

Machine Learning:

- (1) Supervised (2) Unsupervised (3) Reinforcement learning.

1/2 Regression

sklearn → many ml algorithms already implemented for us.
also provides some data sets.

- (1) Linear Regression: linear relation b/w x and y . continuous values

$$y = m_1 x_1 + m_2 x_2 + \dots + m_n x_n + b$$

We need to find m_1, m_2, \dots, m_n, b ($n+1$ parameters)

Goodness of fit / R-squared / coefficient of determination

$$= 1 - \frac{\sum (y - y_{\text{pred}})^2}{\sum (y - \bar{y})^2}$$

sklearn has function called ~~score~~ score to find this.

$$\text{let the line } y = mx + c$$

$$\text{error for pt } i = y_i - (mx_i + c)$$

$$\text{Total error} = \sum |y_i - (mx_i + c)|$$

$$\text{so we take } \sum (y_i - (mx_i + c))^2$$

this value will be same for 2 lines.

so we need find m and c for cost function to have minimum value.

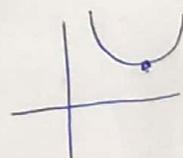
$$\text{cost } f(m, c) = \sum_i (y_i - (m_1 x_1 + m_2 x_2 + \dots + m_n x_n + c))^2$$

where $\{x_1, x_2, x_3, \dots, x_n\}$ are ~~data~~ features.

$$= f(m, c) \quad (\text{assuming 1D}).$$

$$\frac{\partial f(m, c)}{\partial m} = 0 \quad \frac{\partial f(m, c)}{\partial c} = 0$$

solve for m and c .



$$\frac{\partial f(m, c)}{\partial m} = \sum_i (y_i - (m x_i + c)) (x_i) = 0$$

$$\Rightarrow \sum_i \frac{x_i y_i}{N} - m \sum_i \frac{x_i^2}{N} + c \sum_i \frac{x_i}{N} = 0$$

mean

$$\Rightarrow (x * y). \text{mean}() - m (x^2). \text{mean}() - (x. \text{mean})^2 = 0 \quad \text{--- (1)}$$

$$\text{for } c \quad \sum_i (y_i - (m x_i + c)) = 0$$

$$\sum c = \sum y_i - m \sum x_i$$

$$c = \frac{\sum y_i}{N} - m \frac{\sum x_i}{N}$$

$$c = y. \text{mean}() - m x. \text{mean}()$$

--- (2)

$$(i) \Rightarrow (x^T y) \cdot \text{mean}() - m \cdot (x^2) \cdot \text{mean}() - c^T x \cdot \text{mean}() = 0$$

$$(ii) \Rightarrow c = y \cdot \text{mean}() - m \cdot x \cdot \text{mean}().$$

$$= (x^T y) \cdot \text{mean}() - m(x^2) \cdot \text{mean}() - y \cdot \text{mean}()^T x \cdot \text{mean}() \\ + m x \cdot \text{mean}() x \cdot \text{mean}() = 0$$

$$m = \frac{(x^T y) \cdot \text{mean}() - y \cdot \text{mean}()^T x \cdot \text{mean}()}{(x^2) \cdot \text{mean}() - x \cdot \text{mean}()^T x \cdot \text{mean}()}$$

$$c = y \cdot \text{mean}() - m^T x \cdot \text{mean}().$$

$$y_p = mx + c$$

Finding complex decision boundary using LR.

Let insert a dummy variable $x^1 \rightarrow x^2$

$$\text{Now we can do LR } y = m_1 x + m_2 x' + c$$

\hookrightarrow eq of plane.

Also x and x' are dependent on each other, we can do dimensionality reduction and we can get parabola or decision line instead of st. line.

So we can do many stuff $\rightarrow \log x, x^3, \dots, x_1^a x_2, x_1 x_2^a x_3$

So we can add many dummy columns looking at which are important and may do dimensionality reduction later.

But adding many features may lead to overfitting.

$$y = \underbrace{m_1 x_1 + m_2 x_2 + \dots + m_n x_n}_{n+1} + m_{n+1} \theta$$

$$x \rightarrow \begin{bmatrix} & & & & & \\ & \ddots & & & & \\ & & n+1 & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix} \quad \theta = \begin{bmatrix} m_1 \\ m_2 \\ \vdots \\ m_{n+1} \end{bmatrix} \quad Y_p = x^T \theta$$

$$\text{To find } \theta = \boxed{(X^T X)^{-1} X^T Y_{\text{actual}}} = \theta \quad (\text{Normal equation})$$

~~Complexity~~ Complexity $\rightarrow O(n^3)$ as we slow it n is large calculate $X^T X$.

$$\text{or } (n+1)^m \times m^{(n+1)} \quad (n+1)^m \quad m^{m+1}$$

Gradient Decent

$$\text{cost fun} = f = \frac{1}{N} \sum (y_i - (m x_i + c))^2$$

$$\begin{aligned} m' &= m - \alpha \text{ slope}_m \\ c' &= c - \alpha \text{ slope}_c \end{aligned}$$

start with random value of
m and c
 $\alpha \rightarrow$ learning rate

where

$$\begin{aligned} \text{slope}_m &= \frac{\partial f}{\partial m} = \frac{1}{N} \sum (y_i - m x_i - c) (-x_i) \\ \text{slope}_c &= \frac{\partial f}{\partial c} = \frac{1}{N} \sum (y_i - m x_i - c) (-1) \end{aligned}$$

use adaptive α , if cost value increase, decrease α .

