

PCA 2 dimension
T-SNE reduction

KNN - LSH, KDTree

K-fold / cross val

Local Outlier Factor

Classification metrics

ROC - AUC curve

Log-loss < binary
multiclass

R-square score

Naive Bayes + Bernoulli

Gaussian Naive Bayes

Logistic regression

Sigmoid squashing

L1, L2 regularization

→ loss minimization

Grid, random search

Non-linearly separable

Linear Regression

SGD, Batch-GD

PCA derivation

Bag of words

TF-IDF + word2vec

SVM + hard margin
soft margin

dual form - SVM

Kernel SVM

→ RBF
→ poly kernel

nu-SVM

SV regressor

(2) (3)

(4)

(6)

(7)

(12) (13)

(15)

(15) (16)

(16)

(17)

(18)

(21)

(23)

(25)

(27)

(28)

(28)

(29) (1)

(31)

(32)

(34)

(35)

(37)

(39)

(42)

(43)

(45)

(45)

Decision Tree

DT Regression

Ensemble tech

Random forest

Extremely rf

Boosting < residual
loss min²

Gradient Boosting

Adaboost, stacking
XGBoost

Kmeans

Lloyd's, Kmeans++

K-medoids

Agglomerative cluster

DBSCAN clustering

Grid search

DB Scan

Recommender System

Matrix Factorization

Truncated SVD

L1 vs L2 regularizn

SVD + L2 reg

stacking, cascading

pseudo-res vs residual

DT-Regression

Kernel - SVM

(46)

(49)

(50)

(51)

(52)

(53)

(55)

(57)

(56)

(59)

(60)

(63)

(64)

(66)

(28)

(66)

(70)

(71)

(76)

(78)

(79)

(81)

(82)

(84)

(85)

Data Normalization

→ convert all data to $0 \leq x_i < 1$

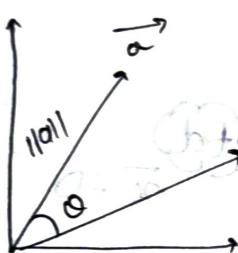
$$x_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

$$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

$$x = x_{\max} = 1$$

$$x = x_{\min} = 0$$

Dot product



$$\rightarrow a \cdot b = [a_1 b_1 + a_2 b_2 + \dots]$$

$$a \cdot b = a^T \cdot b$$

$$\rightarrow a \cdot b = ||a|| ||b|| \cos \theta$$

$$||a|| = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$$

$||a|| \rightarrow$ distance from $(0,0)$

* if $\theta = 90^\circ$,

$$\rightarrow a \cdot b = 0$$

* if \vec{a} is unit vector

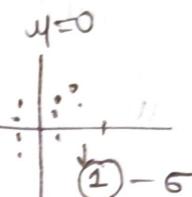
$$\vec{a} \cdot \vec{a}^T = 1 \quad \text{if } \vec{a} \text{ is a unit vector}$$

Data standardization

→ apply Z-score

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

$$\mu = 0, \sigma = 1$$



Covariance

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} f_1(x_1) \\ f_2(x_1) \\ \vdots \\ f_n(x_1) \end{bmatrix}$$

$$\text{Cov}(f_1, f_2) = \frac{1}{n} \sum (x_i - \mu_1)(y_i - \mu_2)$$

$$\text{Cov}(f_1, f_2) = \frac{1}{n} \sum (x_{i1} - \mu_1)(x_{i2} - \mu_2)$$

* if data is std. $\rightarrow \mu = 0, \sigma = 1$

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum x_i y_i$$

$$S_{dxd} = \frac{1}{n} f_1^T f_2$$

$$S_{dxd} = X_{d \times n}^T \cdot X_{n \times d}$$

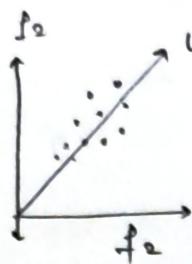
$S_{d \times d}$ — covariance matrix

$$0 = \mu \cdot \mu$$

$$n = 1 \cdot 1$$

②

Principal Component Analysis



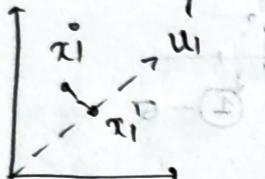
* find direction of a unit-vector such that variance is maximized.

* the variance can be explained by fewer features.

d=2

unit vector x_i : data point

x_i' → projection of x_i on \vec{u}_1



$$x_i' = \text{proj}_{\vec{u}_1} x_i = \frac{\vec{u}_1 \cdot x_i}{\|\vec{u}_1\|} = \vec{u}_1 \cdot \vec{x}_i$$

$$x_i' = \vec{u}_1^T \vec{x}_i$$

→ find \vec{u}_1 such that $\text{var}(x_i')$ is maximized

$$\text{var}(x_i') = \frac{1}{n} \sum_{i=1}^n (x_i^T \vec{u}_1 - \vec{u}_1^T \vec{x}_i)^2$$

as $\frac{\partial \text{var}}{\partial \vec{u}_1} = 0$

∴ maximize $\frac{1}{n} \sum_{i=1}^n (\vec{u}_1^T x_i)^2$

such that

$$\vec{u}_1^T \vec{u}_1 = 1$$

Eigen values

$S_{d \times d}$ (covariance matrix)

for any matrix,

($d \times d$) mat

(d) eigen
vector

we can prove that, $\vec{u}_1 = \vec{v}_1$ → largest eigen vector

let $d=2$, λ_1, λ_2

$\lambda_1 = 3, \lambda_2 = 2$

λ_1 can explain $\frac{3}{5} \times 100\%$ of variance

$$\lambda_1 = 3$$

$$\lambda_2 = 0$$

$$\therefore \text{var exp by } \lambda_1 = \frac{3}{3} = 100\%$$

$$\begin{aligned} \lambda_1 v_1 &= S v_1, \\ \lambda_n v_n &= S v_n \\ \forall i & \quad i=1 \rightarrow n \end{aligned}$$

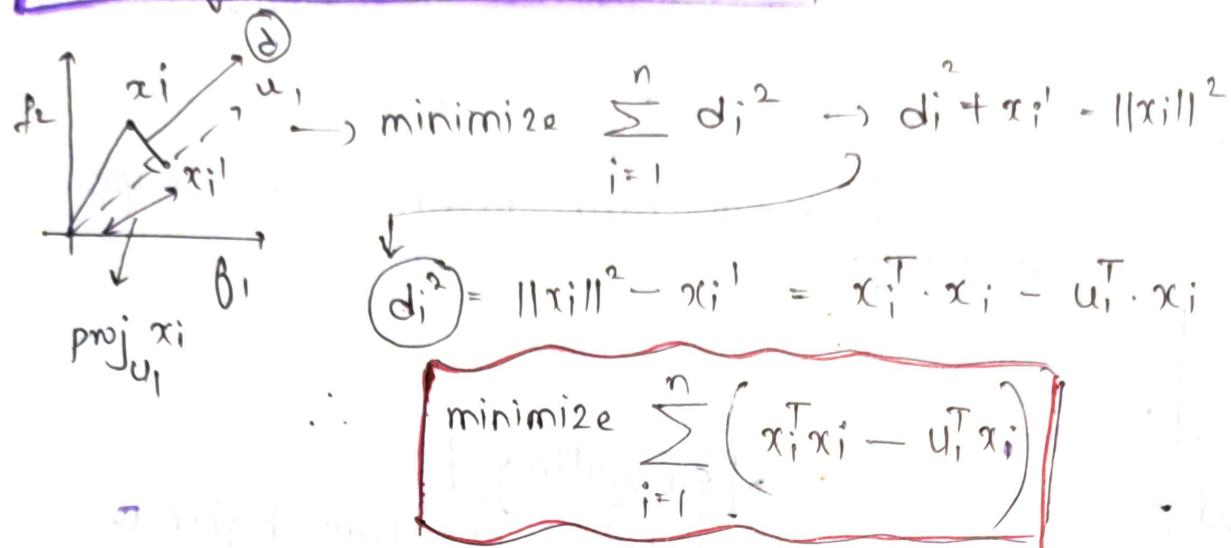
also, $v_i \perp v_j$

$$v_i \cdot v_j = 0$$

λ_i can explain $\frac{\lambda_i}{\sum \lambda_n} \cdot 100\%$ of var

PCA using distance minimization

(3)



PCA for dim reduction

Consider $X_{15k \times 784}$ → top 200 largest Eigen vectors

$$X_{15k \times 784} \times V_{784 \times 200} \rightarrow X_{15k \times 200}$$

$V = [v_1 \ v_2 \ v_3 \ \dots \ v_{200}]$

$v_i = (200 \times 1)$ size

Workflow

① standardize data among dim

② covariance mat $\rightarrow X^T \cdot X \leftarrow$

③ calculate eigen values and vectors for

④ pick top n eigen values vector for n-dimensions

Code

```
→ from sklearn.decomposition import PCA
→ pca = PCA(n_components = 200)
→ p-data = pca.fit_transform(data)
```

$$X_{n \times d} \longrightarrow X_{n \times 200}$$

plot % var

$$\text{var} = \frac{\text{pca.explained_var_sum}}{\text{sumc}}$$

$$\text{cumsum} = \text{np.cumsum}(\text{var})$$

(4) T-SNE

T-distributed stochastic Neighborhood Embedding

T-SNE → can preserve local and global shape.

↳ PCA can't preserve local shape



Neighborhood

↳ points which are nearer to given point

Embedding

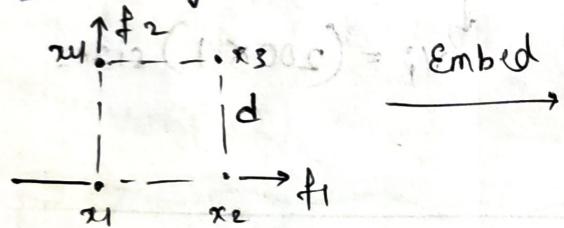
→ place pts from higher to lower dimensional space

stochastic

→ probabilistic

↳ results may vary each time

Crowding problem

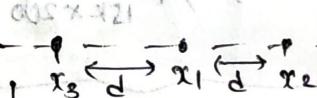


$$N(x_1) = x_2, x_4$$

$$N(x_2) = x_3, x_1$$

$$N(x_3) = x_4, x_2$$

$$N(x_4) = x_3, x_1$$



→ you can't place x_{10} on 1d plane and preserve neighborhood with both x_3 and x_1

→ only x_1 can be preserved

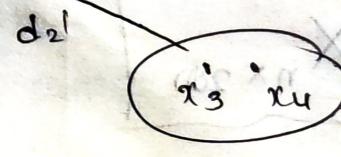
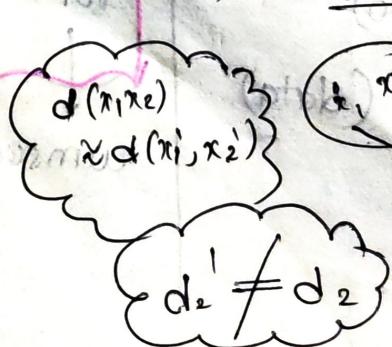
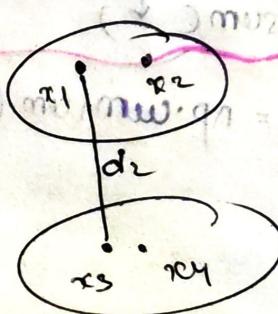
Geometric Intuition

→ intra cluster nbrhood is preserved

→ inter cluster nbrhood may not be preserved

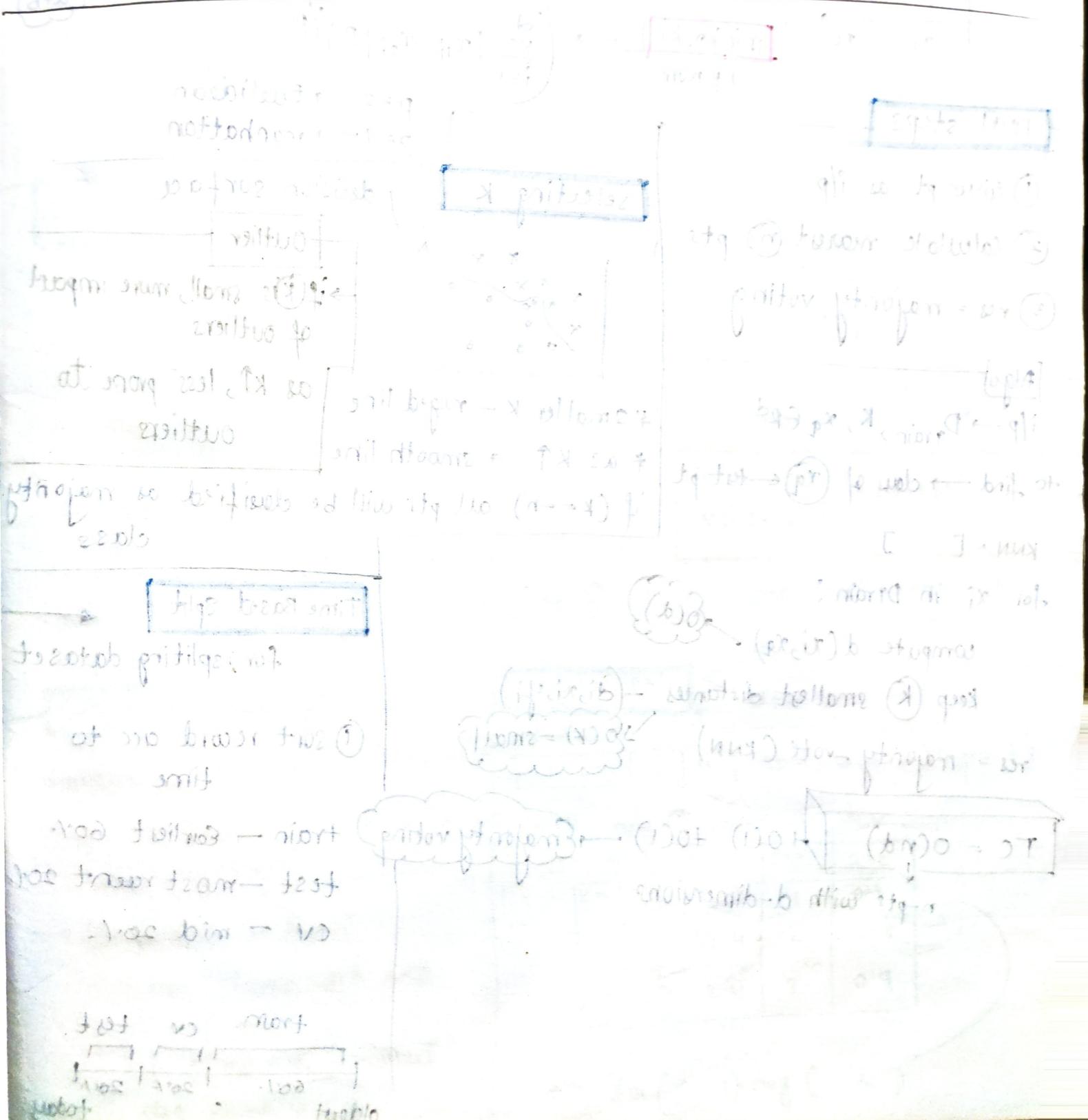
$$d - \text{dim}$$

$$d' - \text{dim}$$



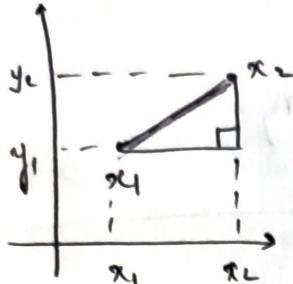
5) t-sne tips

- ① always check multiple params → iteration → perplexity → epsilon → (homrate)
- ② change it until shape stabilize
- ③ $2 \leq \text{perplexity} \leq n$
- ④ cluster size and intercluster distances may vary



dist - 2 pts
norm - 2 vectors

Distances



KNN

hyper-p = K

no of dimension p.

6

L2

Euclidean

$$\rightarrow \sqrt{\sum (x_{2i} - x_{1i})^2 + (y_{2i} - y_{1i})^2} \rightarrow$$

Manhattan

L1

Minkowski

Lp norm

$$\rightarrow \sqrt{\sum |y_2 - y_1| + |x_2 - x_1|} \rightarrow$$

absolute

a+b

KNN steps

- ① given pt as i/p
- ② calculate nearest n pts
- ③ res = majority voting

algo

i/p $\rightarrow D_{train}, K, x_q \in \mathbb{R}^d$

to find \rightarrow class of $(x_q) \leftarrow$ test pt

KNN = []

for x_i in D_{train} :

compute $d(x_i, x_q)$

(d)

keep K smallest distances

(d_i, x_i, y_i)

res = majority vote (KNN)

$\downarrow O(K) = \text{small}$

majority voting

$T_C = O(nd)$

$+ O(1) + O(1) \rightarrow$

n pts with d-dimensions

$p=2 \rightarrow$ Euclidean

$p=1 \rightarrow$ Manhattan

Selecting K



decision surface

Outlier

\rightarrow if K is small, more impact of outliers

as $K \uparrow$, less prone to outliers

* smaller K - rigid line

* as $K \uparrow \rightarrow$ smooth line

if ($K = n$) all pts will be classified as majority class

Time Based Split

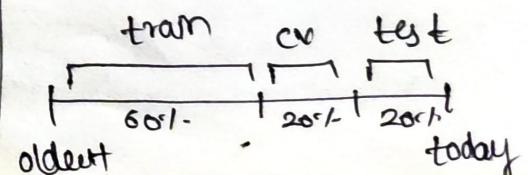
for splitting dataset

- ① sort record acc to time

train - earliest 60%

test - most recent 20%

CV - mid 20%



Cosine similarity and distance

(7)

$$\rightarrow \text{assim} \uparrow \{ \text{dist} \downarrow \}$$

angle b/w vectors

$$\cos\text{-sim} = \cos\theta$$

$$\cos\text{-dist} = 1 - \cos\theta$$

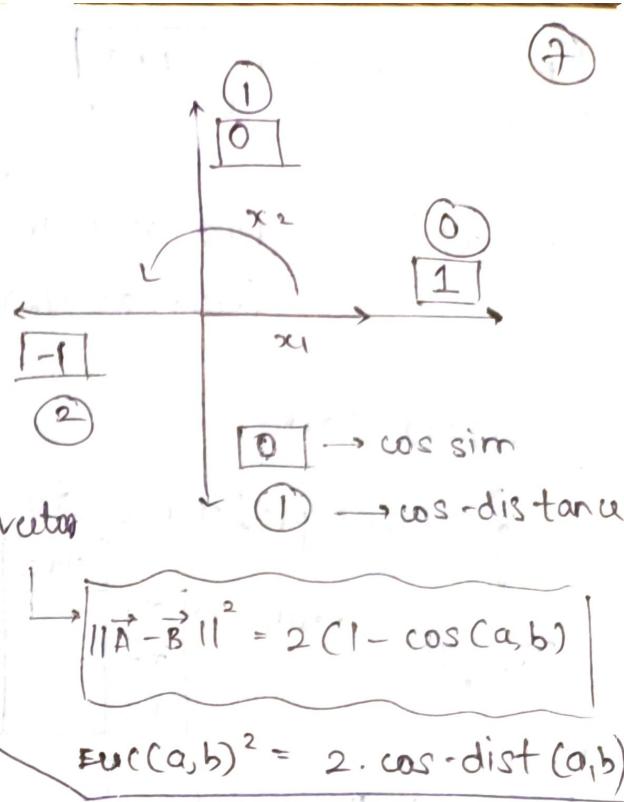
$\cos\text{-sim} = +1$ similar

$\cos\text{-sim} = -1$ different

$$\cos\theta = \frac{\vec{x}_1 \cdot \vec{x}_2}{\|\vec{x}_1\| \|\vec{x}_2\|}$$

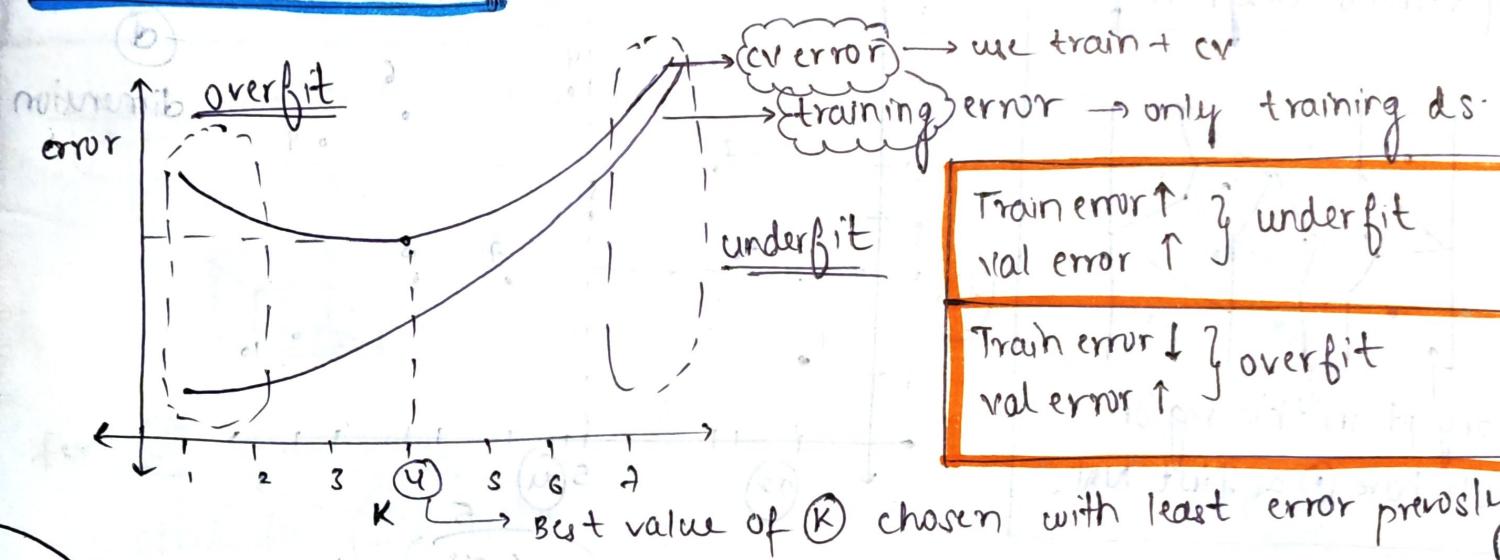
① \vec{x}_1, \vec{x}_2 are unit vectors

$$\cos\theta = \vec{x}_1 \cdot \vec{x}_2$$



$$\text{Euc}(a, b)^2 = 2 \cdot \cos\text{-dist}(a, b)$$

Overfit and underfit



Train error ↑ } underfit
val error ↑ }

Train error ↓ } overfit
val error ↑ }

Choose val - KNN

dataset → dtrain:
(input 100)
dtest:
(output 100)

→ for (x_i, y_i) in dcv:
use dtrain and find
 $f(x_i)$
if $f(x_i) = y_i \rightarrow$ correct
else → wrong

K-fold CV

d-train

1	2	3	4
---	---	---	---

split dtrn into parts

→ pick any ① part as dcv.

Ex: K=1	train	cv	acc
	1 2 3	4	0.9
	1 3 4	2	0.8
	1 2 4	3	0.9

acc for $(K=1) = \arg(\text{acc})$

Regression using KNN

→ majority vote → classification

⑧

inp x_i → find n-nearest neighbours → mean ($y_1, y_2 \dots y_n$)
 $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$
 or
 median ($y_1, y_2 \dots y_n$)

Weighted KNN

give more preference to more nearer points

3NN →

	class	distance	weight	
x_1	+	0.1	1	
x_2	-	0.9	$\frac{1}{\text{dist}}$	
x_3	-	0.10		

$$TC(KNN) = O(n)$$

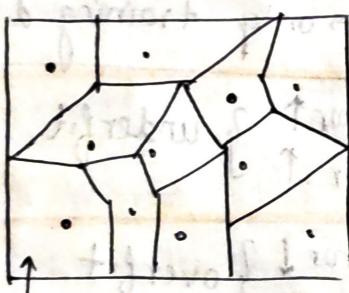
using Kd tree
 $\rightarrow O(\log n)$

Kd tree

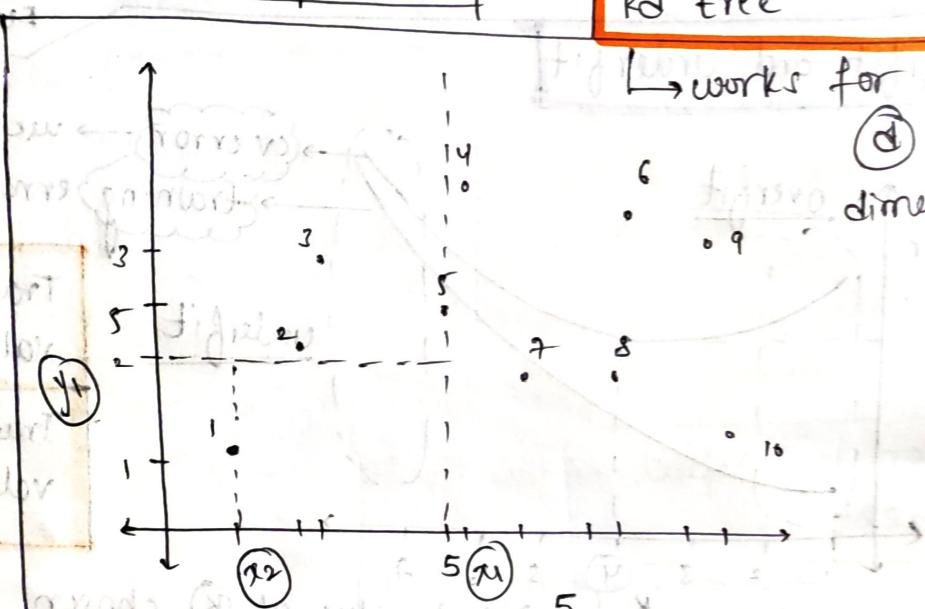
↳ works for small dimension

④

dimension



any pt in this region
 will have ① as first NN.



