

Count plot

```
sns.countplot(df['survived'])
```

```
df['survived'].value_counts().plot(kind='bar')
```

Pie chart

```
df['survived'].value_counts().plot(kind='pie',  
auto_norm=True)
```

Numerical data

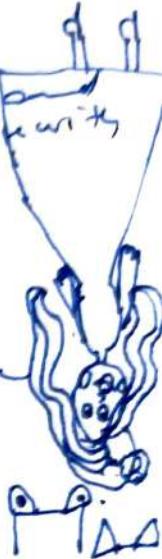
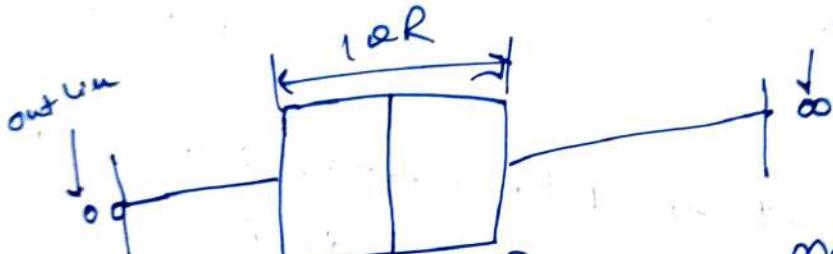
Import matplotlib.pyplot as plt

```
plt.hist(df['Age'], bins=5)
```

we can change the bins value
complexity of bars increase or decrease with
respect to the bins.

Distr Plot

```
sns.distplot(df['Age'])
```



Bow plot

```
sns.bowplot(df['fare'])
```

```
df['Age'].skew()
```

Scatter plot

```
sns.scatterplot(tips['total_bill'], tips['tip'])
```

Bar plot

Two (Numerical - categorical)

```
sns.barplot(df['Pclass'], df['Age'])
```

```
hue = titanic['sex']
```

Box plot (Numerical - categorical)

sns.boxplot(titanic['sex'], titanic['survived'])
 $\lambda = \frac{x-1}{n}$

hue = titanic['survived']
 $\lambda = \frac{x-1}{n}$

Distplot (Numerical - categorical)

sns.distplot(titanic['Age'])

sns.distplot(titanic[titanic['survived'] == 0]['Age'],
 hist = True)

sns.distplot(titanic[titanic['survived'] == 1]['Age'],
 hist = True)

Heatmap (Categorical - categorical)

pd.crosstab(titanic['Pclass'], titanic['survived'])

sns.heatmap(pd.crosstab(titanic['Pclass'],
 titanic['survived']))

titanic.groupby('Pclass').mean()['survived']
 $\lambda = \frac{x-1}{n}$

clustermap (Categorical - categorical)

pd.crosstab(titanic['sibsp'], titanic['survived'])

sns.clustermap(pd.crosstab(titanic['sibsp'], titanic['survived']))

Pairplot sns.heat()

sns.pairplot(sns.heat())

sns.Pairplot(sns.heat(), hue = 'survived')

Lineplot (survived - survived)

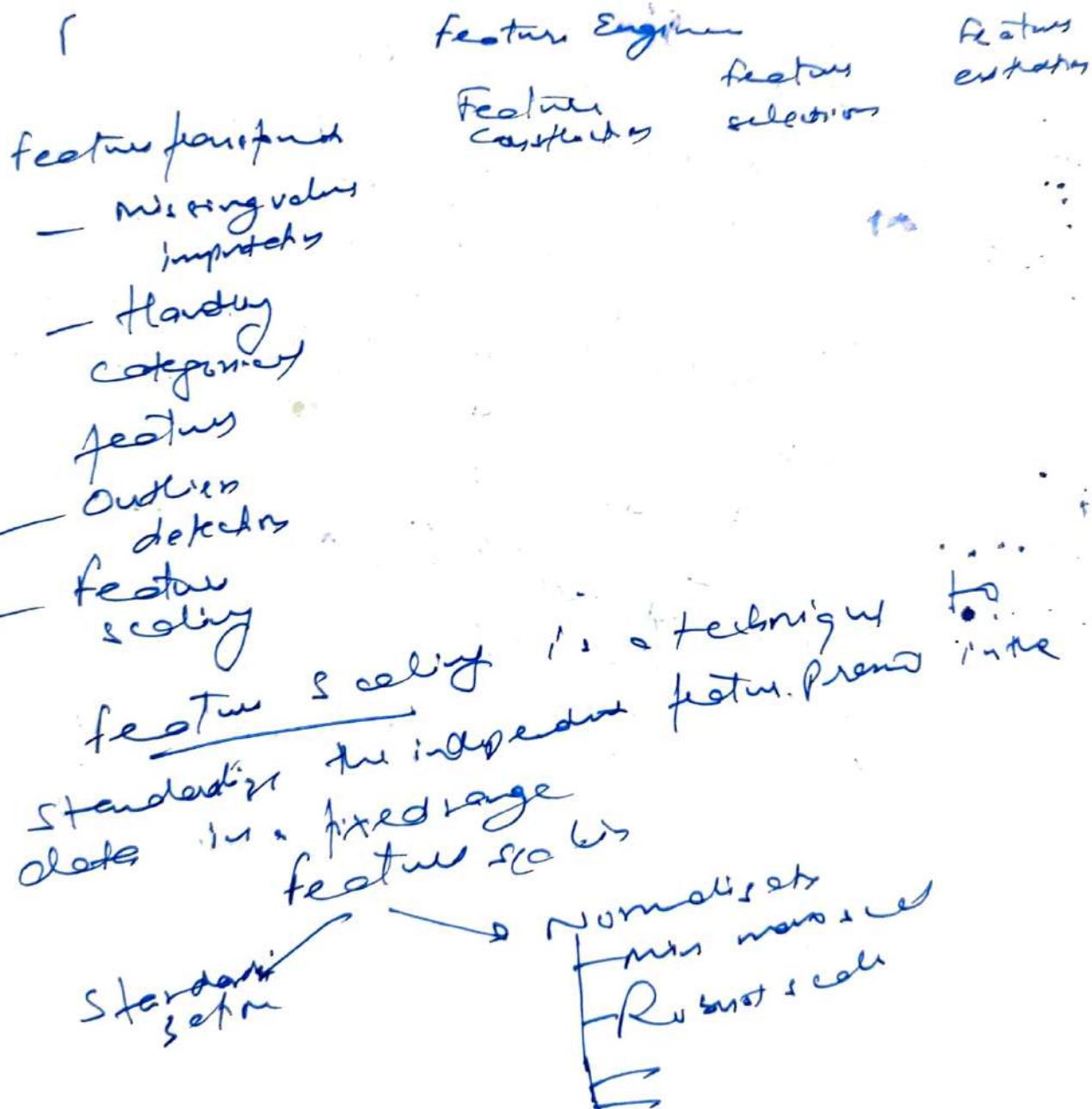
Pandas profilers

1 Pip install pandas-profiling

```
from pandas_profiling import ProfileReport  
prof = ProfileReport(df)
```

```
prof.to_file(output_file='Output.html')
```

feature engineering is the process of using domain knowledge to extract feature from raw data. These features can be used to improve the performance of ML algorithm.



Also called Z-score normalization

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

when to we standardize it



k-means



k-nearest neighbors



PCA



A ~ A'



Gradient descent

the is no need for standardization in Decision tree
Random forest, Gradient boost neighbors

Random forest

- Normalization is a technique often applied as part of data prep. for machine learning. The goal of normalization is to change the values of numerical columns in the dataset to a common scale, without distorting differences in the ranges of values or losing information.



min max scaling

measurements &

max scale was absolute

Robust scaling

min max scaling

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

weight

180

67

81

Order ~~transform source when different, include a
function name, contact, certificate of origin~~

Invoiced transform. transform Column Transform

Pipelines chain together multiple steps so that
the output of each step is used as input to the
next step

- easy to apply the same preprocessing

to have first

Box-Cox transforms the exponent here
is a variable called λ that varies over
the range of -5 to 5 and in the process
of searching we examine all values of λ
finally we can choose the optimal value

value λ for your variable

$$\lambda \geq 0$$

Yeo-Johnson transform is somewhat of
an adjustment of Box-Cox transform, by
which we can apply it to -ve numbers

Box-Cox

$$x_\lambda = \frac{x^\lambda - 1}{\lambda}$$

$$x_\lambda = \frac{x^\lambda - 1}{\lambda}$$

Yeo-Johnson

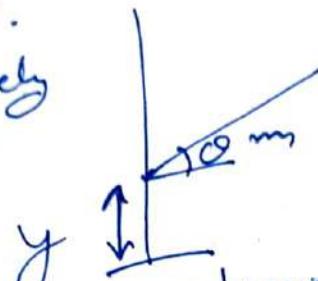
$$\psi(y, \lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \log y & \lambda = 0 \end{cases}$$

Simple linear regression model
 completely linear / sort of linear easy to learn
 Line equation $y = mx + b$ How to find m & b

I'm sort of linear draw a line and called that line best fit line.
 i.e. line which closely pass through all the points

$m = \text{slope}$

$b = y\text{ intercept}$



$x = \text{cgpa}$

$y = \text{package}$

plt. scatter(df['cgpa'], df['package'])

$x = df.iloc[:, 0]$

$y = df.iloc[:, -1]$

- lr. predict(x-test).iloc[0].values.reshape(1, 1))
 . importance of b is that in case value of m, b

direct

closed form
soln

non closed form by

(Gradient descent)
G.D

OLS ordinary least sq.

$$b = \bar{y} - m \bar{x}$$

$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

\bar{x}, \bar{y} mean

why modulus is not used instead of square

$$E = |a_1| + |a_2| + |a_3| + \dots + |a_n|$$

$$E = |a_1|^2 + |a_2|^2 + |a_3|^2 + \dots + |a_n|^2$$

$$E = \sum_{i=1}^n |a_i|^2 \text{ or } \sum_{i=1}^n d_i^2$$

take m , and b such that value of E is minimum

$$d_i = (y_i - \hat{y}_i)$$

$$E = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{ Average error}$$

$$E = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{ Total average}$$

$$\hat{y}_i = mx_i + b$$

$$E(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2$$

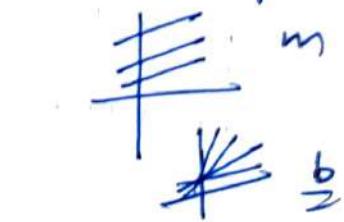
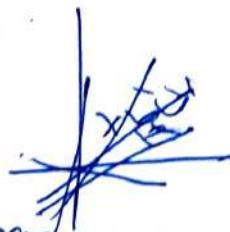
To find the values of m , and b which give

$$\text{minimum } E(m, b)$$

$$y = f(x)$$

$$E(m) = \sum_{i=1}^n (y_i - mx_i)^2$$

~~b - below the chay~~



$$\frac{\partial E}{\partial b} = \frac{\partial}{\partial b} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0$$

$$\sum \frac{\partial}{\partial b} (y_i - mx_i - b)^2 = 0$$

made

~~at origin~~
~~denied~~

use chain rule

$$\sum -2(y_i - mx_i - b) = 0$$

divide w.r.t. y

$$\sum (y_i - mx_i - b) = 0$$

$$\begin{cases} \sum y_i - \sum mx_i - \sum b = 0 \\ \text{divide w.r.t. } b \end{cases}$$
$$\frac{\sum y_i}{n} = \frac{\sum mx_i}{n} - \frac{\sum b}{n} = 0$$
$$\bar{y} = \bar{mx} - \frac{\bar{b}}{n} = 0$$
$$\bar{y} - \bar{mx} = b$$
$$b = \bar{y} - \bar{mx}$$

$$E = \sum (y_i - mx_i - \bar{y} + \bar{mx})^2$$

$$\frac{\partial E}{\partial m} = \sum \frac{\partial}{\partial m} (y_i - mx_i - \bar{y} + \bar{mx})^2 = 0$$

$$\Rightarrow \sum (y_i - mx_i - \bar{y} + \bar{mx}) (-x_i + \bar{x}) = 0$$

$$= \sum (y_i - mx_i - \bar{y} + \bar{mx})(x_i - \bar{x}) = 0$$

divide by -2

$$\sum (y_i - mx_i - \bar{y} + \bar{mx})(x_i - \bar{x}) = 0$$

$$\cancel{\sum} \cdot \sum (y_i - \bar{y})(x_i - \bar{x})(x_i - \bar{x}) = 0$$

$$= \sum [(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2] = 0$$

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

$$\begin{cases} m = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ \sum (x_i - \bar{x})^2 \end{cases}$$

maxima minima
error function

$$\frac{\partial E}{\partial x} = 0$$

$$\frac{\partial E}{\partial m} = \frac{\partial E}{\partial b} = 0$$

Regression Metrics

- (1) MAE
- (2) MSE
- (3) RMSE
- (4) R² score
- (5) Adjusted R² score

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$RMSE = \sqrt{MSE}$$

$$R^2 \text{ score}$$

~~$$1 - \frac{SSR}{SSM}$$~~

MAE ~~different function's~~ not differentiable

advantage at zero

disadvantage same curve

robust to outliers ~~not differentiable~~

SSR - sum of sq. error in regression

SSM - sum of sq. errors in mean line

~~Shows R² without regard to the extent of impact of features~~

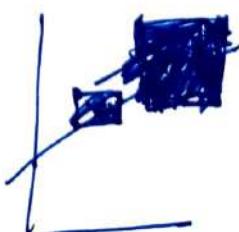
$$\text{Adjusted } R^2_{\text{sq}} = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_{\text{Reg}}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right]_{\text{mean}}}$$

Adjusted R² sq

$$R^2_{\text{Adj}} = 1 - \left[\frac{(1-R^2)(n-1)}{(n-k-1)} \right]$$

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

~~Adjusted loss function's difference~~
diff. $\propto (Y - \hat{Y})^2$ ~~intuitively~~



Multiple linear regression

SLR y^R Mult. P_{LR}
 Polynomiol

$x_1, \{x_2\} x_3 | y$ CGPA | LPA.

CGPA / Grade is LPA

Just extension
of linear
regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$LD = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

$$y = \beta_0 + \beta_1 x_1$$

$$3D: y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\begin{aligned} x_1 &= CGPA \\ x_2 &= IGP \\ y &= LPA \end{aligned}$$

$$LPA = \beta_0 + \beta_1 \times CGPA + \beta_2 \times IGP$$

these coefficients are like weights
then value + become

if β_2 is very small
very small of IGP.

Lecture 5: Multiple Linear Regression Part 2

How find $\beta_0, \beta_1, \beta_2, \dots, \beta_n$?

eg

c gpa	19	gender / female	
x_1	x_2	x_3	y actual
40 dots			

$$y = mx + b$$

$$\hat{y} = \beta_0 + \beta_1 x \quad \text{Predicted}$$

100 student

$$y_1 = \beta_0 + \beta_1 x_{11}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_{100} \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \beta_3 x_{23} \\ \vdots \\ \beta_0 + \beta_1 x_{1001} + \beta_2 x_{1002} + \beta_3 x_{1003} \end{bmatrix}$$

assume 100 rows $\rightarrow n$ rows y column \rightarrow col $\rightarrow m$ cols

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} + \dots + \beta_m x_{1m} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \beta_3 x_{23} + \dots + \beta_m x_{2m} \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \beta_3 x_{n3} + \dots + \beta_m x_{nm} \end{bmatrix}$$

rule

xwrite matrix

$$= \boxed{\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & & & & \\ \vdots & & & & \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix}} \quad \boxed{\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}}$$

$$A\beta = AB$$

$$\hat{Y} = X\beta \quad \cancel{\text{if}}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad e = Y - \hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$e = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

single LR

$$e = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$E = e^T e \quad \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix} = e$$

$$[(y_1 - \hat{y}_1) (y_2 - \hat{y}_2) (y_3 - \hat{y}_3) \dots (y_n - \hat{y}_n)]$$

$$\begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ y_3 - \hat{y}_3 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

$$(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + \dots + (y_n - \hat{y}_n)^2$$

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

$$E = e^T e = (y - \hat{y})^T (y - \hat{y})$$

$$= (y^T - \hat{y}^T)(y - \hat{y})$$

$$= [y^T \quad (\alpha \beta)^T] (y - \hat{y})$$

$$E = y^T y - \underbrace{y^T x \alpha}_{=} - \underbrace{(\alpha \beta)^T y}_{=} + (\alpha \beta)^T x \beta$$

$$\text{Detour } y^T x \alpha = (x \alpha)^T y \quad \text{Proof}$$

$$y = A \xrightarrow{\cancel{x}} \alpha \beta = \beta$$

$$A^T \alpha = \cancel{\beta} B^T A$$

$$\begin{aligned} (A+B)^T &= A^T + B^T \\ (A-B)^T &= A^T - B^T \end{aligned}$$

$$(A \ B)^T = B^T A^T, \quad (A + B)^T = B^T A^T$$

$$A^T B = R^T A \quad \xrightarrow{\text{①}} \quad T \text{ is bijective} \\ A^T B = (A^T D)^T \rightarrow A^T B = C \quad \boxed{C = C^T}$$

$$Y = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad Y^T = (1 \times 4)$$