$$m, c = \text{grad_de}(\text{data}, \alpha, \text{iterations})$$

For general Gradient Decent:

$$\begin{aligned} \text{cost fun} = f &= \frac{1}{n} \sum \left(y_i - (m_1 x_i^1 + m_2 x_i^2 + \dots + m_n x_i^n + c) \right)^2 \\ \frac{\partial f}{\partial m_j} &= -\frac{1}{n} \sum \left(y_i - (m_1 x_i^1 + m_2 x_i^2 + \dots + m_{j-1} x_i^{j-1} + m_{j+1} x_i^{j+1} + \dots + m_n x_i^n + c) \right) (x_i^j) \\ m_j' &= m_j - \alpha \frac{\partial f}{\partial m_j} \end{aligned}$$

Variations of gradient Decent

Batch Gradient: Consider whole dataset in 1 iteration

stochastic Gr.D.: Consider one point by point during iteration and keep updating row.

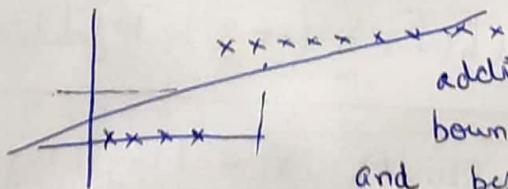
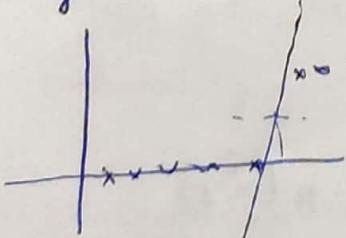
mini-batch Gr.D.:

Classification discrete output.

why not linear Regression?

value is continuous instead! we want discrete.

let us assume a cutoff at $y = 0.5$
but let us have few more data points.



adding few extremes the boundary changes a lot and best fit line would be spoiled.

Also we get large values for extreme x
but we only need {0, 1}?

Logistic Regression

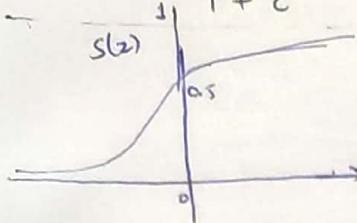
Sigmoid fⁿ.

$$h(x) = \frac{1}{1 + e^{-g(x)}}$$

$$S(z) = \frac{1}{1 + e^{-z}}$$

$$g(x) = mx + c \quad \leftarrow \text{(basic LR line)}$$

$$= m^T x$$



hypothesis

$$h(x) = \frac{1}{1 + e^{-m^T x}}$$

$$h(x) > 0.5 \Rightarrow y_{\text{pred}} = 1$$

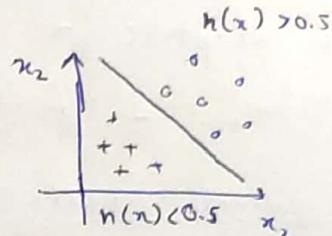
$$m = [m_1, m_2, \dots, m_{n+1}]$$

$$x = [x_1, x_2, \dots, x_n, 1]$$

threshold = 0.5

we know is $z > 0 \Rightarrow S(z) > 0.5$

so if $m^T x > 0$	$h(x) > 0.5$	$y_{\text{pred}} = 1$
if by $m^T x \leq 0$	$h(x) \leq 0.5$	$y_{\text{pred}} = 0$



* we will have a linear decision boundary

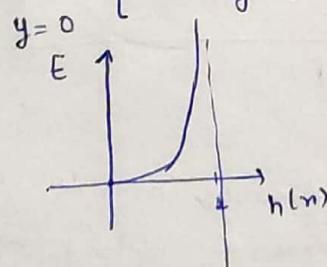
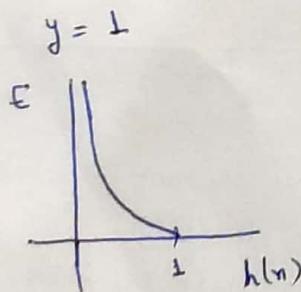
We need to find m.

Cost function

$$\text{in LR is } J(w) = \sum (y^i - h(x^i))^2$$

in Logistic Regression we can use this because it has so many local minima. (Not a convex function)

$$\text{so } E(h(x), y) = \begin{cases} -\log(h(x)) & y = 1 \\ -\log(1 - h(x)) & y = 0 \end{cases}$$



$$E(h(x^i), y^i) = -y^i \log(h(x^i)) - (1 - y^i) \log(1 - h(x^i)).$$

$$E(h(x), y) = \frac{1}{m} \sum_{i=1}^m (-y^i \log(h(x^i)) - (1 - y^i) \log(1 - h(x^i)))$$

$$h(x) = \frac{1}{1 + e^{-m^T x}}$$

We need to find optimal values of m.