① project line at x axis

$x_1 = 5$ → median

$x \leq x_1$

No

② pick y axis

$y_1 = 2$

$y < y_1$

No

③ pick x axis

$x = 1$

$x > x_1$

leaf

Kdtree

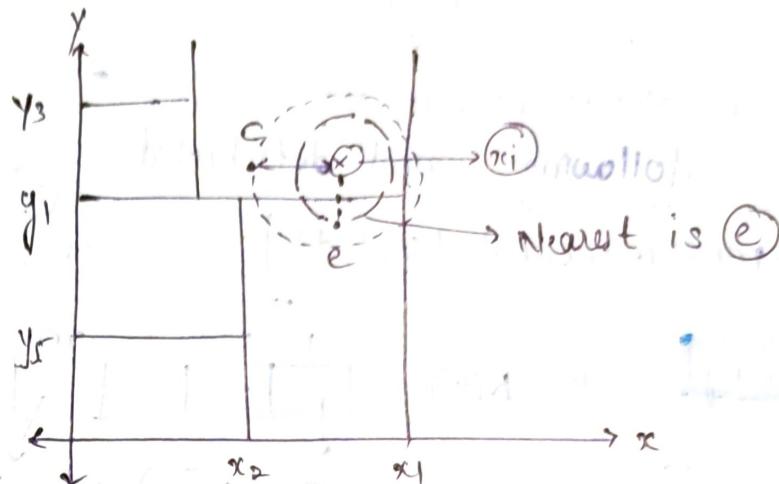
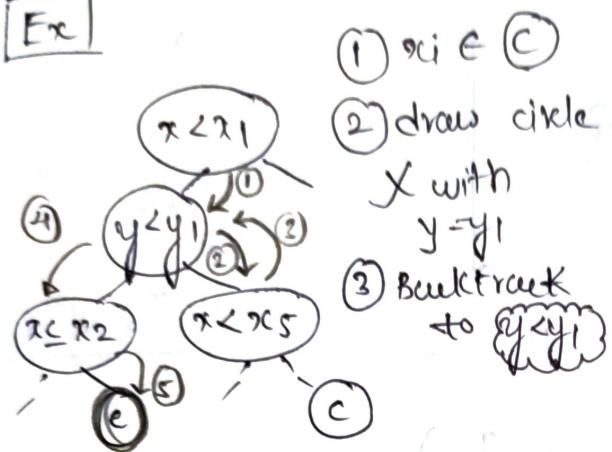
Best case - $O(\log n)$
 Worst case - $O(n)$

Find nn using kd-tree

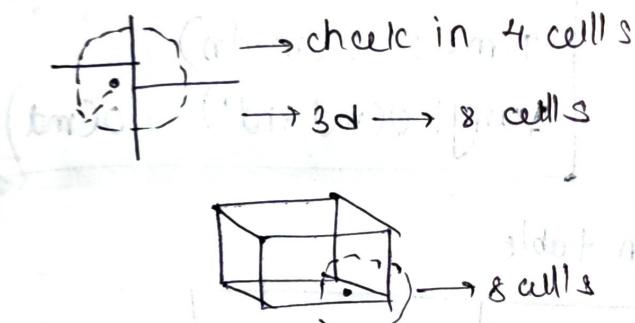
⑨

- ① use kd-tree to find nearest neighbor
- ② draw a circle around that pt and check for any other nearest pt.

Ex



Limitations of kd tree



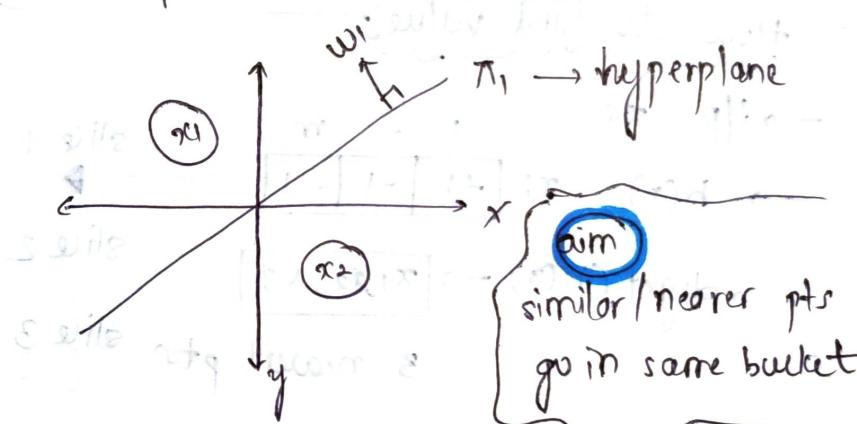
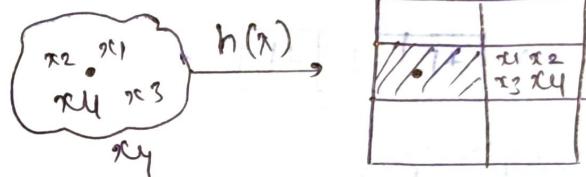
for d dimensions,
 check for 2^d cells

when d is large,

$$T.C \rightarrow O(d \log n)$$

$$2^d \log n \approx n \log n$$

Locality Sensitive Hashing, LSH



$$w_1 \cdot x_1 \geq 0$$

$$w_1 \cdot x_2 \leq 0$$

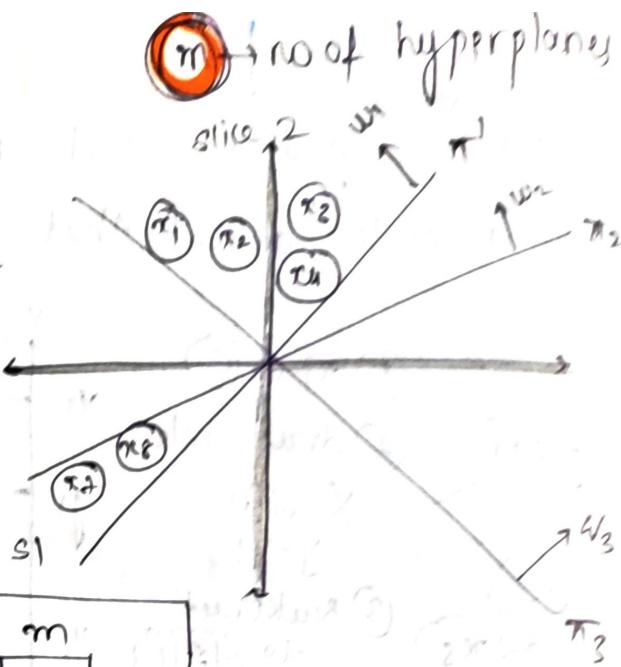
LST for cos-sim

① Generate random hyperplane

$\vec{w} \rightarrow$  d-dim vector

Each value is random no
following normal distributed

`np.random.normal(0, 1, d)`



Ex: slice ^{s1}

$$h(x) = \boxed{\quad | \quad | \quad | \quad | \quad | \quad} \quad \begin{matrix} 1 & 2 & \dots & m \end{matrix}$$

$$\text{sign}(\omega^T x) \rightarrow \text{sign}(\omega_m^T x)$$

$$(x_7, \tau_8) \rightarrow$$

$$(x_7, x_8) \rightarrow \begin{array}{|c|c|c|} \hline +1 & -1 & -1 \\ \hline \end{array}$$

slice ②

$$\begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 \end{pmatrix} \rightarrow \boxed{\begin{array}{|c|c|c|} \hline +1 & +1 & +1 \\ \hline \end{array}}$$

24

space :- $O(nd)$

Time: $O(m \cdot d \cdot n)$

Query: $O(md + nd')$ $\rightarrow O(md)$

How to find value

$\rightarrow i/p \quad x_i$

$$\rightarrow h(x) = x_i \begin{array}{|c|c|c|} \hline +1 & -1 & -1 \\ \hline \end{array}$$

$\rightarrow \text{d.get}(\text{h}(x)) \rightarrow [x_1, x_2, x_3]$

-- m

slice 1

slice 2

3 nearest pts

slice 3

$\begin{array}{ c c c } \hline & 1 & 1 \\ \hline 1 & & \\ \hline \end{array}$	$\{x_2, x_1\}$
$\begin{array}{ c c c } \hline -1 & +1 & 1 \\ \hline & & \\ \hline \end{array}$	$\{x_1, x_3, x_2\}$
$\begin{array}{ c c c } \hline -1 & -1 & -1 \\ \hline & & \\ \hline \end{array}$	$\{x_1, x_2\}$

problems

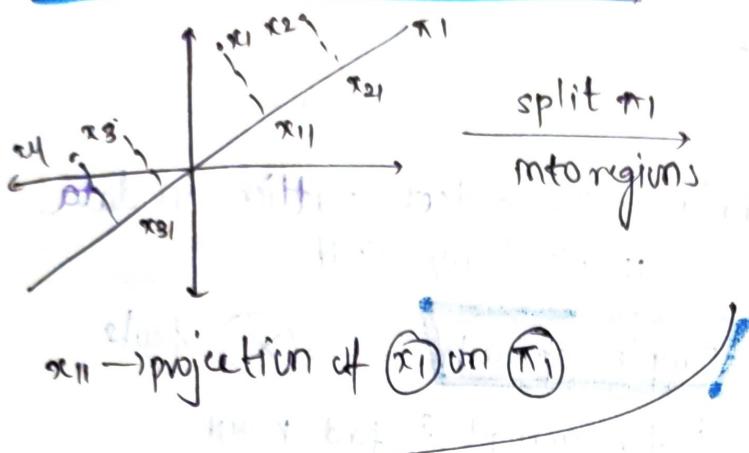
 you could miss nearest neighbor (a date) to
nearest cosine-similarity ②

501-5

have different hashtable
hyperplanes

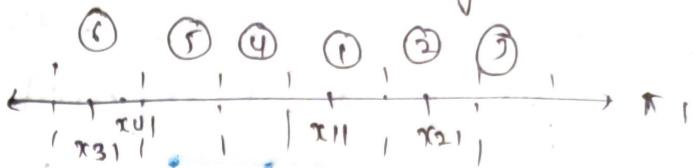
with different

LSH for Euclidean values



* pts falling in some bucket are closer

divide π into ① regions



$$h(\pi) = \boxed{\quad \quad \quad m}$$

region in which pt falls in
on π_1 (for π_1)

$$h(x_{11}) \rightarrow \boxed{1 \quad - \quad -}$$

$$h(x_{21}) \rightarrow \boxed{6 \quad 1 \quad - \quad -}$$

KNN probabilistic class label

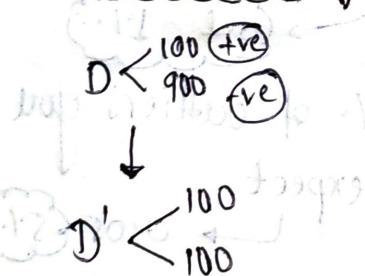
5-NN $(+, +, +, +, +) \rightarrow$ +ve } not fair
 $(+, +, +, -, -) \rightarrow$ +ve

$$\begin{aligned} P(x_i = +) &= 1 \\ P(x_i = -) &= 3/5 \end{aligned}$$

Classification Metrics

imbalanced dataset issue

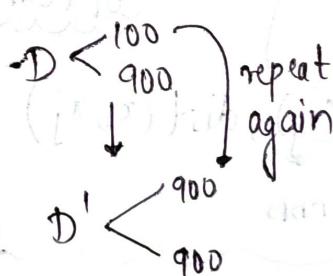
① Undersampling



eliminate 800 randomly

from -ve class

② Oversampling



① repeat

② create synthetics
pts - Extrapolation

③ Class weights

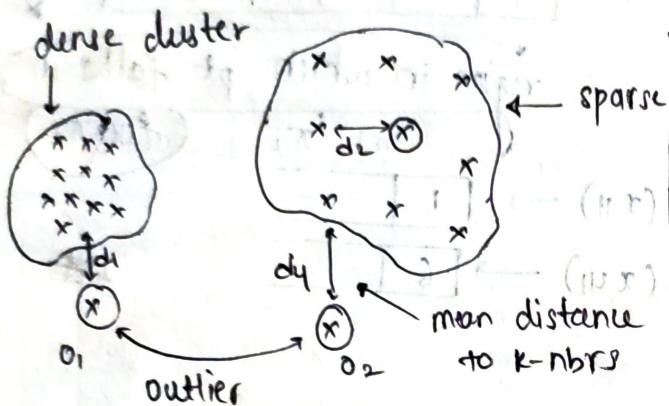
- a) minority class weight
- b) majority class weight

Impact of outliers

- * smaller K more affected by outliers
- * larger K doesn't get affected

local outlier factor LOF

→ technique to detect outliers in data inspired by KNN



Terms

- * k -distance (x_i): distance of x_i from k th nearest neighbor of x_i

Neighborhood (x_i)

$$N_{k=5}(x_i) = \{5 \text{ nearest nbrs}\}$$

Reachability distance

$$rd(x_i, x_j) = \max(k\text{-distance}(x_j), \text{dist}(x_i, x_j))$$

↓
farthest k -nbrs

Ex: $k=5$

if $x_i \in N_{k=5}(x_j)$

$$\text{reach-dist} = 5\text{-dist}(x_i)$$

else

$$\text{reach-dist} = d(x_i, x_j)$$

Simple approach

① for each pt x find K -NN

② calculate mean distance, \bar{d} → KNN

③ sort all to mean dist

large mean dist = outlier

fails: O_2 detected as outlier → $d_1 > \bar{d}$
 O_1 not detected

MAIN SOLN

LOF → code

scikit-learn.nbrs.LocalOutlierFactor

$n_nbrs = n, \dots$

wontamination = ct

def = 1.1

% of outliers you expect

→ 0.05 5%

LOF Flow

① calculate LOF & $x_i \in \text{Dataset}$

② pick top - N (highest)

LOF pts

→ outliers

focal reachability density LRD

density $\rightarrow \times$ per unit area

(13)

$$LRD(x_i) = \frac{1}{\sum_{x_j \in N(x_i)} \text{reach-dist}(x_i, x_j)} \cdot \frac{1}{|N(x_i)|}$$

size of set $\approx K$

\rightarrow average reach-dist of x_i from its neighbors

$$\frac{1}{\text{avg reach-dist}}$$

$$LRD(x_i) \approx \frac{|N(x_i)|}{\sum \text{reach-d}(x_i, x_j)}$$

(No. of pts) / distance

LOF

$$LOF(x_i) = \frac{\sum_{x_j \in N(x_i)} LRD(x_j)}{|N(x_i)|} * \frac{1}{LRD(x_i)}$$

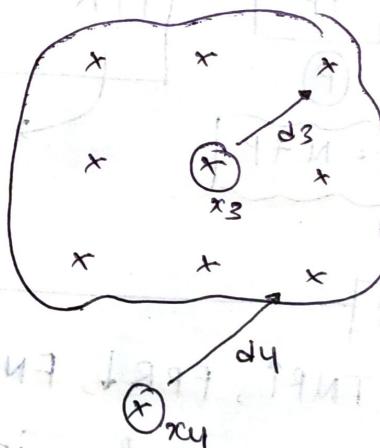
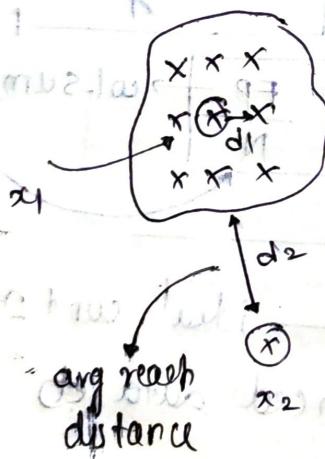
density in nbrs of x_i / avg LRD of pts in $N(x_i)$

LOF $\downarrow \rightarrow$ inlier

LOF $\uparrow \rightarrow$ outlier [a \uparrow b \downarrow]

$LRD(x_i)$ is small compared to nbrs
density around x_i

Example: 947



approximate
 * reach-dist: $d_1 < d_2 < d_3 < d_4$
 * Lrd $\rightarrow 1/\text{avg (r-d)}$
 $\frac{1}{d_1} > \frac{1}{d_2} > \frac{1}{d_3} > \frac{1}{d_4}$
 density(x_i) large

density $\propto \frac{1}{\text{distance}}$

① x_1

$$LRD(x_1) \approx LRD(x_1 \text{ en } N(x_1))$$

\rightarrow LOF small

③ x_3

$$LRD(x_3) \approx LRD(\text{nbrs of } x_3)$$

\rightarrow LOF small

④ x_4

\rightarrow both are small

② x_2

$$LRD(x_2) \ll LRD(x_2 \text{ en } N(x_2))$$

\rightarrow large LOF

$$LRD(x_4) \ll LRD(\text{nbrs of } x_4)$$

\rightarrow LOF is large

Class II metrics

Accuracy

problems

$$acc = \frac{\# \text{ correct preds}}{\text{total preds}}$$

① Imbalanced data

Let $D \leftarrow$
 80% +ve
 20% -ve

problem ②

it can't deal with prob scores

$$acc = \frac{80}{100} = 80\% \rightarrow \text{misleading}$$

y	M_1	M_2
0	0.9	0.6
0	0.8	0.56
1	0.1	0.3
1	0.2	0.4

→ it can't deal
with prob

only way → round off

$M_1 \rightarrow$	1	1	0	0
$M_2 \rightarrow$	1	1	1	0

acc to acc $M_1 = M_2$

but M_1 is better than M_2
in reality

Confusion matrix

actual		0		1	
		0	1	0	1
0		TN	FN	FP	TP
1		FP	TP	TP	FP

was 0, predicted as 1

0		1
0	TN	FN
1	FP	TP
		N P
total = N + P		

$TP \text{ rate} = \frac{TP}{P}$	$TNR = \frac{TN}{N}$
$FPR = \frac{FP}{N}$	$FNR = \frac{FN}{P}$
$FPR = \frac{FP}{N}$ → col-sum	

$$FPR = \frac{FP}{N}$$

F P

↓ what was the prediction
were u correct

uses

$TPR \uparrow, TNR \uparrow, FPR \downarrow, FNR \downarrow \leftarrow$ ideal cond
can work even for imbalanced data set

Ex: (Cancer)

very high TPR

very low FNR

issues

* you need to check 4 values

FP	Type I error
FN	Type II error

Precision

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

* of all **+ve** predictions, * of all the actually **+ve** points
how many were actually **+ve** how many were predicted properly

Recall

$$\text{Recall} = \text{TPR} = \frac{\text{TP}}{\text{P}}$$

15

F1 score

$$\text{F1-score} = 2 * \frac{\text{Precision} * \text{recall}}{\text{Precision} + \text{recall}}$$

$$\approx [\text{avg}(\text{recall}^{-1}, \text{precision}^{-1})]^{-1}$$

ideal \rightarrow F1↑, precision↑, recall↑

Binary classifier?

ROC curve

Receiver operating characteristic curve

can be applied to binary class

Step 1) sort acc to prob class-label (decreasing)

x	y	\hat{y}
x1	1	0.95
x2	1	0.8
x3	0	0.7
x4	0	0.6

\hat{y}	$\hat{y}_{\tau_1} = 0.95$	$\hat{y}_{\tau_2} = 0.8$
1	1	1
1	0	0
0	0	0
0	0	0

select
choose
 τ_i
by
itself

Step 2) Thresholding (τ_i)

if $\hat{y} > \tau_i \rightarrow +ve$
else -ve

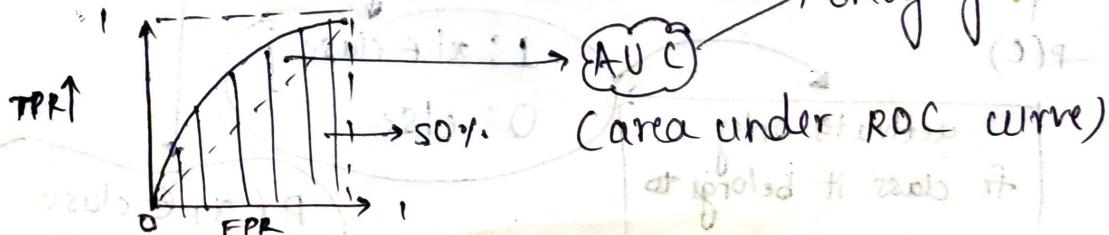
- 1) For each threshold calculate accuracy
- 2) select threshold with max accuracy

Step 4) Calculate FPR and TPR for each $\tau_1, \tau_2, \dots, \tau_n$

To calculate FPR/TPR u need 0,1 got by thresholding

Step 4) plot

only good for binary class



1 → best
0 → worst

AUC properties

- ① prone to imbalanced datasets
- ② only depends on order of \hat{y} and not actual value of \hat{y}
- ③ AUC of random model = 0.5 → model which gives 0 or 1 without logic random
- ④ AUC < 0.5: worse than random

Log-Loss

→ uses probability scores

x	y	\hat{y}
x_1	1	0.9
x_2	1	0.6
x_3	0	0.1
x_4	0	0.4

$$\begin{aligned} x_1 &\rightarrow -\log(0.9) = 0.0457 \rightarrow \text{closer to } 0 \\ x_2 &\rightarrow -\log(0.6) = 0.22 \rightarrow \text{far from } 1 \\ x_3 &\rightarrow -\log(0.1) = 0.0457 \rightarrow \text{closer to } 0 \\ x_4 &\rightarrow -\log(0.4) = 0.22 \end{aligned}$$

$$0 \rightarrow \infty$$

Log-Loss

lesser = better

$$\text{Log-Loss} = -\frac{1}{n} \sum_{i=1}^n (\log(P_i) * y_i) + (1-y_i) (\log(1-P_i))$$

activated when

class = +ve/1

class = -ve/0

English defn

Log-Loss = average of negative of $\log(P(\text{correct class label}))$
 $= \text{Arg of neg of log}(P(\text{correct class}))$

multiclass

(x)

(y)

(\hat{y})

$$(0-n) \quad \begin{matrix} 1 & 0 & \dots & 1 & 0 \end{matrix}$$

$$P(0) \dots P(C)$$

c classes

rows = n

$$\text{Log-Loss}_j = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m y_{ij} P \log(P_{ij})$$

activates only
for class it belongs to

-1: $x_i \in \text{class } j$
0: else

$P(x_i \in \text{class } j)$

R^2 score

only for regression

(17)

error

$$SS_{\text{residual}} = \sum_{i=1}^n (y_i - \hat{y})^2$$

diff from mean

$$SS_{\text{total}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$R^2 = \frac{1 - SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

case 1: No error $\rightarrow SS_{\text{res}} = 0 \rightarrow R^2 = 1 - 0 = 1$ but

case 2: $SS_{\text{res}} < SS_{\text{tot}} \rightarrow R^2 = 0 \text{ to } 1$

case 3: $SS_{\text{res}} = SS_{\text{tot}} \rightarrow R^2 = 0$ (ur model is simple mean model)

case 4: $SS_{\text{res}} > SS_{\text{tot}} \rightarrow R^2 < 0$ (your model is worse than mean model)

issues → can't handle outliers Ex: $R^2 \propto SS_{\text{res}}$

{ if any one error
↓ error is huge SS_{res} gets
affected }

MAD → median absolute deviation

mean std → median mad → not prone to outliers

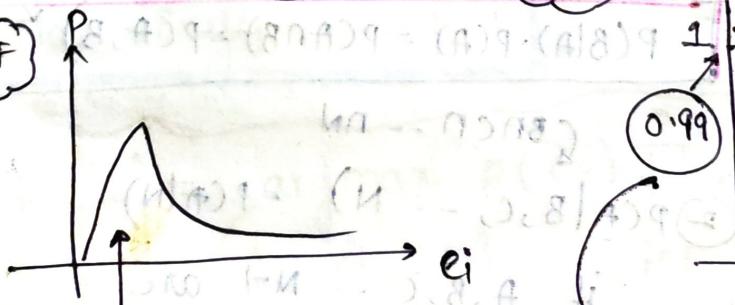
MAD → lower the better
median → errors are small assuming
only few outliers have large error

distribution of errors

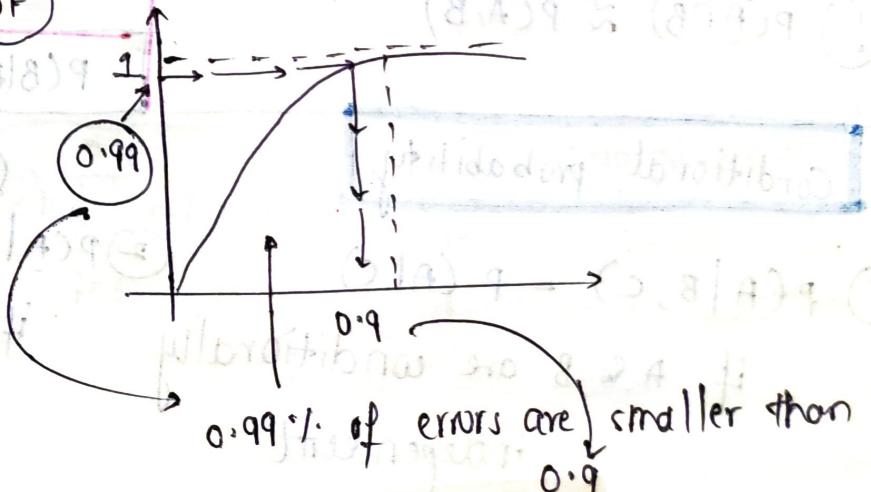
Errors $\rightarrow e_i$

CDF

PDF



goal: most e_i 's are small



Naive Bayes

Chiper param = α

(18)

* conditional prob

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

; $P(B) \neq 0$

probability of occurrence of (A) given (B) already occurred

Events

* Independent event

A, B : events

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

* mutually exclusive

$$P(A|B) = P(B|A) = 0$$

$\therefore A, B$ r mutually exclusive events

also $P(ANB) = 0$

Ex: A: 3 on Dice 1

B: 2 on Dice 1

$$P(A|B)$$

3 has occurred on (B)

(B) → it can't be (2)

as already (3)