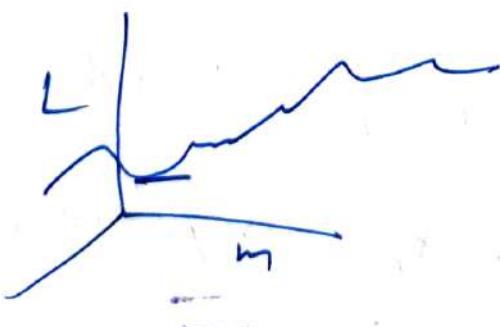
~~-1 x (m/s)~~

$$m \times 1 \rightarrow (1 \times n) \left(\cup \times (m+1) \right) [(m+1) \times 1]$$

$$1 \times (m+1) \quad (m+1) \times 1$$

$$1 \times 1 = [7]^T = [7]$$

$$E = \gamma^\top \gamma - 2\gamma^\top x_\beta + \beta^\top x^+ x \beta$$



$$\frac{d e}{d \beta} = \frac{d}{d \beta} \left[\underbrace{\gamma^T \gamma - 2 \gamma^T x \beta + \beta^T x^T x \beta}_{\geq 0} \right]$$

$$= 0 - 2y^T x + \frac{d}{d\beta} \left[\beta^T x^T P \right] = 0$$

matrix differentiation

$$Y = A^T x A \quad \frac{dy}{d\beta} = 2x A^T$$

$$= 0 - 2y^T x + 2x^T x \beta^T = 0$$

$$= \cancel{x^T x} \beta^T = \cancel{2} y^T x$$

$$\beta^T = \frac{1}{y^T x} \quad \beta^T = y^T x (x^T x)^{-1}$$

$$(\beta^T)^T = [y^T x (x^T x)^{-1}]^T$$

$$\beta = [(x^T x)^{-1}]^T (y^T x)^T$$

$$\beta = [(x^T x)^{-1}]^T x^T y \quad \begin{cases} x = n \times (n+1) \\ (m+1) \times m / m \times (m+1) \end{cases}$$

$$\beta = (x^T x)^{-1} x^T y \quad (m+1) \times (m+1)$$

$$(A \cdot B)^T = B^T A^T, \quad (A^T B)^T = \underline{\underline{B^T A}}$$

$$A^T B = B^T A \quad \text{to be proved}$$

$$A^T B = (A^T B)^T$$

$$\rightarrow A^T B = C$$

$$b = n \times m$$

$$(y^T x \beta)^T = y^T x \underline{\underline{\beta}}$$

$$\boxed{\underline{\underline{C}}^T = C}$$

$$m \times 1 \rightarrow (\cancel{1} \times \cancel{m}) (\cancel{1} \times (m+1)) \underbrace{[(m+1) \times 1]}_{((m+1) \times 1)}$$

$$\underbrace{[1 \times (m+1)]}_{((m+1) \times 1)}$$

Agrahi L-13 Vanishing gradient

Sum comp. Margin maximizing hyperplane

max(Θ)

$$w^T x + b$$

\uparrow

π

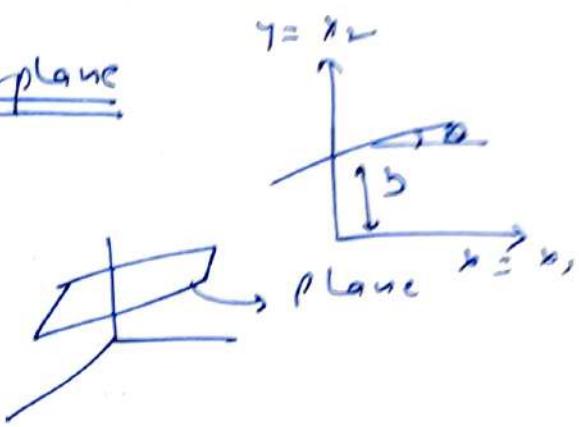
1. very robust to outliers

2. It can also work on non-linear data by using Kernel

3. Applicable for both classification and regression problem

Equation of a Hyperplane

4D, 5D, nD
Hyperplane



for 2D :- $y = mx + b$

↓ general form

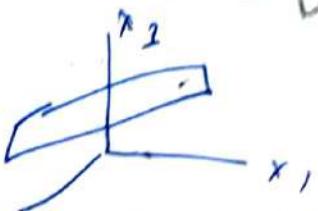
$$ax_1 + bx_2 + c = 0$$

$$w_1x_1 + w_2x_2 + w_0 = 0$$

$$\left| \begin{array}{l} y = -\frac{a}{b}x - \frac{c}{b} \\ m = -\frac{a}{b}, b = -\frac{c}{b} \end{array} \right.$$

My

for 2D



$$w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_0 = 0$$

for Hyperplane in dimension:-

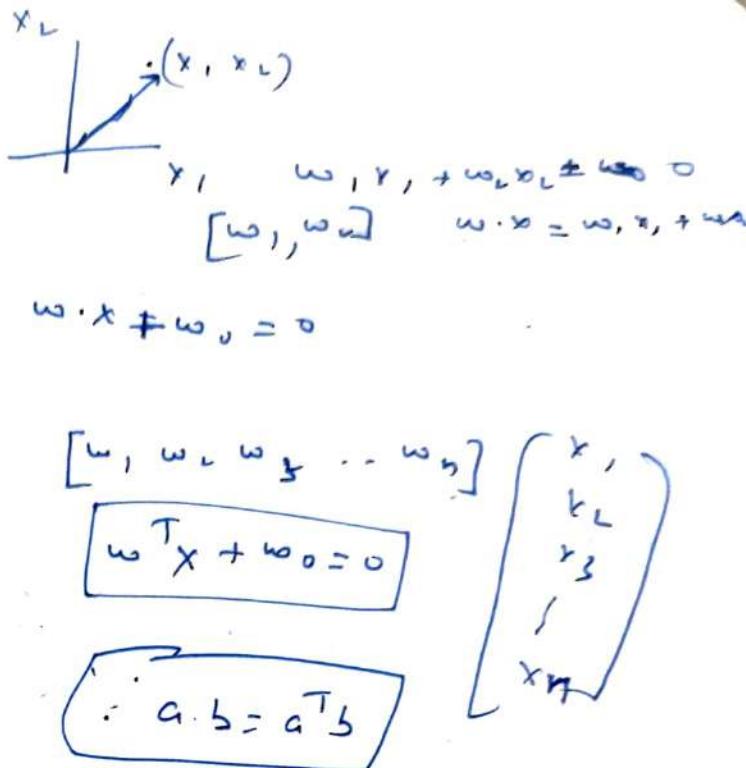
$$(w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + \dots + w_nx_n + w_0 = 0)$$

π

$$\omega = [\omega_1, \omega_2, \dots, \omega_n]$$

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

$$\begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$



$$\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0$$

$$b = -\frac{\omega_0}{\omega_2}$$

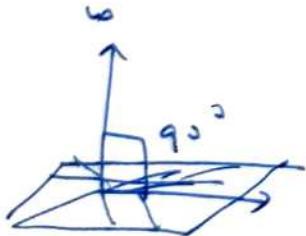
$$\boxed{\omega^T x = 0}$$

$$\omega_1 x_1 + \omega_2 x_2 = 0$$

$$\omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n = 0$$

$$\omega^T x = \omega \cdot x = \|\omega\| \|\omega\| \cos \theta = 0$$

$$\text{when } \theta = 90^\circ$$



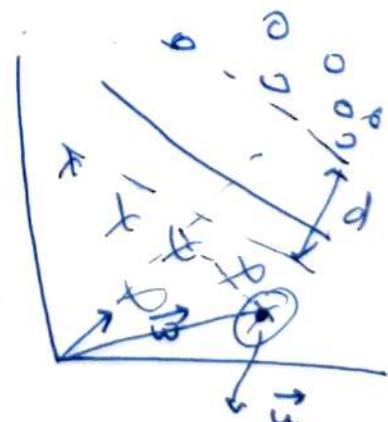
Mathematics for SVM

Decision Rule

$$\vec{\omega} \cdot \vec{u} \geq c$$

$$\vec{\omega} \cdot \vec{u} - c \geq 0$$

$$\vec{\omega} \cdot \vec{u} + b \geq 0$$



+ve/-ve point

$$\vec{\omega} = (2, 3)$$

$$2x + 3y + 3 = 0$$

$$2(3) + 3(-1) + 3 \geq 0$$

$$\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0$$

$$\vec{\omega} = (-2, 0)$$

$$\vec{\omega} = (2, 1) \quad \vec{u} = (3, 2)$$

$$2(-1) + 3(0) + 1 \leq 0$$

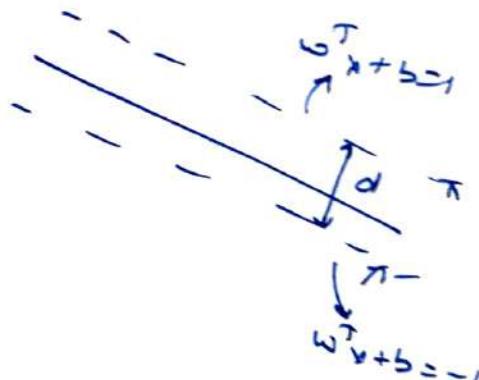
$$-1 \leq 0$$

for any \vec{x}_i :

$$\gamma_i = \begin{cases} +1 & \text{if } \vec{\omega} \cdot \vec{x}_i + b \geq 0 \\ -1 & \text{if } \vec{\omega} \cdot \vec{x}_i + b < 0 \end{cases}$$

1. why?
2. why equal
3. Some other number

constraints



$$\gamma_i (\vec{\omega} \cdot \vec{x}_i + b) \geq 1$$

for support vector

$$\gamma_i (\vec{\omega} \cdot \vec{x}_i + b) = 1$$

$$d = (\vec{x}_v - \vec{x}_i) \cdot \frac{\vec{\omega}}{\|\vec{\omega}\|}$$

$$= \frac{\vec{x}_v \vec{\omega} - \vec{x}_i \vec{\omega}}{\|\vec{\omega}\|}$$

$$\gamma_i (\vec{\omega} \cdot \vec{x}_i + b) = 1$$

$$1 (\vec{\omega} \cdot \vec{x}_i + b) = 1$$

$$\vec{w} \cdot \vec{x}_i = (1 - b) \text{ or } -1$$

$$-1 (\vec{w} \cdot \vec{x}_i + b) = 1$$

$$\vec{w} \cdot \vec{x}_i = -b - 1$$

$$\frac{(1 - b) - (-b - 1)}{\|\vec{w}\|}$$

$$= \frac{1 - b + b + 1}{\|\vec{w}\|}$$

argmax $\frac{1}{2} \frac{\|\vec{w}\|^2}{(\vec{w}^*, b^*)}$ such that $y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$

Hard margin SVM
HARD HORN

Soft Margin SVM

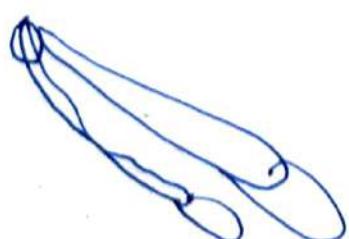
$$\text{argmin}_{(\vec{w}^*, b^*)} \frac{\|\vec{w}\|^2}{2} + C \sum_{i=1}^n \xi_i \quad \begin{matrix} \xrightarrow{\text{zeta}} \\ \text{hinge loss} \end{matrix} \|\vec{w}\|$$

margin error + classification error

logistic loss $+ \lambda \|\vec{w}\|$

$\xi_i = 0$
for correctly classified pt.

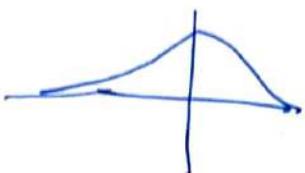
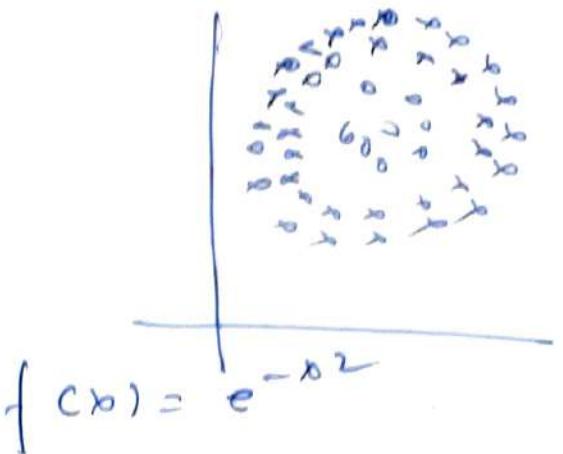
$$C \propto \frac{1}{\lambda}$$



Kernel tricks sum

— 0 0 0 0 0 0 0 0 0 0 0 0 —

$\left\{ \begin{array}{l} RBF \\ \text{Polynomial} \\ \text{Sigmoid} \end{array} \right.$



Grid search —

Ensemble learning

$$L = \sum$$

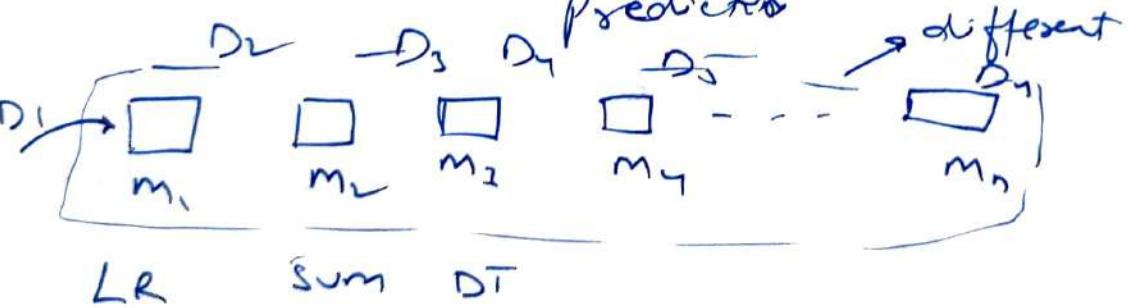
some diff
 same diff
 same abs
 diff abs
 diff abs

Wisdom of crowd

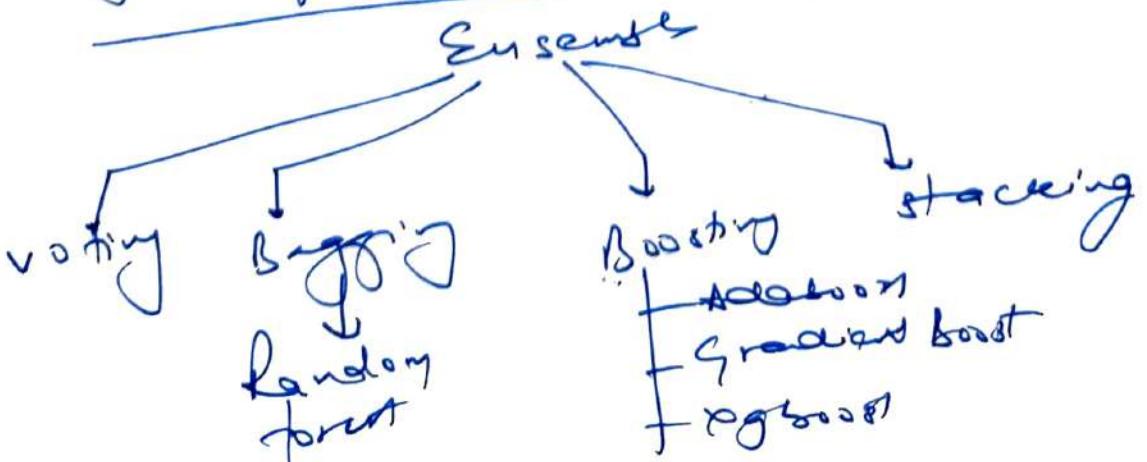
Core idea

Training

Predicts

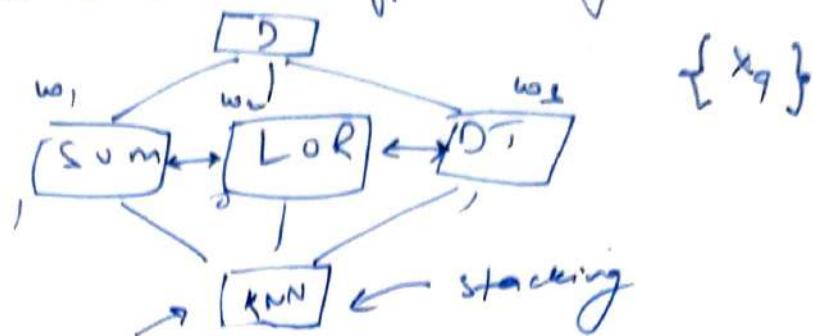


Types of Ensemble learning

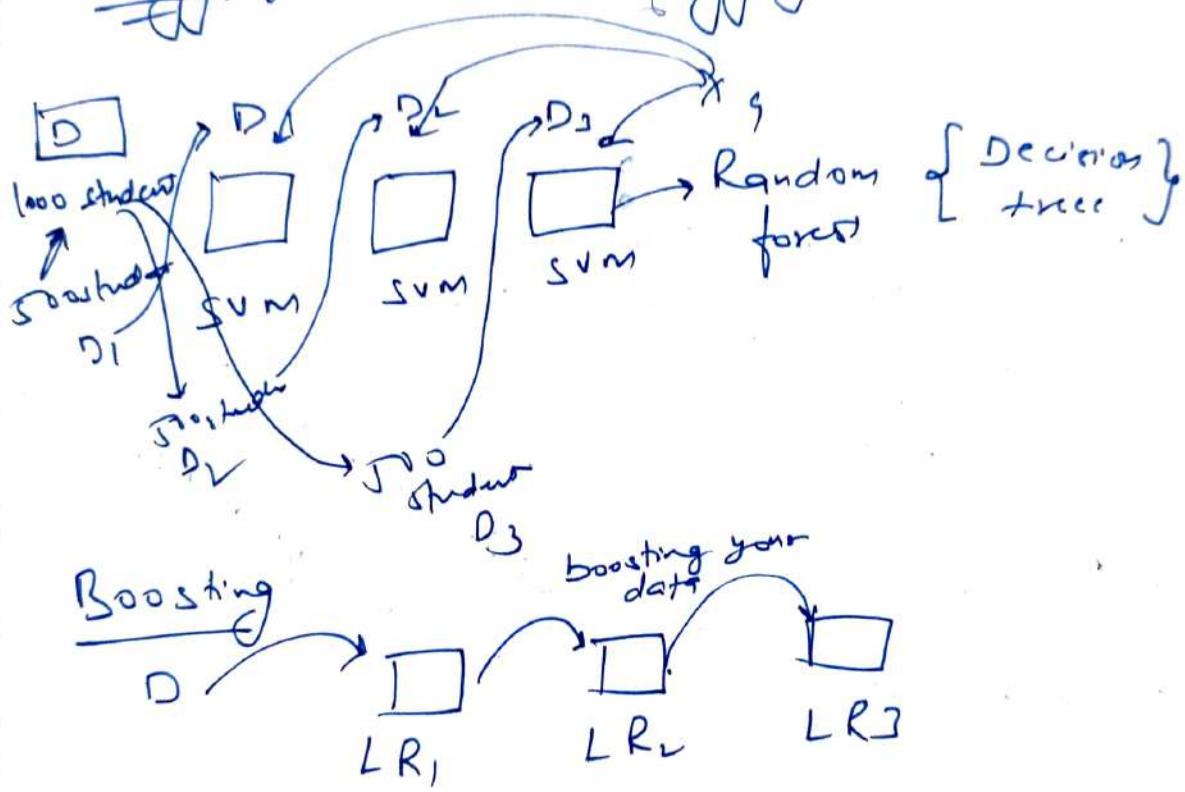


Voting

→ Base model are different algorithms



Bagging → Bootstrap aggregation



- computation cost expensive

Adv. 1. Improvement in performance

2. Bias - Variance reduce

3. $\sqrt{\text{low Bias} + \text{low variance}}$

3. Robustness.

(4) when to use :- Always use ensemble learning

When to use
Always

Voting Ensemble

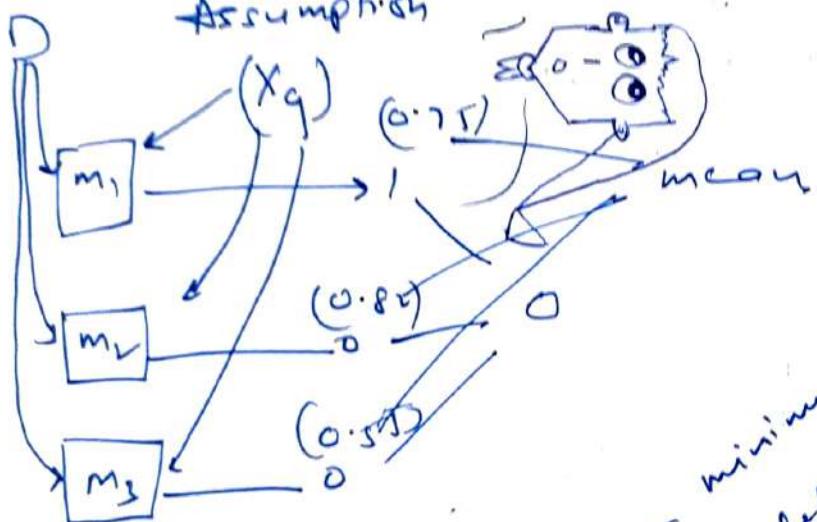
- core / idea / intuition
- Voting classifier
- Voting Regressor