We will use gradient descent

$$m_j^i := m_j - \alpha \frac{\partial E}{\partial m_j}$$

On substituting values of $h(x)$ in error function.

$$E_i(m) = m^T x_i + y^i (-m^T x_i) + \log(1 + e^{-m^T x_i})$$

$$E(m) = \log(1 + e^{m^T x_i}) - \underbrace{y^i (m^T x_i)}$$

$$\frac{\partial E(m)}{\partial m_j} = \frac{1}{1 + e^{m^T x_i}} x_j - y^i x_j$$

$$= -\left(y^i - \frac{1}{1 + e^{-m^T x_i}}\right) x_j$$

$$= -(y^i - h(x_i)) x_j$$

$$\frac{\partial E}{\partial m_j} = -\frac{1}{m} \sum (y^i - h(x_i)) x_j$$

Summary $h(x) = \frac{1}{1 + e^{-m^T x}}$

$$E(h(x), y) = -\frac{1}{m} \sum [y \log(h(x)) + (1-y) \log(1-h(x))]$$

$$m_j := m_j - \alpha \frac{\partial E}{\partial m_j}$$

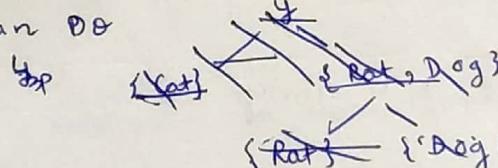
$h(x)$ is other form of probability (P).

$$h(x) \Rightarrow P(y = 1)$$

$$\frac{\partial E}{\partial m_j} = -\frac{1}{m} \sum (y^i - h(x_i)) x_j$$

Multiclass let $y_p = \{\text{cat}, \text{Rad}, \text{Dog}\}$

(one vs rest) we can do



$$y_1 = [\{\text{Cat}\}, \{\text{Rad}, \text{Dog}\}]$$

$$y_2 = [\{\text{Rad}\}, \{\text{Cat}, \text{Dog}\}]$$

$$y_3 = [\{\text{Dog}\}, \{\text{Cat}, \text{Rad}\}]$$

so we get $P(y_1), P(y_2), P(y_3)$

Probability estimates, 3 models are independent.

Multinomial assume we have $\{1, 2, \dots, k\}$ classes.

$$P(y=j) = \frac{e^{m_j \cdot x}}{\sum_{i=1}^k e^{m_i \cdot x}}$$

We are training n^k parameters

function in sklearn model.

`fit()`, `predict()`, `score()`, `predict_proba()`
(mean accuracy)

Why we can add dummy features to get complex boundary.

But overfitting, F!

Regularisation: "I'm gonna save you"

$$\text{cost}' = \text{cost} + \lambda \sum |m_i|$$

Also - high $\lambda \rightarrow$ underfitting

$$L_1 = \text{cost} + \lambda \sum |m_i| \quad \underline{\text{lasso}}$$

$$L_2 = \text{cost} + \lambda \sum m_i^2 \quad \underline{\text{ridge}}$$

* No regularisation on intercept.

Performance measures in classification

(1) Confusion Matrix

sklearn.metrics \rightarrow confusion matrix

		Actual	
		T	F
Predicted	T	TP	FP
	F	FN	TN

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

For a class (n)

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

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN} = TPR = \text{sensitivity}$$

$$\text{sensitivity} = \frac{TP}{TP + FN} = TPR \quad (\text{ability to detect +ve class})$$

$$\text{Specificity} = \frac{TN}{TN + FP} = 1 - FPR$$

Decision Tree! (You know Already)

How do we split? May be Accuracy, Gain ratio etc.

Greedy Algo.

Gini index, Information gain

→ If node is pure : make it leaf.

→ If no feature is left : output majority.

→ Find best feature to split on

→ Make recursive call.

(i) Accuracy → find accuracy of Parent, try splitting on every & valid possible feature and get children's accuracy.

→ choose best feature if \sum children's accuracy for that feature > Parent's accuracy

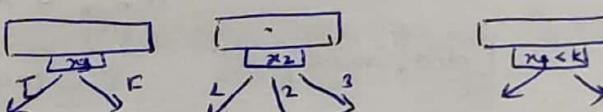
→ If none, make parent a leaf.

* continuous value features:

We decide some decision line say $x_2 \leq 5$ $x_2 \geq 5$

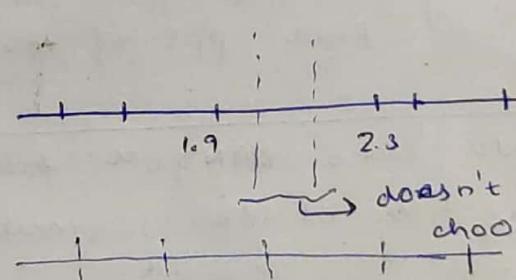
Say data is:

x_1	x_2	x_3	
continuous	discrete	boolean	


How we decide k?

→ sort the data

→ get the mid points



→ find accuracy or other pos matrix for all mid points.

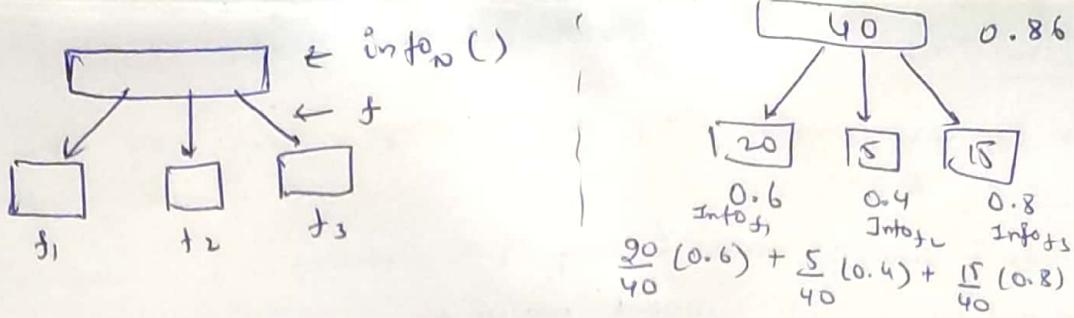
and pick a mid point to make it as decision boundary.