Bayes Theorem

Derivation

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

$$\rightarrow P(A|B) = \frac{P(ANB)}{P(B)} \rightarrow ①$$

$$\rightarrow P(ANB) = P(BNA) \rightarrow ③$$

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

$$\rightarrow P(BNA) = P(B|A) P(A) \rightarrow ②$$

sub ② in ① using prop ③

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

$$P(B|A) \cdot P(A) = P(ANB) = P(A, B)$$

Conditional probability

$$① P(A|B, C) = P(A|C)$$

if A & B are conditionally independent

$$② P(A|B, C, \dots, N) = P(A|N)$$

if A, B, C, ..., N-1 are

conditionally independent

Naive assumption :- all features are conditionally independent

(19)

point to be classified : $X = \langle x_1, x_2, \dots, x_n \rangle$ → features

* NO of class → K

* c_k - class K

$$P(A) P(B|A) = P(A, B)$$

! conditⁿ probability

$$P(c_k | X) = \frac{P(c_k) P(X|c_k)}{P(X)}$$

$P(X)$ → doesn't affect much as it is same for $k=1-n$

scaling term

$$P(c_k | X) \propto P(X|c_k)$$

chain rule

$$P(c_k, x_1, x_2, \dots, x_n) = P(x_1, x_2, \dots, x_n, c_k)$$

$$= P(x_1 | x_2, x_3, \dots, x_n, c_k) \cdot P(x_2, x_3 | c_k)$$

$$= P(x_1 | x_2, \dots, x_n, c_k) \cdot P(x_2 | x_3, x_4, \dots, x_n, c_k) \cdot P(x_3, x_4, \dots, x_n, c_k)$$

$$= P(x_1 | x_2, \dots, x_n, c_k) \cdots P(x_{n-1} | x_n, c_k) P(c_k) \rightarrow ①$$

assumption x_1, x_2, \dots, x_n are conditionally independent

$$P(x_1 | x_2, \dots, x_n, c_k) = P(x_1 | c_k)$$

$$\textcircled{1} \Rightarrow P(x_1 | c_k) P(x_2 | c_k) \cdots P(x_{n-1} | c_k) \cdot P(c_k)$$

$$\therefore P(c_k | x_1, x_2, \dots, x_n) \propto \left(\prod_{i=1}^n P(x_i | c_k) \right) \cdot P(c_k)$$

product

Naive-Bayes-classifier

$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} \prod_{i=1}^n P(x_i | c_k)$$

→ NO of features

→ max value

Naive Bayes with text

(x_i): sentences $\xrightarrow{\text{process}} \text{Bag of words} \rightarrow w_1, w_2, \dots, w_n$
 tokenizing
 stemming
 stop-word removal

(y_i) \rightarrow class [Ex. 0-1]

words	class
—	0
—	0
—	1
—	1

steps

- to find $p(\text{yes} | \langle w_1, w_2, \dots, w_n \rangle) = p(\text{yes}) \cdot p(w_1 | \text{yes}) \cdot p(w_2 | \text{yes}) \dots$
 - $p(w_i | \text{yes}) = \# \text{ points which have class=yes}$
 and have (w_i) word
- $\# \text{ points which have class=yes}$

Laplace Smoothing

- test-data = $\langle w_1, w_2, w_3, w_4 \rangle / w' \rightarrow \text{word which was never seen in training data}$
- $p(\text{yes} | \text{td}) = p(\text{yes}) \cdot p(w_1) \dots p(w_4)$ → you can't calculate $p(w')$

solution: ways

ignore i.e. $p=1 \rightarrow$ if you ignore implies $p=1$ false

$p=0$ whole $p(\text{yes} | \text{td}) = 0$

smoothing

$$p(w' | \text{yes}) = \frac{\# \text{ with } w'}{\# \text{ pts } y = \text{yes}}$$

$= \frac{0}{n_{\text{yes}}} \rightarrow \text{can't be } 0$

$\alpha \rightarrow$ randomly chosen

as $\alpha \rightarrow \infty$

$$p(w' | \text{yes}) \rightarrow 1/2$$

\therefore if $\alpha \uparrow \rightarrow$ uniform probability

K \rightarrow no. distinct values = 2

either word can be there or not

for run-time data \rightarrow no. of values it can take $E.g.$ Outlook \rightarrow 3 rainy, windy, sunny

$$\therefore P(w_i | y=1) = \frac{\# \text{ pts with } (w_i) \& y=1 + \alpha}{\# \text{ pts with } y=1 + \alpha K}$$

(21)

→ no of distinct values
→ (2) for text

Log-probabilities

→ provide numerical stability

$$\text{probm} := p(yes|x) = p(yes) \cdot p(w_1|yes) \cdot p(w_2|yes)$$

↓
small nos ex: $(0.0001) (0.002)$ $\log ab = \log a + \log b$

→ product can quickly become too small

→ Ex: $0.0 \dots 0009$

after ~16 digits after ① python starts rounding off

$$\text{soln}: \log(p) = \log(p(yes)) + \sum_{i=1}^n \log(p(w_i|yes))$$

② impacts more from minority class

* Feature importance: ① w_i having highest $P(w_i|y=\text{class})$ is more imp

* Imbalanced dataset

$$n \begin{cases} n_+ \\ n_- \end{cases} \downarrow \quad n_+ \gg n_- \implies p(\text{class} = +) \gg p(\text{class} = -)$$

will cause problem

majority class is favored by Laplace smoothing.

Soln: ① ignore $p(yes), p(\text{No})$ term ② choose diff value of α for each class ex: α_+, α_-
③ upsampling, downsampling

Numerical features

Gaussian NB

				+
				+
				+
				+
				-
				-
				-

$$\rightarrow x_{ii} = (f_i \text{ col of } x_i) \in \mathbb{R}^d$$

$D \setminus D^+$: data with $y=0$
 D^+ : data with $y=1$

$$P(yes|x_i) = P(yes) \prod_{j=1}^n P(x_j|y=1)$$

Q: 220b

$$\prod_{i=1}^n \prod_{j=1}^{n-i+1} \dots$$

? → ?

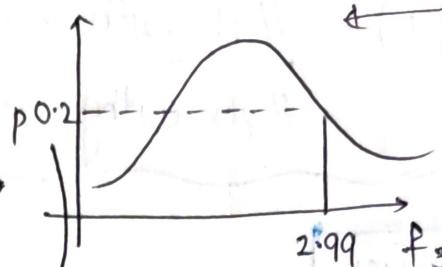
solution :-

$$\text{exp}(x_{i3} | y_{\text{es}})$$

	f_1	f_2	f_3
x_i			

→ plot pdf of f_3 for +ve/-ve

$$\rightarrow x_{i3} = 2.99 \quad p(x_{i3} | \text{yes}) = 0.2$$



assume gaussian/normal
distr

pdf of f_3 on D^{+ve}

step 3

① plot pdf of f_j on $D^{\text{class=val}}$

② assume pdf is gaussian/Normal distr $\rightarrow N(\mu, \sigma^2)$
and find p

③ $p(x_{ij} | \text{yes}) = \text{pdf}(x_i) \text{ from pdf of } f_j$

* NB can't work directly with distance matrix

$$w^T x_i = w^T X$$

$w \cdot X$ dot product

$$\sigma(w^T x_i) = p(y=1 | x)$$

predict point using Log P

$x_q \rightarrow$ find y by $w^T x_q \rightarrow \sigma(y)$

Probability Interpret
of LR

cost fn : for minimization

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -y_i \log(p_i) - (1-y_i) \log(1-p_i)$$

$$p_i = \sigma(w^T x_i)$$

activated
if class = 1 ie $y_i = 1$

activated for
class = 0

$$p(y_i = 1 | x_i) = \hat{y}_i = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$z = x_i^T \bar{w} = w^T w$$

Logistic Regression

(2)

used for classification.

Assumption :- pts are linearly or almost linearly separable

Task :- to find plane $\pi(w, b)$ which best separates the points

where eqn of plane

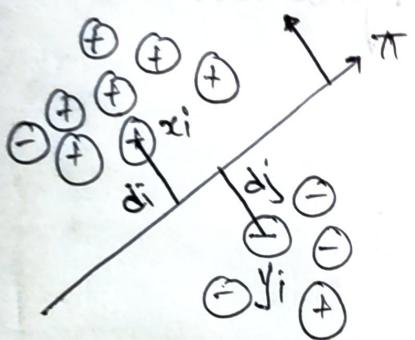
$w^T x + b = 0$ Intercept

$w^T x + b = 0$

$w^T x + b = 0$ normal to plane π

Overall Idea

* let (x_i, y_i) be training dataset and let $y_i \rightarrow +1 \rightarrow \text{positive}$
 $\rightarrow -1 \rightarrow \text{negative}$



→ assume π passes through origin

$\therefore w^T x = 0$ as π is unit vector $\|w\| = 1$

$$\rightarrow d_i = \frac{w^T x}{\|w\|} \quad \therefore d_i = w^T x_i \\ d_j = w^T x_j$$

→ $d_i > 0$; if x_i lies same side of π
 $d_j < 0$; if x_i lies on other side

classifier

: if $w^T x_i > 0 \rightarrow +1 \text{ point}$ $x_i, y_i \rightarrow \text{test point}$
 $w^T x_i < 0 \rightarrow -1 \text{ point}$

take term
predict?

$$y_i \cdot w^T x_i = y_i \cdot (y_i \rightarrow \text{prediction})$$

actual class

case 1 : $y_i = +1$, $w^T x_i > 0$

$$(w^T x_i)(y_i) > 0$$

correct

classify

case 3 : $y_i = +1$, $w^T x_i < 0$

$$(y_i \cdot w^T x_i) < 0$$

case 2 : $y_i = -1$, $w^T x_i < 0$

$$(w^T x_i)(y_i) > 0$$

case 4 :

$y_i = -1$, $w^T x_i > 0$

$$(y_i \cdot w^T x_i) < 0$$

∴ general idea is to have more correct classifications,
 ie $(w^T x_i)(y_i)$ must be as ~~as~~ as possible

general aim \rightarrow

$$\max_w \sum_{i=1}^n (y_i w^\top x_i)$$

$x_i, y_i \rightarrow$ constant
 w^\top variable

$$w^* = \operatorname{argmax}_w \left(\sum_{i=1}^n y_i w^\top x_i \right)$$

optimal w

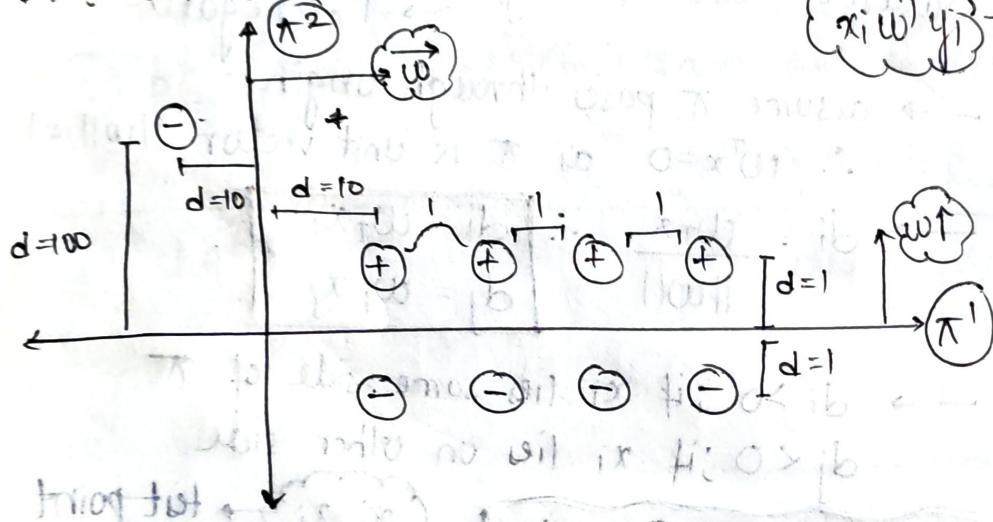
optimization

problem $\rightarrow \max_w$

strong all

Not stable as one outlier can distort optimization

Ex:- one outlier



$x_i w^\top y_i \rightarrow$ called as signed distance

case 1 π_1

$$\sum w^\top y_i x_i = 1 + 1 + 1 + 1 + 1 + 1 - 100 = -92$$

$$\boxed{\sum w^\top y_i x_i = -92}$$

case 2 π_2

$$\begin{aligned} \sum x_i y_i w^\top &= 1 + 1 + 1 + 1 - 1 - 1 - 1 \\ &\quad + (-100)(-1) \\ &= 100 + [10 + 11 + 12 + 13] \\ &= 10 - 11 - 12 - 13 \end{aligned}$$

$$\boxed{\sum x_i y_i w^\top = 100}$$

As to optimization,

π_1 is better than π_2 as ($100 >> -92$)

$$\text{accuracy}(\pi_1) = \frac{8}{9}$$

but logically (π_1) is better

$$\text{accuracy}(\pi_2) = \frac{5}{9}$$

$$\text{acc}(\pi_1) >> \text{acc}(\pi_2)$$

\therefore we need to modify

$$\boxed{\max \sum w^\top x_i y_i}$$

Sigmoid Squashing

(squash) $0, \infty \rightarrow [0, 1]$

(25)

* we found that if $d_i = w^T x_i$ was very large, it impacted model

solution

① use $f(w^T x_i y_i)$

↓
sigmoid-fn ($\sigma(x)$)

why ① tapering

- doesn't ↑ exponentially

for $d \uparrow$

② probabilistic

$w^T x_i = d$ which states that it lies on ℓ

if $d=0 \Rightarrow w^T x_i = 0$

then $\sigma(w^T x_i y_i) = \sigma(0) = 0.5$

signed-dist ($w^T x_i y_i$)

σ (signed-dist)

→ tapering

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

New optimizatn,

$$w^* = \operatorname{argmax}_w \sum_{i=1}^n [\sigma(y_i \cdot w^T x_i)]$$

* $f(x)$ is monotone: f'

i) if $f(x) \uparrow \text{as } x \uparrow$

ii) $a < b; f(a) < f(b)$

if $g(x)$ is nonmonotone function

$$\operatorname{argmax}[g(f(x))] = \operatorname{argmax}[f(x)]$$

$$w^* = \operatorname{argmax}_w \sum_{i=1}^n \frac{1}{1+y_i w^T x_i}$$

$$w^* = \operatorname{argmax}_w \sum_{i=1}^n \frac{1}{1 + \exp(-y_i w^T x_i)}$$

lets take $g(x) = \log x \rightarrow$ monotonous function

$$w^* = \operatorname{argmax}_w \sum_{i=1}^n \log \left(\frac{1}{1 + \exp(-y_i w^T x_i)} \right)$$

$$\text{wkt. } \log\left(\frac{1}{x}\right) = \log(x^{-1}) = -\log x$$

$$\operatorname{argmax}[f(x)] = \operatorname{argmin}[-f(x)]$$

not convex fr
 $\log(x)$ is
convex functn

$$0 = (1) \rho \ell = (\bar{x} + 1) \rho \ell$$

$$\textcircled{1} \text{ becomes, } w^* = \arg \max_w \sum -\log(1 + \exp(-y_i w^T x_i)) \quad (25)$$

$$\rightarrow w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i \cdot w^T \cdot x_i)) \quad \left\{ \begin{array}{l} y_i = +1 \\ y_i = -1 \end{array} \right.$$

geometry
final eqn to optimize which is not prone to outliers

probability based

weight vector

$$w^* = \arg \min_w \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i) \quad \left\{ \begin{array}{l} y_i = 1 \\ y_i = 0 \end{array} \right.$$

$$p_i = \sigma(w^T x_i)$$

$$w^T w = 1$$

\vec{w} is unit vector

$$w = \langle w_1, w_2, w_3, \dots, w_d \rangle$$

$d = \text{no. of feature}$

$$\rightarrow \text{let } z_i = +y_i w^T x_i$$

$$\log(1) = 0$$

$$e^{-x}$$

as $x \rightarrow \infty$
 $f(x) \rightarrow 0$

$$w^* = \arg \min_w \sum \log(1 + \exp(-z)) \geq 0$$

as $\exp(-z) \rightarrow \infty$
 $\log(1 + e^{-z}) \rightarrow 0$

always ≥ 0 because $\log(1) = 0$; $\log(1+x) \geq 0$

minimum value it can take is 0

it occurs when $z_i \rightarrow \infty$ & i

$$z = x_i w^T x_i$$

constant from Dmax

\rightarrow we need to pick w st

$\textcircled{1}$ all pts are correctly classified ie $w > 0$

$\textcircled{2}$ $z_i \rightarrow \infty$

\rightarrow if $z_i \rightarrow \infty$ then $w \rightarrow \infty$ (+ve classes)
 $w \rightarrow -\infty$ (-ve classes)

as $w \rightarrow +\infty / -\infty$ we can say that $w \uparrow$ and cause overfit

$$\textcircled{2} \quad x=0 \quad \textcircled{1} \quad e^x \rightarrow \infty$$

$\textcircled{1m}$ ideally we need loss = 0 for loss = 0 we need $\log(1+x) = \log(1) \Rightarrow$
 \therefore if $e^x \rightarrow \infty, \log(e^x) \rightarrow \infty$

12 regularization add to optimiz^e pblm to prevent $w \rightarrow \infty$

(27)

$$\rightarrow w^* = \underset{w}{\operatorname{argmin}} \sum \log(1 + \exp(-z)) + (\lambda w^T w) \rightarrow w^T w = \sum_{j=1}^d w_j^2$$

minimize both

$$w^* = \underset{w}{\operatorname{argmin}} \left\{ \sum \log(1 + e^{-z}) + \lambda \|w\|_2^2 \right\}$$

tag of war b/w ① and ②

* to min ① $w \rightarrow \infty$, but ② becomes very very large (∞^2)

① → loss term

② → regularization term

if $\lambda = 0$; minimize only ① causes overfit

hyperparam = λ

if $\lambda \uparrow$; min depends on ②, which doesn't depend on y_i .

$$w^* = \underset{w}{\operatorname{argmin}} (\text{logistic loss} + \text{regu})$$

underfit

$$\|w\|_2^2 = \sum w_j^2$$

L1 regularization

$$w^* = \underset{w}{\operatorname{argmin}} \left[\sum \log(1 + \exp(-z)) + \lambda \|w\|_1 \right]$$

$$\|w\|_1 = \sum |w_i|$$

Sparsity :- w/soln to LR is sparse if many w's in $w = \langle w_i \rightarrow 0 \rangle$ are 0's. if $w_i = 0$, then that f_i is less important

$\therefore L_2 \rightarrow w_i$ for unimp f's is small but not 0

$L_1 \rightarrow w_i = 0$; for f_i is unimportant

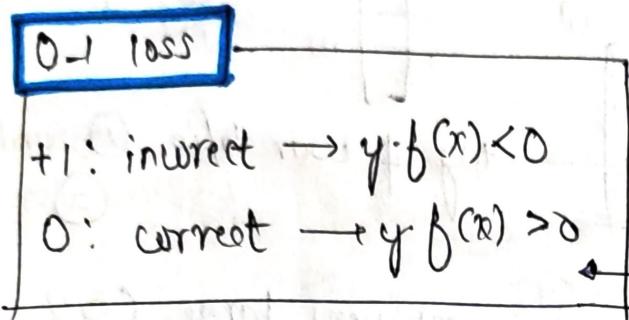
Elastic-Net

$$w^* = \underset{w}{\operatorname{argmin}} \left[\sum \log(1 + \exp(-z)) + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2 \right]$$

\rightarrow as $\lambda \rightarrow \infty$, $e^z \rightarrow 0 \therefore \log(1+0) = \log(1) = 0$ ideal loss nudged

loss-minimization derivation

$w^* = \operatorname{argmin}_w$ (no. of incorrectly classified points)



0-1 loss $f(x)$

+1: incorrect class

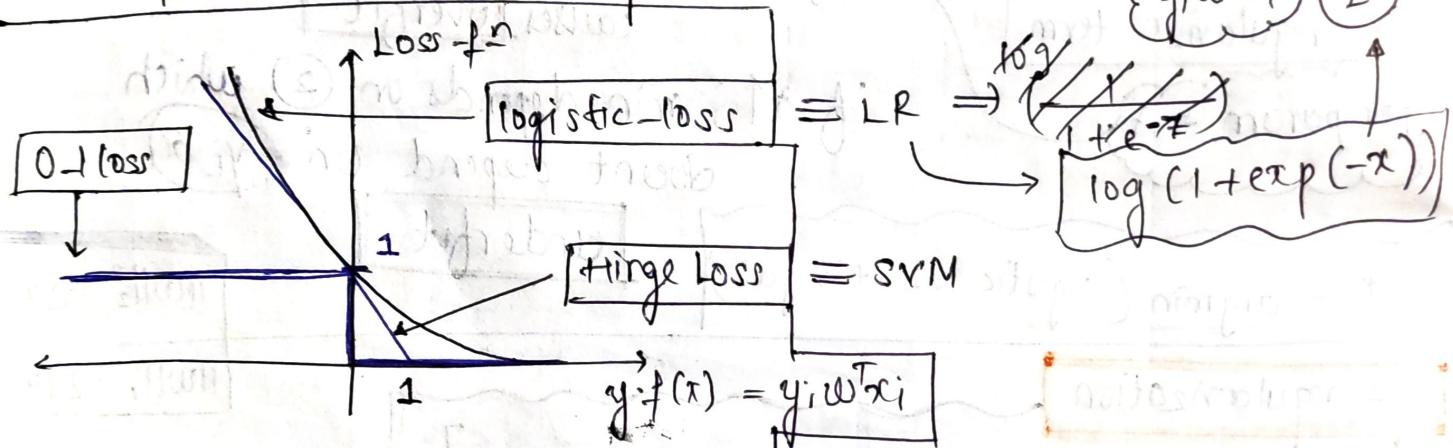
0: correct class

~~y \cdot f(x) > 0~~

$y \cdot f(x) < 0$

* 0-1 loss is not continuous \rightarrow can't differentiate \rightarrow can't optimiz.

\rightarrow optimize 0-1 loss \Rightarrow approximate 0-1 loss

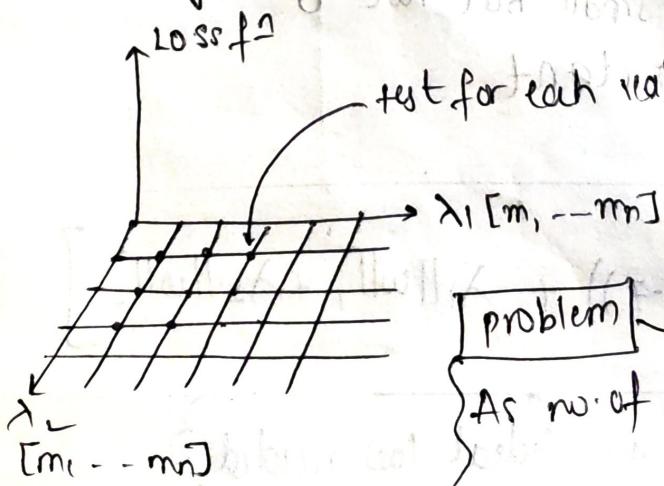


Hyper-param tuning

GridSearch

Ex: ElasticNet

hyper params λ_1, λ_2

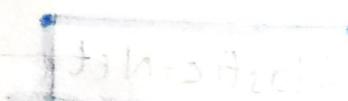


Random Search

* test values randomly picking

$\lambda_1 - [m, --mn]$

$\lambda_2 - [m, --mn]$



As no. of hyper params increase no. of pts

increase exponentially

$$f(x) = w^T x_i$$

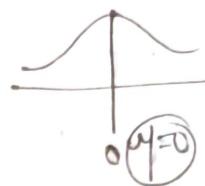
* Logistic regression needs standardization

(29)

	f_1	f_2	f_3	f_4
x_1				
x_2		x_{ij}		
x_3				

$$x_{ij} = \frac{(x_{ij} - \bar{x}_j)}{\sigma_j} \rightarrow \text{centering}$$

↓
scaling [0-1]



Feature Importance: more $|w_j|$ is better affects $(w_j x_i)$

only when features are not collinear

* collinear $\rightarrow f_1 = \alpha f_2 + \beta$ [f_1, f_2 r collinear]

* multicollinear $\rightarrow f_1 = \alpha_1 + \alpha_2 f_2 + \alpha_3 f_3$ [f_1, f_2, f_3 are collinear]

Example

$$\mathcal{D} = \langle x_i, y_i \rangle^n$$

Query point

$$w^* = \langle 1, 2, 3 \rangle; x_q = \langle x_{q1}, x_{q2}, x_{q3} \rangle$$

$$\rightarrow w^T x_q = x_{q1} + 2x_{q2} + 3x_{q3}$$

let assume f_1 and f_2 r collinear

$$f_2 = 1.5 f_1$$

$$\rightarrow w^T x_q = x_{q1} + 3x_{q1} + 3x_{q3}$$

$$= 4x_{q1} + 3x_{q3} = \langle 4, 0, 3 \rangle$$

check collinearity - perturbate

① calculate $w^* = \langle \dots \rangle$

② add small noise $N(0, 0.1)$

$$x_{ij} + \epsilon$$

③ calculate new $\tilde{w} = \langle \dots \rangle$

if $(\tilde{w} \text{ and } w^*)$ differ drastically
then collinear

gives same classifier as $w^* = \langle 1, 2, 3 \rangle$
as collinearity exists

Time Complexity

* Train LR: optimize w^* * run-time complexity

using SGD

$O(nd)$

① space - $O(d)$ [u need to store only w^*
 $= \langle w_1, w_2, \dots, w_d \rangle$]

② time - $O(d)$

d multiplication
+ d addition

you need to find only

$$w^T \cdot x_q$$

$> 0 \rightarrow +$
 $< 0 \rightarrow -$

+ low latency

* memory efficient

(30) what if d is large?