Part 1 - Introduction

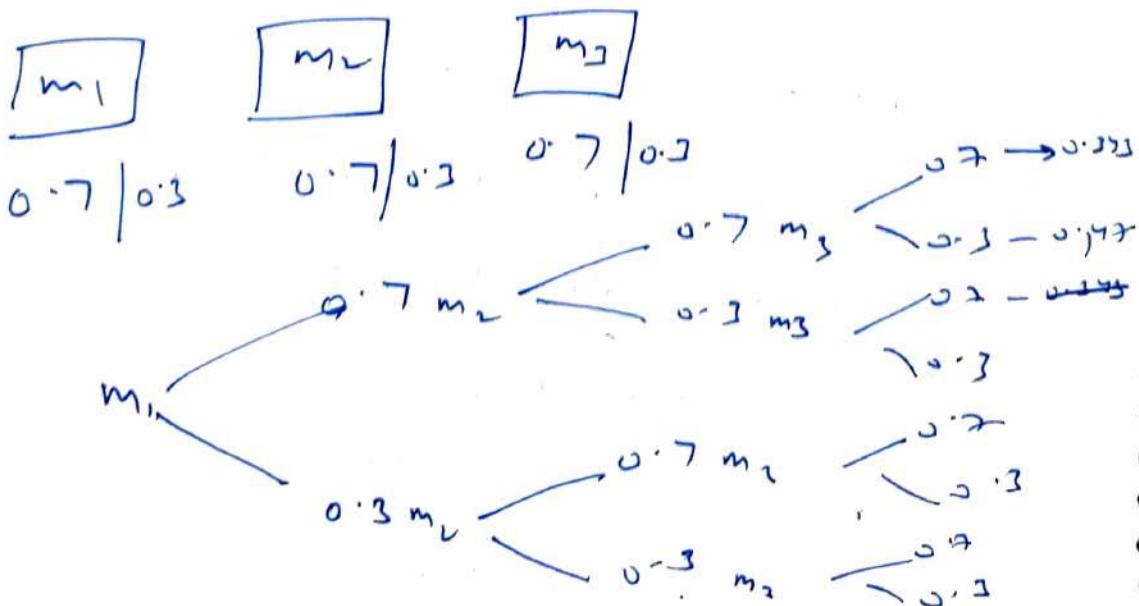
core idea / intuition

why voting works

Assumption



why voting works

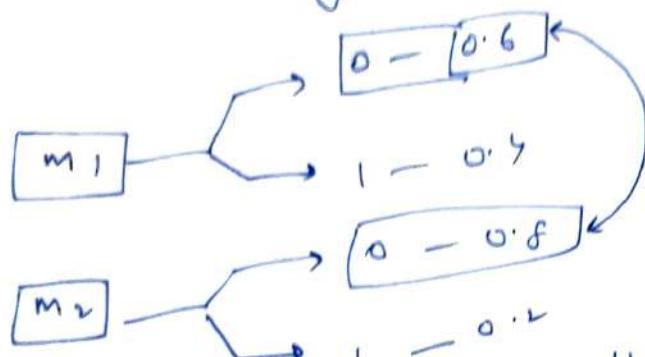


Classifier

- core idea / Demo
- Hard voting vs soft voting

code

- Hyperparameters



0 - 0.7
1 - 0.3

Average calculation

* run both hard and soft voting and select one which is better

Part 3 - Voting Regressor

- core idea + Demo

- code

- Hyperparameters

Bagging Ensemble learning

core idea /

why use Bagging

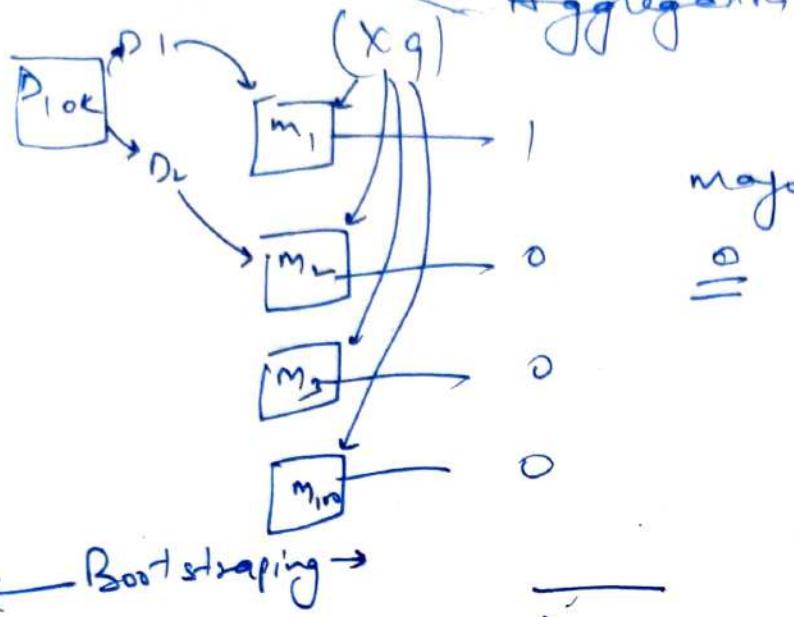
when to use Bagging

code demo

Types of Bagging

Bagging → Bootstrap

Aggregation



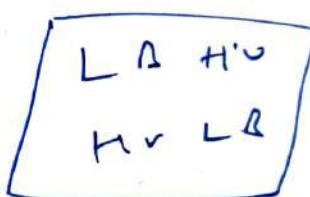
+ same data
different algos
- different data
same algo

majority vote

\oplus

change data
with different
subset
either with
replacement or
without replacement

low bias low variance
LB LU



consistent
result

when to use? Every time

- Random forest
- DT - easy to understand

types

- ① Plasting (without replacement row sampling)
- ② Random subspace (do column sampling with or without replacement)
- ③ Random patches - cut both column and row up but

Part 2 Bagging classifier

OOD score
 (out-of-Bagging) local / n-pp

63% in 37% rows
 are not used

Part 3 Regressor

- Bagging generally give better results than pasting

Random forest

- Decision tree, bagging technique

Random forest

↑ ↗ gr. of trees

Bagging → Bootstrapped Aggregates

↳ sampling (Random)

$D = 1000 \text{ rows} \rightarrow$

$D_{11} \quad D_{12} \quad D_{13} \dots D_{1n}$

+ Row sampling
 + Column sampling
 + combination

Baggit
 Random forest
 is like Bagging
 but only difference
 is that in RF we
 use Decision tree

with replacement
 without replacement

Bootstrap
 Aggregation (Most frequent)

geom

- How Random forest performs so well?

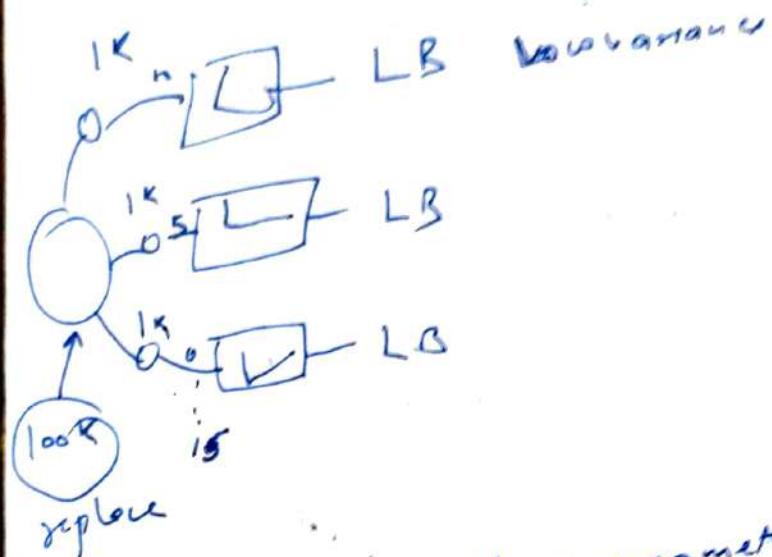
• variance trade-off

$\boxed{LB \quad LV}$

$\beta \propto \frac{1}{v}$

Bias

Random forest [LB+V] connected to
[LB+V]



Random forest Hyperparameters

sklearn.ensemble.RandomForestClassifier

max_samples

max_features

(CART) Decision tree hyperparameters

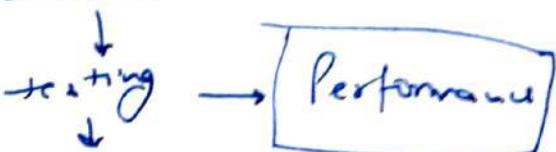
HPT Hyperparameter Tuning RF using GridSearchCV
and RandomizedSearchCV

oob Evaluation out of Bag

Sampling with replacement

out of bag sample :- which are not selected
in model training

oob hidden



377 rows
oob sample
mean validation
sample

feature importance using R and DT. How features
importance calculated?
feature column importance?

Feature importance

$$n_i = \frac{N-t}{2} \left[\text{impurity} - \left(\frac{N-t-\sigma}{N-t} \times \text{right-impurity} \right) \right. \\ \left. - \left(\frac{N-t-L}{N-t} \times \text{left-impurity} \right) \right]$$

How AdaBoost classifier works?

1. weak learners :- Just over 50% Accuracy
 2. Decision Stumps : max depth 1
 3. +1 and -1

stage wise,
additive
method]

$$A = \sin(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_3 h_3(x))$$

Ada Boost - step by step

$$\alpha = \frac{1}{K} \ln \left(\frac{1 - \text{error}}{\text{error}} \right)$$

$$\alpha_1 = \frac{1}{2} \ln \left(\frac{-0.6}{0.4} \right) = 0.20$$

\rightarrow new wt
of miss class
- pred pt
 \rightarrow new wt of
correct class
- pred pt

For misclassified

$$\text{new_wt} = \text{curr_wt} \times e^{\alpha}$$

For correctly classified

$$\text{new_wt} = \text{curr_wt} \times e^{-\alpha}$$

normalized the updated weights

x_1, x_2, y new-wt - range

Bagging Vs Boosting

- 1) Type of model used
- 2) Sequential vs Parallel
- 3) weightage of base learners

5 random

no. of wts used

0.13

0.43

0.62

0.50

0.8

Bagging

df.duplicated().sum() ① df.shape @ df.describe

⑤ df.info(), ④ df.sample(5), ③ df.isnull().sum(),
⑥ df.duplicated().sum() ⑦ df.corr()['Survived'] =

Univariate analysis) ① Categorical Data:

① Countplot sns.countplot(df['Embarked'])

② Piechart
df['sex'].valuecount().plot(kind='pie', autopct='%.2f')

③ Numerical data:

④ Histogram:- import matplotlib.pyplot as plt
plt.hist(df['Age'], bins=5)

⑤ Distplot sns.distplot(df['Age'])

⑥ Boxplot sns.boxplot(df['Age'])

df['Age'].min(); df['Age'].max()
df['Age'].mean(); df['Age'].skew()

Bivariate analysis) scatterplot(Numerical - Numerical)

sns.scatterplot(tips['total_bill'], tips['tip'],
hue=df['sex'], style=df['smoker'], size=df['size'])

Bar plot (Numerical - Categorical)

sns.barplot(titanic['Pclass'], titanic['Age'],
hue=titanic['sex'])

3. Box plot (Numerical - categorical)
sus. boxplot (titanic['sex'], titanic['Age'],
hue = titanic['survived'])
4. Distplot (Numerical - categorical)
sus. distplot (titanic[titanic['survived'] == 0]['Age'],
hist = False)
- sus. distplot (titanic[titanic['survived'] == 1]['Age'],
hist = False)
5. Heatmap (Categorical - Categorical)
sus. heatmap (pd.crosstab(titanic['Pclass'],
titanic['survived']))
- (titanic.groupby('Embarked').mean()['survived'])
- ⑥ Clustermap (Categorical - categorical)
sus. clustermap pd.crosstab(titanic['Pclass'],
titanic['survived'])
7. Pair Plot sus.pairplot (iris, hue = 'species')
- ⑧ Lineplot (Numerical - Numerical)
new = flights.groupby('year').sum().reset_index()
sus.lineplot (new['year'], new['passenger'])
- sus.clustermap (flights.pivot_table(values = 'passenger',
index = 'month', columns = 'year'))

```
from sklearn.preprocessing import OrdinalEncoder
```

```
o_e = OrdinalEncoder(categories=[['no', 'Acc', 'Low',  
                                 'School', '4G', 'PLT']])
```

```
X_train = o_e.transform(X_train)
```

(1) OneHotEncoding using pandas

```
Pd.get_dummies(df, columns=['fuel'], drop_first=True)
```

(2) K-1 OneHotEncoder

```
Pd.get_dummies(df, columns=['fuel'], drop_first=True)
```

```
header=1, cols=1, squeeze=True, skiprows=[0, 1],  
nrows=100, encoding parameter, encoding='Latin-1'  
skip bad line:- , sep='.', encoding='Latin-1',  
error_bad_line=False  
dtype parameter, dtype={'target': int})
```

Handling Data parse_date=['date'])

Converters: def rename(name):

```
def rename(name):  
    if name == "Fogel Chilling":  
        return "FCB"  
    else: return name
```

```
else: return name
```

```
converter = {'Iteam1': rename}
```

15. na_value parameters
, na_values = ['Male']

16. Loading a huge dataset in chunks
(chunksize > 5000)

for chunk in df:
print(chunk.shape)

JSON / SQL

Java script on notation pd.read_json()

working with SQL

- Download file → Xampp control panel

! pip install mysql-connector

import mysql.connector

How big data is?

How does the data look like?

What is the data type of cols?

Are there all missing values

How does the data looks mathematically

Are there duplicate values?