* Information Gain:

$$\text{Entropy} / \text{Info req} : - \sum_{\text{classes}} p_i \log p_i$$

We need minimum randomness / Entropy / Info req to reduce impurities

→ we chose feature which decrease entropy by max.



$$\text{Info gain} = \text{Info}_N() - \text{Info}_F()$$

$$\text{Info}_F() = \sum_i \frac{|D_i|}{|D|} \text{Info}_{D_i}()$$

Pick feature which gives better info gain

* better see an example.

We will be back, a little confusion in this topic

Feature scaling / Data Preprocessing

↳ ~~feature~~ scaling different feature (with variety of range) under some range.

$$(1) \text{ min_max scaling } [0, 1] = \frac{x - x_{\min}}{x_{\max} - x_{\min}} = p$$

$$[min, max] = p(max - min) + min$$

(2)

$$\text{make mean} = 0$$

$$\text{variance} = 1$$

Read ppt of minmax

Random Forest: we take RANDOM training data and features and build many trees (FOREST), result given by majority of trees is taken in account.

② Bagging → Bootstrap Aggregation Algo.

If we have m data points, we select k points with replacement. $\sum k_i = m$
means multiple occurrence of single point is there.

③ Feature selection: k feature from replacement $\overset{n}{\nwarrow}$ feature without replacement (generally $k = \sqrt{n}$)

we can solve regression problems using decision tree,
split on basis of which feature reduces the MSE
and leaf node contains mean value.

Bayes Theorem:

Naïve Bayes:

Bayes Theorem:

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

$y = a_1, a_2, \dots, a_k$ (output classes)

then $\hat{y} = \max_{i=1}^k \left(P(y=a_i | x=x) \right)$

$$\hat{x} = \max_{i=1}^k \left(\frac{P(x=x | y=a_i) P(y=a_i)}{P(x=x)} \right)$$

$$\hat{Y} = \max_{i=1}^k \left(P(x=x | y=a_i) P(y=a_i) \right)$$

↓ output class

$$P(y=a_i) = \frac{|a_i|}{|y|} = \frac{|a_i|}{|a_1| + |a_2| + \dots + |a_k|}$$

$$P(x=x | y=a_i)$$

Independent Events: $P(A \cap B) = P(A) * P(B)$

Given input $x \rightarrow \{x_1, x_2, \dots, x_k\}$ where x_i is value corresponding to feature f_i .

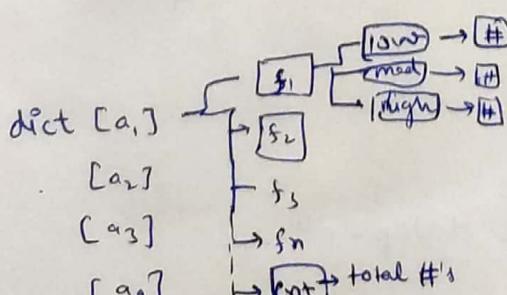
$$P(x=x | y=a_i) = P((x_1, x_2, x_3, \dots, x_k) | y=a_i)$$

* Assumption \rightarrow All features are independent.

$$= P(f_1=x_1 | y=a_i) * P(f_2=x_2 | y=a_i) * \dots * P(f_k=x_k | y=a_i)$$

$$P(x=x | y=a_i) = \prod_{j=1}^k P(f_j=x_j | y=a_i)$$

$$P(f_j=x_j | y=a_i) = \frac{\#(\text{F}_j=x_j \text{ & } y=a_i)}{\#(y=a_i)}$$



$$P(f_j=x_j | y=a_i) = \frac{\text{dict}[a_i][f_j][x_j]}{\text{dict}[a_i][\text{count}]}$$

Given data x-train Y-train

f ₁ f ₂ f ₃ ... f _k	label
low	a ₁
med	a ₂
high	a ₁
low	a ₅
high	a _n

laplace correction: $P(a_i)$ can be 0 if there are no value (x_j^i) present in feature (F_j^i) given a_i

so we modify formula

$$= \frac{\text{dict}[a_j^i][F_j^i][x_j^i] + 1}{\text{dict}[a_j^i][\text{ent}] + |F_j^i|} \rightarrow \text{values } F_j^i \text{ can take}$$

continuous data?

$$p(f_j^i = x_j^i | y = a_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_j^i - \mu)^2}{2\sigma^2}}$$

Naive Bayes can be used as ^{test} species collection

KNN → k-Nearst Neighbours Do feature scaling before KNN
you know what it is!

Cross validation:

low value of k leads to overfitting

high value of k leads to underfitting

Make sure we have relevant data and independent.

→ Assign weights to feature (w_i^i for f_i)

$$\sum_{i=1}^n w_i^i (x_1^i - x_2^i)^2$$

→ Feature selection (backward elimination).

Handling labeled data:

→ use 0,1 in binary data. (2 values).

→ we can't assign 0,1,2,3... to discrete data bcz distance would be effected.