+ if $d = 1000$,

run-time $t_e = 1000$ multiplication + 1000 addition

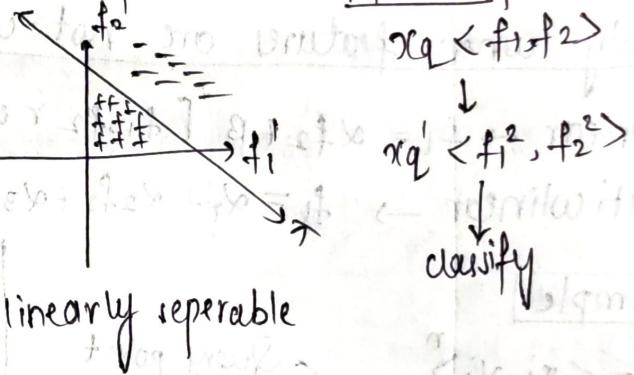
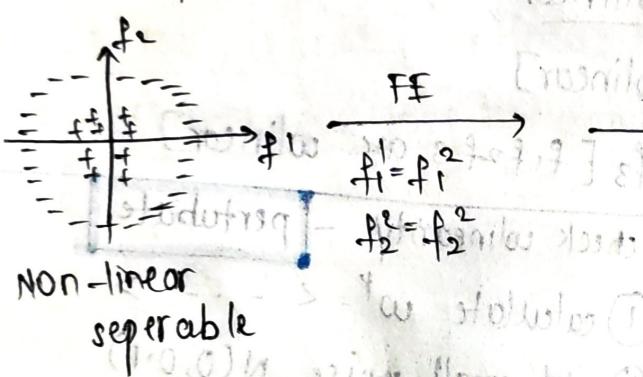
- ① L1 regularization
- ② as $\lambda \uparrow$ sparsity \uparrow .
↳ this can cause overfit \uparrow
latency \downarrow

Non-linearly separable data with LR

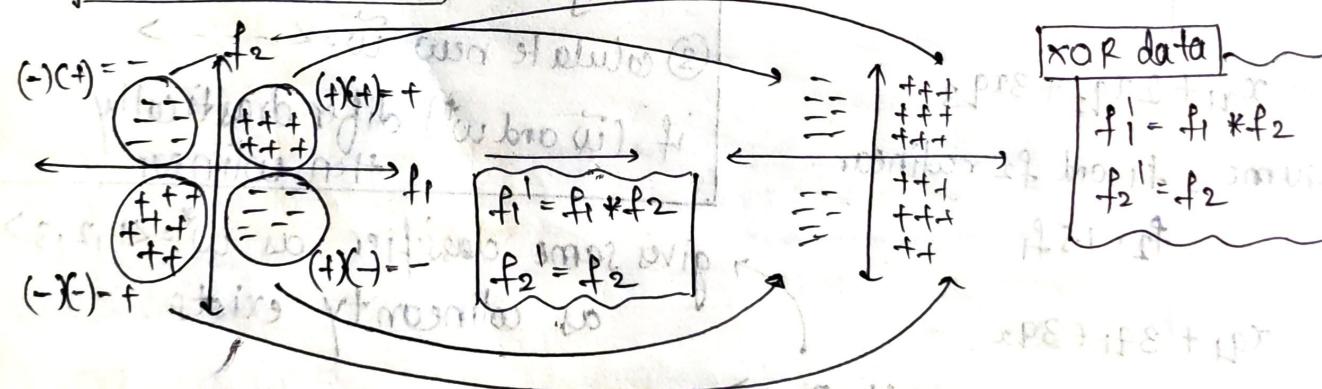
use feature engineering

→ let 2 features $\langle x_{i1}, x_{i2} \rangle$

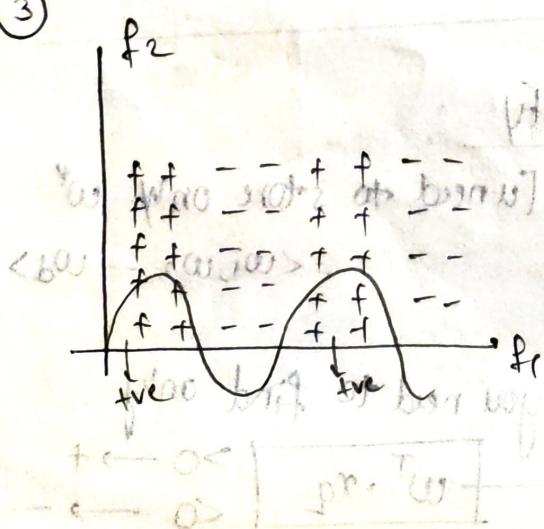
predict



(2) LR for XOR data



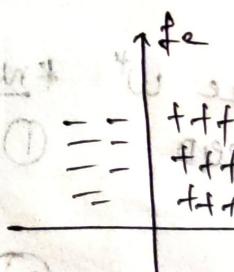
(3)



feature engineering

(b) FE \rightarrow FT

$f_1' = \sin(f_1)$
 $f_2' = f_2$



feature engineering

Typical transforms

NLS - LS

31

- ① $f_1 \cdot f_2, f_1^2, f_2^2, f_2^3 \dots \rightarrow$ polynomial features
 - ② $\sin(f_1), \cos(f_1) \rightarrow$ trigonometric / (trig + poly)
 $\sin(f_1) * \cos(f_2)$
 - ③ OR, AND, XOR \rightarrow Boolean
 - ④ Other $\rightarrow \log(b_i), e(b_i), s(f_i)$

Linear Regression

apply \rightarrow L2, L1 regularization

" ↗ elastic-net

Aim :- find best-fit line/plane (23d)

loss : $y_i - \hat{y}_i \rightarrow$ predicted value
↓
actual value.

$$\text{Loss function} := \sum_{i=1}^n (y_i - \hat{y})^2$$

→ find best fit plane ($w^T x + w_0 = 0$) which best fits (least loss) data

if 2d data $\rightarrow w = \begin{bmatrix} x \\ y \end{bmatrix}$, $w^T x \rightarrow w_1 x + w_0 = 0$
 $ax + b = 0$

$$w^*, w_0^* = \underset{w, w_0}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \underset{w, w_0}{\operatorname{argmin}} \sum (y_i - (w^T x_i + w_0))^2$$

w^*, w_0^*

interpret vector scalar

$f(x) = w^T x_i + w_0$

loss-minimization way

\uparrow Loss fn $\rightarrow (y_i - \hat{y})$

→ hyperbola ($y = x^2$)

feature importance

→ same as logistic regression

→ same regularization

outlier → can't handle directly

- ① D_{Tram}
 - ② remove pts very far away
 - ③ D_{Tram} = outliers
 - loop ④ trim

(32)

Gradient Descent

→ finding maxima/minima is tough → $\frac{df}{dx} = 0, \nabla_x f = 0$
 so we use gd (iterative algorithm)

consider

$$L(w) = \sum_{i=1}^n (y_i - w^T x_i)^2$$

assumptions
 ① No regularization
 ② No intercept: $w^T x_i + w_0 = 0$

$$\therefore \nabla_w L = \sum_{i=1}^n 2(y_i - w^T x_i)(-x_i) \quad \langle x_i, y_i \rangle \text{ are constants}$$

$$\nabla_w L = \sum_{i=1}^n -2(x_i)(y_i - w^T x_i) = 0 \rightarrow \text{find soln}$$

$$D_m = \frac{2}{N} \sum (y_i - \hat{y}) x_i$$

$$D_c = -\frac{2}{N} \sum (y_i - \hat{y})$$

gradient descent

① pick random vector $w_0 = (w_1, w_2, \dots, w_n)$

$$② w_1 = w_0 - r \cdot \nabla_w L$$

update
fn

$$w_j = w_{j-1} - r \cdot \sum_{i=1}^n -2x_i (y_i - w_{j-1}^T x_i)$$

$$w_1 = w_0 - r \cdot \sum_{i=1}^n -2x_i (y_i - w_0^T x_i)$$

$$w_2 = w_1 - r \sum_{i=1}^n -2x_i (y_i - w_1^T x_i)$$

③ if $w_k - w_{k-1}$ is a very small vector

then optimal $w^* = w_{k-1}$

problem

$$\nabla_w L = \sum_{i=1}^n -2x_i (y_i - w^T x_i)$$

as it calculates $i=1 \rightarrow n$
 very difficult to calculate

∴ $w_1 = w_0 - r \nabla_w L$ is computationally expensive

GD for linear regression in 2d plane

33

$$\text{loss} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$y_i = mx_i + c \rightarrow \begin{matrix} \text{w}_{\text{eff}} \\ \text{slope} \end{matrix}$$

$$L = \frac{1}{n} \sum (y_i - (mx_i + c))^2$$

$$\frac{dL}{dm} = \frac{1}{n} \sum 2(y_i - (mx_i + c))(-x_i)$$

$$\frac{dL}{dc} = \frac{1}{n} \sum 2(y_i - (mx_i + c))(-1)$$

$$m_{i+1} = m_i - L \cdot D_m$$

$$c_{i+1} = c_i - L \cdot D_c$$

$$D_m = -\frac{2}{n} \sum (y_i - \hat{y}_i)$$

$$D_c = -\frac{2}{n} \sum (y_i - \hat{y})$$

∴ step ② in GD

$$m_{i+1} = m_i - L \cdot D_m$$

$$c_{i+1} = c_i - L \cdot D_c$$

learning rate

Ridge Regression: $L = \sum (y - \hat{y})^2 + \lambda \|w\|_2^2$

Lasso regression: $L = \sum (y - \hat{y})^2 + \lambda \|w\|_1$

L1

Stochastic Gradient Descent

stochastic = probabilistic

Idea: Instead of all n points, pick any K no. of random points in each iteration.

$$w_j = w_{j-1} - r \sum_{i=1}^K -2x_i (y_i - w_{j-1}^T x_i)$$

K → called batch-size

$$w_{\text{step}}^* = w_{\text{end}}^*$$

Cases

$1 < K < n$

→ if $n = 1\text{ million}$, takes 100 iterations to converge to w^*

if $(K=1)$

if $1 < K < n$

Batch stochastic gradient

BGD

then,
 $K=100$ may take 500 it.
 $K=10$ may take 1000 it.

(34)

PCA constrained optimization

11/11 - 1

PCA optimization prob : $\max_{\mathbf{u}} \sum_{i=1}^n (\mathbf{u}^T \cdot \mathbf{x}_i)^2$ such that $\mathbf{u}^T \cdot \mathbf{u} = 1$

$$\rightarrow \frac{1}{n} \sum_{i=1}^n \mathbf{u}^T \cdot \mathbf{x}_i \Rightarrow \boxed{\sum \mathbf{u}^T S \mathbf{u}}$$

s.t. $\mathbf{u}^T \cdot \mathbf{u} = 1$ and

$$S = \frac{1}{n} \sum \mathbf{x}_i \mathbf{x}_i^T$$

covariance matrix
of $X = X X^T$

$$\therefore \max \sum_{i=1}^n \mathbf{u}^T S \mathbf{u}$$
 wrt $\mathbf{u}^T \cdot \mathbf{u} = 1$

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T S \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1) \rightarrow \text{Lagrange multiplier}$$

$$\frac{\partial L}{\partial \mathbf{u}} = 0 \Rightarrow \frac{\partial}{\partial \mathbf{u}} (\mathbf{u}^T S \mathbf{u} - \lambda (\mathbf{u}^T \mathbf{u} - 1)) = \frac{\partial}{\partial \mathbf{u}} (\mathbf{u}^T S \mathbf{u} - \lambda \mathbf{u}^T \mathbf{u} - \lambda) \\ \text{Hess} + (\mathbf{P} \mathbf{P})^{-1} = 1 \rightarrow \text{minimize} \\ \mathbf{u}^T S \mathbf{u} - \lambda \mathbf{u}^T \mathbf{u} = 0 \quad \boxed{\mathbf{u}^T S \mathbf{u} = \lambda \mathbf{u}^T \mathbf{u}}$$

Eigenvalue / vector

logistic regression with regularization

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \text{ (logistic-loss)} + \lambda (\mathbf{w}^T \mathbf{w})$$

similar

$$\rightarrow \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{min}} \text{ (loss)} \text{ such that } \mathbf{w} \perp \pi, \boxed{\mathbf{w}^T \cdot \mathbf{w} = 1}$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \text{ (log-r-loss)} - \lambda (\mathbf{I} - \mathbf{w}^T \mathbf{w})$$

lagrangian mult

$$\boxed{\mathbf{w}^* = \min \text{ (LR-loss)} - \lambda + \lambda \mathbf{w}^T \mathbf{w}}$$

\therefore Regularization \equiv equality constraint on optimization prob

NLP Basics

Text to vector

BOW
TF-IDF
Word2Vec

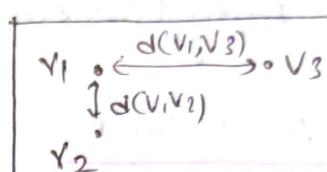
Corpus → set of all documents

38

* text → d-dim vector → fit a plane

consider reviews r_1, r_2, r_3 whose d-dim vector forms are v_1, v_2, v_3

→ $\text{similarity}(r_1, r_2) > \text{sim}(r_1, r_3)$
then, $\text{dist}(v_1, v_2) < \text{dist}(v_1, v_3)$



i.e. more similar ≡ more closer in plane

Bag of words

step ① :- construct a dictionary of size - (d) where 'd' is number of unique words in corpus.

step ② :- for each document - Ex:- $r_1 = \text{this is good this}$

$v_i = [0 \ 1 \ 2 \ \dots \ d]$ # of occurrences of word
 $v_i = [\ \ \ 2 \ \ \ \dots \ 1]$ dict[d] = dict[i]
sparse vector

Calculate distance

$$v_1 - v_2 = \|v_1 - v_2\|$$

this is not good

Ex:- This is good

$v_1 = [1 \ 1 \ 1 \ 0 \ 1 \ 1]$ r_1

This is not good

$v_2 = [1 \ 1 \ 1 \ 1 \ 1 \ 1]$ r_2

(not good)

$$\|v_1 - v_2\|_n = \sqrt{(1-1)^2 + (1-1)^2 + (0-1)^2 + (1-1)^2} = \sqrt{1} = 1$$

score is near but ① word [not] makes difference in real world

Binary BOW

v_i

$v_i = [0 \ - \ - \ - \ 1 \ 1]$

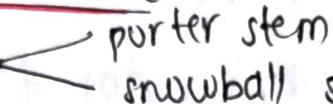
0: if word !present

1: present 1 or more times

36 words → vector flow

① Removing stop words — (this, is, or ...)

② Convert to lowercase

③ stemming 

Ex: tasty, tasteful, taste : taste (stem word)

④ semantic meanings — word2vec

↳ delicious = tasty

⑤ Lemmatization : breaking up sentences into parts

This is New York

uni	grams
bi	BOW

r₁: This is very tasty
r₂: This is not tasty

by default

stop words: not, very

unigram → 1-d for each word

This is very tasty not

--	--	--	--

bi-gram → 1 d for each pair

this	not	very	tasty

n-grams

unigram + BOW : destroys sequential info

bigram + BOW : retains some info
(very tasty)

∴ uni < bi < tri

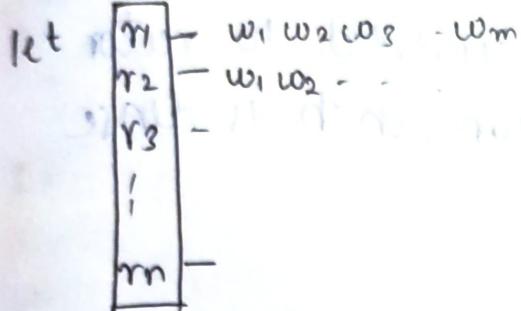
as uni → bi → tri → n

d → increases

∴ increases dimensionality

② TF-IDF Term freq - Inverse doc freq → used in information retrieval

(37)



IDF

$r_i \rightarrow$ records

$D_c = \text{corpus} = \{\text{all records } \{r_i \rightarrow r_n\}\}$

	TF	IDF
frequent word	high	Low
rare word	Low	High

$$TF = TF(W_i, r_j) = \frac{\# \text{ } w_i \text{ occurs in } r_j}{\text{total # words in } r_j}$$

→ it can be interpreted as
prob of word (w_j) present in r_j
 $TF \rightarrow [0 \leftrightarrow 1] \text{ as } dr \geq hr$

$$IDF = IDF(W_i, D_c)$$

$$= \log \left(\frac{N}{n_i} \right) \quad \begin{matrix} \# \text{docs} \\ \# \text{docs} \\ \text{which have } w_i \end{matrix}$$

$$IDF = \log \left(\frac{N}{n_i} \right)$$

$IDF \geq 0$

$\therefore n_i \rightarrow \text{no of records in which } w_i \text{ occurs}$

where

$N \rightarrow \text{no of docs/records}$

$n_i \rightarrow \text{no of docs which contain word } (w_i)$

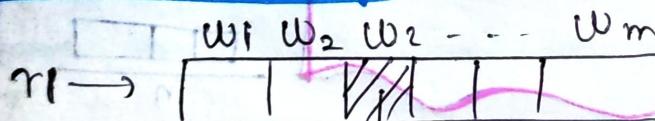
as $n_i \uparrow \downarrow \text{IDF} \downarrow$

$n_i \uparrow \text{IDF} \uparrow$

i.e. rare words: high IDF

freq words: low IDF

TF-IDF format



$r_1 \rightarrow$

$$TF(w_2, r_1) * IDF(w_2, D_c)$$

$\therefore r_i \rightarrow \text{each } w_j \text{ cell} =$

$$TF(w_j, r_i) * IDF(w_j, D_c)$$

drawback → doesn't focus on semantic meaning.

prioritize frequent words in record

prioritize rare words in Corpus

38

Word 2 Vec

deep learning

large corpus text → $w2v$ → each word is converted to vector of d -dimension, which is dense

comes

① similarity

If w_1, w_2 r similar
 v_1, v_2 r similar

nearby

② relationship

$$(v_{\text{queen}} - v_{\text{knight}}) \parallel (v_{\text{man}} - v_{\text{woman}})$$

Avg $w2v$

$w2v$: word → vector

aim: sentence(r_i) → value?

How to convert sentence to vector using $w2v$?

$$\rightarrow r_i: w_1, w_2, w_3, \dots, w_n,$$

$$r_i = \frac{1}{n} w2v(w_1) + w2v(w_2) + \dots + w2v(w_n)$$

$$\rightarrow \text{tf-idf}(r_i) = [t_1 \ t_2 \ t_3 \ t_4 \ t_5 \ \dots]$$

word → vec(d -dim)

tf-weighted $w2v$

→ tf-idf value of (w_i) in r_i

$$\text{tf-idf } w2v(r_i) = \sum_{i: \text{words}} (t_i * w2v(w_i))$$

$$\sum_{i: \text{words}} t_i$$

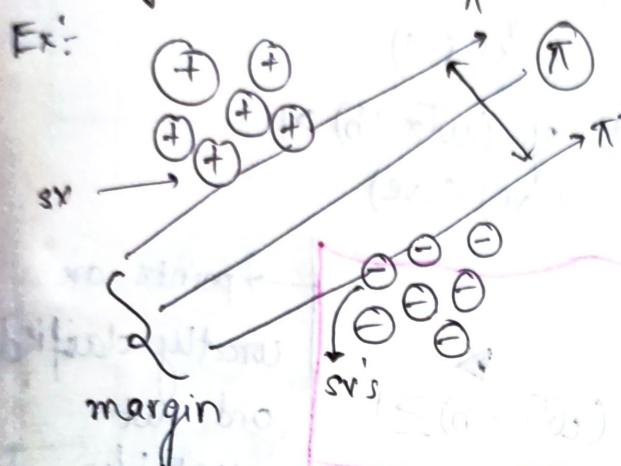
if $(t_i = 1)$

$$\text{tf-idf-}w2v = \text{arg } w2v$$

SVM

(39)

aim: construct optimal hyperplane which maximizes margin



aim: find π such that $d(\pi^+, \pi^-)$ is maximized

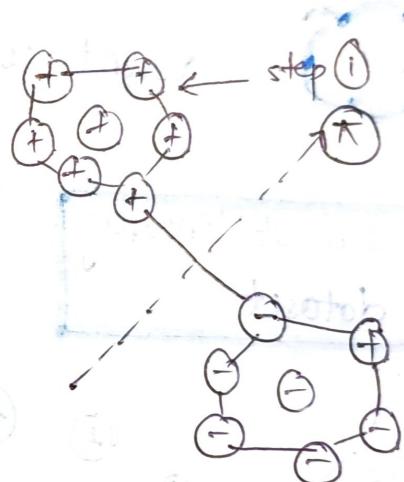
* points on π^+ and π^- are called support vectors (sv)

$$\pi^+ \parallel \pi \parallel \pi^-$$

* while making predictions only π is considered

Alternate Intuition

- ① Construct convex hull around classes of points
- ② Pick 2 points one from each class when joined connect \oplus_c and \ominus_c and shortest distance
- ③ Bisect the line to give π



Mathematical formulation

π : margin maximizing hyperplane

let $\pi: w^T x + b = 0$ \rightarrow geometrically we can get

$$\pi^+ \cdot w^T x + b = 1$$

$$\pi^- \cdot w^T x + b = -1$$

$$d(\pi^+, \pi^-) = \frac{2}{\|w\|}$$

(normal w is \perp to all π, π^+ & π^-)

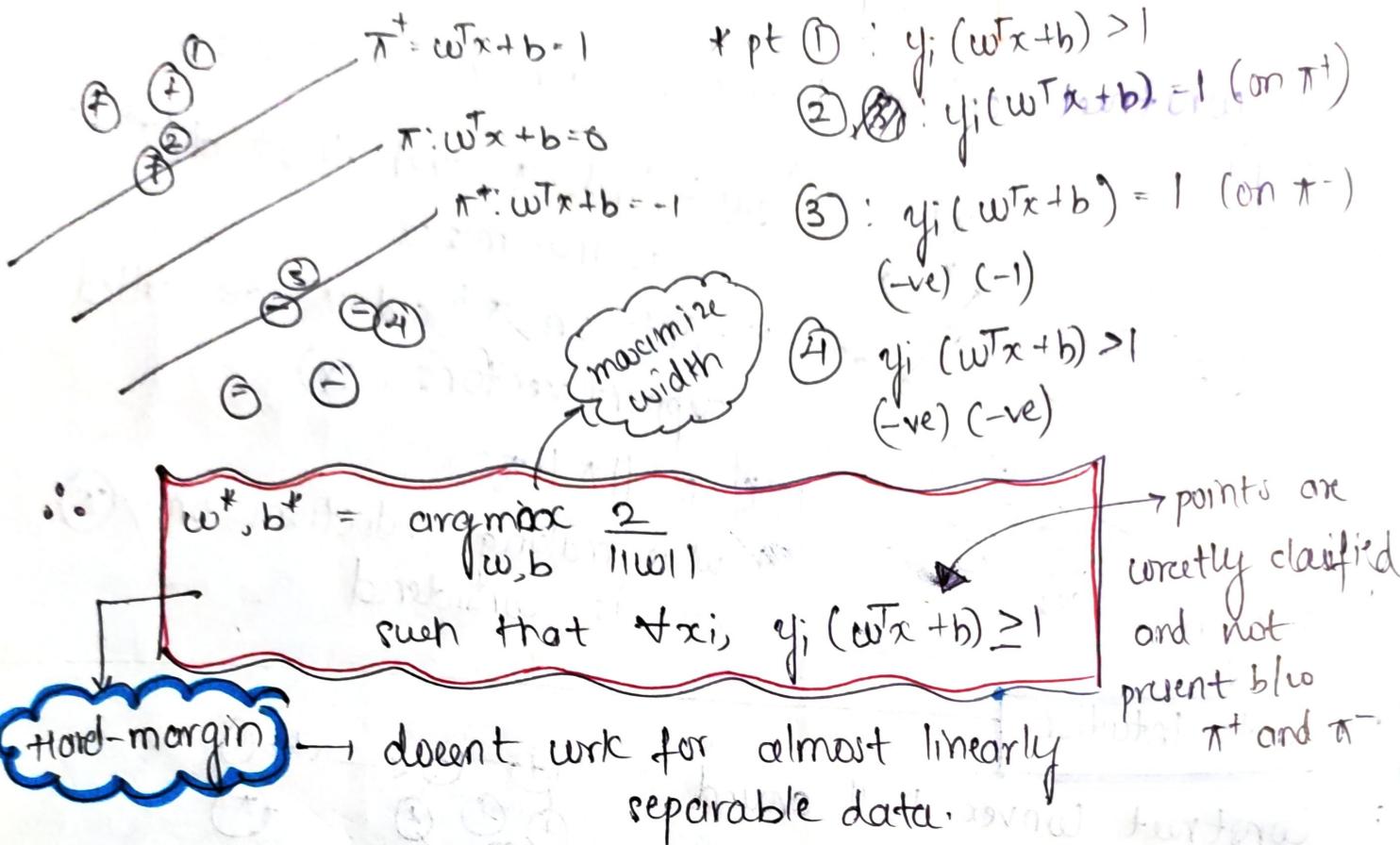
$$\pi^+ \cdot \pi^- = w^T x_1 + b = 1$$

$$- w^T x_2 + b = -1$$

$$w^T(x_1 - x_2) = 2$$

$$x_1 - x_2 = \frac{2}{w^T}$$

(40)



optm for almost linearly
separable dataset

$$\begin{cases} y(w^T x + b) = 1.5 \\ 1 - (0.5) \end{cases}$$

$\beta \rightarrow \text{beta}$

SOFT
MARGIN
SVM

$$\begin{cases} y(w^T x + b) = -1.5 \\ 1 - (2.5) \end{cases}$$

$$\begin{cases} y(w^T x + b) = 0.5 \\ 1 - (0.5) \end{cases}$$

$$\pi^+(w^T x + b = 1)$$

$$\pi^-(w^T x + b = -1)$$

$$\begin{cases} y(w^T x + b) = -0.5 \\ = 1 - (1.5) \end{cases}$$

$$\begin{cases} +1 - 1.5 (>+) \\ y_i(w^T x + b) = -1.5 \\ = 1 - (2.5) \end{cases}$$

from using variable (ξ) to convert values in the form

(ii)

$1 - (\xi)$ we get

$\forall x_i$, if $\xi = 0$; $y_i(w^T x_i + b) \geq 1$ (Correctly classified)

if $\xi > 0$; point is some units distance away from the correct hyperplane in incorrect direction.

$$\therefore w^* b^* = \operatorname{argmax} \frac{2}{\|w\|} \rightarrow \operatorname{argmin} \frac{\|w\|}{2}$$

hyper param.

$$(w^*, b^*) = \operatorname{argmin}_{w, b} \left\{ \frac{\|w\|}{2} + C \left(\frac{1}{n} \sum_{i=1}^n \xi_i \right) \right\}$$

loss fn

less $\sum(\xi_i)$

less misclassif.

such that

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\& \xi \geq 0$$

Note in suki team

$$C = \frac{1}{\lambda} \rightarrow \text{hp of Logistic regression}$$

$C \uparrow$ = overfit

$C \downarrow$ = underfit

why $\pi^+ : +1$ and $\pi^- : -1$

geometric way

$$\pi^+: w^T x + b = +K$$

$$\pi^-: w^T x + b = -K$$

$$w^T x + b = K(x) \leftarrow$$

$$\left(\frac{w^T}{K} \right) x + \left(\frac{b}{K} \right) = 1$$

why same $\frac{1}{K}$

$\rightarrow \pi^+, \pi^-$ are equally

far from π but in

opposite direction

and for ease $K=1$

$$\sum_{i=1}^n w^T x_i + b = 1$$

we don't say
that $\|w\| = 1$

(42)

dual form

primal form

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

st,
 $y_i(w^T x_i + b) \geq 1 - \xi_i$
 $\xi_i \geq 0$

dual
form

Linear SVM $x_i^T x_j$ Kernel SVM : $K(x_i, x_j)$

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

st,
 $\alpha_i \geq 0$
 $\sum \alpha_i y_i = 0$

Note

① x_i is replaced by α_i

② x_i only occur in pairs $\rightarrow x_i^T x_j \Leftrightarrow x_i \cdot x_j = \cos\text{-sim}(x_i, x_j)$

③ primal : $f(x) = (w^T x_q + b)$

predict $f(x)$ given $\|x\| = \|x_q\| =$

$$\text{dual} = f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$$

IMP

Only depends on
support vectors as
 $\alpha_i > 0$

④ $\alpha_i > 0 \rightarrow$ support vectors

$\alpha_i = 0 \rightarrow$ non-support vectors

⑤ $x_i \cdot x_j$ or $x_i^T x_j \equiv \cos\text{-sim}$

$\therefore x_i \cdot x_j$ can be replaced by any similarity

$\rightarrow K(x_i, x_j) = \text{Kernel function}$

dual
form

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \alpha_j y_i y_j) (x_i^T x_j)$$

 $K(x_i, x_j)$

$$f(x) = \sum_{i=1}^n \alpha_i y_i (x_i^T x_q) + b = \sum_{i=1}^n K(x_i, x_q) \alpha_i y_i$$

Kernel SVM used for non-linearly separable dataset

43

$$\rightarrow K(x_i, x_j) = x_i^T x_j \rightarrow K(x_i, x_j) : \text{similarity b/w } (x_i, x_j)$$

Polynomial Kernel

$$K(x_1, x_2) = (x_1^T x_2 + c)^d$$

Quadratic Kernel $c=1, d=2$

$$K(x_1, x_2) = (1 + x_1^T x_2)^2$$

let $x_1 = \langle x_{11}, x_{12} \rangle$ apply Quad-kernel

$$x_2 = \langle x_{21}, x_{22} \rangle$$

$$K(x_1, x_2) = (1 + x_{11}x_{21} + x_{12}x_{22})^2$$

$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11}x_{21} + 2x_{12}x_{22} + 2x_{11}x_{21}x_{12}x_{22}$$

this can be written/rearranged as

$$(1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}) \rightarrow x_1^1 \rightarrow 6d$$

$$*(1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}) \rightarrow x_2^2$$

↓↓

$$(x_1^1)^T (x_2^1) \text{ or } x_1^1 \cdot x_2^1$$

Feature transform

$$\begin{array}{c|c|c} x_1 & x_2 & d \rightarrow d' \\ \downarrow & \downarrow & d' > d \\ x_1^1 & x_2^1 & \end{array}$$

∴ Kernelization ≈ feature transform

$$(f_1, f_2) \xrightarrow[\text{kernel fn}]{\text{kernel trick}} (f_1^1, f_2^1)$$

d -dim

d' -dim

Mercer's theorem

d
(non-linearly
sep)

Kernel trick $\rightarrow (d' - (d' > d))$

(linearly separable)

44

RBF Kernel

Radial Basis Function kernel

$$K(x_1, x_2) = \exp\left(-\frac{d_{12}^2}{2\sigma^2}\right)$$

where,

$$d_{12} = \|x_1 - x_2\| \in \text{dist}(x_1, x_2)$$

σ - hyper parameter

Obs : i) as $d \uparrow$, $K \downarrow$

ie as distance \uparrow , similarity decreases.