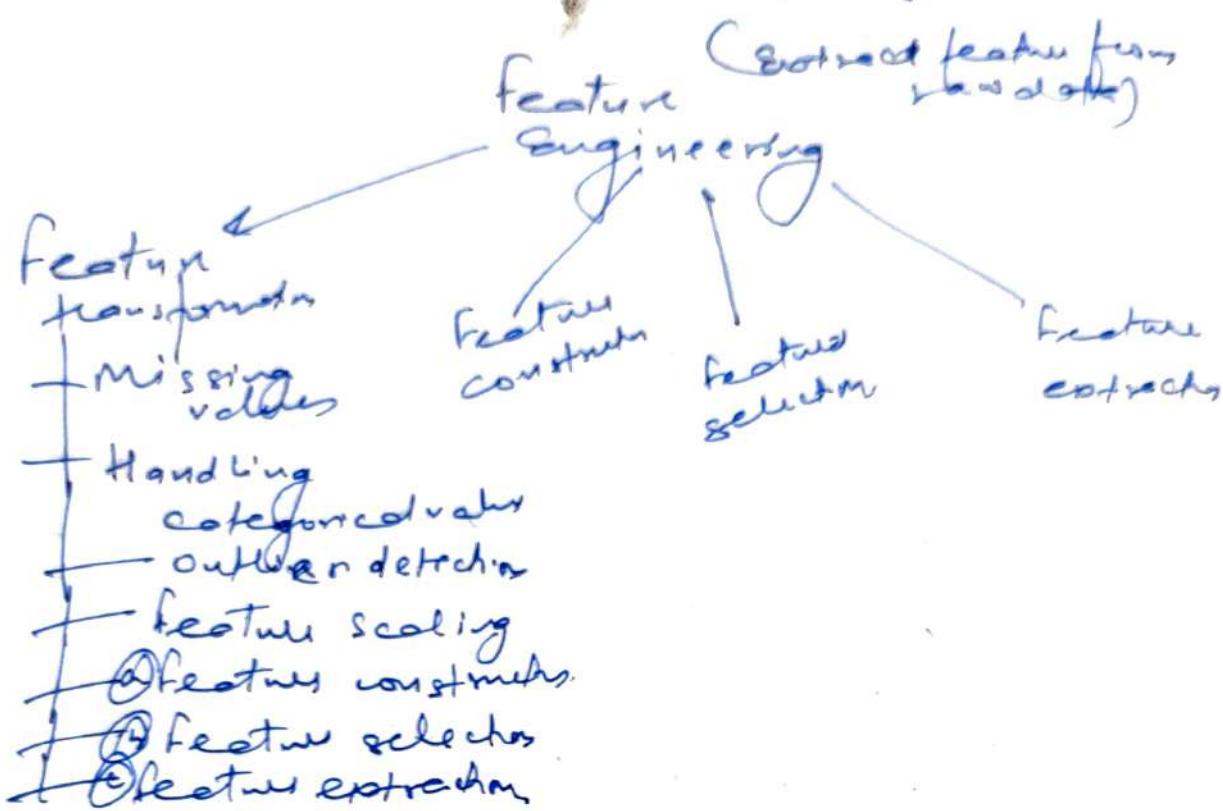
df.duplicated().sum()

How is the correlation b/w w/b

df.corr()['survived']

Multivariate Analysis

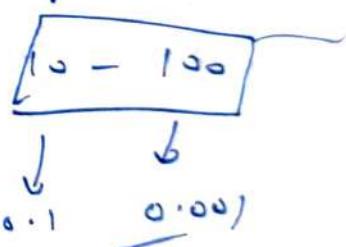
`hue = df['sex'], style = df['smoker'], size = df['size']`
`pd.crosstab(titanic['Pclass'], titanic['survived'])`
`sns.pairplot(titanic, hue = 'species')`



8837785080

feature scaling to standardize the independent feature present in the data in a fixed range

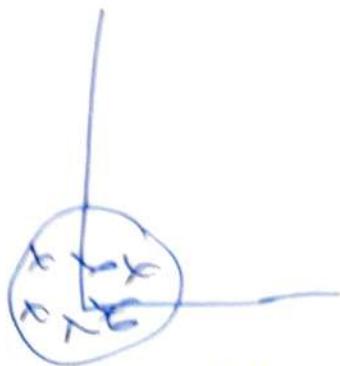
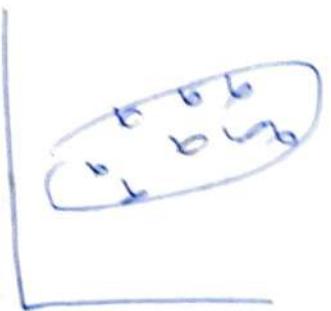
Similar range e.g.



feature scaling

standardization
(Z-score normalization)
minmax scaling
Robust scaling

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



- ④ Denses
- ⑤ Nearest neighbour
- ⑥ PCA
- ⑦ A ∈ ℝ
- ⑧ Gradient Descent

Normalization / Min-Max Scaling

Technique often applied as part of data prep. for ML. The goal of normalization is to change the value of the numeric columns in the data set to use a common scale, without distorting differences in the ranges of values or losing information.

- + min max scaling
- + mean normalization
- + max absolute
- + robust scaling

weights

130

67

81

01

22

57

11

normalization
min-max scaling

mean centering

$$x_i = \frac{y_i - \bar{y}_{\text{mean}}}{\text{range} - \text{range}}$$

range [0, 1]

$$y_i = \frac{y_i - \bar{y}_{\text{mean}}}{y_{\text{max}} - y_{\text{min}}}$$

Max absolute $\|x\|_1 = \sum |x_i|$



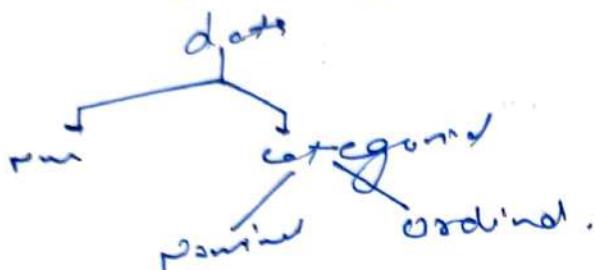
Robust scaling $x'_i = \frac{x_i - \text{mean}}{\text{std}}$

Normalizations vs standardization IQR $\left[71^m, 25^m\right]$

1) Robust location

ordinal Encoding | Label Encoding

Encoding categorical variables



Label encoder use for only target values not for input values



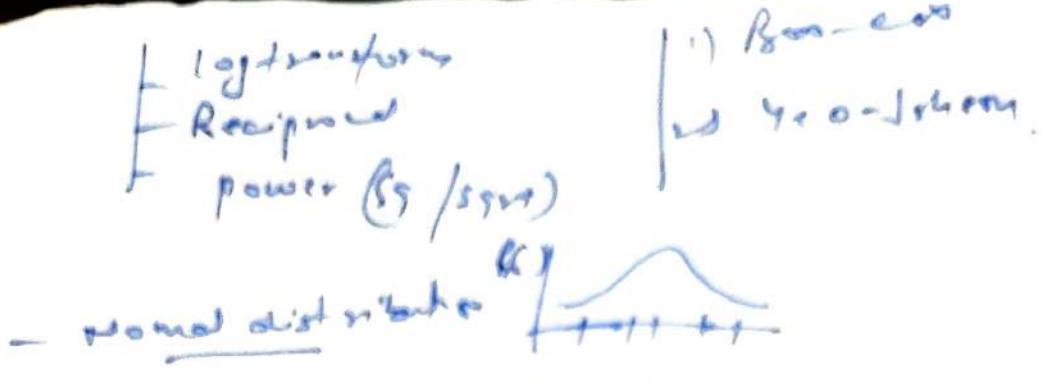
One Hot Encoding

Dummy variables use $(n-1)$ instead of (n)

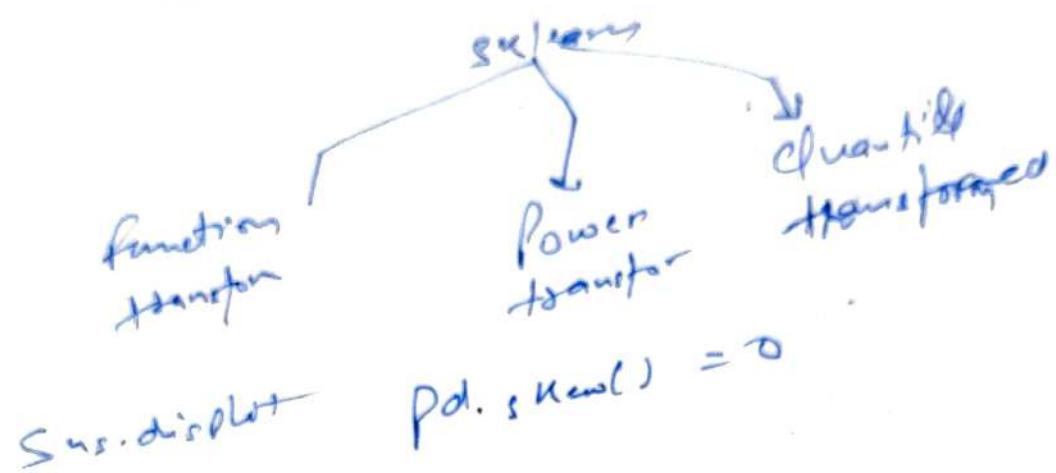
p.d. get_dummies (df, columns=[], drop=True)

Pipelines chains together multiple steps so that the output of each step is used as input to the next step

Pipelines make it easy to apply the same preprocessing to train and test



$$\begin{cases} 1) \text{ Box-Cox} \\ 2) \text{ Yeo-Johnson} \end{cases}$$



Box-Cox transform

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

$n < 0$

The exponent here is a variable called (λ) that varies over the range of -5 to 5, and the process of searching we examine all values of λ . Finally we choose the optimal value (resulting in the best approximation to the normal distribution) for your variable.

Yeo-Johnson transform

$$x_i^{(\lambda)} = \begin{cases} [(x_{i+1})^{\lambda-1}]^\lambda & \lambda \\ \ln(x_i) + 1 & -[(c-x_{i+1})^2 - 1] \end{cases}$$

~~log~~ :- big data, Right skewed data

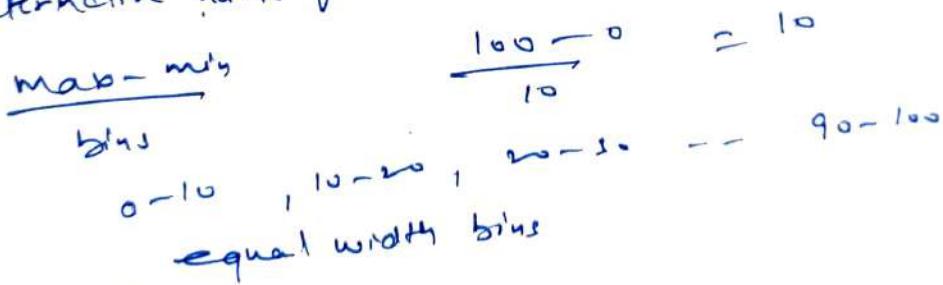
Reciprocal $\frac{1}{x}$	$\sqrt{x^4}$	\sqrt{x}	\log
Left skewed skewed			Right skewed

$$X_{ij} = \begin{cases} \frac{x_i - 1}{x} & \text{if } \lambda \neq 0 \\ \ln(x_i) & \text{if } \lambda = 0 \end{cases}$$

The exponent here is variable called Lambda (λ) that varies over the range -5 to 5, and in the process of searching we examine all values of λ . Finally we choose the optimal value for your variable.

good

~~Discretization~~ i.e. the process of transforming continuous variables into discrete variables by creating a set of contiguous intervals that span the range of the variable values. Discretization is also binning, where bin is an alternative name for interval



- 1) outlier
- 2) spread of data

⑤ equal frequency / quantile Binning

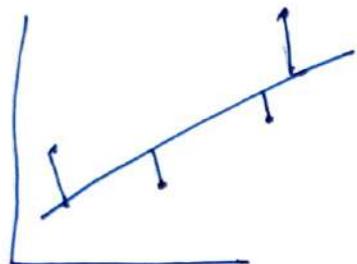
0-16, 16-22, 22-28, 28-34, 34-40

— outlier
— value spread uniform
— Binning clustering

⑥ K means

Gradient descent

- Linear Regression
- Logistic Regression
- Deep learning ...



To find best fit line

Context in Brief

$$\begin{array}{c|c} \text{2 columns} & \text{4 Rows} \\ \text{cgpa} & | \quad L^T \\ \hline = & = \\ = & = \\ = & = \end{array}$$

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

$$g_i = mx_i + b$$

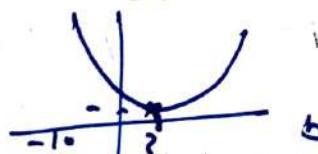
$$L = \sum_{i=1}^n (y_i - (mx_i + b))^2$$

$L(m, b)$ $\not\rightarrow$ depend on m and b

$$m = \frac{78.25}{4} \text{ we know}$$

$$L = \sum_{i=1}^n (y_i - 78.25x_i - b)^2$$

$$L \rightarrow b^2$$

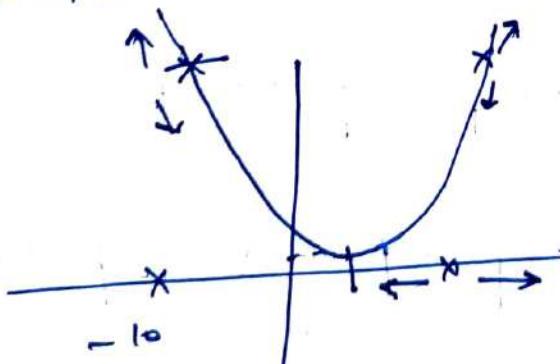


OLS

Step 1/1: Select a Random b

$$\text{Let } b = -10$$

$$L(b)$$



* Find out the slope at a particular point.

Put by derivative and then put value of x .

$$\begin{aligned} y &= x^2 + 2 & x &= 5 \\ \frac{dy}{dx} &= 2x & = & 2(5) \\ & & & = 10 \end{aligned}$$

- If slope is -ve b increased
- If slope is +ve b decreased

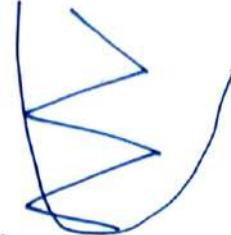
$$b_{\text{new}} = b_{\text{old}} - \eta \text{ slope}$$

η = learning rate

$b_{\text{new}} = b_{\text{old}} - \eta \text{ slope}$

$\eta = 0.01, \text{ or } 0.1, 0.001, 0.0001$

If η is higher or 1 than
this type of change



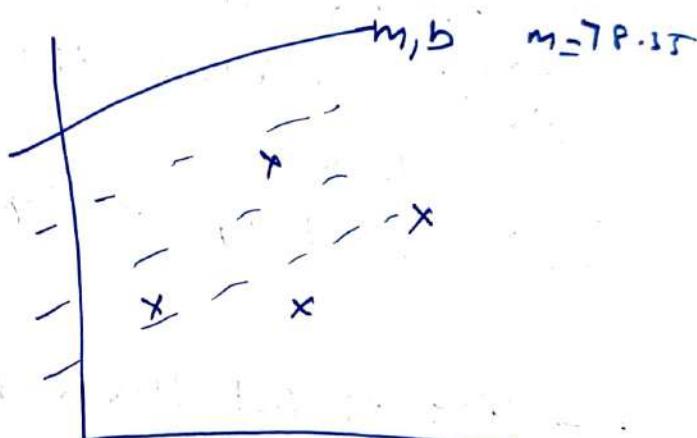
When to stop :-

$$b_{\text{new}} - b_{\text{old}} = 0.0001$$

1) > 0.0001

2) Iterations = 1000, 100
epochs →

Mathematical formulations



Step 1 :- is always
start with random values
if $b = b$
for i in epochs → 1000

$$b_{\text{new}} = b_{\text{old}} - \eta \text{ slope}$$

$\boxed{\eta = 0.01}$

$$L = \sum_{i=1}^n (y_i - \bar{y}_i)^2 \quad \text{Find slope } b = b ?$$

$$b = 0$$

$$\frac{dL}{db} = \frac{d}{db} \left(\sum_{i=1}^n (y_i - \bar{y}_i)^2 \right) =$$

$$\frac{d}{db} \sum_{i=1}^n (y_i - mx_i - b)^2 = -2 \sum_{i=1}^n (y_i - mx_i - b)$$



$$b = 0$$

$$= -2 \sum_{i=1}^n (y_i - 78.25 * x_i - 0)$$

slope($b = 0$)

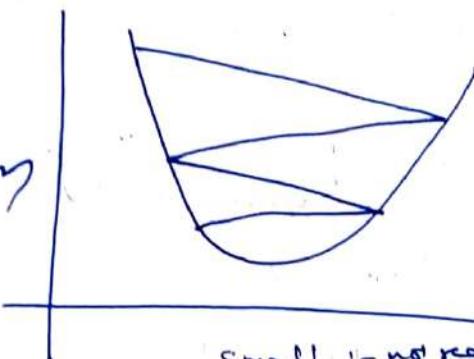
$$b_{\text{new}} = b_{\text{old}} - \eta \text{ slope}_{b=\text{old}}$$

$$\circlearrowleft b_{\text{new}} \quad i=1$$

Visualisation

↓ learning rate η

↑ epochs.