	Red	Blue	Green
Red	0	0	1
Blue	0	0	1
Green	1	0	0
	0	1	0

→ Other KNNS

→ Brute force → compare with all datapoint

→ KD Trees something like BST

→ Ball Trees

Pros of KNN → Easy to understand and code (lol)

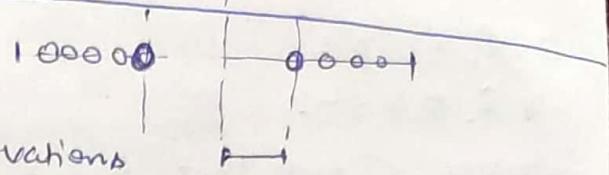
(2) Works for multi-classification

abc. CV - result

Cons of KNN:

- (1) Testing time is huge.
- (2) It can be biased, if training data is biased.
- (3) Curse of dimensionality.

SVM



The shortest distance b/w the observations and the threshold is called margin

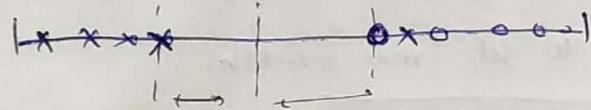
Maximal margin classifier: (MMC) threshold that gives us the largest margin.

MMC is not optimal (think of outliers)

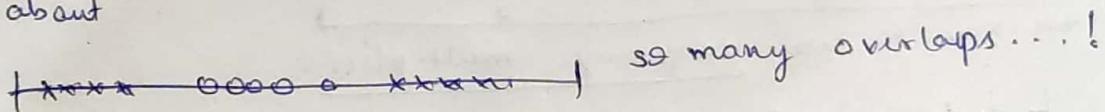
We do cross validation to find the best soft Margin aka Support Vector Classifier.

The support vector classifier is a hyperplane in n-dimension

so SVC can handle outliers nicely using cross validation



But what about



Support Vector machine: We add extra dimensions to points and

then try to fit hyperplane, ~~sup~~ SVC.

→ How to transform the data? SVM uses kernel function

We can find good value of (d) dimension via cross validation

Other kernel is RBF (finds SVM in infinite direction).

Kernels just calculate the sim b/w points in high dimension without actually transforming them.

$$\text{let } \gamma = \frac{1}{2} d^2$$

$$\begin{aligned} \text{Polynomial kernel} &= (a \times b + \gamma)^d \\ &= a^d b^d + a^{d-2} b^2 \gamma^2 + \dots \end{aligned}$$

$$= \underbrace{(a, a^2, \dots, a^d)}_{\gamma} \underbrace{(b, b^2, \dots, b^d)}_{\gamma} \underbrace{\gamma^d}_{\gamma^d}$$

\vec{w} is \perp to margin

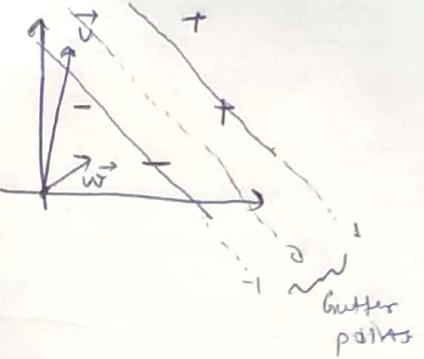
\vec{w} be unknown sample point

$\vec{w} \cdot \vec{x} \geq c$

i.e. $\vec{w} \cdot \vec{x} + b \geq 0$ for + sample

but we still don't know \vec{w} and b .

b is constant



$$\vec{w} \cdot \vec{x}_+ + b \geq 1 \dots \textcircled{1}$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1 \dots \textcircled{2}$$

introduce y_i such that $y_i = 1$ \rightarrow +ve sample
 $y_i = -1$ \rightarrow -ve sample

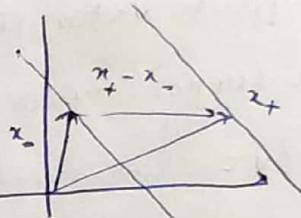
eq 1 and 2 become

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq \pm 1$$

$$\boxed{\begin{aligned} y_i (\vec{x}_i \cdot \vec{w} + b) - 1 &\geq 0 \\ y_i (\vec{x}_i \cdot \vec{w} + b) + 1 &= 0 \end{aligned} \rightarrow \text{For Gutter Points.}}$$

$$\text{width} = (x_+ - x_-) \cdot \frac{\|\vec{w}\|}{\|\vec{w}\|}$$

$$= \frac{\vec{x}_+ \cdot \vec{w}}{\|\vec{w}\|} - \frac{\vec{x}_- \cdot \vec{w}}{\|\vec{w}\|} = \frac{1-b + 1+b}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$$



We are maximising $\frac{1}{\|\vec{w}\|}$ to get max width.

$$\text{minimise } \|\vec{w}\| \sim \text{minimise } \frac{1}{2} \|\vec{w}\|^2$$

We have to use Lagrange multipliers.

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1] \rightarrow \textcircled{3}$$

$$\frac{\partial L}{\partial \vec{w}} = \cancel{\vec{w}} - \sum \alpha_i y_i \vec{x}_i = 0 \Rightarrow \boxed{\vec{w} = \sum_i \alpha_i y_i \vec{x}_i}$$

$$\frac{\partial L}{\partial b} = \boxed{\sum \alpha_i y_i = 0}$$

Substitute in eq 3.