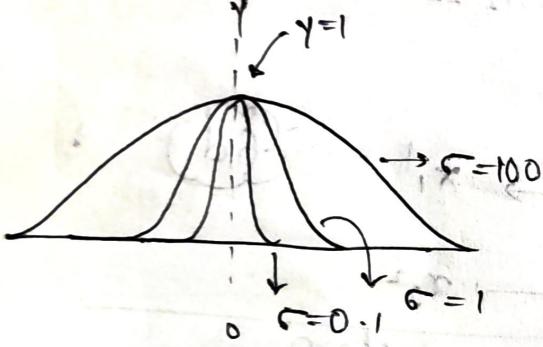
ii) effect of σ

plot $\rightarrow e^{-\frac{x^2}{2\sigma^2}}$

$$\sigma = 0.1 \text{ (wide)} \quad \sigma = 1 \text{ (medium)} \quad \sigma = 100 \text{ (narrow)}$$

$$\sigma = 1 \quad \text{--- (****)}$$

$$\sigma = 100 \quad \text{--- (*****)}$$



\therefore as $\sigma \uparrow$, no of pts considered for similarity increases

σ is similar to (R) in KNN

Time complexity

Train

$$O(n^2)$$

i) SMO

ii) sequential minimal opt

(SMO)

$\rightarrow O(n^2)$: kernel SVM

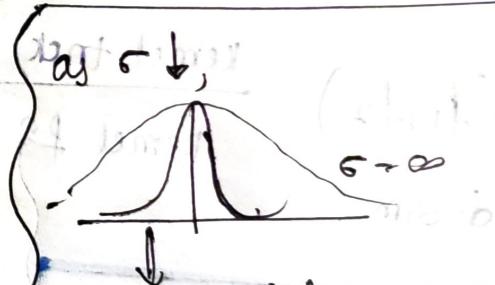
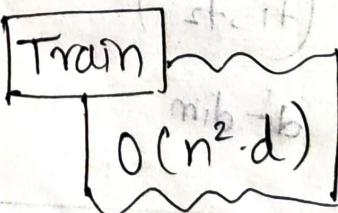
$\rightarrow O(nd^2)$: if $d < n$

\therefore if n is very large,
dont prefer SVM

Runtime $f(x_a) = \sum \alpha_i y_i K(\cdot) + b$

$$O(kd)$$

$\rightarrow k$: no of support vectors



only consider
nearest points

$$\exp\left(-\frac{d^2}{2\sigma^2}\right) \cdot \sigma^2 = \frac{1}{2\sigma^2} \equiv \exp(-\sigma^2 \cdot d^2)$$

nu-SVM

1) default SVM - C-SVM

$$\text{hyperp} = C \ (C \geq 0)$$

→ hp: $\begin{cases} \text{nu} \rightarrow 0 \leq \text{nu} \leq 1 \\ \text{nu} \geq \text{frac}^n \text{ of errors} \\ \text{nu} \leq \text{frac}^n \text{ of SVectors} \end{cases}$

usage

(45)

nu - max error which controls

① want error < 1%

$$\text{nu} = 0.01$$

Ex: if nu = 0.01

1% errors $\leq 1\%$ → upper bound

support vectors $\geq 10\%$ → lower bound

nu-SVM can control errors

but can't control SV's → Prod depends on only SV's (ie $\alpha_i > 0$)

Support Vector Regression (SVR)

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

s.t.,

$$\begin{aligned} y_i - \hat{y}_i &\leq \epsilon \\ \hat{y}_i - y_i &\leq \epsilon \end{aligned}$$

$$\hat{y}_i = w^T x_i + b$$

ϵ is similar to minimum error
ie. $\epsilon = 0.1$

max error is $\rightarrow 0.1$

$$|\hat{y} - y| \leq \epsilon : y - \hat{y} \leq 0.1; \hat{y} - y \leq 0.1$$

case

if $\epsilon = 0$:

overfit chances \rightarrow minimize error

if $\epsilon = 100$:

underfit \rightarrow error to minimize is

overfit ($\epsilon = 0$)

underfit ($\epsilon = 100 \dots$)

best fit

→ data (x_i, y_i)

→ train

find optim al w, b

→ predict

$$\hat{y}_i = w^T x_i + b$$

↳ prediction

46

SVM - key conclusions

Best → u have kernel
 Worst → n is large
 sv are large

* Outliers → little impact as only support vectors matter

Exception :- RBF kernel with small $K \approx 6$

* Bias-variance = $C \uparrow$ - overfit - var↑

$C \downarrow$ - underfit - bias↑

* Large dimensions → good for SVM as it tries to $d \rightarrow d^*$ ($d^* > d$)

[Entropy] → degree of randomness

Decision Tree

$$\log_2 = \lg$$

$$\text{Entropy } H(Y) = -\sum_{i=1}^K P(Y_i) \log_2(P(Y_i))$$

K: no of classes

binary → yes/no [K=2]

$$D(14) \begin{cases} +9 \\ -5 \end{cases} \quad H(Y) = -\frac{9}{14} \lg\left(\frac{9}{14}\right) - \frac{5}{14} \lg\left(\frac{5}{14}\right)$$

consider 2 state var-y

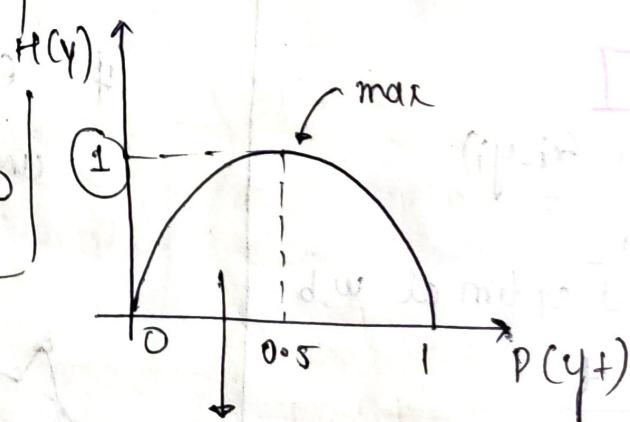
$$\begin{array}{l} ① \\ D \begin{cases} 5 (+) \\ 10 (-) \end{cases} \quad H(Y) = 1 \quad (\text{maximum}) \end{array}$$

$$\begin{array}{l} ② \\ D \begin{cases} 0 (+) \\ 10 (-) \end{cases} \quad H(Y) = 0 \end{array}$$

Results

i) $H(Y) = 1$; if equiprobable

ii) if one class dominates fully $\rightarrow H(Y) = 0$

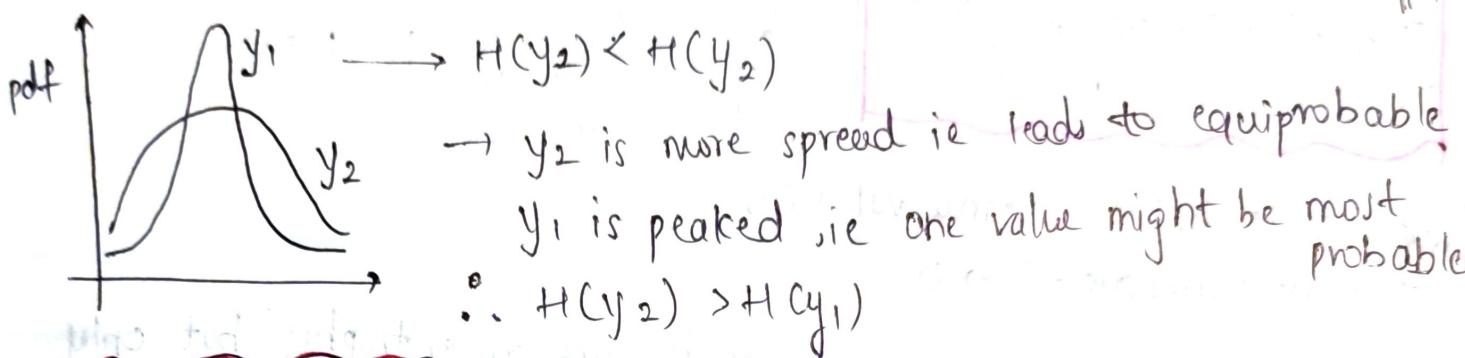
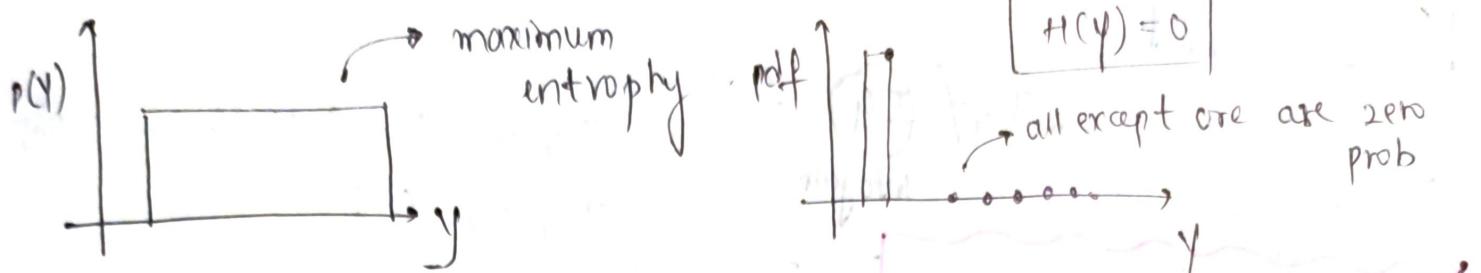


symmetric because $P(Y+) = 1 - P(Y-)$

TC $\propto (n \log n \cdot d)$ training

\propto (depth) predict

PDF based Interpretation



peak \uparrow entropy \downarrow

$$IG = \text{Ent}(\text{parent}) - \text{weighted entropy}(\text{child})$$

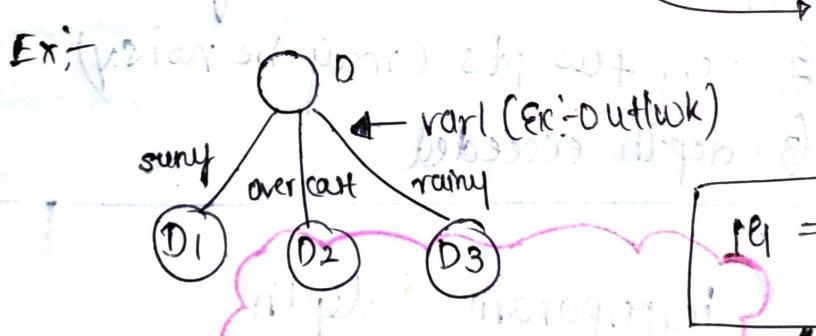
info-gain

categories

$$IG(Y, \text{var}) = \sum_{i=1}^K \frac{|D_i|}{|D|} \times H_{D_i}(y) + H_D(y)$$

Entropy before splitting

weighted entropy after splitting



$$IG = -\text{weighted entropy} + \text{entropy (previous)}$$

$$IG(D, \text{var1}) = \left[\frac{|D_1|}{|D|} H_{D_1}(y) + \dots \right] + H_D(y)$$

* $IG \uparrow$ better
 * weighted ent $\downarrow \rightarrow IG \uparrow$

Gini Impurity

$$y < \frac{y^+ + y^-}{2}$$

$y \rightarrow y_1, y_2, \dots, y_k$
(k -categories)

$$I_q(y) = \sum_{i=1}^k p(y_i)^2$$

case 1

$$\text{if } y < 5$$

$$I_q(y) = 1 - (0.25 + 0.25) = 0.5$$

$$I_q(y) = 0.5$$

$$I_q(y) = 1 - \sum_{i=1}^k (p(y_i))^2$$

gini impurity
→ equi-probable → max-val = 0.5
→ one dominate → 0

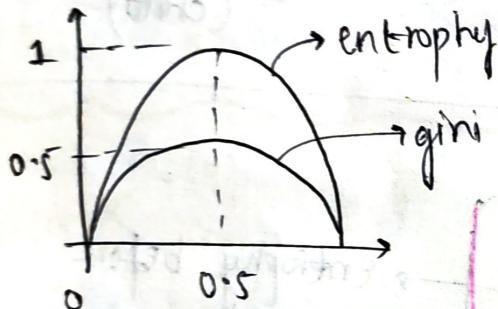
case 2

$$y < 0 \quad I_q(y) = 1 - 1 = 0$$

$$I_q = 0$$

Entropy vs Gini

gini is same as entropy but only less scaled ie max value is 0.5



why gini > entropy?

→ no log-term to calculate

Initial split

$$H(y+, y-)$$

when to stop

- ① calculate $H(y)$
- ② find weighted-entropy for each variable $I_q(y, \text{outlook/sunny/l-})$
- ③ for first split

hyperparam = depth

$d \uparrow \rightarrow \text{overfit}$

$d \downarrow \rightarrow \text{underfit}$

for ③ u can use entropy / gini impurity

Numerical features

$$f_1 - (0 \leftarrow 10000)$$

then

$$\begin{cases} f_1 < \tau_1 \\ f_1 < \tau_2 \\ \vdots \\ f_1 < \tau_n \end{cases} \quad \begin{cases} \text{pick best information} \\ \text{gain yielding} \\ \text{threshold} \end{cases}$$

$$z_i \rightarrow \text{thresholds}$$

steps

$$f_1 - (f_{11}, f_{12}, f_{13}, \dots, f_m)$$

(1) sort f_1

(2) pick thresholds

$(\tau_i = f_{1i})$ and split

(3) τ_i which maximizes

Categorical ft with large categories

Ex: pinode - may have 1000's of pinodes

→ (solution) convert to numerical feature

consider - class = $f_1 - (0, 1)$

→ pinode = $\{P_1, P_2, P_3, \dots, P_n\}$

$$tp = P_i = P(y=1 | P_i)$$

probability that ($y=1$) given (pinode = P_i)

Mean Encoding

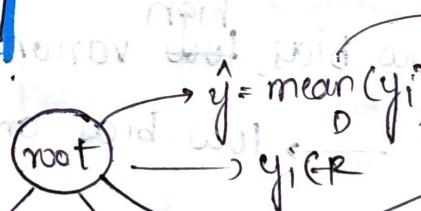
$$Ex: P_2 = P(y=1 | P_2)$$

$$= \frac{\#(y_i=1 \& P_j=P_2)}{\# P_2}$$

Decision Tree regression

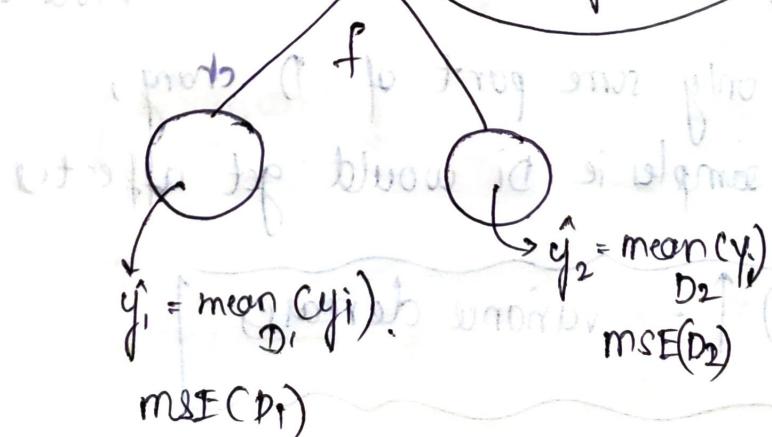
\leftarrow MSE \Leftrightarrow Entropy

Joint is wide



$$\frac{1}{m} \sum (y_i - y)^2$$

MSE = Entropy



(5)

before split

→ after split

$$\therefore \text{IG} = \text{MSE}(D) - \text{weighted MSE}(D_1, D_2)$$

variance $\uparrow \equiv$ overfit

classif	reduce entropy - ideal = 0
regressn	reduce · MSE / MAD - ideal = 0

Ensemble Techniques

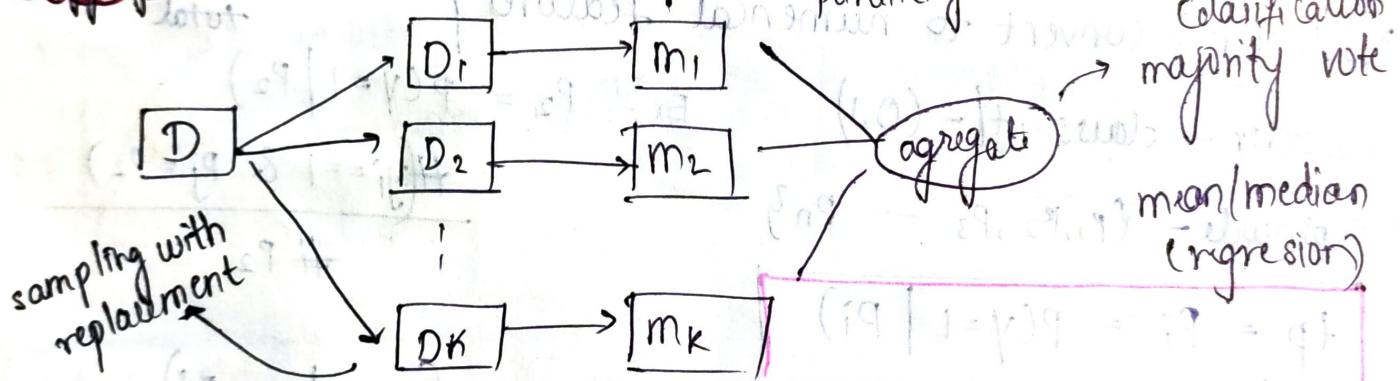
Random Forest

- based on Bagging
(bootstrap aggregation)

Bagging

reduce variance

Bagging
Boosting
Stacking
Cascading



variance: how much model changes when data changes
→ lower the better

* if each M_i have low bias, ^{high} variance

then after bagging → low bias and low variance

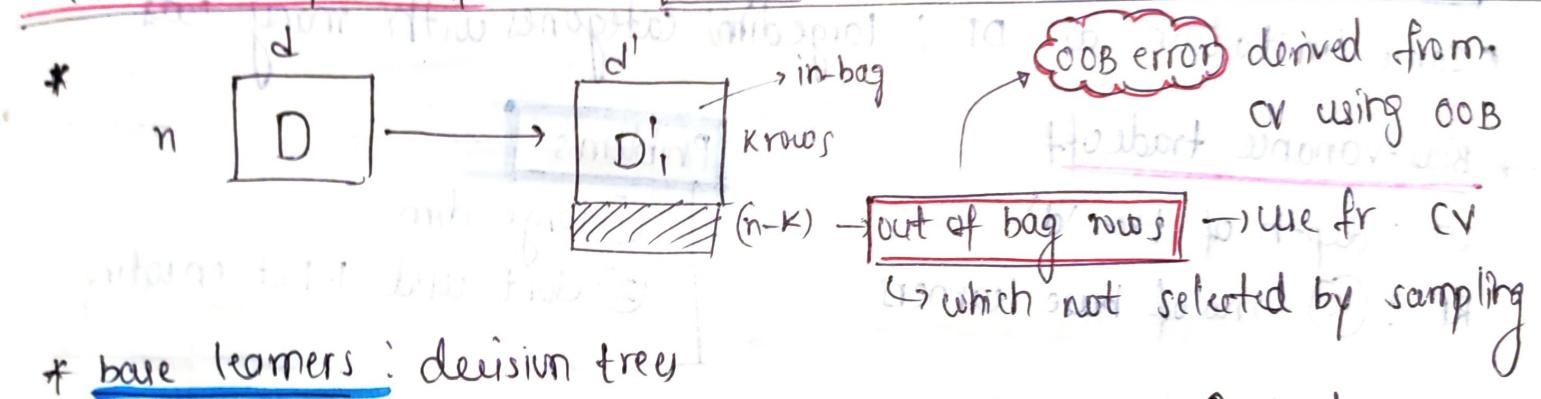
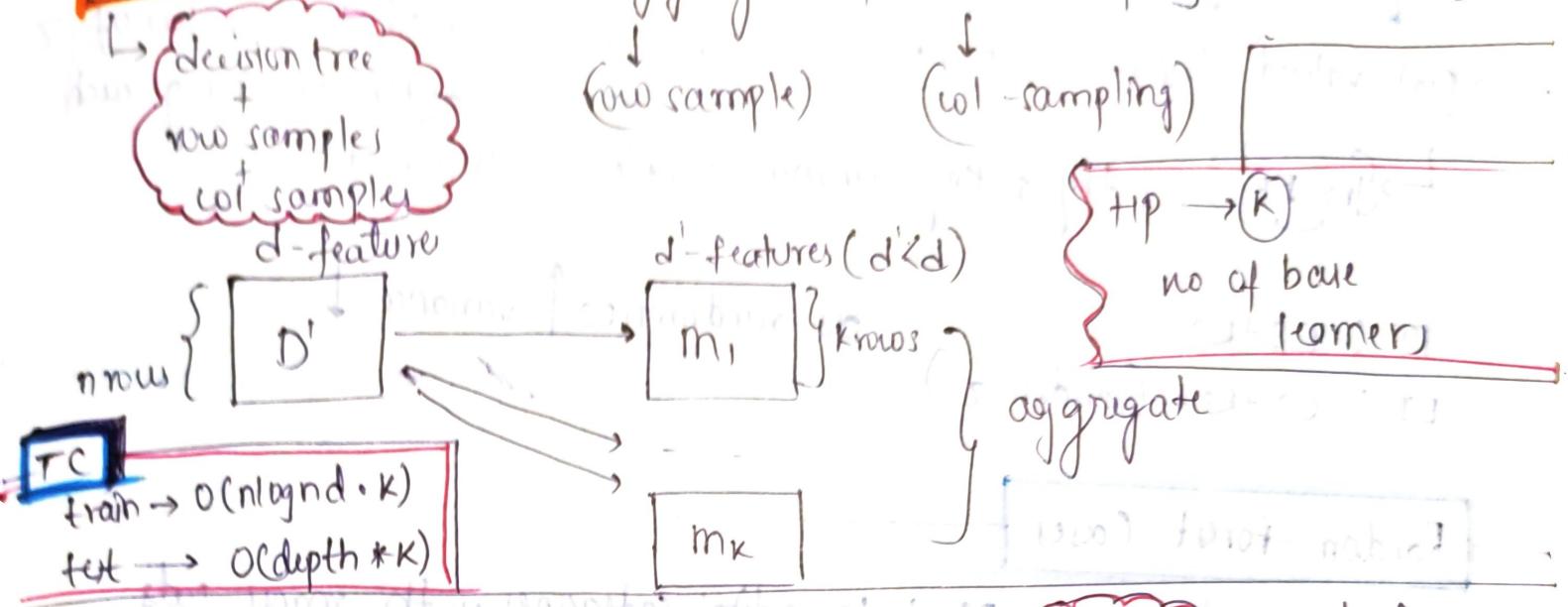
→ low or high

reason: if D changes only some part of D changes

∴ only few samples i.e. D_i would get affected.

as K (no of samples) ↑ → variance decreases

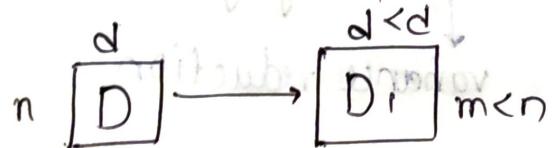
Random Forest: DT + Bagging + feature sampling (51)



* base learners: decision tree

→ DT have large depth bcoz if depth ↑ variance ↑ and we can reduce variance using bagging

Hyperparam CSR RSR



↳ new/col sampling rate

$$\text{CSR} = \frac{d'}{d}$$

$$\text{RSR} = \frac{m}{n}$$

↳ 50% CSR

↳ pick 50% of cols in D' from D

→ as CSR / RSR ≈ variance ↓

General way

↳ fix CSR and RSR → try best value of K → choose best K using CV

⑤2 Extremely Random Tree

→ Real valued features: RT / DT → sort values → pick best value of? after checking each
↳ (in ERT) → try picking random sample of thresholds

RT: CS + RS

ET: CS + RS + (random z)

as randomness ↑ variance ↓

Random-forest Cases

* same limitation as DT: large dim categories with many cats

→ Bias-variance tradeoff

DT: depth of tree 'd'

RF: K → no of base learners

Problems

① large dim

② dont use i-hot encoding

Feature Importance

DT: more imp to feature which reduces Entropy at level

RF: ↪ which reduces entropy in all K model.

