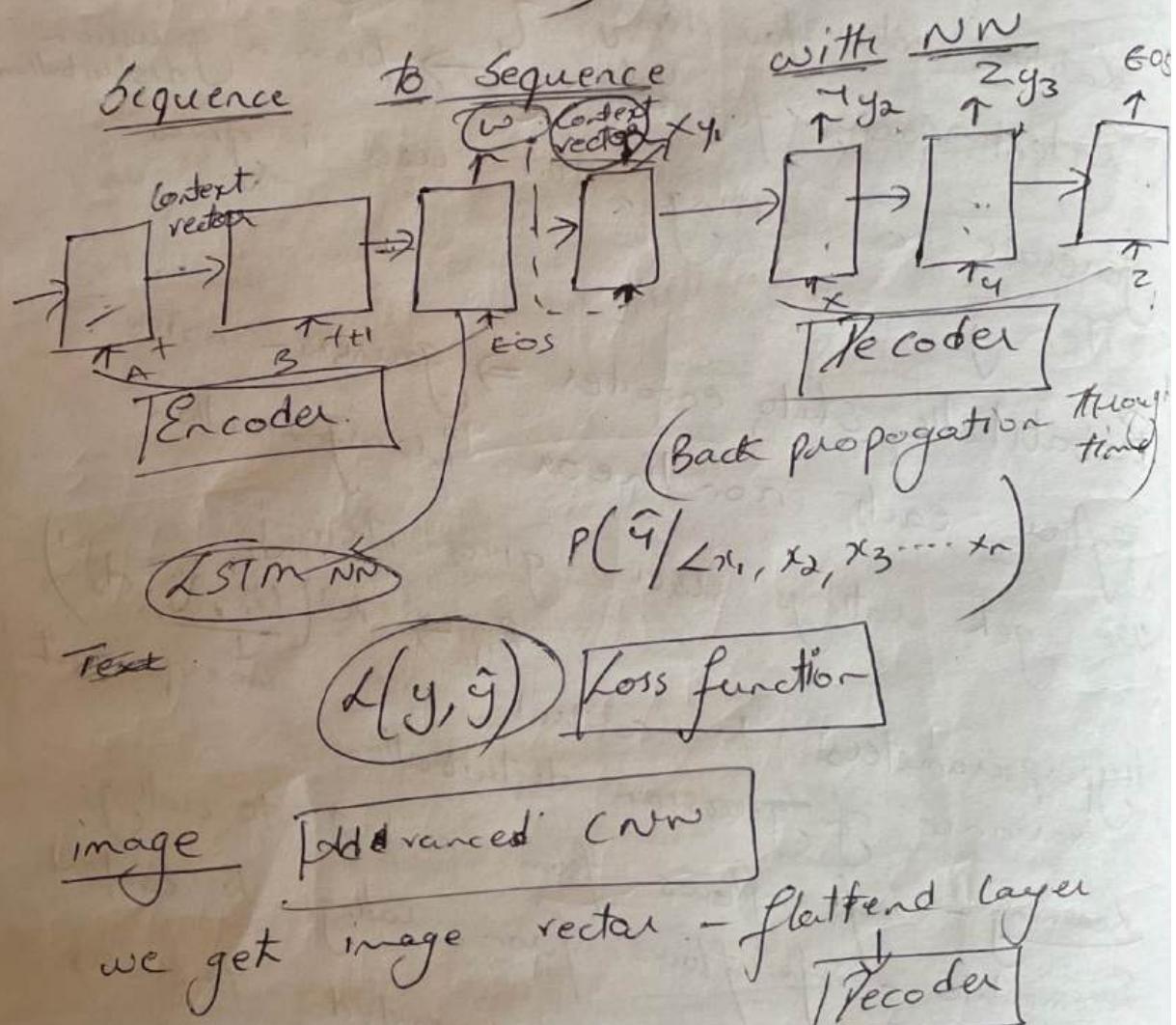
CNN

- feed forward NN using filters & pooling
- Size of i/p & o/p is ~~not~~ fixed.
- spatial data (image)
Image recognition, classification, medical analysis etc

RNN

- decoupling network that feeds the result back into the network
- Size of input & output ~~not~~ may vary
- Good for temporal data / sequential
- Test translation, sentiment analysis (NLP tasks)
Speech analysis.

Max pooling is a pooling operation that calculates the max value of a patch in feature map.
 Usually used after each convolutional layer
 [Translating the image by small amount does not change affect the values of most pooled o/p]



Disadvantages

Long sentence, the accuracy decreases

↓
 Self attention (Bi-directional LSTM)
 architecture.

LSTM & GRU

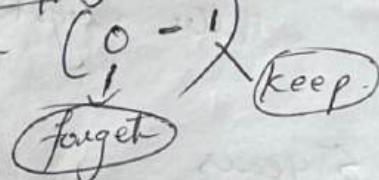
long short term memory, Gated Recurrent
~~Vanishing gradient problem (Recurrent NN) units~~

new gradient value becomes small with
 P time & it doesn't contribute much to
 learning.

internal gates that regulate flow of info
 which data to keep / throw away]

$$\text{Tanh} \Rightarrow -1 \text{ to } 1$$

gates - forget gate, ip gate, op gate
 Sigmoid activation



GRU

Reset gate, update gate

↓
 How much past info to forget)