$$L \Rightarrow \frac{1}{2} \left(\sum_i \alpha_i y_i \vec{x}_i \right) \cdot \left(\sum_j \alpha_j y_j \vec{x}_j \right) = \sum \alpha_i y_i$$

$$L \Rightarrow \frac{1}{2} \left(\sum_i \alpha_i y_i \vec{x}_i \right) \left(\sum_j \alpha_j y_j \vec{x}_j \right) - \left(\sum \alpha_i y_i \vec{x}_i \right) \left(\sum \alpha_j y_j \vec{x}_j \right) = \sum \alpha_i y_i b + \sum \alpha_i$$

$$L = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j (\alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j) \rightarrow \text{Max of this.} \quad \textcircled{4}$$

depend on dot product.

$$\sum a_i y_i \vec{w} \cdot \vec{x}_i + b \geq 0 \quad \text{+ve sample.}$$

dot product

Won't work with non linear decision boundary.

transformer: $\phi(\vec{x}) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$ to max
 $\phi(x_i) \cdot \phi(x_j)$

$$K(x_i, x_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

\hookrightarrow kernel function. we actually don't have to do ~~the~~ transform so ~~$\phi(\vec{x})$~~ , all depends on dot product.

some kernels: (1) $(\vec{w} \cdot \vec{x} + b)^d$ Polynomial Kernel.

$$(2) e^{-\frac{\|\vec{x}_i - \vec{x}_j\|}{c}} \quad \text{Radial basis kernel}$$

Neural Network

Activation Function \rightarrow sigmoid f^n , ReLU,

We get x-axis coordinate Soft Plus.

from the weight s and use activation function to get Y.

weights come from Activation Function.

Act. (weight + biases) \rightarrow new perception.

ReLU



$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

$$f(x) = \max(0, x).$$

Main PPT's after mid sem'.

KNN → pattern recognition

→ While scaling we can add weight to the important features. w_k
we can learn w_k using cross-validation.

Ez:)

Kmeans Algo: (Unsupervised learning)

StadQuest

PPT ↗

- (1) Select a K (clusters to identify) ($k=3$)
- (2) Randomly select K distinct dp. ~~(Centroids)~~
- (3) Measure distance of remaining points to clusters and classify them to nearest cluster.
- (4) Calculate mean of each cluster \leftarrow (centroid)
- (5) Repeat. (measure distance along the mean this time)
- (6) ~~Reset~~ When it does not change its result and calculate the variance.

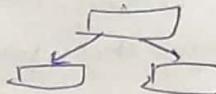
centroid can be real or imaginary dp

Repeat whole with different K and diff. starting points.
The best is where we have a gradual change in variance.

Add Boost → contains stump

resource: StadQuest

PPT ↗



The tree with just one node and two leaves is stump.

A tree uses all features to make decision, but a stump would take one feature to make decision.
So 'stumps' are weak learners.

In RF, each tree has equal vote to classification
In Adaboost, each stump has different vote to classification.

In RF, each tree is independent.

In Adaboost, order is important, the error of first stump influences how second stump is made and soon.

Steps

(1) Make stumps of

(2) Assign sample weight (this will tell which dp is important)
(initially for all = $1/\text{no. of samples}$)

(3) Make stumps using each feature and calculate gini-index.
The one with lowest would be the first stump.

(4) Now we find "vote"

$$\text{Amount of vote} \stackrel{\text{say}}{=} \frac{1}{\log(1 - \frac{\text{Total Error}}{\text{Total Error}})} \quad \left\{ \begin{array}{l} \# \text{ of incorrectly} \\ \text{classified with their} \\ \text{sample weight.} \end{array} \right.$$

(5) Now we modify sample weight.

We increase the sample weight of incorrectly classified weight of that feature, so that next time we try to classify it correctly and decrease all other sample weight.

$$\uparrow \text{Net sample weight} = \text{sample weight} \times e^{\text{amount of vote.}}$$

$$\downarrow \text{and calculate new sample weight.}$$

and find rest of stamps.

Cross Ensemble Technique → combining multiple models. Need to talk from Main

→ Bagging: Random Forest
(Bootstrap aggregation)

→ Boosting: ADA BOOST, GRADIENT BOOSTING, XG BOOST.

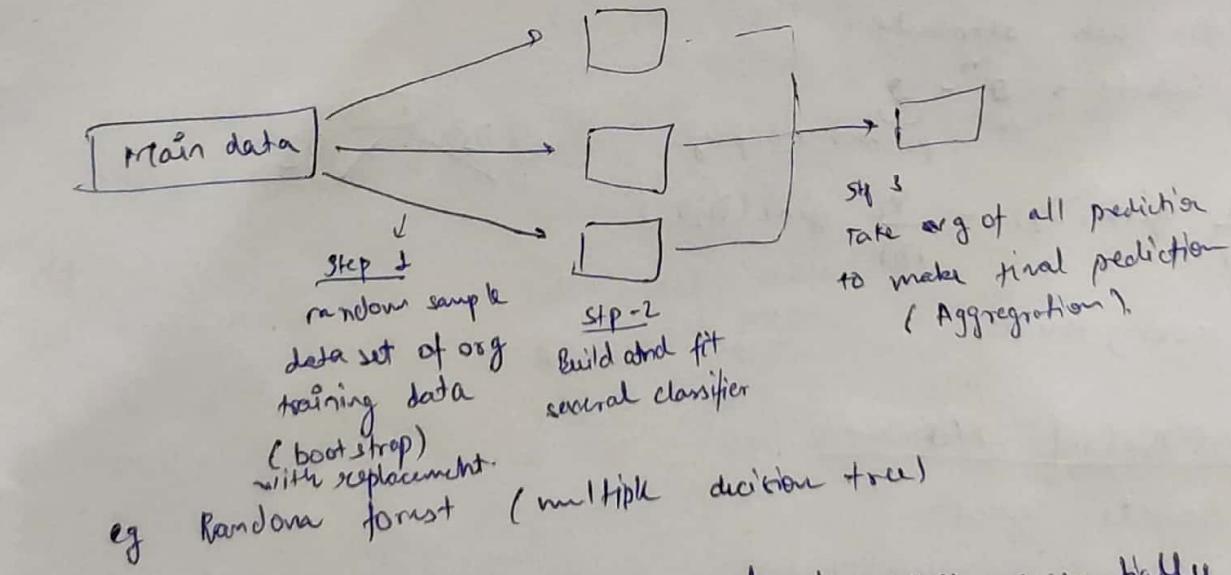
Error → Bias + variance.