Bagging

= high variance, low bias

randomiz + aggregation

base learners

↓ variance reduction

Generally →

$$\text{Error} = \text{bias}^2 + \text{varf } \epsilon$$

→ too much

↓ irreducible error

↓ random

Boosting

Gradient Boosting (GBDT)

AdaBoost (adaptive boost)

(53)

Boosting = low-var, high-bias + additive combine
bare model reduce bias

Core Algo

low bias = DT with d↑
high var

low var → shallow DT
high bias (d↓)

Bare model

→ Step 0
* $D_{\text{train}} (x_i, y_i)_i^n \rightarrow M_0 \quad y = h_0(x) \rightarrow$ will have high error as D_{train} is low var high bias
error_i = $y_i - h_0(x_i)$ $y_i - \hat{y}_i = AE$

$\{x_i, y_i, \text{error}_i\}_i^n$

→ Step 1
 $\{x_i, \text{error}_i\}_i^n \rightarrow M_1 \quad y = h_1(x) \rightarrow$ it tries to predict $\text{error}_i = y_i - h_0(x_i)$
 $\Rightarrow F_1(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x)$ model at end of stage 1

Step 2

$\{x_i, \text{error}_i\}_i^n \rightarrow M_2 \quad y = h_2(x) \Rightarrow F_2 = \alpha_0 h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x)$
 $\downarrow, y_i - F_1(x_i)$

Step K

$\Rightarrow F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$ trained to fit residual error @ prev step
additive weighted model $h_i(x)$ is trained on $\{x_i, \text{error}_i\}$
err at end of stage ($i \rightarrow$)

as residual error $\downarrow \Rightarrow$ train error $\downarrow \Rightarrow$ bias \downarrow

$\text{err}_i = y_i - F_{i-1}(x_i)$

5A Residual

* residual at stage K = $\text{err}_i = y_i - F_K(x_i)$ used in building model F_{K+1}

loss minimization

let $L(y_i, F_K(x_i)) = (y_i - F_K(x_i))^2 \rightarrow$ squared loss

let $F_K(x_i) = z_i$

$$\frac{dL}{dz_i} = 2(y_i - z_i) \leftarrow -2(y_i - z_i)$$

$$\Rightarrow -\frac{\partial L}{\partial z_i} = 2(y_i - z_i)$$

$$-\frac{\partial L}{\partial z_i} \approx (y_i - F_K(x_i))$$

negative gradient

called pseudo-residual

idea

* (x_i, err_i) is used for training

$\text{err}_i \rightarrow$ residual;

$\text{err}_i \rightarrow$ pseudo-residual

$$-\frac{\partial L}{\partial F_{i-1}(x)}$$

stage i

(x_i, err_i)

pseudo residual

not residual
 $y_i - F_K(x_i)$

can have any loss function, only constraint it should be differentiable

L

∴ Gradient Boosting uses pseudo-residuals

Hence, you can use any loss function (L)

i.e. minimize any diff'ble loss function

residuals keep on getting smaller

$$\arg \min \alpha_1 x_1 + \alpha_2 x_2$$

Gradient Boosting

Input: $\{(x_i, y_i)\}$, $L(y, f(m))$: differentiable, $M \rightarrow$ no of base learners

Algo:

① Initialize model

$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum L(y_i, \gamma)$$

find γ
which minimizes

Here if L is squared-loss

$$\min_{\gamma} (y_i - F_k(x_i))^2 \rightarrow \gamma^2 = \text{mean}$$

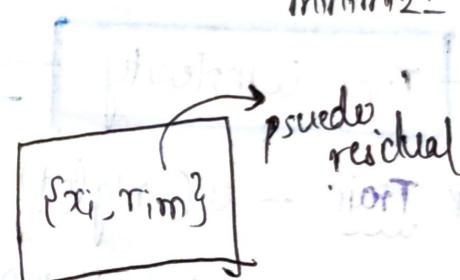
② for $m=1$ to M

① Compute pseudo-residuals

$$\gamma_m = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]$$

$$F(x) = F_{m-1}(x)$$

② fit learner $h_m(x)$ with training data



③ Compute constant γ_m

$$\gamma_m = \underset{\gamma}{\operatorname{argmin}} \sum [L(y_i, F_{m-1}(x_i)) + \gamma h_m(x_i)]$$

④ update

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

easy why?

③ output

$$F_m(x) = F_0(x) + \sum_{i=1}^m \gamma_i h_i(x)$$

only γ is unknown

$$F_M(x) = h_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

$$\sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma_m h_m(x_i))$$

average

(56) Regularization and shrinkage

* Bias-var tradeoff $\approx k/M$

$M \uparrow$ - overfit

$M \downarrow$ - underfit

shrinkage → modifying update rule

$$F_m(x) = F_{m-1}(x) + \gamma_m \cdot h_m(x)$$

$\gamma \rightarrow$ learning rate / shrinkage coefficient

Time Complexity

→ Not easy to parallelize

* Train - $O(n \lg(nd) * M)$

[space] - $O(\text{store each tree} + \gamma m)$

* run - $O(\text{depth} * M)$

→ small \Rightarrow we use shallow trees (good for low latency applications)

XGBoost - best implementation of eBDT

→ eBDT + row sampling + column sampling

added params

* colsample_bytree

* colsample_level

* reg_alpha → L1 regularization

$$\alpha L_1 + \lambda L_2$$

* reg_lambda → L2 regularization

sklearn.ensemble.GradientBoostingClassifier()

hyper-params → MeV

→ cross validation

→ grid search

→ random grid search

minimizing error

(SVD)

more $\gamma \rightarrow O^{(LVB)}$

more importance given to $h_m(x)$

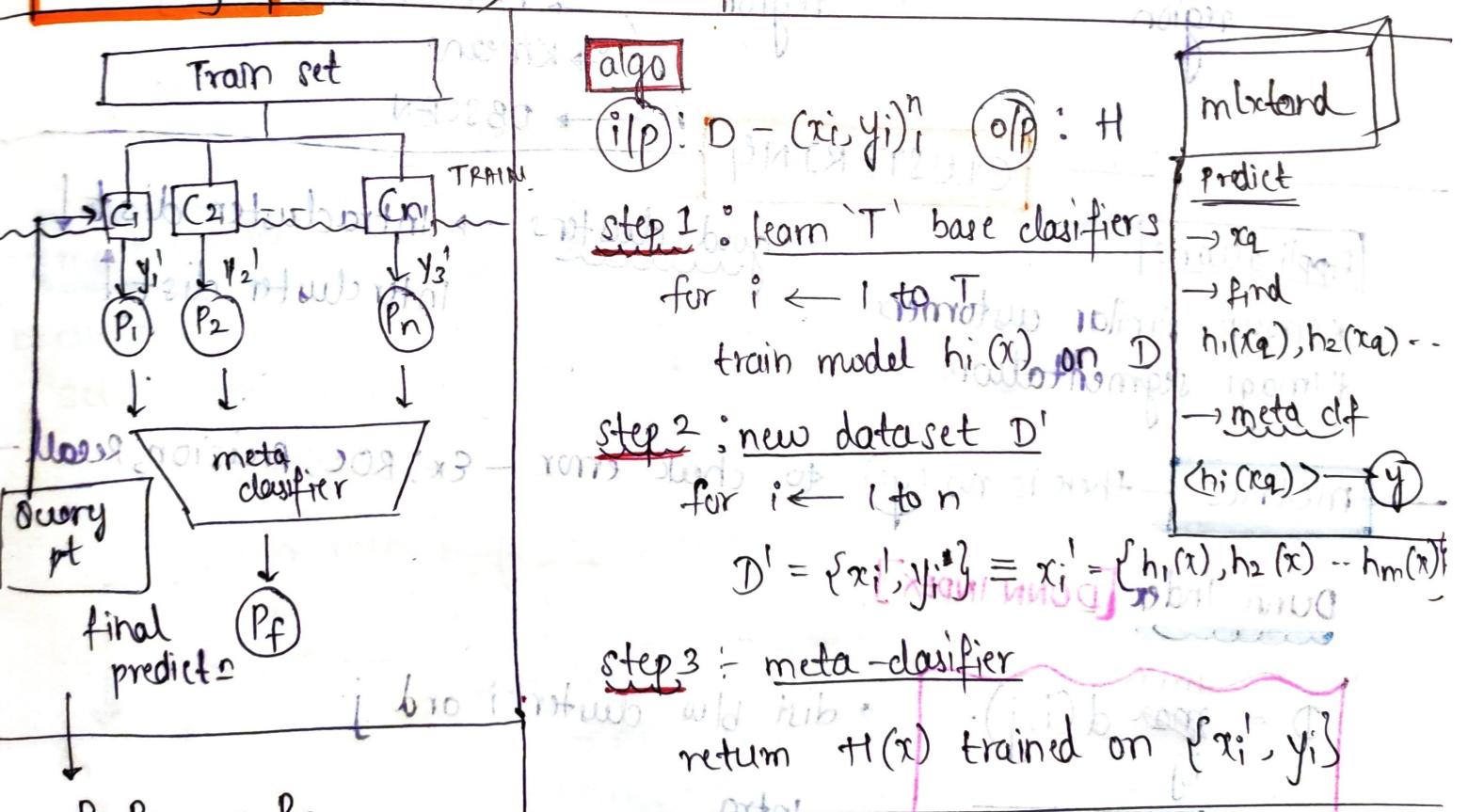
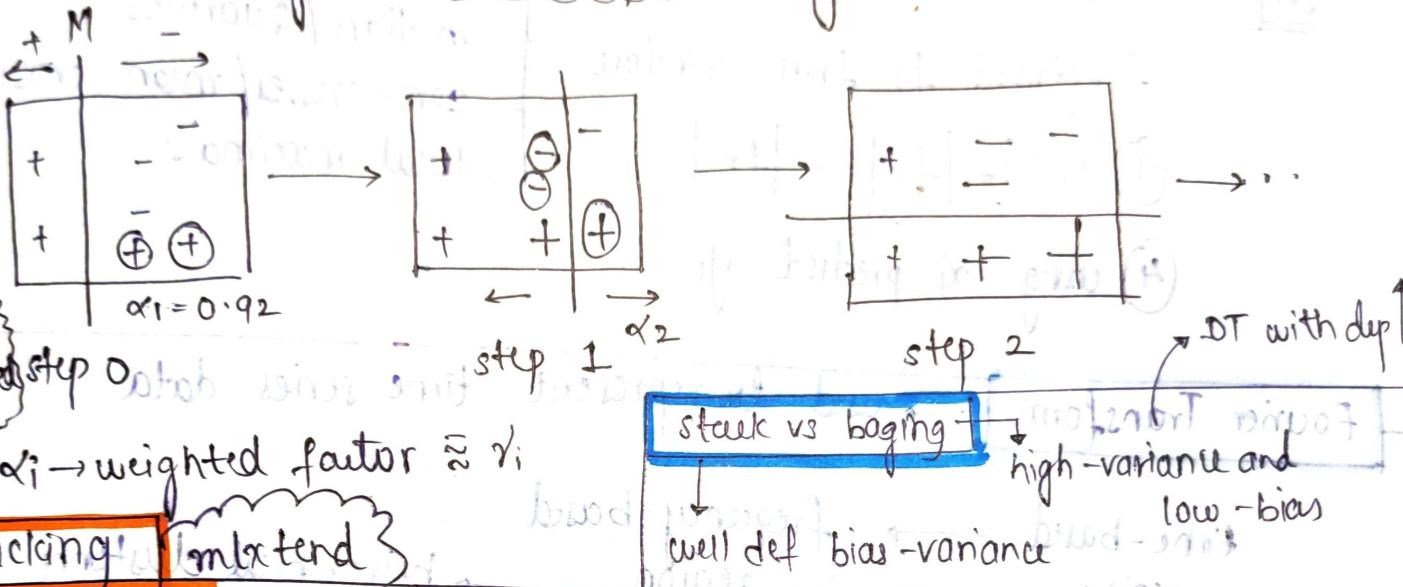
(SVD)

AdaBoost used for face detection

57

main idea: at each step keep on changing weights assigned to pts to incorporate errors.

* weights increase exponentially



Case

- ① classifiers must be as different as possible
- ② useful in Kaggle competitions

58

Featurization

→ page 18 book 2

Moving window

- steps:
- ① fix window width
 - ② extract f_i from window
 - ③ $x_i : [f_1 | f_2 | \dots | f_n]$
 - ④ using x_i predict y_i

max-min

avg

mean/sd

median/quartiles

zero-crosses/mean cross

local maxima's

Fourier Transforms

→ used to represent time series data

time-based
regionfrequency-based
region

hierarchical clustering

K-means

DBSCAN

CLUSTERING

Applications

+ group similar customers

+ image segmentation

good clusters

→ Intracluster dist ↓

intercluster dist ↑

metrics

there is no y_i to check error - ex: ROC, precision, recall

Dunn Index [DUNN INDEX]

$$D = \frac{\min_{i,j} d(i,j)}{\max_{i} d(i)}$$

dist between cluster i and j

$$\max_{i} d(i) \rightarrow \text{intercluster distance}$$

no of clusters - k

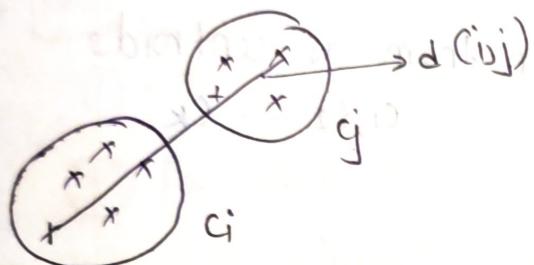
 $D \uparrow$ → better

→ more inter & min intra

calculating d and d'

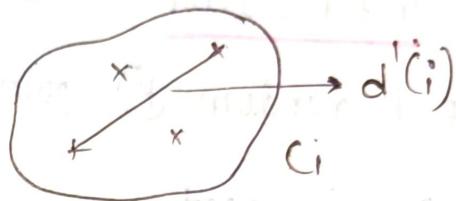
$d(i,j) \rightarrow$ inter

↳ distance b/w 2 farthest pts in c_i and c_j



$d'(i) \rightarrow$ intra

↳ farthest dist b/w 2 pts in cluster i



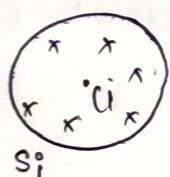
Centroid based clustering

K-means

hyper param - K : no of clusters

Centroid : geometrically a mean

Ex:-



$$c_i = \frac{1}{n} \sum_{j=1}^n x_j$$

$n \rightarrow$ size of S_i

mathematical formulation

k -centroids — $c_1, c_2 \dots c_K$

sets — $S_1, S_2 \dots S_K$

$$\begin{aligned} & \forall x_i \in S_j \\ & \forall i, j, S_i \neq S_j \end{aligned}$$

NP-hard pb/lm

$$\arg \min_{c_1, c_2, \dots, c_n} \sum_{i=1}^K \left(\sum_{x \in S_i} \|x_i - c_i\|^2 \right)$$

(dist b/w points in set and its centroid)

→ Intracluster dist

$a_s \uparrow c \uparrow$

s.t. that,

$x \in S_i$

$S_i \neq S_j$

approximation

(Lloyd's algorithm)

60

Lloyd's algorithm

$D = \{x_i\}_{i=1}^n$, K -clusters

Step 1 : Initialization

→ pick random (K) points and assign them as centroids
 c_1, c_2, \dots, c_K

Step 2 : Assignment

→ for each pt x_i in D :

→ select nearest c_j :

$$\min(d(x_i, c_j)) \quad \forall j = 1 \text{ to } K$$

→ add x_i to nearest c_j set S_j

Step 3 : - recompute centroids

→ for each centroid $c_j \quad \forall j = 1 \rightarrow K$

$$c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

Step 4 : repeat ② and ③ until convergence

new updated centroids and older ~~older~~ one
 dont change much

$$\text{ie } |c_{\text{old}} - c_{\text{new}}| \rightarrow 0$$

Initializing centroids → how do you assign K pts as centroid 61

- ① try randomly
- ② K-means++

→ Kmeans++ Init $[C_1, C_2 \dots C_K]$ (probabilistic)

① pick C_1 randomly

② for each $x_i \in D$,

calculate $d^2(x_i, \text{Nearest Centroid})$

x_1	d_1
x_2	d_2
x_n	d_n

first case → only 1 centroid $\rightarrow d(x_i, C_1)$

③ → pick a point from $D - \{C \text{ already assigned}\}$

as centroid with probability $\approx d_i$

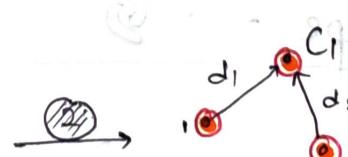
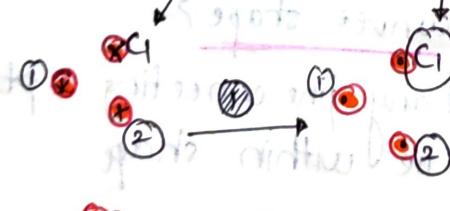
proportional

proportion

$p(x_i \text{ chosen as centroid}) \approx d_i$

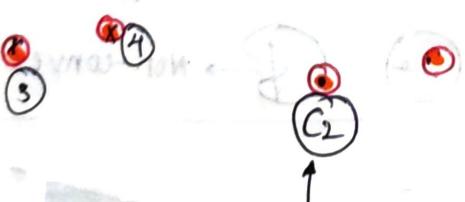
$K=3$

Init



x_1	d_1
x_2	d_2
x_4	d_4

as $d_4 \uparrow$
 $p(x_4) \uparrow$



why $\text{prob}(d)$ and not directly $\text{max}(d_i)$?

→ outliers will always become centroids

biggest outlier
smallest outlier

(62)

Best K

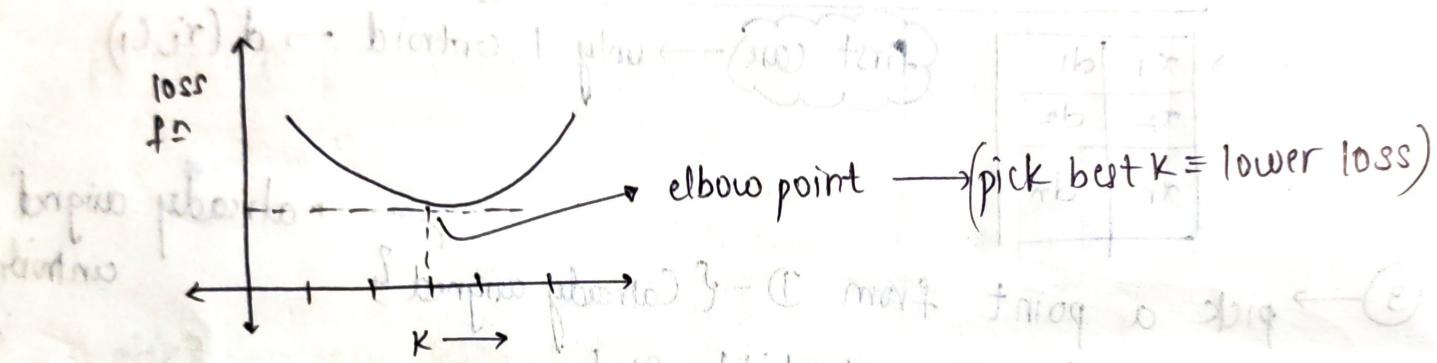
① Domain knowledge: you know exactly how many class

Ex:- $K=2 \rightarrow$ disease pred yes/no

\rightarrow spam - yes/no

② Knee / Elbow method

\rightarrow loss fn : $\sum_{i=1}^K \left(\sum_{x \in S_i} \|x - c_i\|^2 \right)$ [Train model on different values of K]



Limitations

① Different sized clusters :- K-means ideally creates same sized clusters

② Differing density :- it tends to stretch them out

③ Non-globular/conver shape \rightarrow

④ Outliers $(n-K)$ to swap
 $(n-K)$ to calculate loss

NY Taxi

convex shape?

* any line connecting 2 pt lie within shape



Non-convex

Solutions

* Increase no of clusters K
 \rightarrow then recombine (tough)

Kmeans \rightarrow Lloyd

K-medoids \rightarrow PAM

calculate distance to each medoid
swap by each non-medoid
 $O(K*(n-K)^2) \neq O(n^2)$

K-mediods | PAM (partitioning around mediod)

(33)

→ problem with centroids | Each centroid is actually a datapoint

* Not interpretable: $c_i \notin D$, ie centroid may not belong to D , u can get info about it

PAM

Initialization: - K-means + f (probabilistic)

$O(n^2)$

Assignment: - same as K-means

$\forall x_i \in g_j$ if for (x_i) the closest mediod is (c_j)

*** update: a) swap each mediod with non-mediod $\rightarrow (n-k)$ swaps.

b) if loss decreases keep swap else undo swap.

$$\sum_{i=1}^k \sum_{x_i \in S_i} \|x_i - m_j\|^2 \rightarrow \text{dist}(x_i, m_j)$$

can be kernelized

[Kernel matrix]

repeat until converge

use

* At last each mediod is a datapoint

$\Rightarrow \forall c_i \in D$ so easily interpretable.

$$* \sum \sum d(x_i, m_j)$$

less prone

to outliers

can be replaced by kernel trick.

Time Complexity

{LINEAR TC} \rightarrow LLVODS

space

$$O(nkdi)$$

$n = \# \text{pts}$

$K = \# \text{of clusters}$

$d = \# \text{of dimensions}$

$i = \# \text{of iterations}$

$$O(nd + kd)$$

store

store K

centroids

store n pts

(64) **Hierarchical clustering**

- Agglomerative → bottom up
- ~~Divisive~~ Divisive → top down

Agglomerative algorithm → No need to give K as hyper-param

- assume all pts are clusters → keep on reduce no of clusters
- * you need distance b/w pts as well as clusters

Dendrogram: a tree which depicts sequence of merges/splits

Algorithm

dist matrix

- ① Compute proximity matrix
- ② Initially all pts are clusters
- ③ Repeat
 - a) merge the 2-closest clusters
 - b) update proximity matrix

until only one cluster

	P1	P2	P3	...
P1	0			
P2		0		
P3			0	

	P1, P3	P2	P4	
P1, P3	0			
P2		0		
P4			0	

needs update

main issue → how to update prox-mat?

how to define similarity b/w clusters

Def intercluster distance

pick lowest value of each $\min(\text{mat}) \dots$

→ **MIN** $\text{sim}(C_i, C_j) = \min_{P_i \in C_i} \min_{P_j \in C_j} \text{sim}(P_i, P_j)$ [pick smallest line formed by (P_i, P_j) where $P_i \in C_i$ and $P_j \in C_j$]

→ **MAX** pick 2 farthest pts in C_i, C_j

→ **GROUP AVE** $\text{sim}(C_i, C_j) = \frac{\sum \text{sim}(P_i, P_j)}{|C_i| * |C_j|}$

→ less prone to noise
biased towards globule

ward → 12
some as group-avg
but $\sum (d(P_i, P_j))^2$
 $(|C_1||C_2|)$

→ Centroid distance calculate $d(\text{centroid}(c_i), \text{centroid}(c_j))$

(63)

↳ can't kernelize easily? → because c_i, c_j may not be in \mathcal{D}
∴ we don't have $d(c_i, j)$

pros & cons

* MIN — works for nonelliptical shapes

⊗ sensitive to noise and outliers

* MAX — less sensitive to noise

⊗ Breaks large clusters into smaller ones

⊗ biased towards globular shape

MIN — single link

MAX — complete link

Time complexity

vv large

space $O(n^2)$ = cuboid

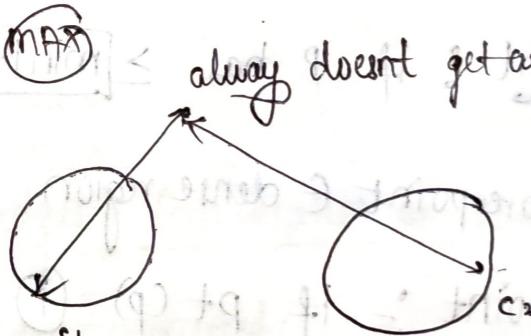
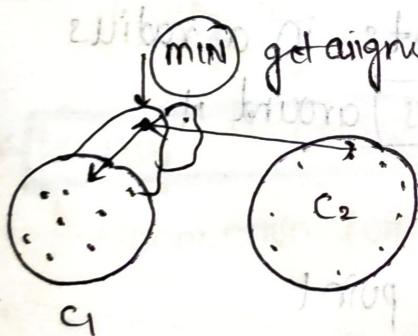
matrix ($n \times n$) \rightarrow proximity matrix

time $= O(n^3) \Rightarrow n$ iterations

↳ In each update proximity matrix which is n^2

MAX vs MIN in presence of noise

$O(n \cdot n^2) = n^3$



∴ min accumulates more and more noise

from $\min \leftarrow \emptyset$

$$E^3 = (0.975)^3$$

etc etc

new edition
aborted etc

(2)

66

DBScan

Density Based clustering

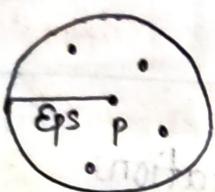
main aim → get sparse regions → noise
dense regions → clusters

measuring density

→ min Points \geq hyperparameters
Eps

① density at point P :- No of points in hypersphere formed by $\text{radius} = \text{Eps}$ around point P.