} less operations
 so little faster
 to train than
 LSTMs

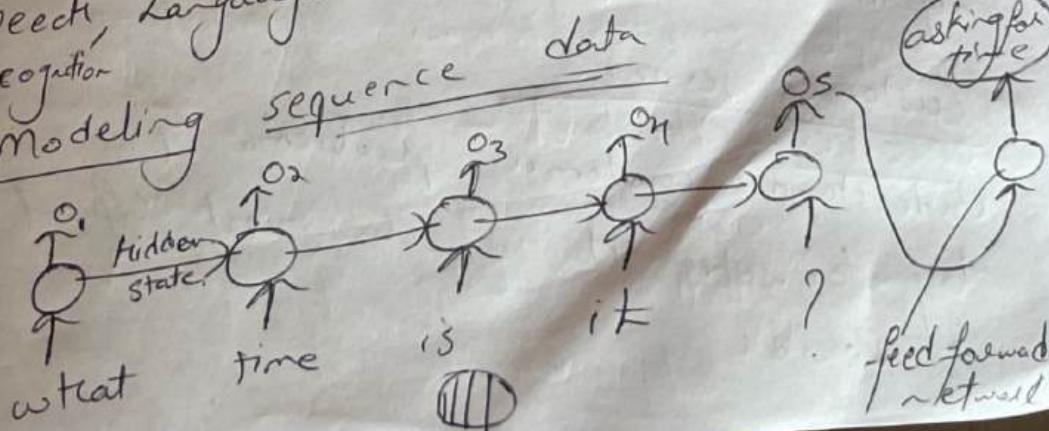
Recurrent N.N

Problem: Speech, Language translation, stock prediction

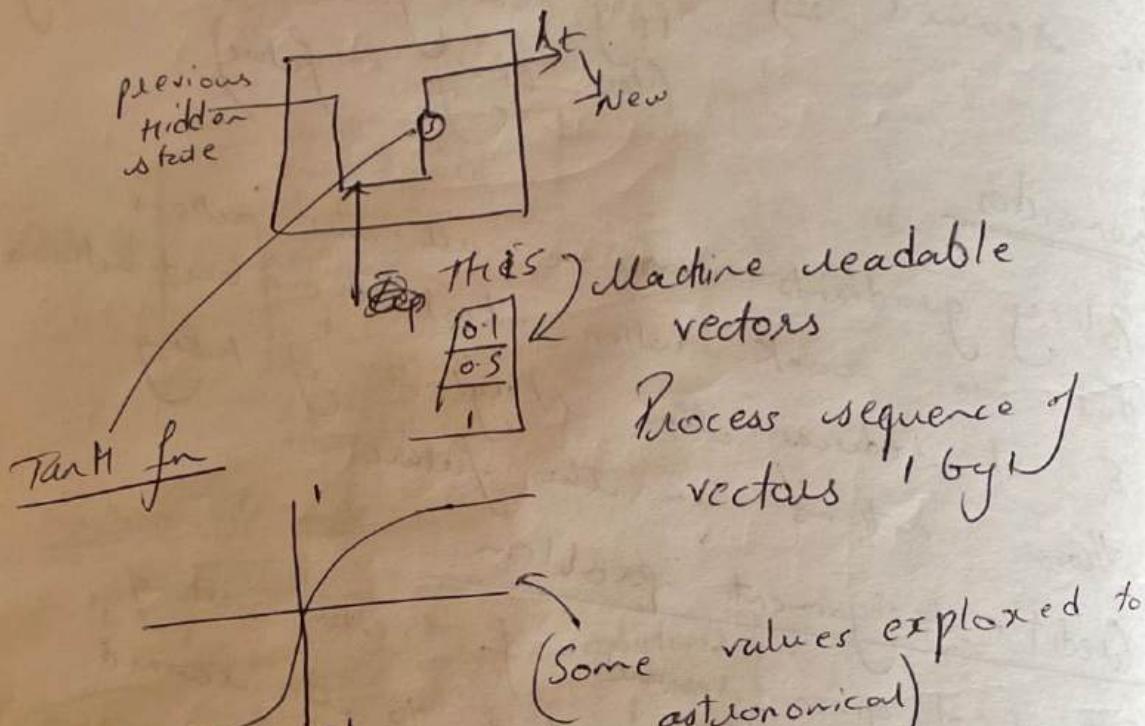
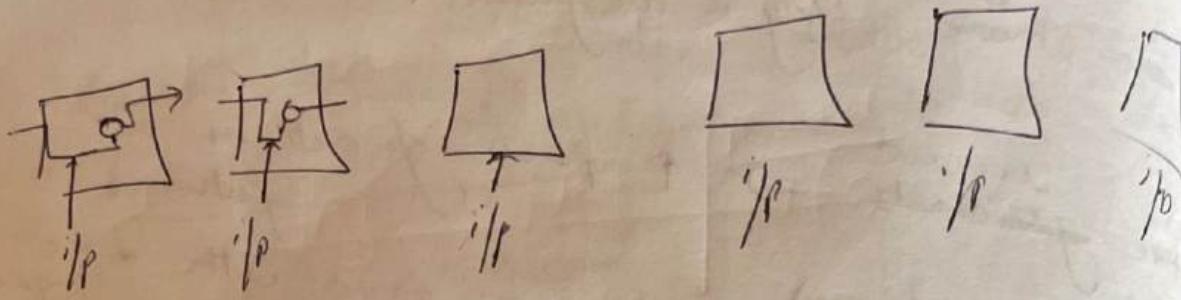
recognition