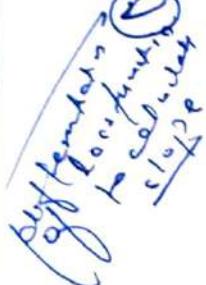


Few discussion

① Learning rate

- small : - not reaches the end
- large : - bypass the convergence

② The Universality of Gradient descent



$$b = 0$$
$$b = b_{\text{old}} - \eta \text{ slope} \quad \frac{\partial L}{\partial b} = \boxed{\boxed{\frac{\sum (y_i - f_i)^2}{LR}}}$$

Adding m into the mix

LOR \rightarrow L family

step 1 :- Initialize random values for m and b

$$m = 1$$

$$b = 0$$

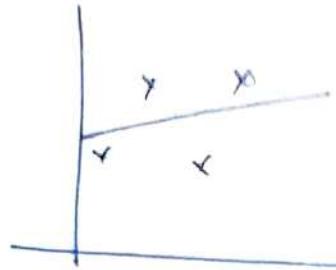
$$\text{Epoch} = 10, \text{LR} = 0.01$$

for i in epochs:

$$b = b - \eta \text{ slope}$$

$$m = m - \eta \text{ slope}$$

step 2

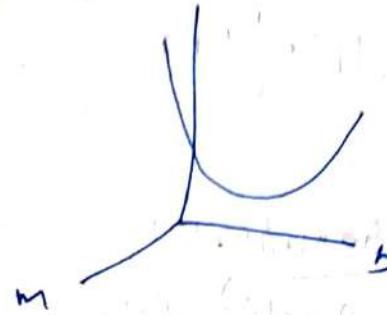


$$L = \sum (y_i - \hat{y}_i)^2$$

$$L = \sum (y_i - mx_i - b)^2$$

$L(m, b)$

$$\frac{\partial L}{\partial b} \quad \frac{\partial L}{\partial m}$$



$$\sum (y_i - mx_i - b)^2$$

$$\frac{\partial L}{\partial b} = \sum (y_i - mx_i - b)$$

$$= -2 \sum (y_i - mx_i - b)$$

$$= \text{slope of } -b \text{ at } b = 0$$

$$\frac{\partial L}{\partial m} = 2 \sum (y_i - mx_i - b)$$

$$\frac{\partial L}{\partial m} = -2 \sum (y_i - mx_i - b) x_i$$

slope = $m = -\frac{1}{2} \sum (y_i - mx_i - b) x_i$

Property
scale-free data

Type of gradient descent

- ① Batch
- ② stochastic
- ③ mini batch

$$m = 1, b = 0$$

$$m_n = m_0 - \eta \times (\text{Slope})_{m=0}$$

$$b_n = b_0 - \eta \times (\text{Slope})_{b=0}$$

$$\frac{\partial L}{\partial m} \quad \frac{\partial L}{\partial b}$$

- use whole data, eg 300 rows monitoring and then give decision
- stochastic: - only on single row, fast, error prone
- minibatch in b/w batch and stochastic
use predefined batch size

$$\text{cgp} \mid \text{LPA} \quad \boxed{\begin{matrix} \text{cgp} \\ \text{LPA} \end{matrix}} \quad (\bar{x}, \bar{y})$$

cgp = 1/2 / gender / LPS

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Mathematical formulation

n-dim dataset 3-columns

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

LPA (cgp) (ig)

cgp	ig	LPS
8.1	93	12
7.5	95	1.5

(2, 2)

① Random Values

$$\beta_0 = 0, \beta_1, \beta_2 = 1.$$

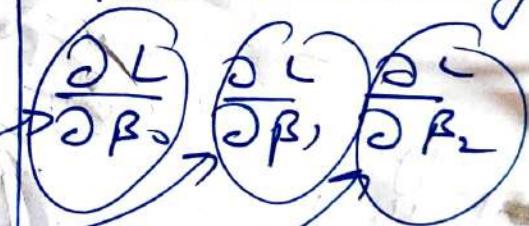
$$\text{epochs} = 100, \Delta = 0.1$$

$$\beta_0 = f_0 - \eta \text{ slope}$$

$$\beta_1 = f_1 - \eta \text{ slope}$$

$$\beta_2 = f_2 - \eta \text{ slope}$$

$L(\beta_0, \beta_1, \beta_2)$
four dimensional graph



n-dimension

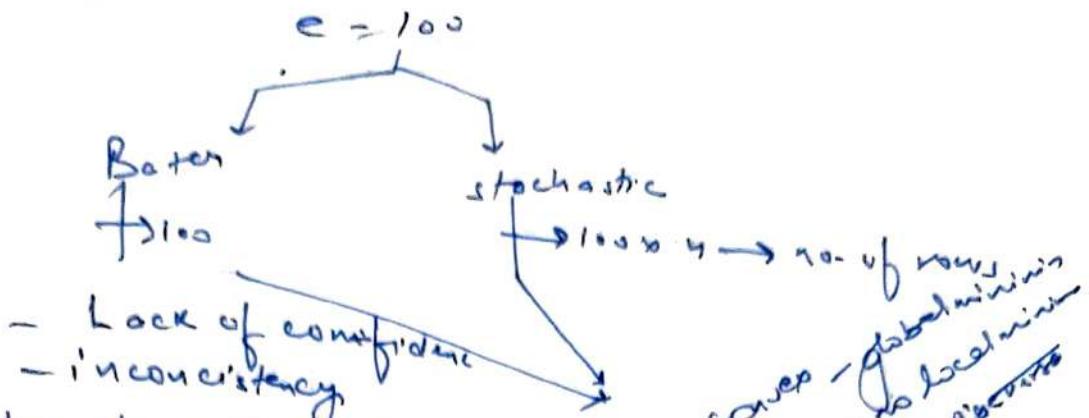
$$(n+1) \beta_0 - f_0$$

$$\overbrace{L}^{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{rows 2, when 2}$$

$$\frac{1}{2} [(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2]$$

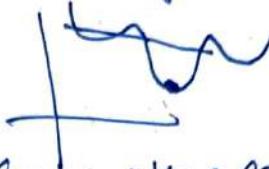


No. of epochs \rightarrow fixed



When to use SGD?

- ① \rightarrow Big data \rightarrow SGD
- ② Non convex function.



Learning schedules

Vary learning rate (α) with epochs

SGD Regression

Mini Batch GD

(How frequently you update)

Batch :- group of rows

batch update
per epoch

$n = 1000 \rightarrow$ batches $\rightarrow 100$ \rightarrow 10 batches

Ridge Regression

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda m^2$$

overfitting where m is slope
model fits better on training
but on testing does not perform

$$L = \sum_{i=1}^n (y_i - m x_i - b)^2 + \lambda m^2$$

gather minimum

$$\frac{\partial L}{\partial b} = 0 \rightarrow b$$

$$b = \bar{y} - m \bar{x}$$

$$\bar{y} = y_{\text{mean}}$$

$$\bar{x} = x_{\text{mean}}$$

m is slope

$$\frac{\partial L}{\partial m} = 0 \rightarrow m$$

$$L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + \lambda m^2$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})(-x_i - \bar{x}) + 2\lambda m = 0$$

$$= -2 \sum_{i=1}^n (y_i - \bar{y} - mx_i + m\bar{x})(x_i - \bar{x}) + 2\lambda m = 0$$

$$= \lambda m - \sum_{i=1}^n (y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

$$= \lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2 = 0$$

$$= \lambda m = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + m \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$= \lambda m + m \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$$

$$= m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

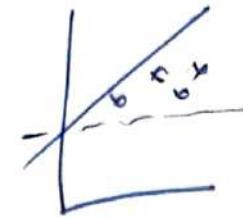
hyperparameter
alpha

$$= m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \lambda =$$



$$y = \frac{mx+b}{\uparrow \uparrow}$$

reduce m value to get
it on ~~wrong~~ right value
bias ↑ variance ↓



$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda m^2$$

$$L = \sum_{i=1}^n (y_i - mx_i - b)^2 + \lambda m^2$$

$$b = \bar{y} - m\bar{x} \quad \bar{y} \rightarrow y_{\text{mean}} \quad \bar{x} \rightarrow x_{\text{mean}}$$

$m = \text{slope}$

$$\frac{\partial L}{\partial b} = 0 \quad \rightarrow b$$

$$\frac{\partial L}{\partial m} = 0 \quad \rightarrow m$$

$m = \text{slope}$

$$L = \sum_{i=1}^n (y_i - mx_i + \bar{y} + m\bar{x})^2 + \lambda m^2$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (y_i - mx_i - \bar{y} - m\bar{x})(-x_i + \bar{x}) + 2\lambda m = 0$$

$$= -2 \sum_{i=1}^n (y_i - \bar{y} - mx_i + m\bar{x})(x_i - \bar{x}) + 2\lambda m = 0$$

$$= \lambda m = \sum_{i=1}^n [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

$$= \lambda m = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - m \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$= \lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + \lambda \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

↑

$$= \lambda w + \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2$$

$$= \lambda w + \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2 = \sum_{i=1}^m (y_i - \bar{y})(w_i - \bar{w})$$

$$= \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})(w_i - \bar{w})$$

$$\sum_{i=1}^m (w_i - \bar{w})^2 + \lambda$$

$\lambda \rightarrow \text{hyperparameter}$

$$w = \frac{\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_n | \mathbf{y}}{\downarrow \quad \downarrow \quad \dots \quad \downarrow} \boxed{w_0 \quad w_1 \quad w_2 \quad \dots \quad w_n}$$

(n+1)

Ridge

$$L = \sum_{i=1}^m (y_i - \bar{y})^2$$

$$= (\mathbf{x}w - \mathbf{y})^T (\mathbf{x}w - \mathbf{y})$$

$$\mathbf{y} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \text{values} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad (n+1)$$

$$\mathbf{x} = \begin{bmatrix} 1 & \mathbf{x}_1^T & \mathbf{x}_2^T & \dots & \mathbf{x}_n^T \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} & \end{bmatrix} \quad \text{size} \quad m \times n$$

Normal linear regression \rightarrow Ridge

$$L = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$\begin{cases} \lambda w_1^2 + \lambda w_2^2 + \lambda w_3^2 + \\ \dots + \lambda w_n^2 \end{cases}$$

$$\lambda (w_0^2 + w_1^2 + w_2^2 + \dots + w_n^2)$$

$$\begin{bmatrix} w_0, w_1, w_2, \dots, w_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \overline{\overline{w}}^T w$$

$$\begin{array}{l} (a-b)^T = a^T - b^T \\ (a+b)^T = b^T a^T \end{array}$$

$$\left(\frac{\partial L}{\partial w} \right) = ?$$

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

$$L = [x^T - y^T] (xw - y) + \lambda w^T w$$

$$= (w^T x^T - y^T) (xw - y) + \lambda w^T w$$

$$= w^T x^T x w - w^T x^T y - y^T x w + y^T y + \lambda w^T w$$

$$L = w^T x^T x w - 2 w^T x^T y + \overline{\overline{y^T y}} + \lambda w^T w$$

$$\frac{dL}{dw} = \cancel{2x^T x} w - \cancel{2x^T y} + 0 + 2\lambda w = 0$$

$$x^T x w + \lambda w = x^T y$$

$$\frac{x^T x + \lambda I}{I} w = \cancel{x^T y}$$

$$w = (x^T x + \lambda I)^{-1} x^T y$$

$$w = x^T x + \lambda I^{-1} x^T y$$

solve for Ridge regression

Verlustfunktion

$$L = \sum_{i=1}^n (\gamma_i - \tilde{\gamma}_i)^2$$

vector form $L = (xw - y)^T (xw - y) + \lambda \|w\|^2$

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

w_0, w_1, \dots, w_m C parameters

$$w_0 = w_0 - \gamma \frac{\partial L}{\partial w_0}, \quad w_1 = w_1 - \gamma \frac{\partial L}{\partial w_1}, \dots, w_m = w_m - \gamma \frac{\partial L}{\partial w_m}$$

$$\Delta w_0 \rightarrow w_{\text{new}} = w_{\text{old}} - \boxed{\gamma \frac{\partial L}{\partial w_0}} \rightarrow \text{Gradient } \frac{\partial L}{\partial w_0}$$

$$\begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_m} \end{bmatrix}$$

$$L; w$$

$$L = \frac{1}{2} (xw - y)^T (xw - y) + \frac{1}{2} \lambda w^T w$$
$$= \frac{1}{2} (w^T x^T - y^T) \overbrace{(xw - y)}^{\frac{1}{2} (xw - y)^T} + \frac{1}{2} \lambda w^T w$$

$$\begin{aligned}
 & \frac{1}{2} [\mathbf{w}^T \mathbf{x}^T \mathbf{x} \mathbf{w} - \cancel{\mathbf{y}^T \mathbf{x} \mathbf{w}} \cancel{\mathbf{w}^T \mathbf{x} \mathbf{y}} - \cancel{\mathbf{y}^T \mathbf{w} \mathbf{w}} + \mathbf{y}^T \mathbf{y}] + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} \\
 & = \frac{1}{2} [\mathbf{w}^T \mathbf{x}^T \mathbf{x} \mathbf{w} - \cancel{\frac{2 \mathbf{w}^T \mathbf{x} \mathbf{y}}{2}} + \mathbf{y}^T \mathbf{y}] + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} \\
 & = \frac{dL}{d\mathbf{w}} = \frac{1}{2} \left[2 \mathbf{x}^T \mathbf{x} \mathbf{w} - \cancel{2 \mathbf{x}^T \mathbf{y}} \right] + \frac{1}{2} \lambda \mathbf{x}^T \mathbf{w} \\
 & = \mathbf{x}^T \mathbf{x} \mathbf{w} - \cancel{\frac{\mathbf{x}^T \mathbf{y}}{2}} + \lambda \mathbf{w} = \frac{dL}{d\mathbf{w}} \left(\frac{\Delta L}{\Delta \mathbf{w}} \right) \\
 & = \mathbf{w} = \begin{bmatrix} w_0 & w_1 & \dots & w_m \\ 0, 1, \dots, m \end{bmatrix} \quad \text{starting}
 \end{aligned}$$

epoch $\mathbf{w} = \mathbf{w} - \gamma \frac{dL}{d\mathbf{w}}$

$\mathbf{w} \rightarrow \underline{\text{final answer}}$

$$\mathbf{w} = \mathbf{w} - \gamma \frac{dL}{d\mathbf{w}}$$

epoch times

$$\boxed{\frac{dL}{d\mathbf{w}} = \mathbf{x}^T \mathbf{x} \mathbf{w} - \mathbf{x}^T \mathbf{y} + \lambda \mathbf{w}}$$

Lasso Regression (\hat{L})

$$L = \text{MSE} + \lambda \|\mathbf{w}\|$$

$$\lambda [|w_1| + |w_2| + \dots + |w_n|]$$

code, implement lars

→ compare Ridge vs Lasso regression

Simple

$$y = mx + b$$

$$b = \bar{y} - m\bar{x}$$

$$\bar{y} = \text{mean}(y)$$

$$\bar{x} = \text{mean}(x)$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

make m is not affected in sum two terms

$$b = \bar{y} - m\bar{x} \quad m = ?$$

$$\text{Loss } L = \sum_{i=1}^n (y_i - \bar{y})^2 + \lambda|m| \quad \text{for methods simple}$$

$$\frac{dL}{dm} = \sum_{i=1}^n (y_i - mx_i - b\bar{y} - m\bar{x})^2 \quad |m|$$

$$\begin{cases} m > 0 \\ |m| = m \end{cases}$$

$$\frac{d}{dm} \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda m$$

$$2 \sum (y_i - mx_i - \bar{y} + m\bar{x}) (-x_i + \bar{x}) + 2\lambda = 0$$

$$-2 \sum [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) + 2\lambda = 0$$

$$- \sum [(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2] + 2\lambda = 0$$

$$- \sum (y_i - \bar{y})(x_i - \bar{x}) + m \sum (x_i - \bar{x})^2 + 2\lambda$$

$$m \sum (x_i - \bar{x})^2 = \sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

Lasso
for $m > 0$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2}$$

for $m = 0$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

for $m < 0$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

Sparcity occurs why?

$$m < 0 \quad m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

$$\lambda > 0$$

$$m = \frac{-100 + \lambda}{50}$$

$$= \frac{-100 - 150}{50}$$

$$= -5$$

$$\begin{aligned} \lambda &= 0, m = -2 \\ \lambda &= 50, m = -1 \\ \lambda &= 100, m = 0 \\ \lambda &= 150, m = 1 \end{aligned}$$

1) $0 \rightarrow$
2) $0 \rightarrow$ stop

why Ridge regression no sparsity

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2 + \lambda}$$

Deno noise

$m = 0$ when numerator is zero

In ridge λ is numerator
Inlasso λ is denominator

$$\lambda = 10000000$$

Elasticnet Regression

it's a combination of Ridge and Lasso

Ridge

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

mse

overfitting
reduces

used where all feature
are important

↓ 100 column.



Lasso

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|_1$$

- feature selection
- if any feature is not important than eliminate that particular feature

If you have large data set and you are unable to decide which one is applied in that case go to elasticnet regression

$$L = \sum (y_i - \hat{y}_i)^2 + \alpha \|w\|^2 + \beta \|w\|$$

$$\lambda = \alpha + \beta$$

$$= \frac{\alpha}{\alpha + \beta}$$

λ and $\ell 1$ -ratio

two hyperparameters

$$\lambda = 1, \ell 1\text{-ratio} = 0.5$$

$$\alpha = 0.5 \quad b = 0.5 \quad \ell 1\text{-ratio} = 0.5$$

$$\uparrow \ell 1\text{-ratio} = 0.9$$

90% ridge regres
10% Lasso regres

$$L1 = \frac{\alpha}{\lambda}$$

$$\alpha = \lambda^{-1} \quad g^b = T^{-1}$$

Input wt. → multicollinearity

$$x_1 | x_2$$

Logistic Regression

Introduction:-

- easy method to learn.
- 2 perspective
 - + geometric
 - + probability

Required
+ linearly separable
or nearly linearly
separable

- If we apply logistic regression on nonlinear data
then we didn't get good result.