2D



$$d(P) = 4$$

② dense region :- a hypersphere with $\text{radius} = \text{Eps}$ that contains atleast min pts no of points

$D = \{x_i\}$, min pts, Eps ← given input

* core point :- if pt P has $\geq \text{min pts}$ no of points in a radius of Eps around it

③ (P)
∴ core point \in dense region

* border point :- if pt (P)

i) p is not a core point
ii) pt Q ∈ Neighbourhood (Q)

$Q \rightarrow \text{core point}$

$\text{dist}(P, Q) \leq \text{Eps}$

Neither core
Nor border \Rightarrow Noise pts

⑤

Ex: let $\text{Eps} = 1$
 $\text{minpts} = 4$

x_1 = core

x_2 = border

x_3 = noise

pts around circle = 3 < 4

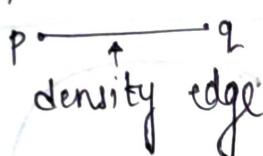
but, $d(x_2, x_1) \leq \text{Eps}$ and x_1 is core pt

* Density Edge ⑥

① p, q are core points

② $\text{dist}(p, q) \leq \text{Eps}$

then



* Density Connected points ⑦

↳ connection formed by density edge

$O(\log n)$ using K-d tree

range Query

$si \rightarrow \text{range-q}(x_i, D, \text{Eps})$

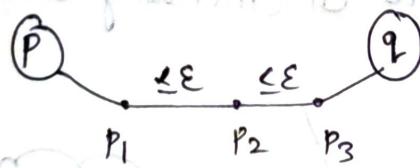
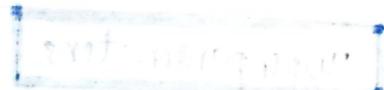
↳ all points within Eps radius of (x_i)

Time Complexity

$O(n \cdot \log n)$ range query using Kd-tree

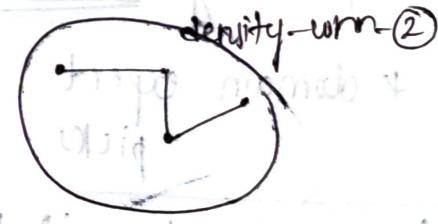
space

$O(n)$



$P_1, P_2, P_3 \rightarrow$ core pts
 $P_1 - P_2$: density edges

density-connected ①



cluster 2

(68)

Algorithm**DBSCAN**

~~range Query~~ \Rightarrow implemented
KD-tree

Input: $\{x_i\} = D$, minPts, Eps

~~but what is (apart from)~~

Step 1: $\forall x_i \in D$, label as core pt, border pt, noise pt

$s; \rightarrow \text{rangeQuery}(x_i, D, \text{Eps}) = \text{return all pts within}$
 (x_i, Eps) ~~properly~~ \rightarrow \rightarrow

Step 2: Remove all noise points

Step 3: for each core point P not assigned to cluster;

i) create new cluster with P

ii) add all pts that are density connected to P
into this new cluster

iteration \leftarrow

Step 4: \forall border point, assign it to nearest core pts cluster

hyperparameters

minPts, Eps

dim

① minpts

* usually $mp > d + 1$

and $mp = 2 * d$

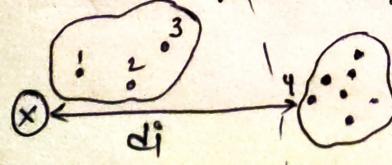
* for noisy data

pick larger MinPts

* domain expert picks

$di \uparrow \rightarrow$ noisy pts if $d(x_i, 4^{\text{th}})$ \uparrow

then to remaining pts will be even min



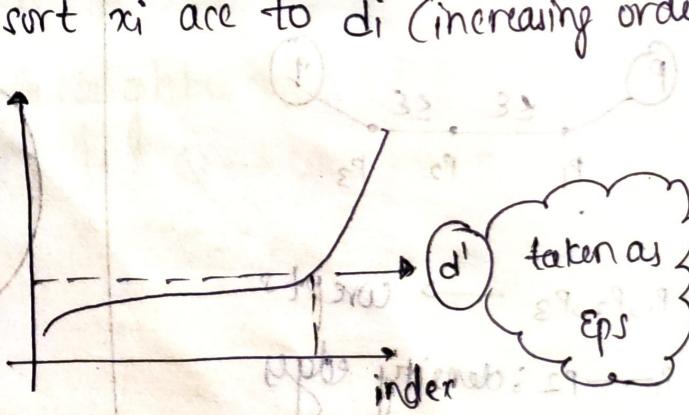
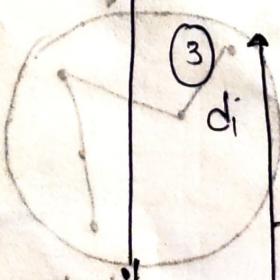
Ex: minPts = 4

② Eps

① $\forall x_i \rightarrow d_i$: dist from x_i to 4th

nearest neighbour of x_i

② sort x_i are to d_i (increasing order)



Advantages

- ⑦ resistance to noise
 - * you don't need to specify no of clusters at start
 - * 2 HP - m^p, Eps
 - * suitable for databases
R*-tree for retrieval
 - * cluster shape may vary

Disadvantages

- * data with varying densities
 - * high-dimensional data
 - * v.v. sensitive to hyperparameters
 - * domain expertise needed to choose b/w m_p, ϵ^*

69

(+) profit will increase

broad fronted

it is not

I am stronger at stock still we
will do the waterribution of wind - I am
still ribbed in my strength we

stage 2 → ratson view ←
WOB to man

$$\boxed{mT} = \boxed{e^{\frac{T}{k}} \cdot T} = \boxed{U}$$

2011-10-10

1912 1912-1913

Recommender Systems

Dataset

Let A be matrix \textcircled{A} \longrightarrow

n: no of users

m: no of items

$A(i,j) \rightarrow$ score given (Netf(1?))
watched/not (youtube)

leave blank if not watched

$$\text{sparsity} = \frac{\# \text{ non-empty cells}}{\# \text{ of total cells}}$$

Task

Ques Given user u_i and I, I_m
suggest In which (u_i) will like

Type

content based

Collaborative filtering (CF)

classif., regression

Core assumption for CF

→ users who agreed / liked in past will also agree in future

* content based

→ doesn't use A_{ij}

→ we meta data to represent u_i, I_j

Ex: Ij - Enne, autor, director, de cade, y or

$u_i \rightarrow$ favorite genre, likes, dislikes

Similarities

→ item-item sim

→ user - user sim

wr-wr

$$U_i^T = \begin{bmatrix} I_1 & I_2 & \dots & I_m \end{bmatrix}$$

4 mm (1) = 10 cm

(S^u_{ii})

→ user vector

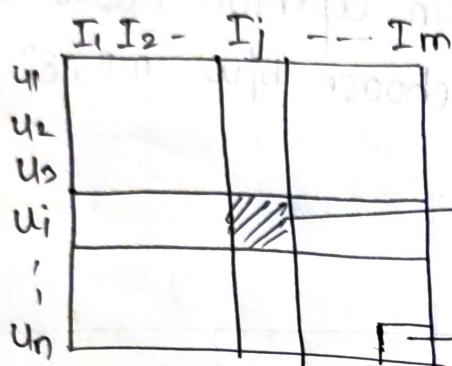
same as BoN

$$\text{sim}(u_i, u_j) = \cos \text{me}(u_i, u_j) = u_i^T u_j$$

Recommender Systems

Dataset

let ds be matrix A \rightarrow n : no of users
 m : no of items



$A(i, j)$ \rightarrow score given (Netflix)
 watched/not (youtube)

leave blank if not watched/rated

$$\text{sparsity} = \frac{\# \text{ non-empty cells}}{\# \text{ of total cells}}$$

Task

given user U_i and I, I_m
 suggest In which (U_i) will like

content based

classification, regression

Type

Collaborative filtering (CF)

Core assumption for CF

\rightarrow users who agreed / liked in past will also agree in future

content based

\rightarrow don't use A_{ij}

\rightarrow we meta data to represent U_i, I_j

Ex: I_j - genre, actor, director, decade, year

U_i \rightarrow favorite genre, likes, dislikes

Similarities

item-item sim

user-user sim

user-user

$$U_i^T = \begin{bmatrix} I_1 & I_2 & \dots & I_m \end{bmatrix}$$

$$\text{sim}(U_i, U_j) = \text{cosine}(U_i, U_j) = \underline{U_i^T U_j}$$

$$(S_{ij}^u)$$

$$\text{null}(U_j)$$

\rightarrow sparse
 same as BoW

(S^u) → user similarity matrix $\Rightarrow S_{ij}$ (similarity b/w user i, j) (7)

Ex: recommend movies to u_n

① find users u_i - u_m similar to u_n

② recommend movies watched by u_i - u_m to u_n

problem with user-user smm

* user preferences change over time

if (ur > items)

use item-item smm sys

Item-Item based

why?

rating generally don't change significantly after initial time

(S^I) → item similarity matrix

Factorization
Matrix Decomposition

PCA

→ centered data

$$S_{d \times d} = \frac{X \cdot X^T}{n-1}$$

PCA or mat decomp

Eigen decomposition PCD

$$S_{d \times d} = X^T \cdot X$$

find top eigen values $\lambda_1, \lambda_2, \dots, \lambda_n$
and eigen vectors w_1, w_2, \dots, w_n
such that $w_i^T \cdot X = 0$

$$\lambda_1 w_1 = S w_1, \\ \lambda_2 w_2 = S w_2, \dots$$

$$S_{d \times d} = W_{d \times d} \cdot \Lambda_{d \times d} \cdot W_{d \times d}^T$$

$S = \text{covariance}(X)$

where,

$$W_{d \times d} = [w_1 \ w_2 \ \dots \ w_n]$$

eigen vectors

matrix of eigen vectors

$$\Lambda_{d \times d} = [\lambda_1 \ \lambda_2 \ \lambda_3 \ \dots \ \lambda_n]$$

diagonal matrix of eigen values

(22)

Singular Value decomposition

SVD

→ SVD → can work when X is rectangular

→ PCA → $S = \text{cov}(X) = X^T \cdot X$ is square matrix + symmetric

$$X_{n \times d} = U_{n \times n} \cdot \Sigma_{n \times d} \cdot V^T_{d \times d}$$

$$U = \begin{bmatrix} u_1 & u_2 & \dots & u_n \end{bmatrix}$$

left singular vector of $X_{n \times d}$

right singular vector of $X_{n \times d}$

$$\Sigma_{n \times d} = \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & \ddots & \\ & & & s_d \\ & 0 & & \\ & & \ddots & \\ & & & s_d \end{bmatrix}$$

so so diagonal matrix

s_i : singular values of X

SVD ↔ PCA

→ u_i = i^{th} eigen vector of $X \cdot X^T$

v_i = i^{th} eigen vector of $X^T \cdot X$

$$\frac{s_i}{\sqrt{n}} = \lambda_i$$

Eigen values

Not. 13 ($n \times d \cdot d \times n = n \times n$)

This is S ($d \times d$)

NMF (Non-negative MF)

SVD

$$X_{n \times m} = U_{n \times k} \cdot \Sigma_{k \times k} \cdot V^T_{k \times m}$$

$$A_{n \times m} = B_{n \times d} \cdot (C^T)_{d \times m}$$

$$B_{ij} \geq 0 \quad C_{ij} \geq 0$$

$$B_{ij} \geq 0 \quad C_{ij} \geq 0$$

Each cell in B and C is non-negative

not sure

$$A_{ij} = B_{i \cdot} \cdot C_j$$

or

$$A_{ij} = B_{ij} \cdot C_i$$

MF for Collaborative Filtering

(73)

$A_{n \times m}$ n-users
m-items, sparse matrix $\rightarrow A_{ij}$

$$r_{ui} = b_i^T c_u$$

let $A_{n \times m} = B_{n \times d} C^{(CT)}_{d \times m}$

where, $d \leq m, d \leq n$

$A_{ij} = B_i \cdot C_j$ represented as col vector

ith row of B jth row of C represent scalar value

$$A_{ij} = B_i^T \cdot C_j = B_j^T \cdot C_i$$

$$A_{ij} = B_i^T \cdot C_j \rightarrow A = B C^T$$

optimization problem

∴ optimization prob =

B_i and C_j
now represented as column vectors

$$\underset{B, C}{\operatorname{argmin}} \sum_{i, j \in (A_{ij} \text{ is non-empty})} (A_{ij} - B_i^T \cdot C_j)^2$$

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

if empty you
can't optimize

don't use empty cells

5 steps

→ Step 1

$$\underset{B, C}{\operatorname{argmin}} \sum (A_{ij} - B_i^T \cdot C_j)^2$$

Step 3 matrix completion

$$A = \begin{bmatrix} & & & & & \\ & \vdots & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix}^m$$

Empty (no rating) cell

→ Step 2

get values of

filling empty cell

= recommendation

$$A_{10,5} = B_{10}^T \cdot C_5$$

→ fill all gaps

B and C

predictive value

now, if u_i had not rated I_i it would be empty

now we have predicted rating $(B_{10}^T \cdot C_5)$

74

K-means as MF

K-means : $D \cdot \mathbb{R}^{n \times n}$

↳ n : no of pts

K : no of centroids

optⁿ prob

$$\min_{C_i} \sum_{i=1}^K \sum_{x \in S_i} \|x - C_i\|^2$$

→ let's take a matrix (Z),

$$Z(k \times n) = \begin{bmatrix} 1 & 2 & \dots & 5 \\ 2 & & & \\ \vdots & & & \vdots \\ K & & & \end{bmatrix} - n \quad \text{where,}$$

$$z_{ij} = \begin{cases} 1, & \text{if } x_j \in S_i \\ 0, & \text{else} \end{cases}$$

$$j_8 = 4$$

$x_5 \in S_7$ and all other ele in col = 0 (each pt \in only 1 cluster)

→ assume matrices C and X

$$X = d \times n \quad \begin{bmatrix} 1 & 2 & \dots & n \\ x_1 & x_2 & \dots & x_n \end{bmatrix}$$

$$C = d \times k \quad \begin{bmatrix} c_1 & c_2 & \dots & c_k \end{bmatrix}$$

Now,

$$\min_{C_i} \sum_{i=1}^K \sum_{x \in S_i} \|x - C_i\|^2$$

$$\min_{Z_j} \sum_{i=1}^K \sum_{j=1}^n z_{ij} \|x_j - C_i\|^2$$

$$X = C \cdot Z$$

using MF find (duempore)

$$X = C \cdot Z$$

$$2^2 + 0^2 + 0^2 = 2^2$$

where

$$\sum_{i=1}^k z_{ij} = 1 \quad (\text{columns of } Z)$$

due to this it is NP-hard

$$\min_{CZ} \|X - CZ\|^2$$

new, optimization

$$X_{d \times n} = C_{d \times k} \cdot Z_{k \times n}$$

hyperparam tuning

→ (d)

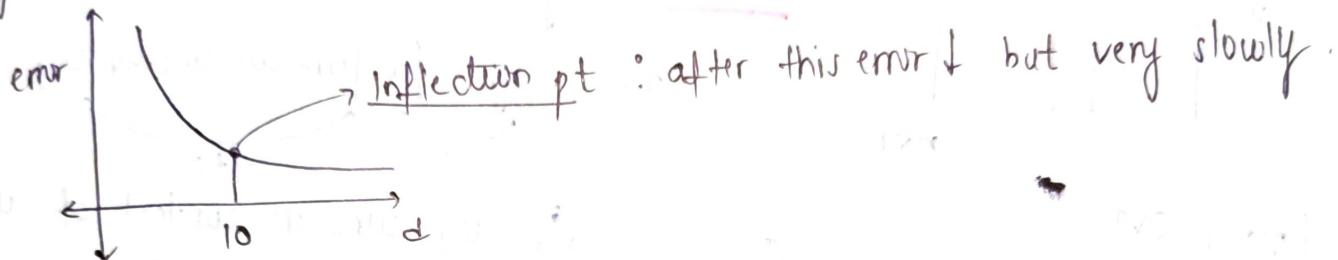
$$A_{n \times m} = \underbrace{B}_{n \times d} \cdot \underbrace{C^T}_{d \times m}$$

75

$$\boxed{\text{error}} : \sum_{ij} (A_{ij} - B_i^T C_j)^2$$

hyperparam = (d) → ① problem-specific

② systematic way : (as $d \uparrow$ error \downarrow)



Netflix prize

input : r_{ui} → rating by user (u) on movie. (i)
 q_i → item vector
 p_u → user vector

$i = \text{no of items/movies}$
 $u = \text{no of users}$
 $r_{ui} \equiv A_{ij}$

opt'n pblm : $\min \sum_{(u,i)} (r_{ui} - q_i^T p_u)^2 + \lambda (\|p\|^2 + \|q\|^2)$

cold-start problem

what happens when new user joins? u don't have its data

(sol) use metadata abt user
→ ip, browser, OS

[ip browser OS]

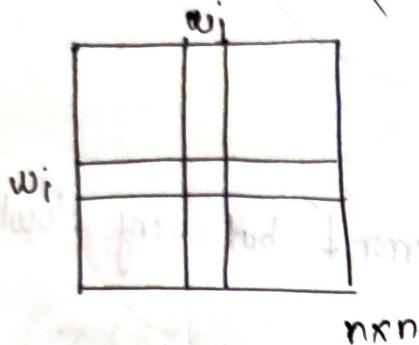
history show don't

Word vectors as MF

$D = \{ \text{reviews}_i \}_{i=1}^n$

$n = \# \text{ of words in corpus}$

co-occurrence matrix (X)



Co-occurrence + truncated SVD

v.v large [so instead of all n words, are $\approx K$ most important words] $K \ll n$

(76)

$X_{ij} = \# \text{ of times } w_j \text{ occurs in context of } w_i$
context: neighborhood ($\text{ex} = \text{neig} = 3$)

$$r_2 = \underbrace{w_j w_2}_{3} \boxed{w_i} \underbrace{w_3 w_2 w_3 w_4}_{4+2} w_j$$

$\therefore w_j$ occurs in context of w_i 1 time in r_2

MF SVD

$$X_{n \times n} = U_{n \times n} \cdot \sum_{k=1}^K \Sigma_{k \times k} \cdot V_{k \times n}^T$$

(using SVD)

$$U_{n \times n} = \begin{bmatrix} u_1 & u_2 & \dots & u_n \end{bmatrix}$$

$$\Sigma_{n \times n} = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_2 & \\ & & & \ddots & \sigma_n \end{bmatrix}$$

$\sigma_i \rightarrow$ singular values

$$V_{n \times n} = \begin{bmatrix} v_1 & v_2 & \dots & v_n \end{bmatrix}$$

Truncated SVD

→ discard last columns → $\Sigma_{n \times n}$ pick top K singular values

$$U_{n \times n} = \hat{U}_{n \times K}$$

$$\Sigma_{n \times n} = \hat{\Sigma}_{K \times K}$$

$$V^T = \hat{V}^T$$

$$1 \leq K \leq n$$

↑ holding fraction

(PCA) - pick top K eigen values

take \hat{U} after truncated SVD

$$U = U_1, U_2, \dots, U_K$$

$U_i = \text{word vector for word } i$

final word vector

Eigenface → PCA on Images

64x64

77

* Dataset → 400 images, each 64×64



$$\textcircled{1} = [n_1 \ n_2 \ \dots \ n_{400}] \Rightarrow A = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ 400 \end{bmatrix}$$

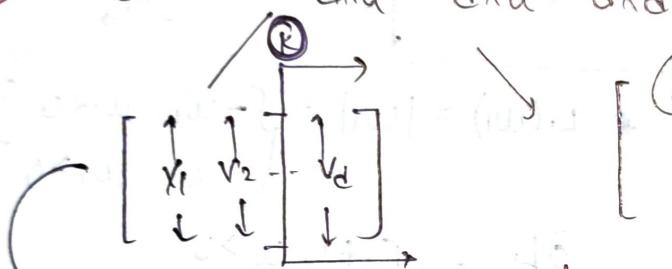
$$A = \begin{bmatrix} 1 & 2 & \dots & 400 \end{bmatrix} \quad n \times d$$

[step 1] compute covariance matrix

4096

$$A_{n \times d} \rightarrow \text{cov}(A) = S_{d \times d}$$

$$\textcircled{2} \quad S_{d \times d} = W_{d \times d} \Lambda_{d \times d}^{1/2} W^T$$



→ take top K eigen vectors corresponding to top K eigen values

$$\textcircled{3} \quad k_{d \times K} = [v_1 \ v_2 \ \dots \ v_K] \quad K \ll d$$

$$A_{n \times d} \cdot k_{d \times K} = A_{n \times K}$$

→ K -dim vector for each image

∴ d -dim → K -dim

PCA

[Note] as ($K \downarrow$) → clarity of photos decreases

as $K \downarrow$ clarity of photos decreases

$$[64 \times 64] \text{ photo} \rightarrow [10 \times 10]$$

but clarity

of a speed limit sign

photo

photo

photo

photo

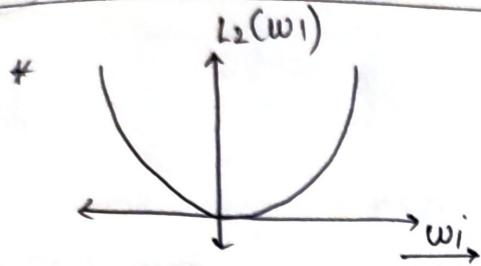
photo

why L₁ reg creates sparsity

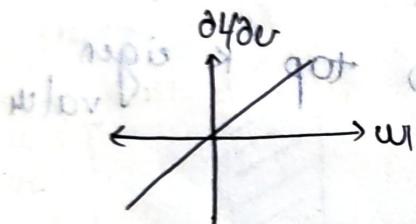
L₂-reg

$$* L_2 = \|w\|_2^2 = w^T w$$

$$= \sum_{j=1}^n w_j^2$$



$$* \frac{\partial L_2}{\partial w_i} = \frac{\partial w_i^2}{\partial w} = 2 \cdot w_i$$



$$* w_{new} = w_{old} - r(2w_i)$$

Observations

as $w_i < 0$, the gradient $\frac{\partial L_2}{\partial w}$

also decreases gradually

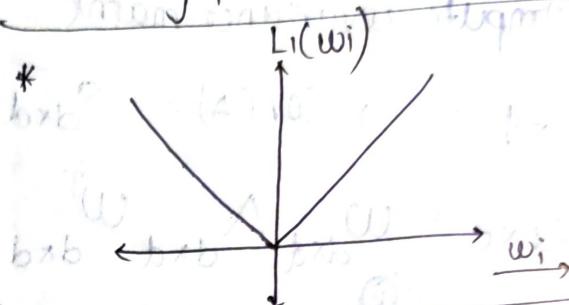
$\therefore w_{new} = w_{old} - A$ — decreases constantly

w_{new} keeps \downarrow slowly,
because as $w_i \downarrow$, A gradually

Sparsity ie ≈ 0

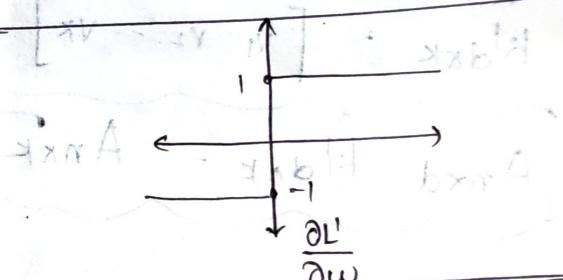
L₁-reg

$$* L_1 = \|w\|_1 = \sum_{j=1}^n |w_j|$$



$$* L_1(w_i) = |w_i| = \begin{cases} -w_i, & w_i \leq 0 \\ w_i, & w_i > 0 \end{cases}$$

$$\frac{\partial L_1}{\partial w_i} = \begin{cases} +1, & w_i > 0 \\ -1, & w_i \leq 0 \end{cases}$$



$$* w_{new} = w_{old} - r(w_i)$$

Observations

as $w_i < 0$, the grad does not decrease, it remains constant

$$\therefore w_{new} = w_{old} - A$$

↓
constant

$\therefore w_{new}$ keeps \downarrow because A is constant (1)

SVD with L2 normalization

79

logistic regression

Loss Fn

$$L = \sum_i -y_i (\log(\sigma(w^T x + b))) - (1-y_i) (\log(1-\sigma(w^T x + b))) + \frac{\lambda}{2} \|w\|^2$$

wkt,

$$\left[\frac{d\sigma(z)}{dz} = \sigma(z)(1-\sigma(z)) \right] \text{ & } \left[(d+x^T w) \sigma(z) = (d\sigma(z)) \right]$$

$$\nabla_w L = \nabla_w \left[-y \log(\sigma(w^T x + b)) - (1-y) \log(\sigma(-w^T x - b)) + \frac{\lambda}{2} \|w\|^2 \right]$$

$$= -y \cdot \frac{\sigma(w^T x + b)(1-\sigma(w^T x + b)) \cdot x}{\sigma(w^T x + b)} + (1-y) \frac{\sigma(-w^T x - b)(1-\sigma(-w^T x - b))}{\sigma(-w^T x - b)} \cdot (-x) + \frac{\lambda}{2} \cdot 2\|w\|^2$$

$$= -y (1-\sigma(w^T x + b)) \cdot x + (1-y) (1-\sigma(-w^T x - b)) \cdot (-x) + \frac{\lambda}{2} \|w\|^2$$

$$\begin{aligned} &= -xy + xy \sigma(w^T x + b) + (x-xy) \sigma(-w^T x - b) + \lambda w \\ &\quad \cancel{-xy + xy \sigma(w^T x + b)} + x \sigma(w^T x + b) - \cancel{xy \sigma(-w^T x - b)} + \cancel{\lambda w} \\ &= -xy + x \sigma(w^T x + b) + \lambda w = -x(y - \sigma(w^T x + b)) + \lambda w \end{aligned}$$

$$\therefore \nabla_w L = \sum_{i=1}^n -x_i (y_i - \sigma(w^T x_i + b)) + \lambda w$$

$$\nabla_w L = - \sum_{i=1}^n x_i (y_i - \sigma(w^T x_i + b)) - \lambda w$$

$$\begin{aligned}
 \textcircled{2} \quad \nabla_b L &= -y \left(\log(\sigma(w^T x + b)) \right) + (1-y) \log(\sigma(-w^T x - b)) + \frac{\lambda \|w\|_2^2}{2} \\
 &= -y \cdot \frac{\sigma(w^T x + b)}{\sigma(w^T x + b)} (1 - \sigma(w^T x + b)) (1) - (1-y) \frac{\sigma(-w^T x - b)}{\sigma(-w^T x - b)} (1 - \sigma(-w^T x - b)) \\
 &= -y (1 - \sigma(w^T x + b)) + (1-y) \sigma(w^T x + b) \\
 &= -y + y \sigma(w^T x + b) - \sigma(w^T x + b) - y \sigma(w^T x + b) \\
 &= -y + \sigma(w^T x + b)
 \end{aligned}$$

$$\begin{aligned}
 \text{Final} \quad \nabla_b L &= \frac{1}{n} \sum_{i=1}^n (-y_i - \sigma(w^T x_i + b)) = ((d+x^T w) \cdot \alpha) p(x) - \frac{1}{n} \sum_{i=1}^n (-y_i - \sigma(w^T x_i + b))
 \end{aligned}$$