prediction

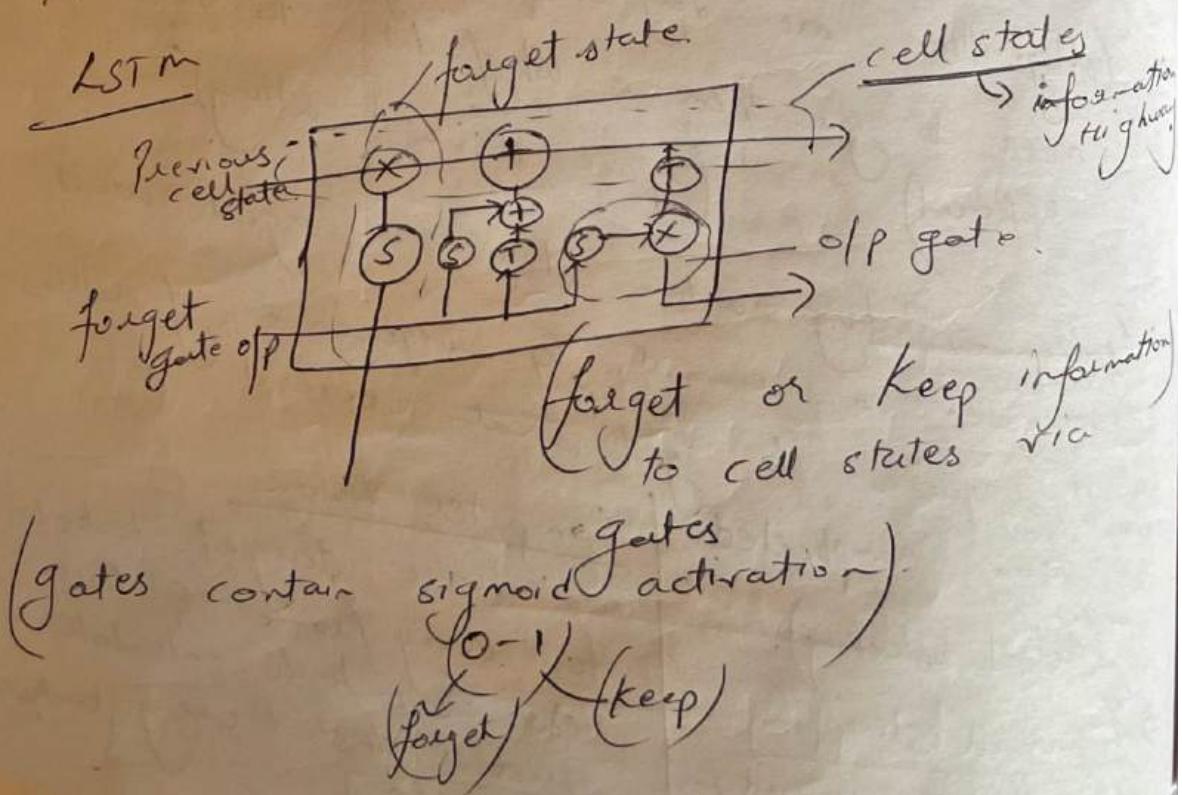
asking for type



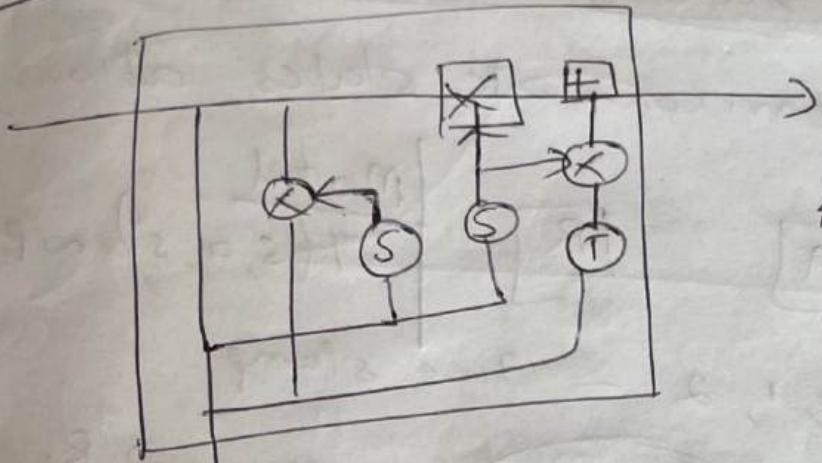
Recurrent Neural Network



(Some values explained to astronomical)
RNN uses lot less computational resources



GRU



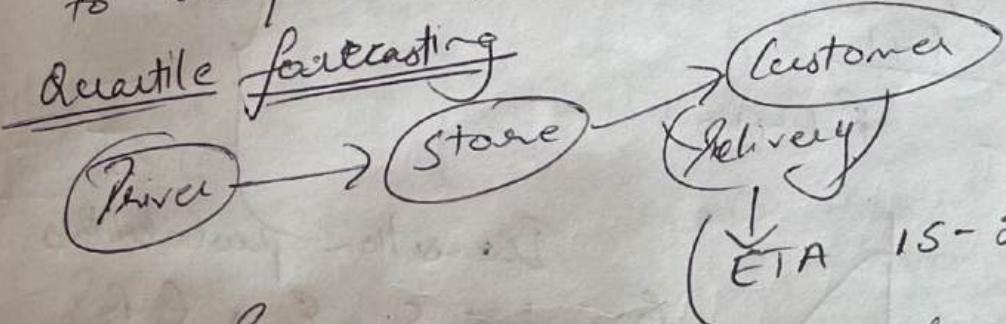
Removed cell state. & used hidden state to transfer info instead

~~How much Reset past iffo to forget gate~~

~~update~~
if forget +
if gate

Mitigate short term memory using gates

→ Gates are now used to regulate flow of info.
to be passed from one step to next.



Quantile Regression
we want to penalize loss if
(τ) percentile is \downarrow & predict \uparrow
($1-\tau$) percentile is \uparrow , predict is \downarrow

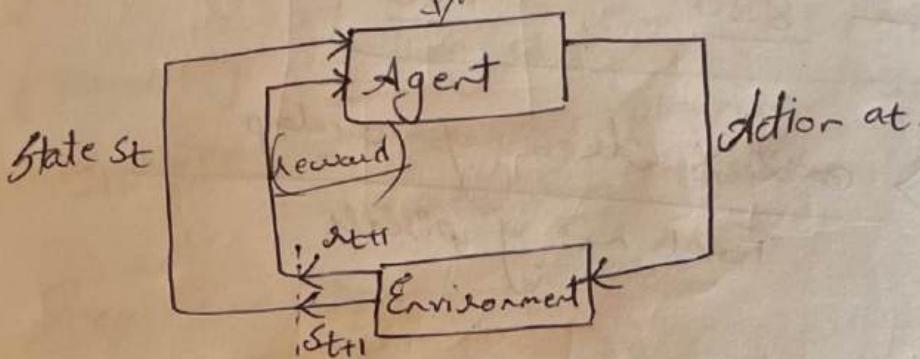
10% - 90% chance
(25 mins - 35 mins) → (Package light)

GBM Regression (objective = "quantile"
criterion = "tau")

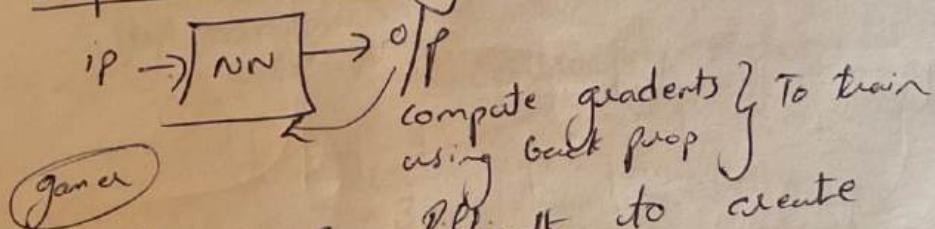
Deep Reinforcement Learning

(Robots)
 current technology has good hardware
 but it lacks the Intelligence to do a task

Software problem.



Supervised Learning



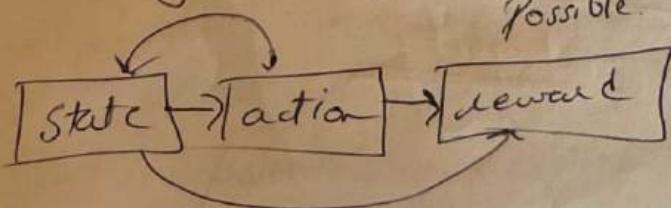
Dataset ← difficult to create
 Agent can't be better than human to
 play that game

In Reinforcement learning → we don't know
 the target label. (Go up/down)

(Policy Network)

$$NN \text{ op} = \begin{matrix} \text{Up / Down} \\ 0.7 \quad 0.3 \end{matrix}$$

Reward → +1]
 Penalty → -1]
 optimize its policy to
 achieve as much reward as
 possible.



compute gradients for the actions that our agent choose
are more likely in the future.

+ve reward \Rightarrow \uparrow prob of actions in future.
-ve reward \Rightarrow Apply some grad with -ve sign
(less likely in future)

Downsides

Policy gradients assume all the actions that we took when lost are bad actions & it reduces the likelihood of taking those actions in the future.