- Perception trick

$$Ax + By + c = 0$$

$$Ax_1 + Bx_2 + c = 0$$

$$Ax_1 + Bx_2 + Cx_3 + d = 0$$

↑ sing dimensions

A, B, C to find these values

Desired
for r_i

start random value of A, B, C

- random student, ask from different random points
/ transformation \rightarrow ~~The~~ ~~The~~ ~~Yon~~

loop \rightarrow randomly pick point \rightarrow shifting the line from
response of randomly pick point

- 1000 epochs

- till convergence

How to identify these
- no epochs

How to label a region

Transformations

C \rightarrow Line more parallel up
- rd and down more

~~A~~ \rightarrow from point Y more

~~B~~ \rightarrow from X more

Not Good X

$$2x + 3y + 5 > 0$$

$$2x + 3y + 5 < 0$$

$$2x + 3y + 5 = 0$$

+ve
-ve
along

$$(4,5) \rightarrow 4, 5, 1$$

$$\begin{array}{r}
 & 2 & 3 & 5 \\
 -4 & 5 & 1 \\
 \hline
 2x - 2y + 4 = 0
 \end{array}$$

when one
more row
region has
~~one~~
subtracted

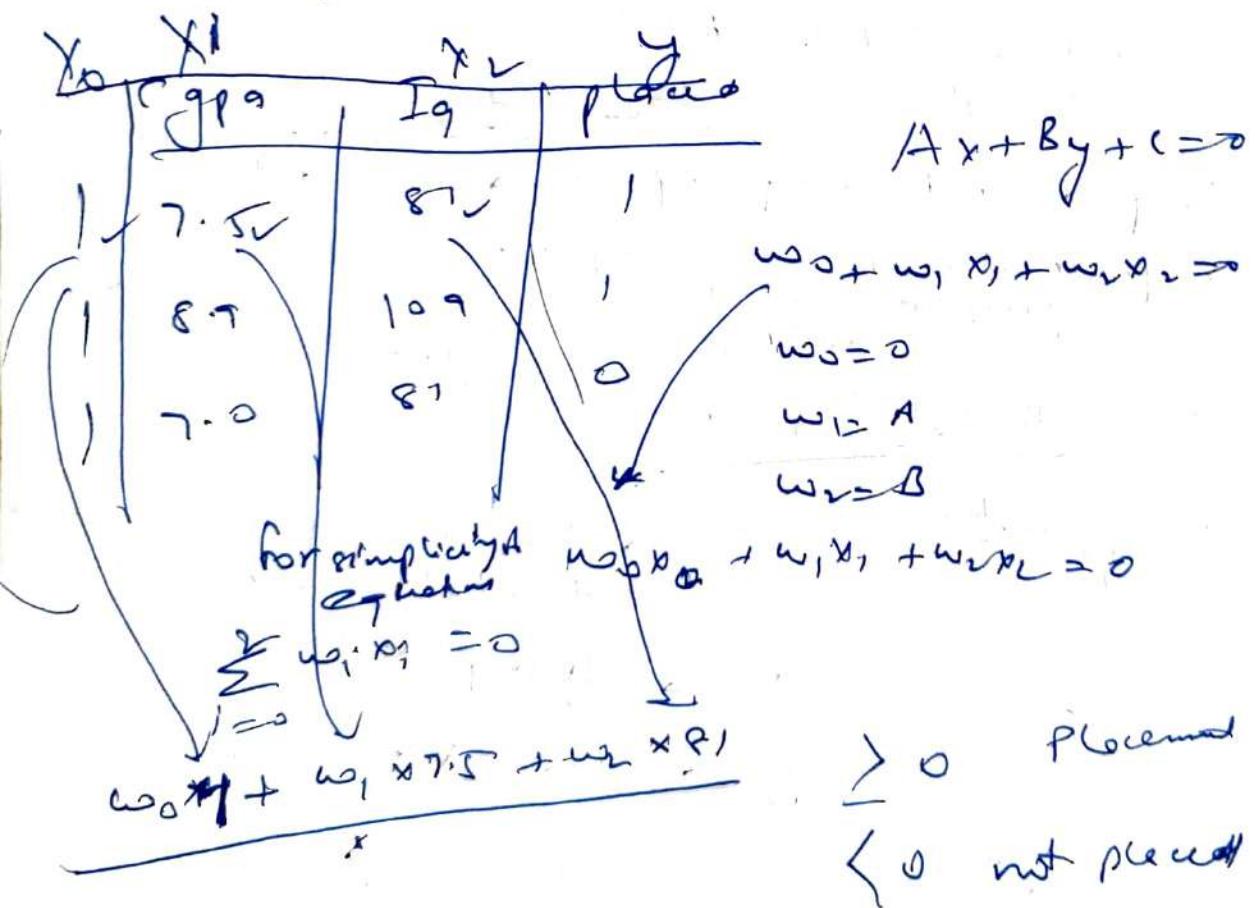
$$(1,3) \rightarrow (1,3,1)$$

$$\begin{array}{l}
 \text{learning rate} \rightarrow 0.01 \\
 \hline
 =
 \end{array}$$

$$\begin{array}{r}
 & 2 & 3 & 5 \\
 +1 & 3 & 1 \\
 \hline
 3x + 6y + 6 = 0
 \end{array}$$

when
less
more
-ve reg's
that
add

coeff = coeff of coordinates



$$\begin{array}{c}
 \text{work with } w_0, w_1, w_2 \\
 \text{work with } x_0, x_1, x_2 \rightarrow \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \\
 \text{work with } y_0, y_1, y_2 \rightarrow \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix} \\
 \text{work with } w_0, w_1, w_2 \rightarrow \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \\
 \text{work with } x_0, x_1, x_2 \rightarrow \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \\
 \text{work with } y_0, y_1, y_2 \rightarrow \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix}
 \end{array}$$

decide epoch $\rightarrow 1000$
 $\eta = 0.01$
 for i in range epochs
 randomly select pt
 or student
 \rightarrow if

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{[x_0, x_1, x_2]}{[x_i]}$$

$$\text{if } x_i \in P \text{ and } \sum_{i=1}^3 w_i x_i < 0$$

$$w_{\text{new}} = w_{\text{old}} + \eta x_i$$

simplified the algo

for i in range 1000 randomly selected student

$$w_n = w_0 + \eta (y_i - \hat{y}_i) x_i$$

$P = \text{positive regions}$

$$y_i \quad \hat{y}_i \quad y_i - \hat{y}_i$$

$$\begin{array}{ccc}
 1 & 1 & 0 \\
 0 & 0 & 0
 \end{array}
 \quad w_n = w_{\text{old}}$$

$$\begin{array}{ccc}
 1 & 0 & 1 \\
 0 & 1 & -1
 \end{array}
 \quad \cancel{w_n = w_0 + \eta x_i}$$

$$\cancel{\text{change the previous approach}} \quad \cancel{w_n = w_0 - \eta x_i}$$

- misclassified - line pull

correctly classified - line push

$$w_n = w_0 + (y_i - \hat{y}_i) x_i$$

$$g = 0^-$$

$$x > z(\infty) \quad \text{sigmoid function} \quad \frac{P}{1+e^{-z}} = \sigma(z)$$

Sigmoid function give probability function

sigmoid function
 give probability function
 define
 probability
 function
 sigmoid function
 give probability function

$$\sigma = 0.5$$

$$P = 0.5$$

$$w_n = w_0 + \eta \sum (y_i - \bar{y}) x_i$$

$$y_i = \sigma(z)$$

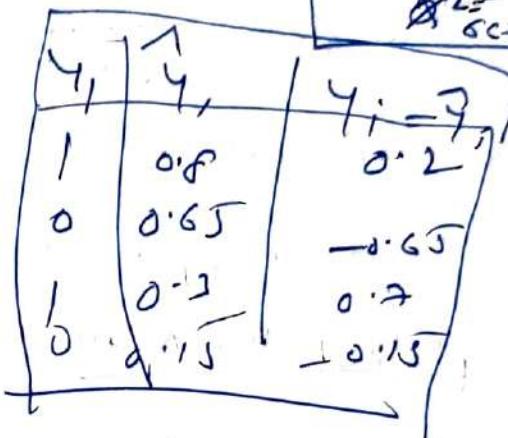
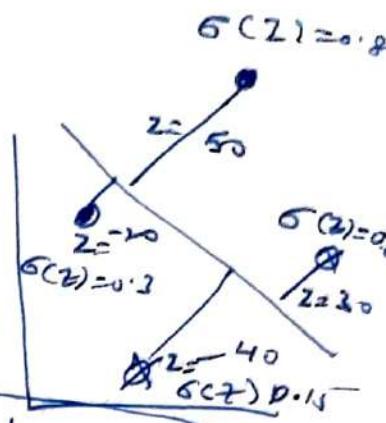
where $z = \sum w_i x_i$

$$w_n = w_0 + \eta * 0.2 * x_1$$

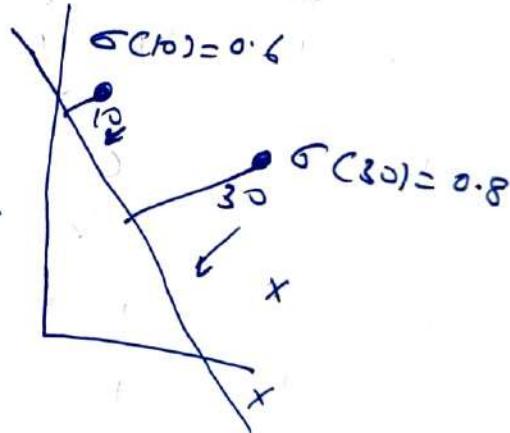
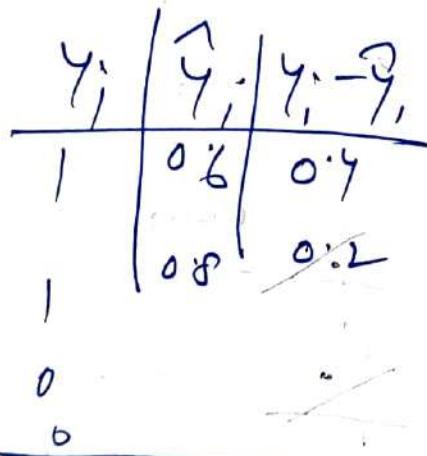
$$w_n = w_0 + \eta * 0.65 * x_1$$

$$w_n = w_0 + \eta * 0.7 * x_1$$

$$w_n = w_0 - \eta * 0.15 * x_1$$



start with
single output
expander



$$w_n = w_0 + \eta * 0.7 * x_1 = x_1$$

$$w_n = w_0 + \eta * 0.2 * x_1 = x_2 \quad x_1 > x_2$$

Maximum likelihood (based on probability)

$$\begin{aligned} \bar{y} &= \sigma(z) \\ z &\equiv \sum w_i x_i \end{aligned}$$

loss function error predictions

methodology is
not true because
of 1000 steps every
time give new result

If target word data not met i.e. it's more than predicted
(is very rarely)

$$\log(0.5) = \log a + \log b + \log(c \cdot d)$$

$$\log \text{max} = \log(0.7) + \log(0.4) + \log(0.8)$$

0-1 = -ve
convert all logs into -ve $\log \frac{0.1}{0.9}$

$$\log(\text{max}) = -\log(0.7) - \log(0.4) - \log(0.8) - \log(0.9)$$

cross entropy minimize

$$-y_i \log(\hat{q}_i) - (1-y_i) \log(1-\hat{q}_i)$$

$$-y_1 \log(\hat{q}_1) = -\log(q_1) \\ = -\log(0.7)$$

$$-y_2 \log(\hat{q}_2) - (1-y_2) \log(1-\hat{q}_2) \\ = -\log(1-\hat{q}_2) = -\log 0.4$$

$$y_2 = 1$$

$$-y_3 \log(\hat{q}_3) = \log(\hat{q}_3) = -\log(0.4)$$

$$y_3 = 0, (1-y_3) \log(1-\hat{q}_3)$$

$$= -\log(1-\hat{q}_3)$$

$$L = \sum_{i=1}^n -y_i \log(\hat{q}_i) - (1-y_i) \log(1-\hat{q}_i) \\ = -\log(0.8)$$

$$L = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

min.
w.r.t. w's

\rightarrow gradient descent.

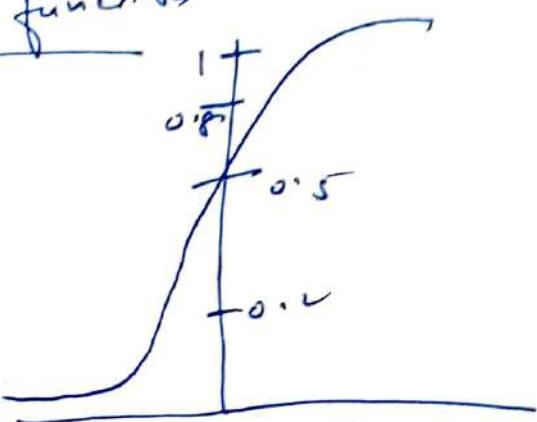
\rightarrow log loss errors

Binary cross entropy
function.

Derivative of sigmoid function

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

differentiation reciprocal
rule



$$\begin{aligned} \frac{d}{dx} \left(\frac{1}{x} \right) &= \frac{d}{dx} (x)^{-1} \\ &= -x^{-2} = -\frac{1}{x^2} \end{aligned}$$

$$\frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] = \frac{d}{dx} \left[(1+e^{-x})^{-1} \right]$$

$$= -\frac{1}{(1+e^{-x})^2} \frac{d}{dx} (1+e^{-x}) \quad \boxed{\text{known } e^{-x} = e^x}$$

$$= \frac{-1}{(1+e^{-x})^2} \frac{d}{dx} (e^{-x}) = \frac{-e^{-x}}{(1+e^{-x})^2} \frac{d}{dx} (-x)$$

$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\frac{1 \cdot e^{-x}}{(1+e^{-x})(1+e^{-x})} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \sigma(x) \left[\frac{e^{-x}}{1+e^{-x}} \right]$$

Sumant@2022

$$= \sigma(x) \left[\frac{1+e^{-x}-1}{1+e^{-x}} \right] = \sigma(x) \left[\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right]$$

$$\sigma(x) \left[1 - \sigma(x) \right] \Rightarrow \sigma^*(x) = \sigma(x) \left[1 - \sigma(x) \right]$$

L-474 Logistic Regression - Gradient descent

Rows = m

columns = n

1	2	3	n	y
x_{11}	x_{12}	x_{13}	x_{1n}	y_1
x_{21}	x_{22}	x_{23}	x_{2n}	y_2
⋮	⋮	⋮	⋮	⋮
x_{m1}	x_{m2}	x_{m3}	x_{mn}	y_n

$$\hat{y} = \left[w_0, w_1, w_2, \dots, w_n \right]^T$$

$$\sigma(w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n) = \hat{y}_1$$

$$\sigma(w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n) = \hat{y}_2$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \begin{bmatrix} \sigma(w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n}) \\ \sigma(w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n}) \\ \vdots \\ \sigma(w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn}) \end{bmatrix}$$

$$\hat{y} = \sigma \left(\begin{bmatrix} w_0 \\ w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n} \\ w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n} \\ \vdots \\ w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn} \end{bmatrix} \right)$$

$$= \sigma \left(\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \right)$$

$$\hat{y} = \sigma(\omega^T w)$$

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$$L = -\frac{1}{m} \left[\sum_{i=1}^m y_i \log(\hat{y}_i) + \sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i) \right]$$

$$\sum_{i=1}^m y_i \log(\hat{y}_i) = y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3 + \dots + y_m \log \hat{y}_m$$

$$[y_1, y_2, y_3, \dots, y_m] \begin{bmatrix} \log \hat{y}_1 \\ \log \hat{y}_2 \\ \vdots \\ \log \hat{y}_m \end{bmatrix}$$

$$[y_1, y_2, y_3, \dots, y_m] \log \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix}$$

$\hat{y} \log \hat{y} = y \log(\sigma(\omega^T w))$

$$L = -\frac{1}{m} [\hat{y} \log \hat{y} + (1-y) \log(1-\hat{y})]$$

where $\hat{y} = \sigma(\omega^T w)$ \uparrow gradient descent
 $[\omega]$ final

Loss function in matrix form

$$L = \frac{1}{m} [y \log(\sigma(w^T x)) + (1-y) \log(1-\sigma(w^T x))]$$

\uparrow minimum

gradient descent

$$w = [\quad]$$

$$\left[\frac{\Delta L}{\Delta w} \right]$$

for i in epoch:

$$w = w - \gamma \frac{\Delta L}{\Delta w}$$

$$w_0 \rightarrow w_n$$

$$\frac{\Delta L}{\Delta w} = \left[\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]^{(n+1)}$$

$$L = -\frac{1}{m} \left[y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

$$\frac{dL}{dw} = \frac{d}{dw} y \log \hat{y} \Rightarrow y \frac{d}{dw} \log \hat{y} = y \frac{d}{d\hat{y}} (\hat{y})$$

$$\Rightarrow \cancel{y} \frac{d}{d\hat{y}} \sigma(wx) \Rightarrow \cancel{y} \sigma(wx) [1 - \sigma(wx)] \frac{d}{d\hat{y}} (\hat{y})$$

$$= \cancel{y} \cancel{x} (1-\hat{y}) x = \boxed{y(1-\hat{y})x}$$

$$\frac{d}{dw} (1-y) \log (1-y) \Rightarrow (1-y) \frac{d}{dw} \log (1-\hat{y}) \Rightarrow \cancel{(1-y)} \frac{d}{d\hat{y}} (\hat{y})$$

$$\Rightarrow \frac{(1-y)}{(1-\hat{y})} \frac{d}{d\hat{y}} (1-\hat{y}) \quad \left\{ \hat{y} = \sigma(wx) \right\}$$

$$= -\frac{(1-y)}{(1-\hat{y})} \frac{d}{d\hat{y}} \sigma(wx) \Rightarrow -\frac{(1-y)}{(1-\hat{y})} [\sigma(wx) [1 - \sigma(wx)]]$$

$$\Rightarrow -\frac{(1-y)}{(1-\hat{y})} \cancel{y(1-\hat{y})x} \boxed{-\hat{y}(1-y)x} \frac{d}{d\hat{y}} (\hat{y})$$

$$\frac{dL}{dw} = -\frac{1}{m} [y(1-\hat{y})x - \hat{y}(1-y)x]$$

$$= -\frac{1}{m} [y(1-\hat{y}) - \hat{y}(1-y)]x$$

$$= -\frac{1}{m} [y - y\hat{y} - \hat{y} + y\hat{y}]x$$

$$\boxed{\frac{\Delta L}{\Delta w} = -\frac{1}{m} (y - \hat{y})x}$$

gradient descent update

$$\boxed{w = w + \gamma \frac{1}{m} (y - \hat{y})x}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \quad (n+1, 1)$$

$$x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \quad (m, n+1) \rightarrow (m, n+1)$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \\ y_{n+1} \end{bmatrix} \quad (n+1, 1)$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \quad (m, 1)$$

$$(n+1, 1) \quad \begin{bmatrix} w \\ \frac{1}{m} \end{bmatrix} \quad (n+1, 1) \quad (m, 1) \quad (m, n+1)$$

$$\cancel{(n+1, 1)} \quad \cancel{\frac{1}{m}} \quad \cancel{(n+1, 1)} \quad \cancel{(m, 1)} \quad \cancel{(m, n+1)}$$

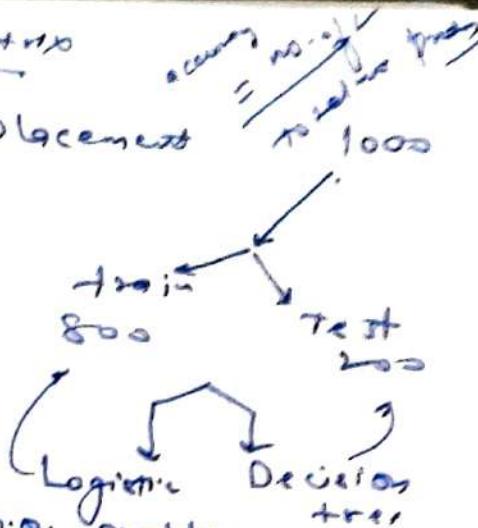
$$\frac{1}{m} (y - \hat{y})x \quad \rightarrow \quad (m, n+1)$$

Accuracy and confusion matrix

Accuracy

no. of correct prediction = accuracy
 total no. of prediction

$$\frac{8}{10} = 0.8 = 80\% \text{ accurate}$$



Accuracy of multi classification problem

How much accuracy is good?