\therefore final update eqn,

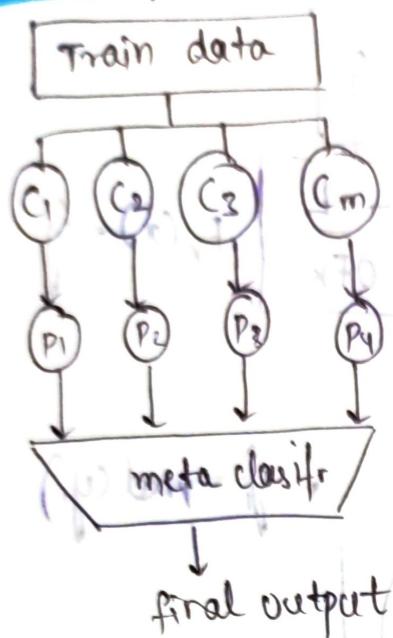
$$\begin{aligned}
 b^{t+1} &\Rightarrow b^t - \alpha \left(-\sum_{i=1}^n y_i - \sigma(w^T x_i + b) \right)
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow b^{t+1} &= b^t + \alpha \left[\sum_{i=1}^n y_i - \sigma(w^T x_i + b) \right] \quad w=w^t, b=b^t \\
 \Rightarrow w^{t+1} &= w^t - \alpha \left[\sum_{i=1}^n x_i (y_i - \sigma(w^T x_i + b)) - \lambda w \right] \quad w=w^t, b=b^t
 \end{aligned}$$

$$w^t + ((d+x^T w) \cdot \alpha - \beta) x = w^t + (d+x^T w) \cdot \alpha x + \beta x = \frac{1}{n} \sum_{i=1}^n x_i (y_i - \sigma(w^T x_i + b)) + \frac{\lambda}{n} w$$

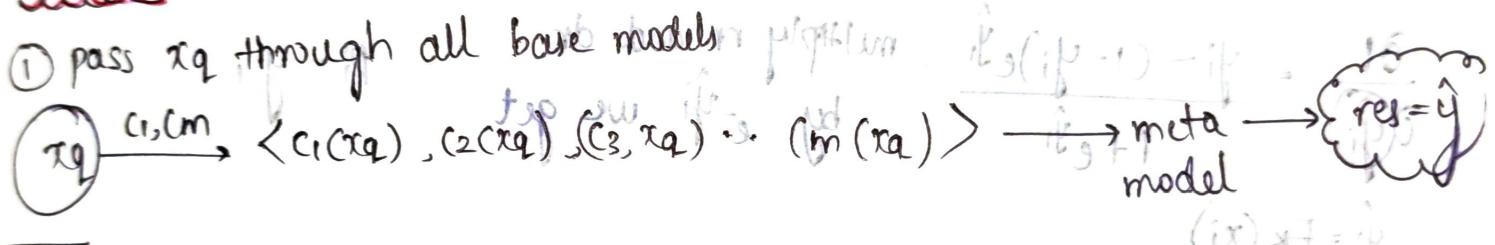
$$\boxed{w^t - ((d+x^T w) \cdot \alpha - \beta) x = \frac{1}{n} \sum_{i=1}^n x_i (y_i - \sigma(w^T x_i + b)) + \frac{\lambda}{n} w}$$

stacking



- Train algorithm: $\{x_i, y_i\}_{i=1}^n$ (for $i = 1, 2, \dots, n$)
- [Step 1] train m models (Base classifiers) using $\{x_i, y_i\}$ training data, C_1, C_2, \dots, C_m
 - [Step 2] New dataset $D' = \{(x'_i, y'_i)\}_{i=1}^{n'}$ (for $i = 1, 2, \dots, n'$)
 - for each data point x'_i , pass it through each classifier
 - $x'_i = \{h_1(x'_i), h_2(x'_i), \dots, h_m(x'_i)\}$
 - [Step 3] train meta classifier using D'

Prediction: x_q

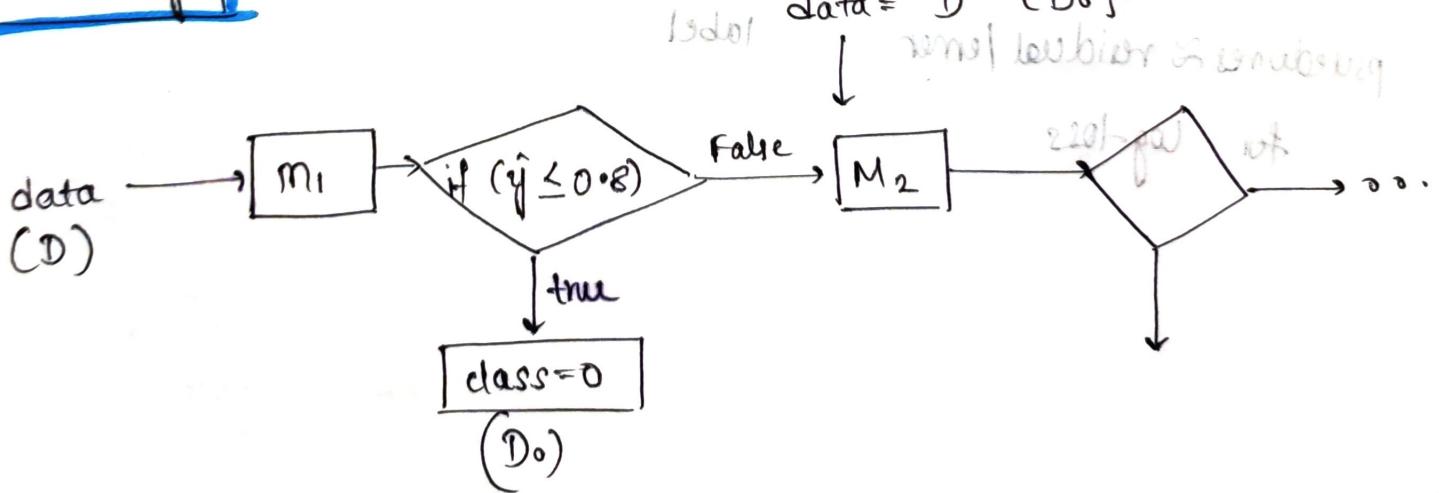


TC

train: t_c to train m base models + 1 meta model

prediction: t_c to run/predict of m base models + meta model

Cascading



(22)

psuedo-res = residual ?

① Squared Loss [refer page (84)]

$$\text{residual} = y - \hat{y}$$

psuedo res ~ residual

$$\text{psuedo-res} = \frac{-\partial L}{\partial F_k(x)}$$

stoch. noise

$$F_k(x) = \hat{y}$$

② Log-loss

$$L = y_i \cdot \log(p_i) + (1-y_i) \log(1-p_i)$$

$$\hat{y}_i = p_i$$

$$\hat{y}_i = \text{pred}(y_i)$$

$$L = y_i \log \left(\frac{1}{1+e^{-\hat{y}_i}} \right) + (1-y_i) \log \left(\frac{e^{-\hat{y}_i}}{1+e^{-\hat{y}_i}} \right)$$

ok : not bled

$$\frac{\partial L}{\partial y_i} = y_i - (1-y_i)e^{\hat{y}_i}$$

multiply numerator and denominator by $e^{\hat{y}_i}$ we get

$$\hat{y}_i = F_k(x_i)$$

$$\frac{-\partial L}{\partial F_k(x)} = p_i - y_i = \frac{1 - y_i}{1 + e^{-(F_k(x))}}$$

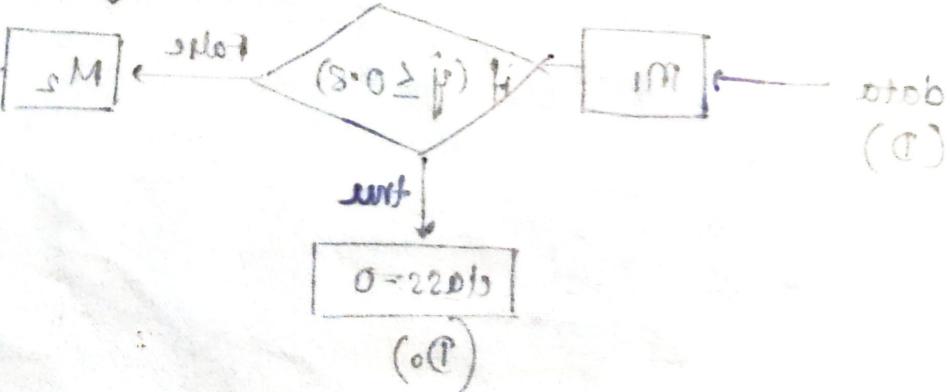
error

actual class label

∴ for bin class

psuedores ~ residual / error

for log-loss



Decision Tree Extras

83

* it doesn't need standardization as it is not distance based.

TC :- Train :- $\sim O(n \log n \cdot d)$ → n : #pts
 d : dimensions
 sorting for numerical features
 run_TC :- max-depth of decision tree
 ↳ very low %, low-latency

∴ DT \propto dim

Best when
 large data
 small dim
 low-latency

Cases

* Imbalanced data can't be used or handled directly

+ as $d \uparrow \rightarrow$ train + c ↑ ∴ avoid one-hot-encoding
 why? $f_i = \text{state}_{\text{one}} \rightarrow f_i = \langle f_{i1}, f_{i2}, f_{i3} \dots f_{in} \rangle$

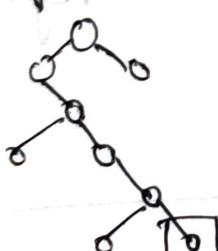
multi-class $d \uparrow \rightarrow$ train + c ↑
 ↳ one vs rest
 ↳ Entropy handles multi-class → at leaf node do majority vote

Outliers

→ as depth ↑ - outlier impact more

why?

unstable tree



may we have two 10 atm
 2 points and they maybe
 outliers

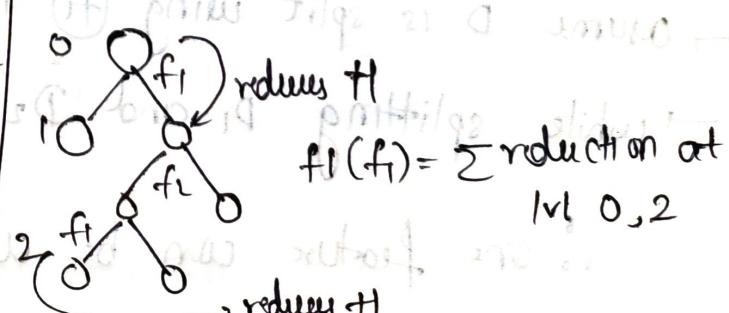
(17)

Interpretability

Feature Importance

fI for (f_i) = reduction in H or I_g

$fI(f_i)$ = normalized sum of reductions in H / I_g due to f_i



$$fI(f_i) = \sum \text{reduction at lvl } 0, 2$$

Decision Tree Regression

→ Impurity = MSE

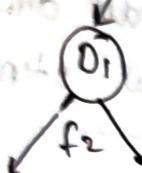
Ex-

D

f₁

split using f₁

$$\text{MSE}(D) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$



split using f₁

$$\text{Split using } f_1 = \text{MSE}(D_1) + \text{MSE}(D_2)$$

$$f_2 = \text{MSE}(D_1) + \text{MSE}(D_2)$$

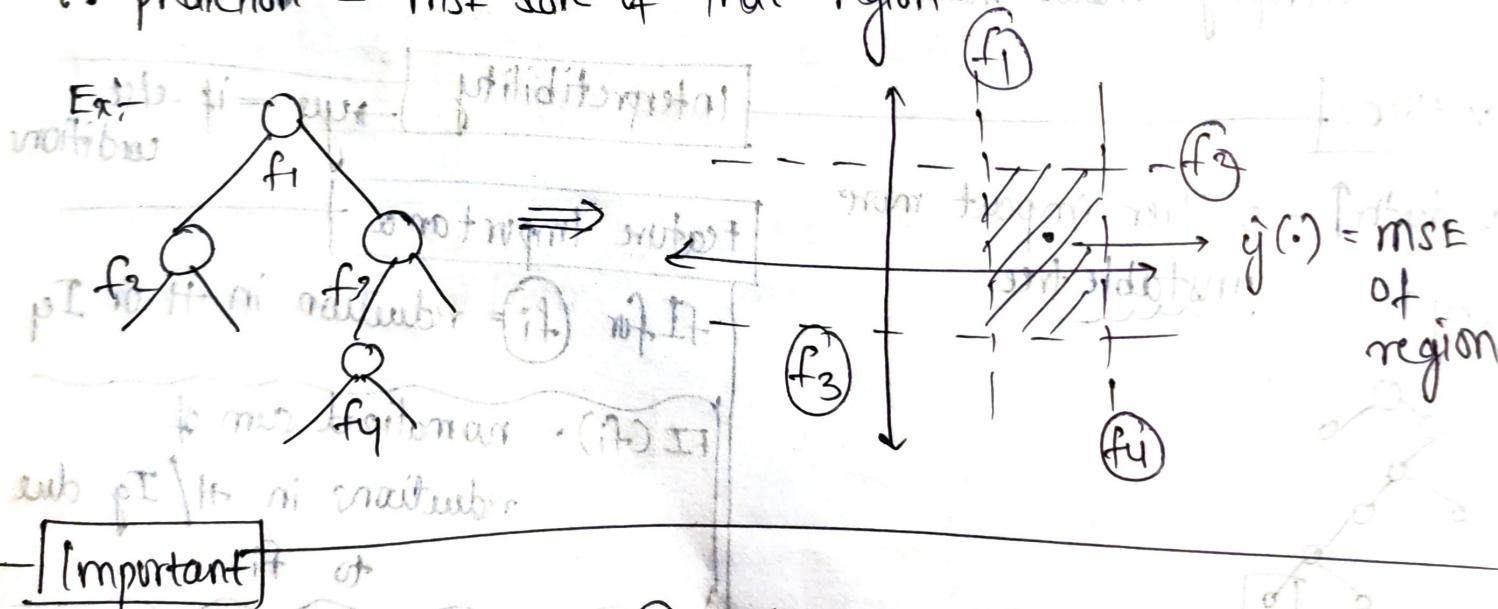
∴ which reduce MSE(D) most is selected as split

How to predict

→ given a point x_q → follow path and reach root node

at leaf node $\left(\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right)_{n=\text{no. of nodes in leaf}} = \hat{y}$

∴ prediction = MSE score of that region



[Important]

→ assume D is split using f₁ into D₁ and D₂

→ while splitting D₁ and D₂ we again consider f₁

∴ one feature can be used for splitting given it reduces left entrop

Kernel SVM

(85)

Normal SVM

→ Train phase and get optimal \vec{w}

→ Test $\text{sign}(\vec{w}^T \vec{x} + b)$ if. + class
-, -ve class

Kernel SVM

Train find best \vec{w}

Test decision fn = $\sum_{i=1}^n \alpha_i y_i K(x_i, x_q)$

$\rightarrow \text{sign}(\sum \alpha_i y_i K(x_i, x_q)) = \text{class label}$

Log-R

train \vec{w}

test: $p = \sigma(\vec{w}^T \vec{x})$
 $\Rightarrow 0.5 = \text{class}$
L, co

Query/test point

Kernel (RBF)

$i=1 \rightarrow n$ is set of all SV's

$\alpha_i y_i$: dual-coefs

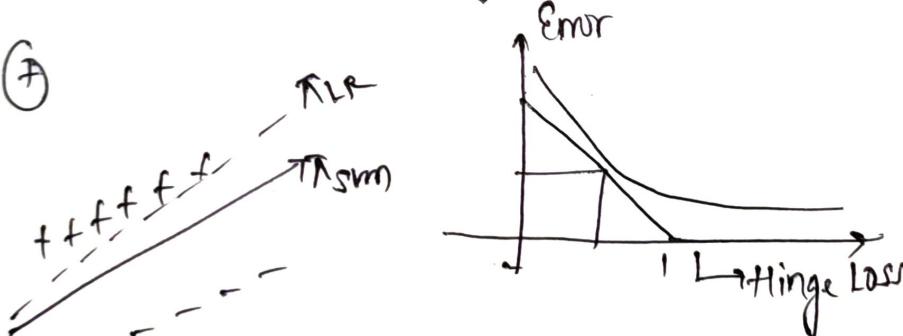
x_i = support-vectors

x_q = test point

$\sigma(-\gamma \cdot \text{euclidean dist})$

why LR is more impacted by outliers and SVM is less?
Ex: you have one correctly classified pt but is very far from π .
In case of LR, π_{lr} changes but π_{svm} doesn't change why?

⑦



→ In hinge loss, if point is correctly classified ie $y_i \vec{w}^T \vec{x} \geq 1$ then error = 0 no matter if > 0

→ but in log-loss,

although error ≈ 0 but not exactly ≈ 0

Conf matrix terms

* Sensitivity : $(TPR) = \frac{TP}{TP+FN}$ [proportional of positives correctly identified as +ve]

* Specificity : $(TNR) = \text{proportional of negative or correctly identified as } \ominus\text{ve}$
 ie out of all $\ominus\text{ve}$ how many $\ominus\text{ve}$ are correctly classified

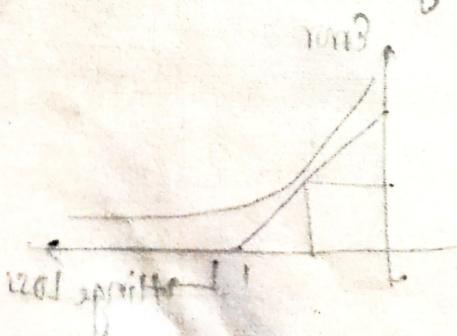
$$\text{sensitivity} = \text{recall} = TPR$$

$$\text{specificity} = 1 - FPR$$

$$* TPR = \frac{TP}{P} = \frac{TP}{TP+FN} = 1 - FNR$$

$$* TNR = \frac{TN}{N} = \frac{TN}{TN+FP} = 1 - FPR$$

first is muz done without ref before going down in q_1 then
 next ref goes to had by bivalve pattern and not up;
 (prior option has been used up so it goes to q_1 to use)



pattern is being +, 2nd option not +
 options left is \leq ref in bivalves
 so \leq if option not +

2nd ref in had

options left are 0 or 1 in had

ALGO DEEP DIVE

(35)

KNN

Train : No train
Time complexity = $O(n^2)$

Test : $O(nd)$
 K is small
* you need large space complexity

* TCF is also very huge \Rightarrow high latency
Space : $O(nd)$
store all train dataset

* forward feature selection
→ Helps to speed up nearest neighbor finding

No of comparisons : Best - $O(\log n) \cdot K$ → check all adjacent boxes
Worst - $O(n) \cdot K$

* problem for dimension 2 → 4 checks [generally d is small size]
* d-dimensions → d^d checks

TC : $O(2d \cdot \log n) = 2O(d \cdot \log n)$ being +

due to adjacent checks
* data must be uniformly distributed
worst case TC $\neq O(\log n)$ more if no of snarls

LSH

* Train TC : $O(m \cdot dn)$ [m - no of hyperplanes]
you need to take each point and check which hyperplane
on which side it lies \rightarrow $\begin{matrix} 0 & 1 & 0 & 1 \\ \pi_1 & \pi_2 & \dots & \pi_m \end{matrix}$

* Query : $O(md + nd)$

ifp : x_q
you need to find

$$\text{sign}(w^T x_q) \quad w \in \mathbb{R}^m$$

if the hash value (bucket has n elements), you need to traverse each to compare
(maybe find k-nearest).

③ Naive-Bayes

IND 9/10 3/10

* Bias-variance in Laplace smooth

* outliers

$\alpha=0$: you give prob to even rare words \rightarrow overfit

① test

new word - laplace smooth

$\alpha \uparrow$: $p \rightarrow 1/2 \rightarrow$ random model

② train

\therefore underfit \rightarrow if word occurs less than a certain threshold (α) ignore

Gives

prob of word given query \rightarrow laplace smooth with high α value

① good for high-dim text

\rightarrow can handle text with high dim.

Assumptions

\rightarrow use log-probabilities

* conditional dependence of features

② good for high-dim text data

+ categorical features

Time complexity

\hookrightarrow used as baseline models

Train $\rightarrow O(n)$

③ Not good for real data

Test $\rightarrow O(nk)$

④ Prone to overfit when Laplace

Space \rightarrow low

smoothing is not used

⑤ highly interpretable

[10 11 10]

(not trivial) interpretation of b_i
of bias, elements in

$(b_{ir} + b_{ir})0$: prev

(p) $\cdot q_i$

logistic-regression

37

- ④ we used distances \therefore standardization is mandatory
- ⑤ if ~~multi~~ multi-collinearity \rightarrow perturbation test
- ↳ can't use L1/L2 as feature importanc
- Time Complexity
- ⑥ Train :- $SE(D) = O(nd)$
- ⑦ Test/run time
- only w space $= O(d)$ is needed
 - time $\leftarrow O(d) \rightarrow \sum_{i=1}^d w_i x_i$
- ⑧ very good for low latency appls
- ↳ Internet companies
- ⑨ if d is large \Rightarrow L1 regularization
- which would make non-useful features = 0

Assumptions

- ① Linearly separable
- ② ! collinearity

Outliers

↳ reduces impact

Handling?

* fit hyperplane

remove very far pts



fit again

Best Case

- * linearly separable data
- * low-latency required
- * low TC

large d \uparrow more chances of finding optimal π^* but as $d \uparrow$ the latency also increases.

$$LR = \text{log-loss} + \lambda \|w\|_2^2$$

- $\lambda \uparrow \rightarrow$ underfit
- $\lambda \downarrow \rightarrow$ overfit

$L1 \text{ reg}$ $\leftarrow \lambda \uparrow \rightarrow$ sparsity $\uparrow \rightarrow$ lat \downarrow

Support vector machines

SVM

Time Complexity

Train $\propto O(n^3)$

Test $\propto O(nd)$ if $d \ll n$
sequential minimal opt (smo)

* Train $\approx O(n^2)$ for kernel SVM

* optimized, $O(nd^2)$ if $d \ll n$

Running time $O(Kd)$ if $K \ll n$
no of support vectors
no margin \Rightarrow no of support vectors

Problems

* No Interpretability esp. Kernel SVM

what is slope of w \Rightarrow

from below which

is most likely to have non

other slopes given \Rightarrow [200 to 8]

higher function \Rightarrow

lower function \Rightarrow

function to work on \Rightarrow best

to do this \Rightarrow limit to

non linear data can't be fit

$(w_1x_1 + w_2x_2 + b) \rightarrow 0$

if $x_1 \leftarrow 1$

if $x_2 \leftarrow 1$

$b \leftarrow b + w_1x_1 + w_2x_2 \leftarrow 1 \Rightarrow \dots$

$O(n^2)$ is large in which few are not used when n is large

so less computation

Outliers

* very less impact — depend only on

$(b) \propto -abs(w)$ \Rightarrow more

* large $n \rightarrow$ high latency

$\propto \frac{1}{n} \rightarrow$ no of support vectors \rightarrow large

more function with no of support vectors

$\propto \frac{1}{n} \rightarrow$ large

Conclusion

longer process \Rightarrow less

processes \Rightarrow less

time

longer process \Rightarrow less

processes \Rightarrow less

longer process \Rightarrow less

longer process \Rightarrow less

Decision Tree

* Scaling is not needed (it is not distance based method)

Categorical with many cat. → mean encoding / pick-top

cat → numerical feature

$$P(y_i=1 | P_j)$$

P_i	y_i
P_1	0
P_2	0
P_3	1

$$P(y_i=1 | P_j) = \# \text{times } y_i=1 \text{ over } \# \text{times } P_j$$

midp of P_i

Overfit / underfit

depth \uparrow → overfit (you can grow tree to fit even noisy points)

how predict?

At each leaf you can have only one point - pure node

multiple nodes - majority vote

Time complexity

Train : $O(n \log n d)$ d - at each node you need to check for all features for info gain

$n \log n$ - mostly when numerical attributes

space :- traindt → nested if else

$O(C(\# internal nodes + \# leaf nodes))$

runtime : $x_q \rightarrow y_q$

$O(\max\text{-depth of any leaf node})$ | generally depth $<$

Cost

good when $\begin{cases} \text{low latency} \\ \text{large data} \\ \text{low dimensions} \end{cases}$

10

Heuristic clustering

first method

space :- $O(n^2)$ time :- $O(n^3) \rightarrow n$ iterations

at each iter you may need to join 2 clusters

Limitations

update similarity matrix

* no clear objective fn / optimization pblm

* MIN - outlier pblm

MAX - breaks large, can't accommodate different sized clusters

* space and time complexity

DBSCAN

Pros

* resistant to noise

* handle diff cluster shapes

* no need to specify K value

Cons

* not entirely deterministic

* depends on distance measures

* cannot cluster datasets with large differences in densities

(K) value

+ user defined
((user for K))

+ K ← user input

> affects clustering (user K → user output)

penalized cost

stab prob

cost

noise cost

new best