Credit assignment problem

Instead of giving reward for entire set of actions, which actions lead to the reward. (Algorithms are very sample inefficient & we need to give them a lot of data to be efficient) (Sparse reward setting)

- Problem \rightarrow Taking random steps, other agent won't ever see a reward & sequence of actions it requires to get that reward is too complicated.

(Robot stacking)

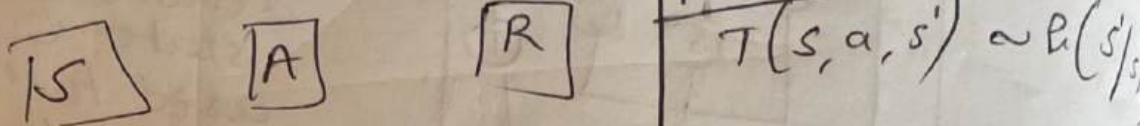
\Rightarrow Reward shaping

Manually design reward for manual that needs to guide agent for desired behaviour

Computer Vision
we have target label for every frame & this lets us to calculate gradient descent using backprop

Markov's Decision Process

Agent, environments, states, actions, i.e.



$t = 0, 1, 2$ — Time stamp.

s_t, a_t

$s_{t+1} \rightarrow R_{t+1}$

$$f(s_t, a_t) = R_{t+1}$$

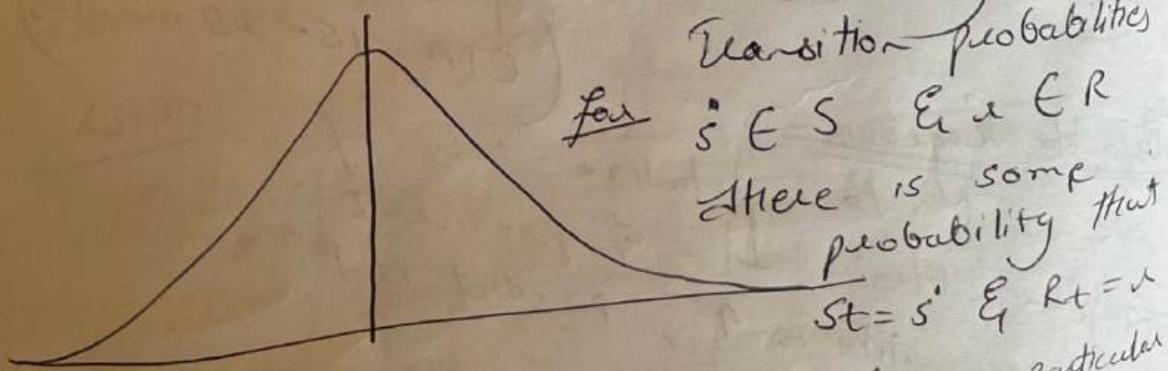
At time t

- Env. is in state s_t .

- Agent selects action a_t .

At time $t+1$

- R_{t+1}



This probability is determined by a particular
values of preceding states $s \in S$ & actions $a \in A$

$a \in A(s)$

Set of actions that can be made
from state s

Only present matter

Policy $\pi(s) \rightarrow a$ action out next action
to take
at any state

$\pi^* \rightarrow$ maximizes the reward over
a long time.

$s_0, s_1, s_2 \rightarrow$ learn for policy

$s, a \rightarrow s, a, r$
If policy π is optimal, it does tell us
what to do in what ever situation
we are in.

$$\pi^* = \arg\max \mathbb{E} \left[\sum \gamma^t R(s_t) \mid \pi \right]$$

Utility $U^\pi(s) = \mathbb{E} \left[\sum \gamma^t R(s_t) \mid \pi, s_0 = s \right]$

Reward gives immediate feedback
(Utility for the state is going to be
term reward) Delayed rewards

optimal policies $\pi^*(s) = \arg\max \mathbb{E} \left[\sum T(s, a, s') U(s') \right]$

$U(s) = U^\pi(s)$

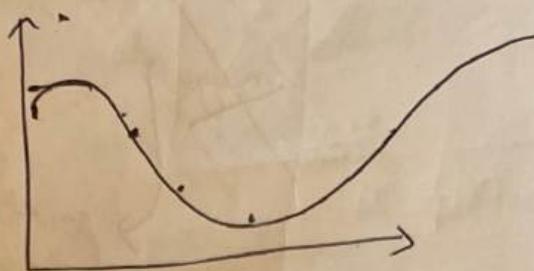
$$V(s) = R(s) + \gamma \max \sum T(s, a, s') U(s')$$

Bellman's Equation Utility Reward Discount of all
rewards that
I get at that
point

optimization

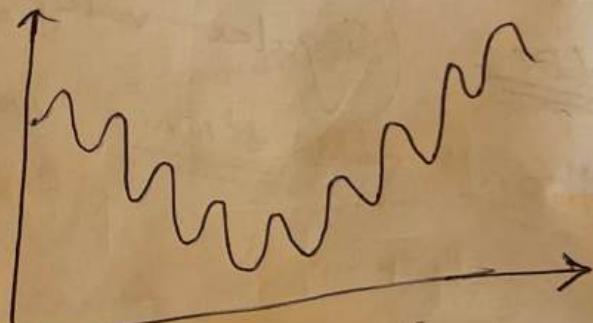
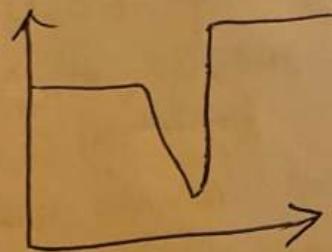
→ Minimize error / loss / cost function by adjusting hyperparameters

Gradient Decent



Exhaustive search
gradient decent
genetic algorithms

if slope is steeper take a big step
if the slope is shallower take small steps
e.g. if the slope is shallow we are interested in global minima



(Race track)

(walking)

Genetic algorithms,

simulated annealing etc

Gradient decent

Exhaustive search

(truck)

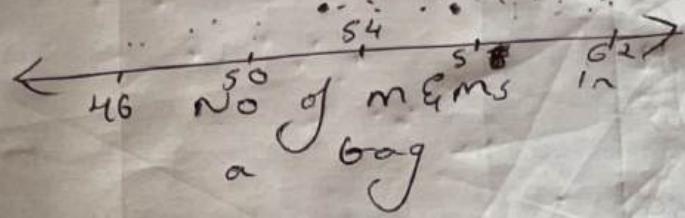
- More assumptions
- Sensitive
- Efficient to compute.