It depends on the problem we are solving

e.g. in medical field, 100% accuracy is needed

even 99% accuracy is supposed to be bad

- self driving car 100% accuracy is required

whereas shopping for food recommendation 80% is good

90% accurate \rightarrow too many wrong.

Accuracy score not tell us about type of mistake (nature of mistake)

Confusion matrix: (predicted nature of mistake)

		Prediction	
		1	0
Actual	1	TP 26	FN 6
	0	Type I error FP 0	TN 29

of correct pred / # of predictions

$$Acc = \frac{TP + TN}{TP + TN + FP + FN}$$

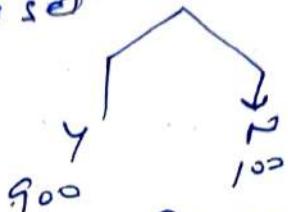
Type 1 Error ^{relative confusion matrix} for multiclass
 Type 2 Error ^{failure} Predict classification problem

Predicted	0	1	2
Actual	0	1	2
0	10	2	1
1	2	7	1
2	1	1	17

Diagonal elements are true predictions
 $\frac{\text{true pred}}{\text{total}} = \text{accuracy}$

when accuracy is misleading

i) imbalanced data set



when ratio is not appropriate

Precision, Recall and F1 score

Precision

$$fP_a > fP_b \quad | \quad fP_a < fP_b$$

what proportion of predicted positive is truly positive
 called precision

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

Recall :- what proportion of actual +ve is correctly classified

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

$$\text{Recall A} = \frac{1000}{1200}$$

F1 score: It's a combination of precision and recall.

Harmonic mean towards lower side

$$F_1 \text{ score} = \frac{2 \cdot P \cdot R}{P + R}$$

+ macro precision or calculate

$$\frac{0.86 + 0.66 + 0.58}{3} = 0.70$$

+ weighted precision

$$\left(\frac{40}{100} \cdot 0.86 \right) + \left(\frac{34}{100} \cdot 0.66 \right) + \left(\frac{24}{100} \cdot 0.58 \right)$$

$$= 0.71$$

If your classes balanced then used macro

If your classes imbalance then used weighted

$$R_D = \frac{25}{40} = 0.62, \quad R_C = \frac{30}{24} = 0.88$$

$$R_R = \frac{37}{20} = 0.58$$

+ macro recall

+ weighted recall

$$F_{1D} = \frac{2 \cdot P_D \cdot R_D}{P_D + R_D}$$

$$= \frac{2 \cdot 0.86 \cdot 0.66}{0.86 + 0.66}$$

softmax Regression (multinomial Regression)

			Yes	No	Optout
cgpa	1.9	164			
7.1	71	0			
8.5	85	1			
9.5	95	2			

{65, 65} → yes, no, optout

softmax \rightarrow general \rightarrow for n classes

softmax function

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$k = \text{no. of classes}$

yes = 1

no = 0
opposite = 1

$$\sigma(z)_i = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z_3) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

values add up to 1
and adding all three equal to 1

Training intuition

cgp = 1/9 | placed?

- - - - 1
- - - - 2

+ transform your data
one hot encode

cgp	1/9	0	0	0	0	0	0	0	0
		z1	z2						
1		1	0	0	0	0	0	0	0
0		0	1	0	0	0	0	0	0
0		0	0	1	0	0	0	0	0
0		0	0	0	1	0	0	0	0
0		0	0	0	0	1	0	0	0
0		0	0	0	0	0	1	0	0
0		0	0	0	0	0	0	1	0
0		0	0	0	0	0	0	0	1



$$\text{cgp} = 1/9 = 0 \\ M_1$$

$$\text{cgp} = 1 \\ M_2$$

$$\text{cgp} = 1/3 = 1/3 \\ M_3$$

$$\boxed{\begin{matrix} 3 \text{ coeff} \\ w_0 & w_1 & w_2 \\ w_1 & w_2^{(0)} & w_0^{(0)} \end{matrix}}$$

$$\boxed{\begin{matrix} y_1 \text{ no} \\ w_1^{(1)}, w_2^{(1)}, w_0^{(1)} \end{matrix}}$$

$$\boxed{\begin{matrix} m_3 \text{ opt out} \\ w_1^{(2)}, w_2^{(2)}, w_0^{(2)} \end{matrix}}$$

$$z_1 = 7 \times w_1^{(1)} + 7 \times w_2^{(2)} + w_0$$

$$z_2 = 7 \times w_1^{(2)} + 7 \times w_2^{(4)} + w_0$$

$$z_3 = 7 \times w_1^{(3)} + 7 \times w_2^{(3)} + w_0^{(3)}$$

$$(0.45) \quad C(\gamma) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$C(0) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$(0.25) \quad C(0) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

Loss functions

$$L = \frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1 - \hat{y}_i) \log(1 - \hat{y}_i)$$

$$\boxed{L = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{y}_k^{(i)})}$$

$$x_1 \quad x_2 \quad y \quad y_{k=0} \quad y_{k=1} \quad y_{k=2}$$

$$\begin{array}{cccccc} x_{11} & x_{12} & 0 & 1 & 0 & 0 \\ x_{21} & x_{22} & 2 & 0 & 0 & 1 \\ y_{21} & x_{22} & 2 & 0 & 0 & 1 \end{array}$$

$$y_0^{(1)} \log(\tilde{y}_1^{(1)}) + y_1^{(1)} \log(\tilde{y}_2^{(1)}) + y_2^{(1)} \log(\tilde{y}_3^{(1)}) \\ + y_0^{(2)} \log(\tilde{y}_1^{(2)}) + y_1^{(2)} \log(\tilde{y}_2^{(2)}) + y_2^{(2)} \log(\tilde{y}_3^{(2)}) \\ + y_0^{(3)} \log(\tilde{y}_1^{(3)}) + y_1^{(3)} \log(\tilde{y}_2^{(3)}) + y_2^{(3)} \log(\tilde{y}_3^{(3)})$$

$$L = y_0^{(1)} \log(\tilde{y}_1^{(1)}) + y_1^{(1)} \log(\tilde{y}_2^{(1)}) + y_2^{(1)} \log(\tilde{y}_3^{(1)})$$

$$\tilde{y}_1^{(1)}, \tilde{y}_2^{(1)}, \tilde{y}_3^{(1)} = ?$$

$\rightarrow s = \text{softmax}$

$$\tilde{y}_1^{(1)} = \sigma(w_0^{(1)}x_{11} + w_1^{(1)}x_{12} + w_2^{(1)}x_{13} + w_3)$$

$$\tilde{y}_2^{(1)} = \sigma(w_0^{(2)}x_{21} + w_1^{(2)}x_{22} + w_2^{(2)}x_{23} + w_3)$$

$$\tilde{y}_3^{(1)} = \sigma(w_0^{(3)}x_{31} + w_1^{(3)}x_{32} + w_2^{(3)}x_{33} + w_3)$$



$$\frac{\partial L}{\partial w_{01}}, \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{21}} = \tilde{y}_{11} - y_{11}$$

gradient loop \rightarrow 10 epochs $i = 1 \dots 10$
 $w_{01}^{(1)} = w_{01}^{(1)} - \eta \frac{\partial L}{\partial w_{01}}$

$$w_{01}^{(1)} = w_{01}^{(1)} - \eta \frac{\partial L}{\partial w_{01}}$$

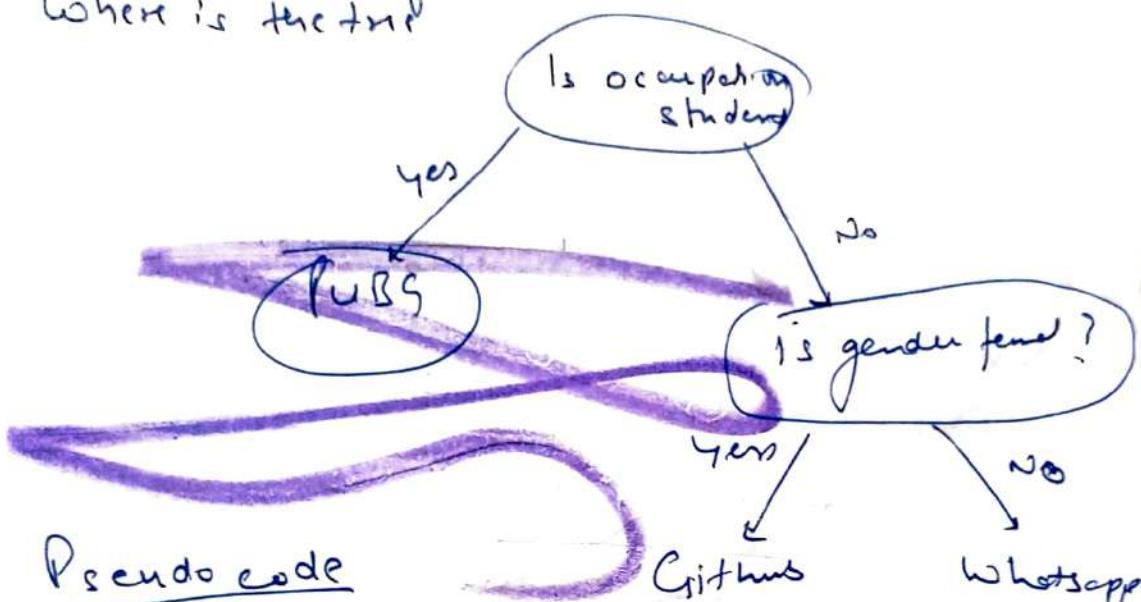
Polynomial Logistic Regression

- non-linear data
- logistic regression
- Polynomial,
degree x_1 , x_2
- useful in non-polynomial lines ↓
data
- convert non-linear to linear by applying poly —
 x_1 , x_2 , x_3 — x_1^0 x_1^1 x_1^2 — x_2^0 x_2^1 x_2^2 — y
- non-lin features. w_0 w_1 , w_2 w_3 in $w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$

Logistic hyperparameters

Decision tree Geometric Intuition

nested if else condition
where is the tree?



Pseudo code

Github

WhatsApp

- Begin with your training dataset, which should have some features variables and classification or regression output
 - Determine the best feature in the dataset to split the data on; more on how we define "best feature" later
 - Split the data into subset that contains the correct values for this best feature. This splitting basically defines a nodes for tree i.e. each node is splitting point based on certain feature from our data
 - Recursively generate a new tree node by using the subset of data created from step 3
- Root node - splitting - Branch (subset) - Decision tree - Leaf node

Conclusion: Programmatically speaking, Decision trees are nothing but a giant structure of nested if-else condition.

Mathematically speaking, Decision tree use hyperplanes which run parallel to any one of the axes to cut your coordinate system into hyper cuboids.

Root node splitting Branch / subtree.

Decision node, Leaf node

Some unanswered questions

How to decide which column should be considered as root node?

How to select subsequent decision nodes?

How to decide splitting criteria in case of numerical columns?

CART - Classification and Regression Trees

The logic of DT can be also be applied to regression problems, hence the name CART

Decision tree Entropy

What's entropy the measure of disorder

or measure of purity / impurity

or measure of entropy?

How to calculate entropy?

$$E(C) = \sum_{i=1}^n P_i \log_2 P_i$$

$$E(C) = P_{\text{yes}} \log_2 (P_{\text{yes}}) + P_{\text{no}} \log_2 (P_{\text{no}})$$

calculating entropy for 3 class

$$H(A) = -P_A \log_2 (P_A) - P_B \log_2 (P_B) - P_C \log_2 (P_C)$$

Observation

- more the uncertainty more is entropy
- for a 2-class problem the min. entropy is 0 and the max. 1,
- for more than 2 classes, the min. entropy is 0 but the max. can be greater than 1,
- both \log_2 or \log_e can be used to calculate entropy

Entropy vs probability

Entropy for continuous probability

Information Gain is a metric ~~equivalent to~~ to train Decision trees. Specifically, this metric measures the quality of a split.

The information gain is based on the decision tree's entropy after a data set is split on an attribute constituting a decision tree is all about finding attribute that returns the highest information gain.

$$\text{InformationGain} = E(\text{parent}) - \left\{ \text{weighted Average} \right. \\ \left. * \in \{\text{chicken}\} \right\}$$

Gini Impurity

$$E = -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$$G_I = 1 - (P_Y^2 + P_N^2)$$

$$\overline{P_Q} = \frac{\text{wt}_Q}{\text{wt}_Q + \text{wt}_N}$$

computation is fast

Handle Numerical data

1. Sort the data on the basis of numerical column.
2. Split the entire data on the basis of every value of user-rating.

(5) $\text{max}[\text{IS}_1, \text{IS}_2, \dots, \text{IS}_n]$

(6) Do this recursively until you get all the leaf nodes.

Overfitting

Perform well on training

$\text{max_depth} = \text{None}$ Perform badly on test

underfitting

- Perform till leaf node

$\text{max_depth} = \text{None}$: only one split

is not achieved

By tuning the hyperparameters toward underfitting and overfitting

hyperparameter overfitting, underfitting

PCA Comps to machine learning

Principal component analysis:

1/1

feature extraction \rightarrow COD curse of Dimensionality
eg photographer taken photo at diffrent angles.
3D movement in 2D image

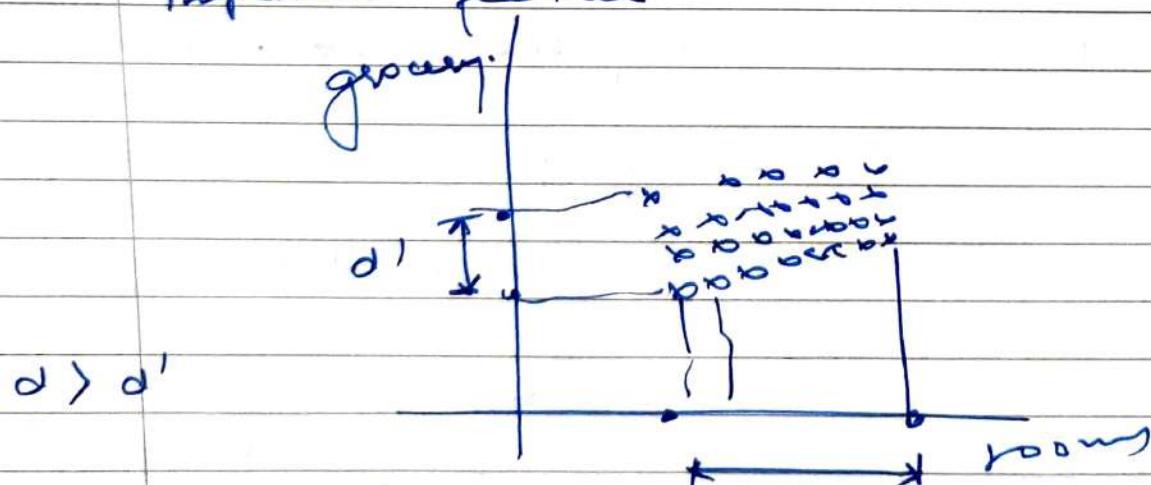
Benefit of PCA:

1. High dimensions to lower dimensions
2. Faster execution of alg
3. Visualization

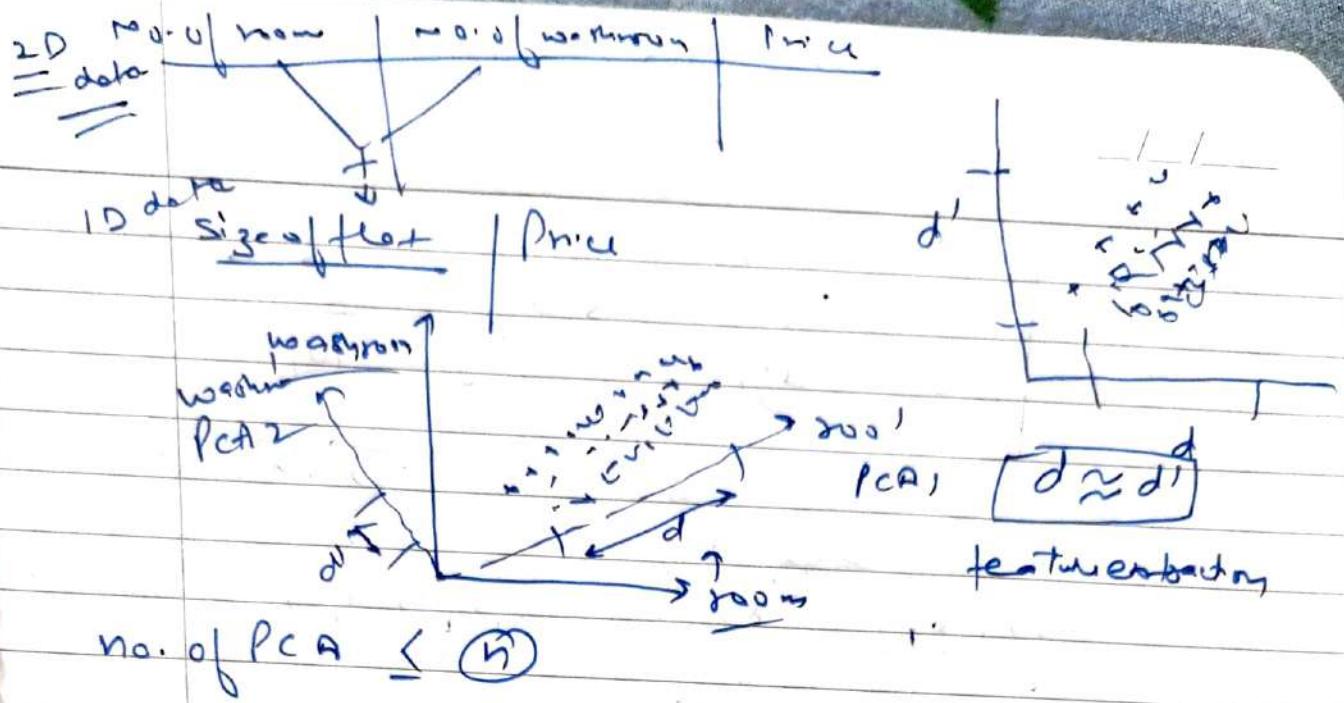
Geometric intuition :-

No. of rooms	No. of grocery shop	Price
3	2	60
4	0	130
5	6	170
7	10	90

feature selection: Reduce the feature only selected important feature

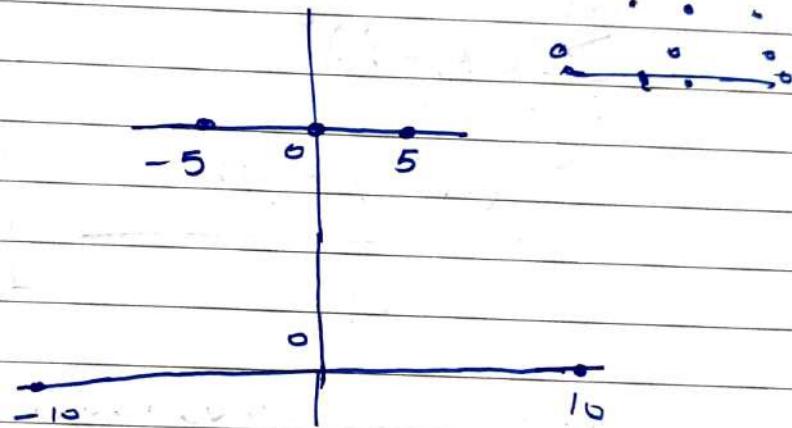


- pick that feature d on which variance is high.
- select the column whose spread is more



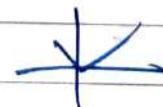
Variance spread \rightarrow imp?

mean



$$\text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{-5 + 0 + 5}{3} = \frac{0}{3}$$

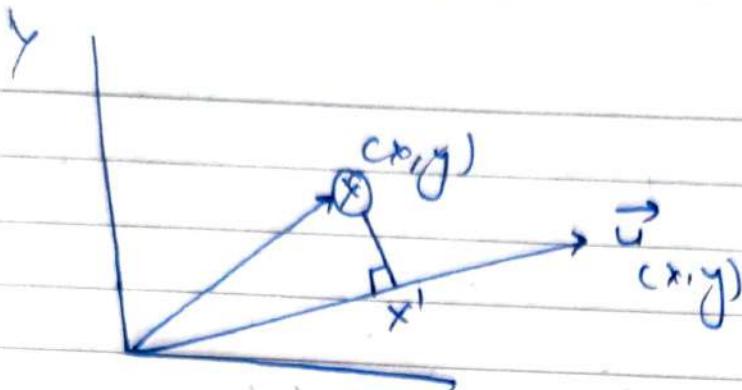
$$\text{Spread} \propto \text{to variance} \quad \frac{10 + 0 + 10}{3} = \frac{20}{3}$$



Why variance is important?