using Bagging we can decrease variance
using Boosting bias

Simplest
sketch

PPT v

Bagging: → bootstrap aggregation reduces variance of an estimate by taking mean of multiple estimates.



Boosting reduces bias by training weak learners sequentially, each trying to correct predecessor.

- eg:
→ ADA BOOST (multiple stamps)
→ Gradient Boosting (GBM)
→ XG BOOST

CART Algo (Classification and Regression Tree)

Gini \rightarrow Gain ratio:

- Binary Tree
- Uses GINI Index

- If a data set D exemplified from n classes,

$$gini(D) = 1 - \sum_{i=1}^n p_i^2$$

- If a data set D is split on A into two subsets D_1 and D_2

$$gini_A(D) \rightarrow \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- * Reduction in impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

The one has minimum Gini index, is selected as the splitting attribute.

- (1) calculate gini of whole dataset

$$gini(D) = 1 - \sum p_i^2$$

- (2) calculate for each attribute

→ Total subsets $\rightarrow 2^n - 2$
power + empty.



$$\frac{D_1}{|D|} gini(D_1) + \frac{D_2}{|D|} gini(D_2)$$

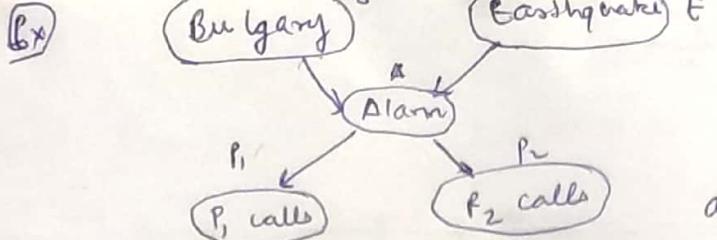
and keep doing this

- Bayesian Belief Network

(1) Directed acyclic graph

(2) conditional Probability Table.

convenient for representing probabilistic relation b/w multiple events.



Given

$$P(B=T) = 0.001$$

$$P(B=F) = 0.999$$

$$P(E=T) = 0.002$$

$$P(E=F) = 0.998$$

and other probabilities

see video 5 minute engineer.

* Rule Based classification.

'IF' 'AND' 'THEN'

Rules are extracted using these three keywords.

e.g.: IF outlook = 'sunny' AND
Humidity = 'High'
THEN Play = 'No'

* Association Rule Mining (Market-Basket Analysis)

$A \Rightarrow B$
If Then

subset
understand!

3 metrics which help in Association

(1) support (2) confidence (3) lift

$$\text{support} = \frac{\text{freq}(A, B)}{N}$$

$$\text{confidence} = \frac{\text{freq}(A, B)}{\text{freq}(A)}$$

To turn rules
into association
rules

How frequent both items are
bought together.

$$\text{lift} = \frac{\text{support}}{\text{supp}(A) \times \text{supp}(B)}$$

where A, B occurs independently

e.g.:

Let transactions be

T_1	A	B	C
T_2	A	C	D
T_3	B	C	D
T_4	A	D	E
T_5	B	C	E

lets create rules.

	support	confidence	lift
$A \rightarrow D$	2/5	2/3	10/9
$C \rightarrow A$	2/5	2/4	5/6
$A \rightarrow C$	2/5	2/3	5/6
$B \& C \rightarrow A$	1/5	1/3	5/9

Apriori Algorithm: uses frequent item set to generate ~~itemsets~~ association rules. It is based on concept that a subset of a frequent itemset must also be frequent itemset.

Frequent item set is an itemset whose support value is greater than a threshold value.

Eg

TID	Item	Itemset	Support	Itemset	Support	Itemset	Support
T1	134	{1,3,4}	3	{1,2,3}	2	{1,3,5,5}	2
T2	235	{2,3,5}	3	{1,3,5}	3	{2,3,5}	2
T3	1235	{1,2,3,5}	4	{2,3,3}	2		
T4	25	{2,5}	1	{2,5,3}	3		
T5	135	{1,3,5}	4	{3,5,3}	3		

↓ 1st Iteration

↓ 2nd Iteration.

We can do some pruning

for eg: we can discard those data set containing {1,2,3}.

$$\text{eg} = \{1,2,3\} \{1,2,3\}, \{1,5\}, \{2,5\}.$$

(Spartan Eng)

so frequent items are {1,3,5}, {2,3,5}

Now check their subset:

eg for $I = \{1,3,5\} \rightarrow \{13\} \{35\} \{15\} \{135\}$.