- Few assumptions
- Robust
- Expensive to compute

(Randomness to other steps & jumps)

Examples of optimizations - Bag of means

Single parameter



$d = \text{actual} - \text{guess}$

$$\text{cost} \Rightarrow d^2 = 1 + 4 + 16 + 64$$

$|d|$

$$\text{loss} = d_1^2 + d_2^2 + \dots + d_n^2$$

$$d = \sum_i (n_i - n_{\text{est}})^2$$

Loss function
(if we take $|d|$ then cost would be median)

$$\frac{d}{dn} \left(\sum_i (n_i - n_{\text{est}})^2 \right) = 0$$

Slope = 0

$$\Rightarrow \sum_i \frac{d}{dn} \left[(n_i - n_{\text{est}})^2 \right]$$

$$0 = \sum_i 2(n_i - n_{\text{est}})$$

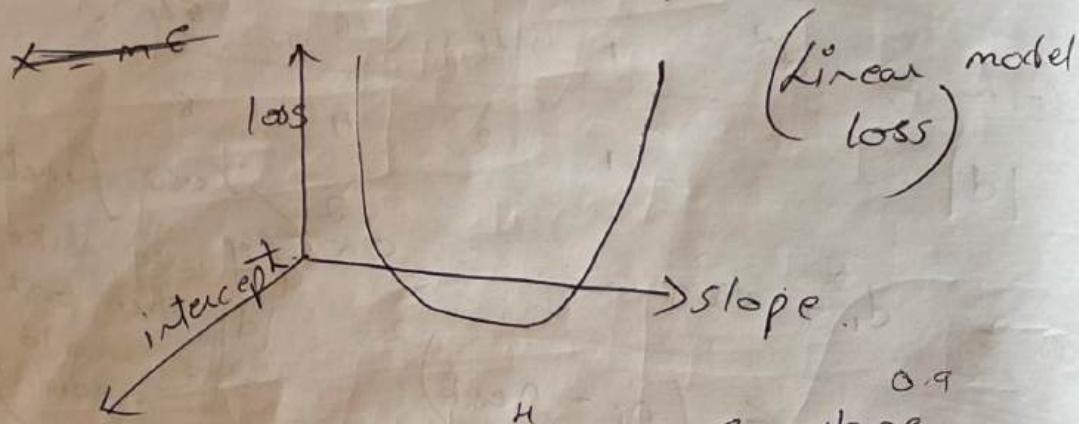
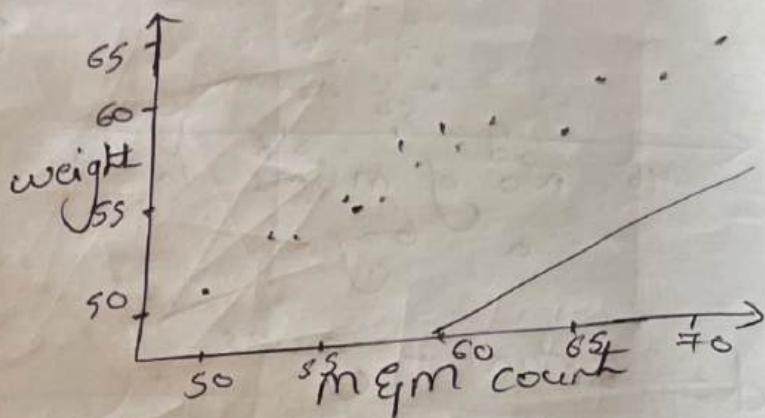
$$\sum_i n_i - \sum_i n_{\text{est}} = 0$$

$$n_{\text{est}} = \frac{\sum_i n_i}{m} = \bar{n}$$

↓
(Average of count)

mem count (vs) bag weight

cost function
is square
deviation
 (d^2)



by changing we minimize the loss function.

(Many parameter model)

gradient decent gives different answers or where we start. Inorder is sensitive to initial guess

- (Non-convex loss function)

Scatter extract approach

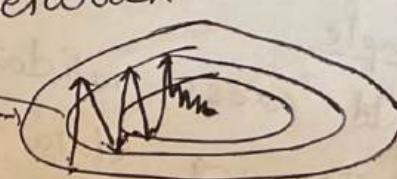
→ ~~Select~~ Several random starting points & running gradient decent on each of them (Better chance of global optima)

Stochastic Gradient Descent with Momentum

It accelerates gradient vectors in right directions, thus leading to fast convergence.

We define a momentum, which is moving averages of our gradients. We then use it to update the weights of the network.

- These up & down oscillations delay GD (α) overshoots the minima.
On iteration t compute $dw \in db$ to current mini batch.



↓ slow learning rate
↔ faster learning rate.

$$vdw = \beta v_{dw} + (1-\beta) dw \quad (\text{Ball rolling down analogy})$$

friction

$$vdb = \beta v_{db} + (1-\beta) db \quad \text{acceleration}$$

velocity

$$w = w - \alpha v_{dw}, b = b - \alpha v_{db}$$

charge weights & bias to smooths out the steps of gradient decent.

$$\beta = 0.9$$

Most common value

α = Learning rate.

Initially $v_{dw} = 0, v_{db} = 0$

Adam Optimization

$$v_{dw} = 0, s_{dw} = 0, v_{db} = 0, s_{db} = 0$$

On iteration t
compute dw, db using current
mini batch.

momentum

$$\begin{cases} v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw \\ v_{db} = \beta_1 v_{db} + (1 - \beta_1) db \end{cases}$$

RMSProp

$$\begin{cases} s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2 \\ s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 \end{cases}$$

$$v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^t) \quad v_{db}^{\text{corrected}} = \frac{v_{db}}{(1 - \beta_1^t)}$$

$$s_{dw}^{\text{corrected}} = s_{dw} / (1 - \beta_2^t) \quad s_{db}^{\text{corrected}} = \frac{s_{db}}{(1 - \beta_2^t)}$$

$$w = w - \alpha \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}}} + \epsilon} \quad b = b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}}} + \epsilon}$$

α = Learning rate (to be tuned)

$$\beta_1 = 0.9 (dw) \quad \beta_2 = 0.999 (dw^2)$$

$$\epsilon = 10^{-8}$$

2) Naive Bayes

$$P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}$$

Posterior probability Predictor

likely hood Prior probability

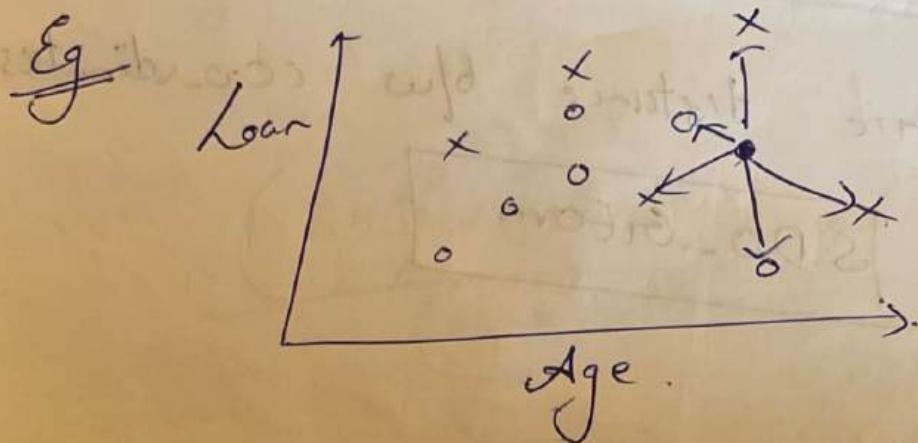
Decision Trees, Random forest &
 Neural networks
 only difference here in classification
 is that the output is discrete
 rather than continuous.

Unsupervised Learning
 Patterns of input data without
 reference to labeled outcomes.
 → Clustering → Dimensionality Reduction.

Clustering ~~KNN, (k means)~~
 grouping of data points
 frequently used in customers segmentation,
 fraud detection, document classification etc
 (k means, Hierarchical, mean shift,
 Density based)

K nearest neighbours (KNN - Classification)

Stores all available cases & classifies new cases based on a similarity measure.



Similarity functions

Manhattan

$$\sum (x_i - y_i)$$

Euclidian,

$$\sqrt{\sum (x_i - y_i)^2}$$

Cosine

$$\left\{ \begin{array}{l} \text{Minkowski} \\ \left(\sum_{i=1}^k (|x_i - y_i|) \right)^{\frac{1}{q}} \end{array} \right.$$

$q = 1$ Manhattan
 $q = 2$ Euclidian

Should k be always odd?

25	40,000
35	60,000
52	100,000

New point

43 30,000

Standardized Euclidian

$$x = \frac{x - \text{min}}{\text{max} - \text{min}}$$

(it is wrong more weightage is given to ~~loan~~ age & age is almost discarded)
So

Number of neighbours

Select most freq neighbours
take the ratio & say it

(Eg 3-Red 2-blue)

(60% it is red, 40% it is blue)

Categorical attributes

male male 0

male Female 1

Closer to 0 they are similar

& away from zero they are dissimilar

As number of dimensions increase
distances become more closer

Kmeans

Clustering learning

Unsupervised

Goal is to find groups in data with number of groups represented by k

algorithm

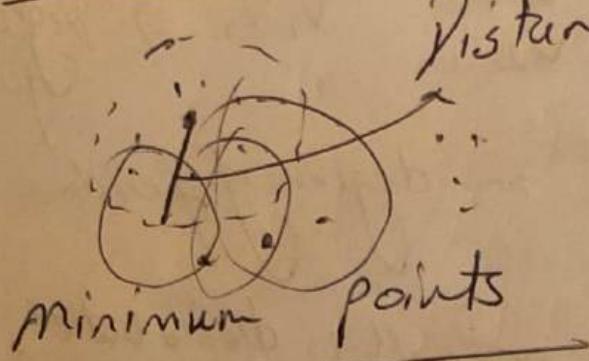
- Randomly assign k centroids (each data point is assigned to its nearest cluster)
 - Update centroids by taking mean of data pts in the cluster
 - Stop when no data point is reassigned to its nearest cluster (will converge)
- Output - \bullet centroids for k clusters.
it can be used to label new data
Labels for training data.

Example grouping inventory by sales

Choosing K - for range of K plot mean distance to the centroid & pick the elbow point where rate of decrease sharply stiff.

DBSCAN

Epsilon



Distance that we set
Minimum number of \sim points that
must exist for a cluster to be
formed. (*including the start*)

Each member of cluster will broadcast
their own radius.

- They look to find new data points
to join the cluster

→ New points again search for new no
when no more data points could be
found, the cluster is finalized.

⇒ We start again by randomly selecting
new starting points; ~~again~~

Linear discriminant Analysis

Seachies linear combinations of variables that best separates 2 classes.

$$Z = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_d x_d$$

$$\beta = C^{-1} (\mu_1 - \mu_2)$$

Pooled covariance matrix $C = \frac{1}{n_1 + n_2} (n_1 C_1 + n_2 C_2)$

Assessment

$$\Delta^2 = \beta^T (\mu_1 - \mu_2) \text{ (Mahalanobis distance)}$$

$\Delta > 3$ means 2 groups are averages differ by 3 standard deviations & overlap b/w 2 groups is small.

Model Evaluation for classification mode

To avoid overfitted & overoptimistic
we use Hold out & cross validation
Hold out ↓
 To evaluate the model.

Hold out

Training set }
 Validation set }
 Test set
 Speed simplicity flexibility
~~disadvantages~~
 Variability in test &
 train data set

Cross validation

Split data into training set
 Eg k fold
~~k = 10~~ (Data split into 10 parts)
 * trained on $\frac{9}{k}$
 & tested on $\frac{1}{k}$

Hyperparameters - Depth,

* Confusion Matrix *

		Actual	
		Yes	No
Model	Yes	TP	FP
	No	FN	TN

True positive

False positive

false negative

True negative

(Sensitivity
 $\frac{TP}{TP+FN}$)

(Specificity
 $\frac{TN}{TN+FP}$)

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

(precision)
Positive predicted value - $\frac{TP}{TP+FP}$

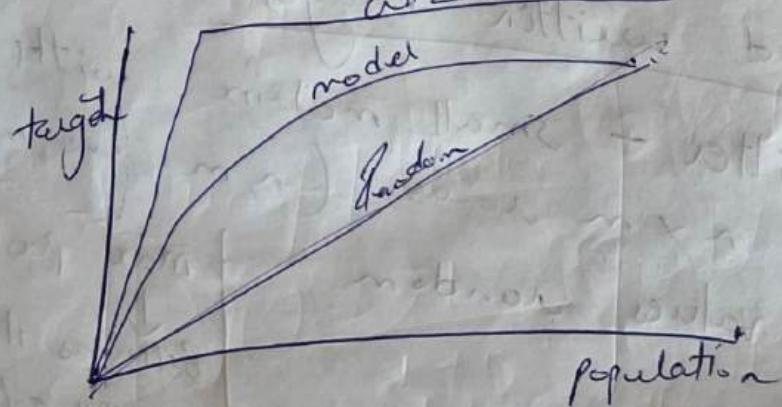
Negative predictive value - $\frac{TN}{TN+FN}$

Sensitivity - Accuracy of model for actual positive cases.

Specificity - Accuracy of model for actual negative cases.

39961

Gini coefficient
Measure of effectiveness of classification model calculated as the relation obtained with & without the model.