For every subsets of I , output the rule:

$s \rightarrow (I-s)$ i.e (s recommends I-s)

if $\text{support}(I) / \text{support}(s) \geq \text{min_conf value}$.

so we check which rule is selected

say $\{1,3,5\} \rightarrow \{1,3,5\} - \{35\} [1 \& 3 \rightarrow 5]$

$\Rightarrow \text{support}(1,3,5) / \text{support}(1,3,5 - \{35\}) = 2/3 > 60\%$

~~135~~

$\{13\} \rightarrow \{1,3,5\} - \{35\} [1 \rightarrow 3 \& 5]$

$\Rightarrow \text{sup}(135) / \text{sup}(1) = 2/3 > 60\%$

$\{5\} \rightarrow \{1,3,5\} - \{35\} [5 \rightarrow 1 \& 3]$

~~X~~ $\text{sup}(1,3,5) / \text{sup}(5) = 2/4 < 60\%$

Frequent Patterns:
 eg {milk, bread}, {computer, Antivirus} →
 type: Itemset, sequential

FP Tree Construction:

min support → 2

TID	Items
1	b d c a
2	e d c
3	a b
4	a c d
5	f g d b

Item	Support
d	4
a	3
b	2
c	2
f	1
g	1
e	1
t	1

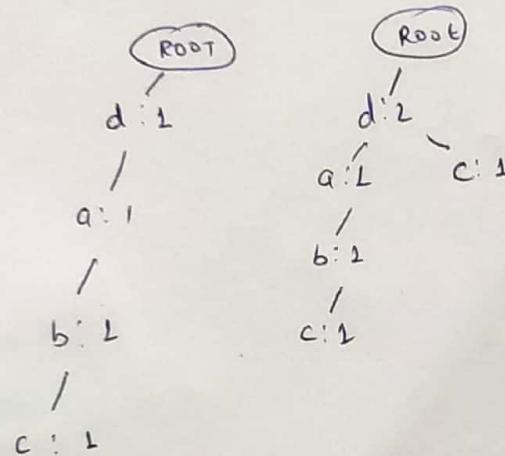
Priority ↓

ordered frequent items

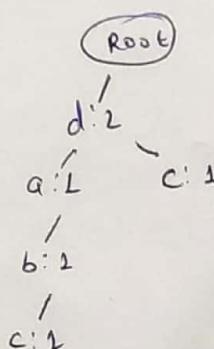
TID	Items
1	d a b c
2	d c a
3	a b
4	d a c
5	d b

construct FP tree

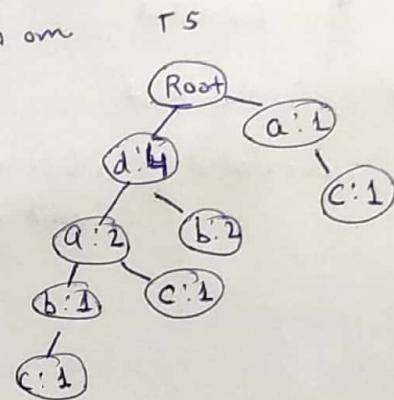
Transition 1



T 2.



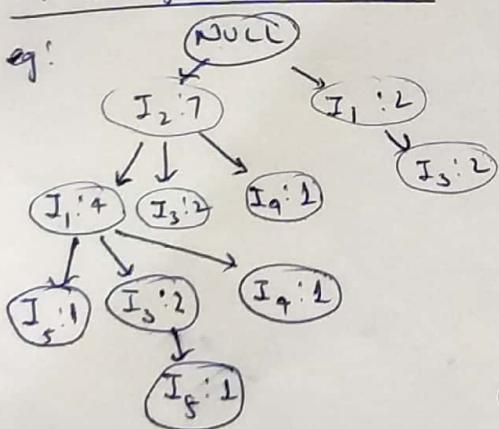
... 30 min



T 5

FP tree growth Algorithm:

eg:



I₅:2

I₄:2

I₃:6

I₁:6

I₂ is ignored as its
connected to root.

conditional Pattern Based:

For I_S:

{I₂, I₁:2} {I₂, I₁:1}

For I_A:

{I₂, I₁:1} {I₂:1}

For I_B:

{I₂, I₁:2} {I₂:2} {I₁:2}

For I_L:

{I₂:4}

conditional FP-tree.

Find common prefix

For I_S:

{I₂:2, I₁:2} {I₂:1}

For I_A:

{I₂:2} {I₂:1}

For I_B:

{I₂:4, I₁:2} {I₁:2}

For I_L:

{I₂:4}

min support = 2

Frequent Pattern Generation!

sets to vapas add kardo

eg $I_5 \nearrow I_2 \rightarrow I_5$ $\rightarrow \{I_2, I_5 : 2\}$
 $I_5 \searrow I_2 \rightarrow \{I_1, I_5 : 2\}$
 $I_5 \rightarrow I_2 \rightarrow \{I_2, I_5 : 2\}$

See video! 5 minute
eigenly.

LIFT and χ^2

$\text{lif}t(A, B) = 1$ A, B are independent
 > 1 positive correlated
 < 1 -ve correlated

$$\text{lif}t(A, B) = \frac{\text{confidence}(AB \rightarrow B)}{\text{support}(B)}$$

$$\chi^2 = \sum \frac{(O \text{ (Observed)} - E \text{ (Expected)})^2}{E}$$

$$= \frac{\text{support}(A \cup B)}{\text{support}(A) \times \text{support}(B)}$$

$$= \frac{P(A \cup B)}{P(A) P(B)}$$

$\chi^2 = 0$: independent

> 0 correlated, either + or -
need additional test