K-S is a measure of degree of separation between positive & negative distributions Best - 100

ROC chart

Sensitivity
(True positive rate)

1-Specificity
(False positive rate)

AUC - Area under ROC curve

~~Provides~~

for ch models

6/w [0.5 - 1] Easy to compare 2 ROC's for model evaluation

values

Covid +ve all
Precision desired
include TN

Receiver operating characteristic

Summarizes
confusion
matrix the
each threshold
produces.

Replace
Precision

CNN Images \rightarrow matrix of pixel
Convolutional Neural Network

Basics of Text Analysis

Character encoding

UTF - 8

7 digit representation of 6 symbols

GREP - globally searched regular expressions

N grams - Unigrams, bigrams
This is, This is, is Rithvik
Rithvik

Creating Text database

Corpus - stores actual data

Tidy text, Tm - (packages)

map reduce algorithm
↓
Mapping text to lookup table which
has words
Reduce it by aggregation

Preprocessing → remove stop words

Eg to, the, and, of etc } Tidy to

They don't add value }

→ Remove punctuations

→ Based on case, remove numbers

→ Word case - lower or upper case.

→ Remove whitespace

stemming

take the root word

e.g. process, processing

(keeps just
process)

Document

Term

narrative

Docs

around

conflict great hard

1

2

0

2

0

1 3

3

0

3

0

1

[counts of words in documents]

Natural Language Processing

How computers carryout instructions

→ How to deal with text data

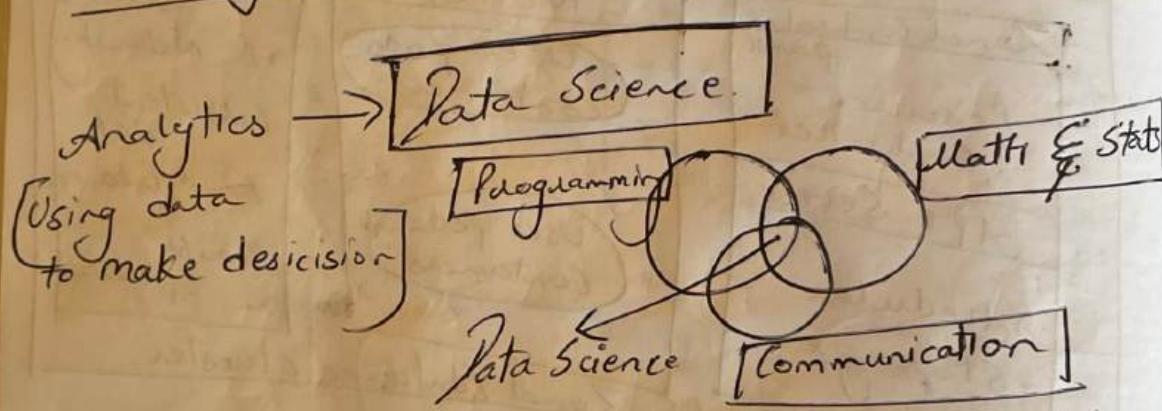


Sentiment analysis

[Text is positive or negative]

Topic modelling - Label the document based on text

Text generation - Replies to messages



Programming - pandas, sklearn, de, nltk, Textblob, gensim

Math & Stats - corpus, document term matrix, sets, sentiment, topic modelling, text generation

Communication - scope, visualize, extract insight, expertise

work flow

Start with question → Get & Clean → EDA → Data

Share insights ← Apply Tech

Data gathering

[How, how much, which data, time range]

Data cleaning

- Corpus

Corpus

DTM

Pandas data frame.

Comedian	Transcript

Document term matrix

Comedian	right	wish	me	back	...
Ali Wong	1	0	1	2	
Dave Chappelle	0	1	3	0	

Clean

- Remove punctuation
- lower case
- remove numbers

ae → python

Tokenize

Unigrams,
bigrams etc

Bag of words

Question

How is Ali Wong different?

Visk learn

Count vectorizer
(to create DTM)

Shift + tab → function details

Exploratory data analysis

→ Top words

→ Vocabulary

→ Amount of Profanity

Visualization

Word cloud

Comedian	unique words	total words	sentences	words/min

Sentiment Analysis

- Input \Rightarrow Corpus \Rightarrow [order matters]
→ Text blob - Python library averages plurality & subjectivity
- Output - Sentiment score
~~objectivity~~ Subjectivity Score [How opinionated they are]

Topic modelling

- Input - DTM & no of topics, no of iterations
 - Library - gensim
(Latent Dirichlet Allocation)
LDA
 - Output - find themes across the comedy routine.
- LDA - Learn top mix in each document,
add word mix in each topic
- Steps → Choose no of topics in corpus
- Randomly assign each word in doc. to a topic
 - Go through every word & its topic assignment in each document
look how often it occurs in topic overall. Based on this assignment info assign word to new topic
- \Rightarrow if topics don't make sense, use nouns only
[linguist manually labelled them]

Analytical modelling

After analysing patterns → Build a model to predict x.

(To get multiple perspectives)

Model is
only as good
as the way
data is composed
from / assembled.

↓
Hypothetical situations
(How extreme data
at model)

(Avoid missing/incorrect/redundant
data as it dears the model)

(Include everything that is necessary
not everything that is possible)

make predictions & draw conclusions → others can
draw their own conclusions

→ No good in knowing it ⇒ one has to
use information to take decisions

Dimensionality Reduction

The process of reducing dimensions of feature set

→ feature elimination

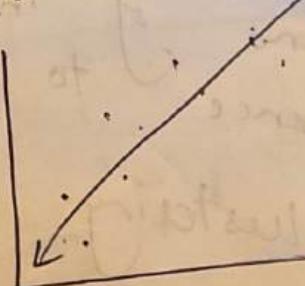
→ feature extraction

Popular method Principal component analysis (PCA) → We get principle components & the variance explained by PC.

Weight & bias (Behaviour when you train)

Bias ← Inability of model to capture the true relationship ~~to~~ is called

Bias (Eg Linear regression)
Straight line can't be curved like



(Behaviour when you Test)

Variance - The difference in fits b/w data sets is variance

(overfitting)

Eg - Model fits perfectly with training set but has high sum of squares in test set

A good ML model should have low Bias (Model true relationship accurately)

& variance

(consistent prediction across different datasets)

Variance bias trade off
we want to ~~fit~~ train the model well but we don't want to overfit
the data is a ~~way that~~
whether how do we optimize our hyperparameters we if we don't touch the test set?
K fold cross validation (~~for~~ very useful for hyperparameter optimization)

Example
→ 10 fold cross validation
split training data into 10 equal sections
we train the algorithm / model on $9/10^{\text{th}}$ of the training set & evaluate its performance on remaining $1/10^{\text{th}}$ of train set
then we train on different $9/10^{\text{th}}$ of train set train set & test on $1/10^{\text{th}}$ of train set
Repeat this process 10 times & we get 10 different measures of model performance
After each iteration, we compare the output results, access the accuracy & adjust the hyperparameters.

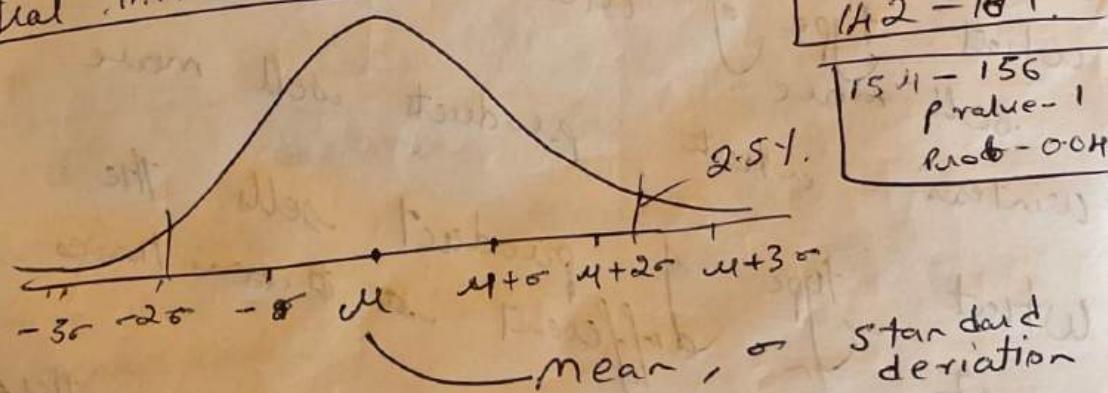
P-value

Probability of obtaining result at least as extreme as the observed result of a statistical hypothesis test, assuming that null hypothesis is correct.

⇒ Probability is coming from distribution (Area under the distribution)

$p < 0.05$ [we reject null hypothesis]

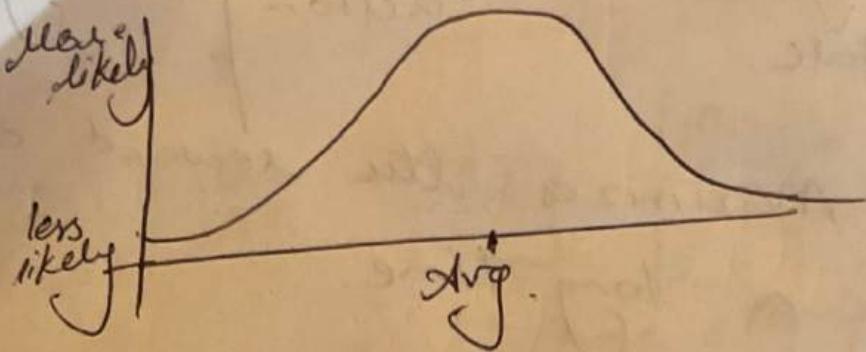
Central limit theorem



P-value is the probability that random chance generated the data or something else that is equal to or greater than what was observed.

prob of getting 5 heads	$HH - 0.25$	$\frac{1}{32} = 0.03125$
prob of getting 5 tails	$HT - 0.5$	$\frac{1}{2} = 0.5$
(Since it's not unusual to get 5 heads in a row)		$\frac{1}{32} + \frac{1}{32} = 0.0625$

Distributions



95% of measurement falls in $\pm 1-2$ standard deviations of around mean width tells us the probability in which measurements like b/w sds

\rightarrow High prob
if $b/w \pm 2$ sd

Central Limit theorem

low prob to lie
 $b/w \pm 1.2$ sd

Means are normally distributed.

of samples

[We don't know what distribution our data comes from. \Rightarrow Sample means are normally distributed]

(Sample size must be atleast 30)

Rule of thumb

Candy dist exception

make confidence intervals, t tests, ANOVA

Difference b/w means of samples

Normalization & Standardization for feature scaling

25 gears
magnitude
Counts

weight
70 kgs

\Rightarrow Normalization (min max scalar)
between 0 - 1.

to scale down features

$$\frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Standardization - Based on standard normal dist

$\frac{1}{\sqrt{6}}$ (standard scalar) $\mu = 0$ & $sd = 1$

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

Euclidean distance / Gradient Descent

\Rightarrow Standardization performed better

images - Normalization
 CNN / ANN
 Helps to learn weights quickly

R² Coefficient of determination

$R^2 \uparrow$ (amount of variation in the data are explained by variable)

0-1 good
But some fields of study have great amounts of unexplainable variance (Human behaviour)
 $R^2 < 50\%$

If you have low R^2 but independent variables are statistically significant; we can still draw important conclusion.

Elastic net (Combination of ridge & lasso regression)

$$\frac{1}{2m} \sum (y - xw)^2 + \alpha (\text{ratio}) \sum_{j=1}^p |w_j| + \alpha (1-\text{ratio}) \sum_{j=1}^p w_j^2$$

If no of predictors (p) \geq no of observations (n)
Lasso will at most pick n predictors as non zero.

If there are 2 or more highly collinear variables Lasso randomly picks one of them which is not good for interpretation.

Inter cluster - distance b/w data points with cluster center

Intra cluster - Distance b/w data point in one cluster to data point in another cluster.

F distribution

$$F = \frac{x_1/\nu_1}{x_2/\nu_2}$$

$x_1, x_2 \Rightarrow$ Random Chi-square variable
 $\nu_1, \nu_2 \Rightarrow$ Degree of freedom

$[0 - \infty]$, Value depends on degree of freedom

(Assumes - Populations are normally distributed,
 Populations are independent of each other)

F ratio $\frac{s_1^2}{s_2^2} = F$

(to find variance of 2 samples)

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

F statistic / F value

To identify means b/w two populations
 are significant or not.

| Ratio of 2 variances
 Used in ANOVA

| t-test single variable diff or not
 F group of variables stat significant or not

To compare more than 1 level
 of independent variable with multiple groups
 we use f distribution.

Rolling a dice & expected payoff

Roll - 1

Average \Rightarrow 3.5 expected payoff

2 Rolls

I won't roll if I get a 4, 5, 6 as my expected payoff is 3.5

However if I roll 1, 2, 3; I roll again

$$\frac{3}{6}[3.5] + \frac{3}{6}[5] \Rightarrow 4.25$$

3 Rolls

I won't roll if I get 5 or 6 as with 2 rolls my expected payoff is 4.25

$$\frac{4}{6}(4.25) + \frac{2}{6}(5.5) = 4.\overline{66}$$

[Theory of optimal stopping for markov chain
assigns a payoff for each state of markov chain
you'd collect if you stop playing
with the chain in the state.]