Problem formulation

2D \rightarrow 1D



$$\frac{\vec{v} \cdot \vec{x}}{|\vec{v}|} = 1$$

$$\vec{v} \cdot \vec{x} = \vec{v}^T \vec{x}$$

$$[x_1, y_1] \quad [x_2, y_2] = [x_1, x_2 + y_1, y_2] \text{ scalar}$$

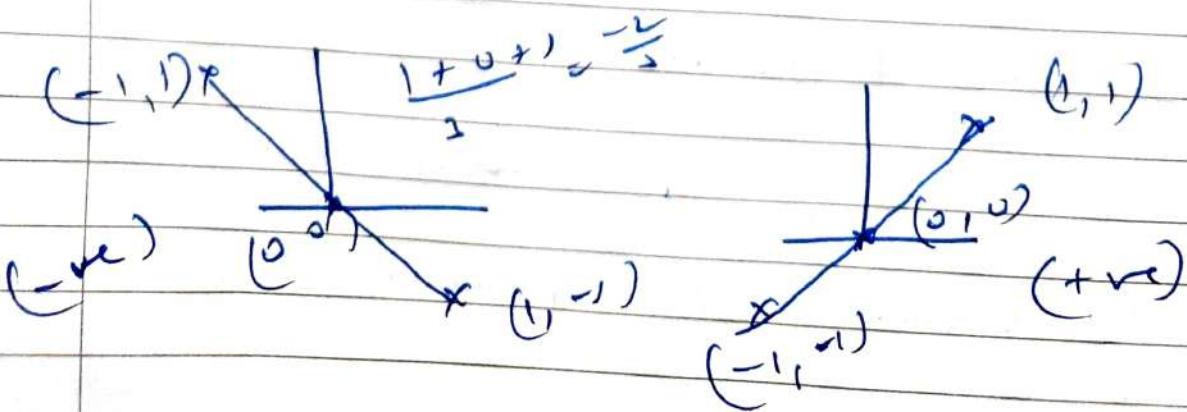
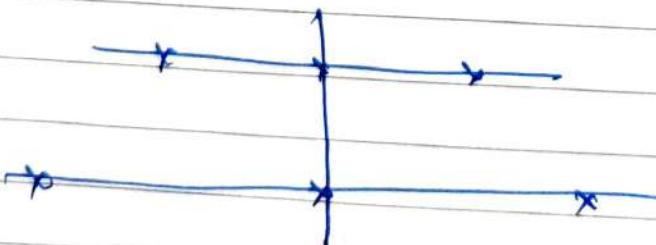
Select unit vector whose variance is minimum

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

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

\rightarrow optimization

covariance and covariance matrix



~~$x_1 \perp x_2$~~ covariance

$$[2x_1]$$

x_1

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

$$\text{cov}(x_1, x_1)$$

$$\downarrow$$

 $\text{cov}(x_1)$

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

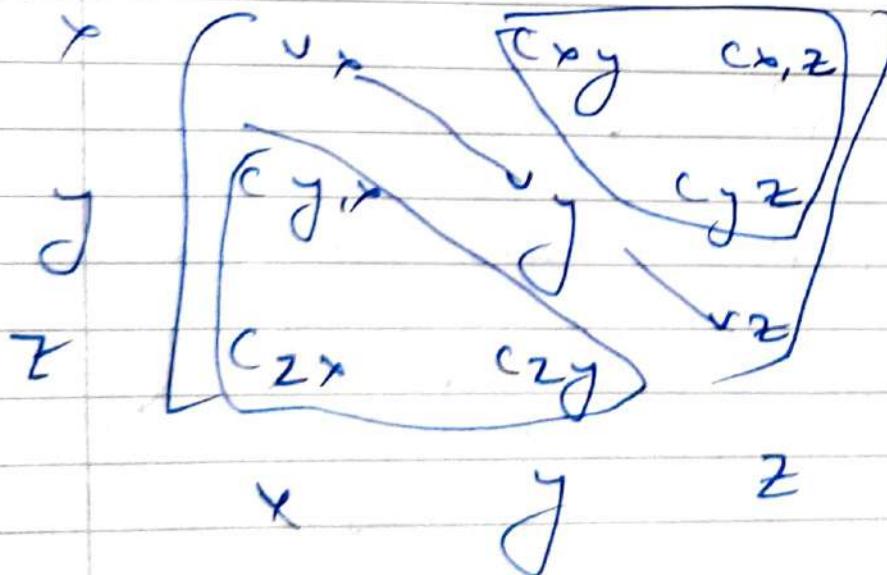
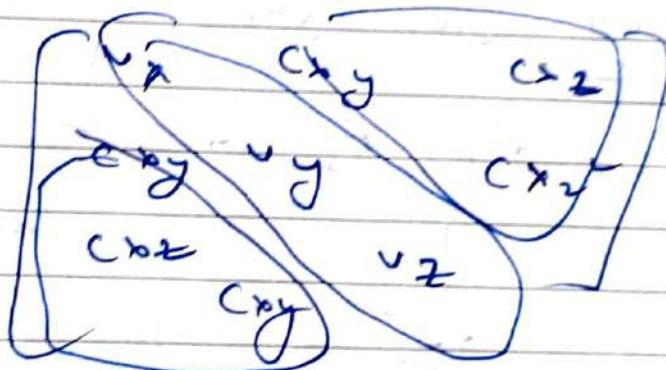
$$\text{cov}(x_2, x_1)$$

$$\text{Var } x_2$$

square

symmetric

$$\text{cov}(a, b) = \text{cov}(b, a)$$



eigen decomposition of covariance matrix

eigen vectors and eigen values.

$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

eigenvector an special vector
whose direction not change but
magnitude change

~~eigenvalues~~ How much eigenvector shrink
~~of transform~~

near by

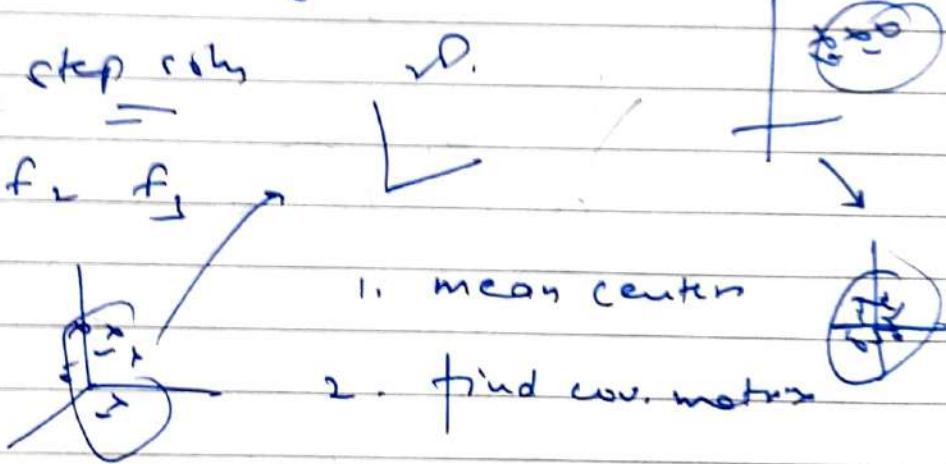
$$A \vec{v} = \lambda \vec{v}$$

↙
eigen values

5. Step by step solve

$$f_1 \ f_2 \ f_3$$

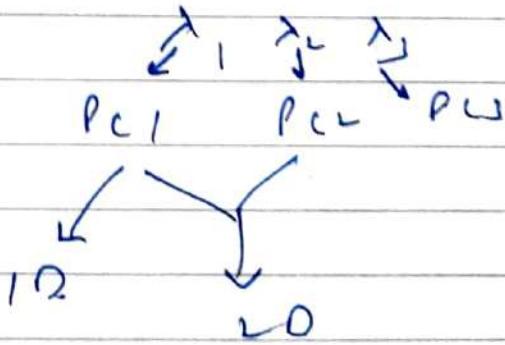
3D



$$t_1 \begin{bmatrix} v_{t_1} & c(t_1 t_2) & c(t_1 t_3) \\ t_2 & c(t_1 t_2) & v_{t_2} & c(t_2 t_3) \\ t_3 & c(t_2 t_3) & v_{t_3} \end{bmatrix}$$

w^v
matrix

3D → 2 vector



How to transform points:

$$+1 \mid t_1 \mid t_2 \mid \text{target}$$

$$U^T x$$

$$(1000, 2)$$

$$(1, 2)$$

$$(2, 1)$$

$$\begin{pmatrix} 1000 \\ 2 \\ 1 \end{pmatrix}$$

$$P_{C1} \mid \text{target}$$

find optimum number of

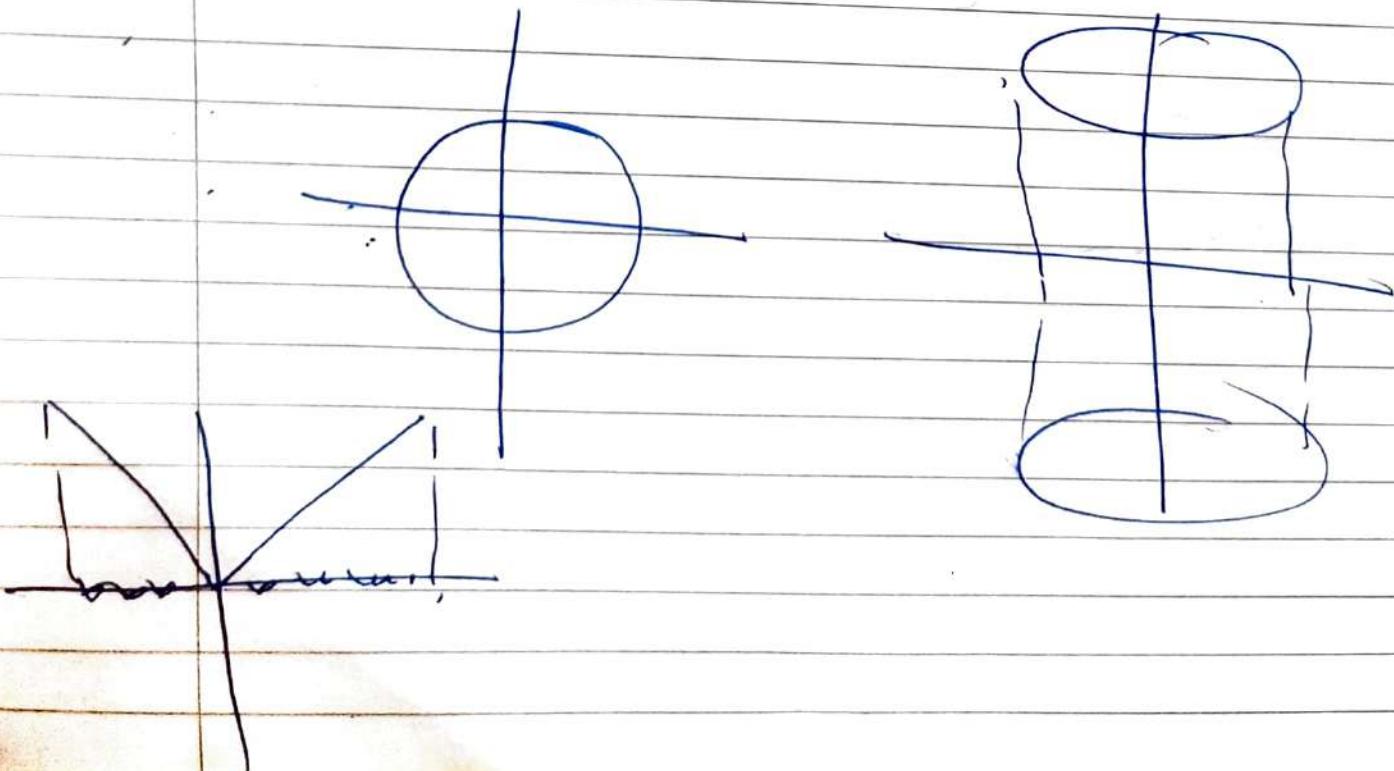
$$e_v = \lambda$$

$$784 = \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_{784}$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_{784}} \times 100 \rightarrow \text{percentage}$$

90%.

Q When PCA do not work?

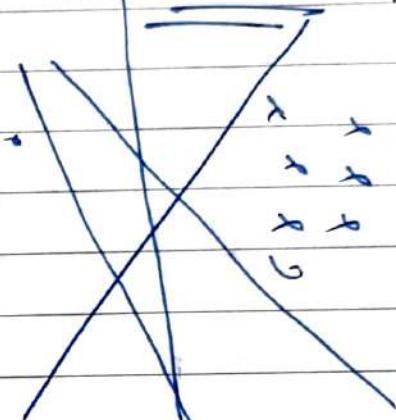


Support vector machine Ahmed Kumar

- concept of margin
- which one to select

- middle of these two datasets
- margin is maximize from both sides of the dataset
- maximum distance from closest point

which one select?



Theory of SVM

Max-margin → objective of SVM

question to be asked

1. How to measure the distance from a pt. to a plane
2. How to formalize a max-margin problem
3. How to solve max-margin problem

— Max margin

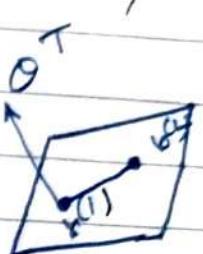
$$g_0 = \mathbf{x}'\mathbf{\theta} + \theta_0$$

\mathcal{F}_+

Theorem Let g_θ be a linear discriminant function
 & let $H \subset \{x | g_\theta(x) = 0\}$ be separating hyperplane
 then

i. $\theta / \| \theta \|_2$ is the normal vector of H i.e.
 for any $x^{(1)}, x^{(2)} \in H$ we have $\theta^T (x^{(1)} - x^{(2)}) = 0$

(ii) for any $x^{(0)}$, the distance b/w $x^{(0)}$ & H is
 $d(x^{(0)}, H) = \frac{g_\theta(x^{(0)})}{\|\theta\|_2}$



Proof: (i) We consider two points $x^{(1)} \in H$, $x^{(2)} \in H$ since both the points are on the hyperplane, we have

$$g_\theta(x^{(1)}) = 0 \quad g_\theta(x^{(2)}) = 0 \quad \text{and hence}$$

$$\theta^T (x^{(1)}) = 0$$

$$\theta^T x^{(1)} + \theta_0 = 0, \quad \theta^T x^{(2)} + \theta_0 = 0$$

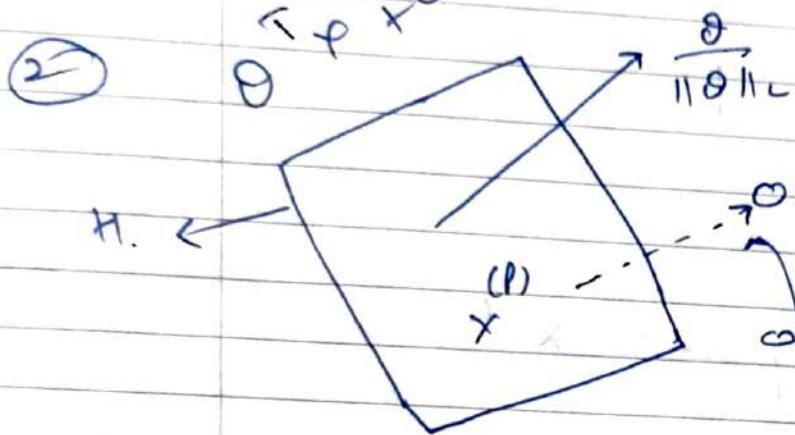
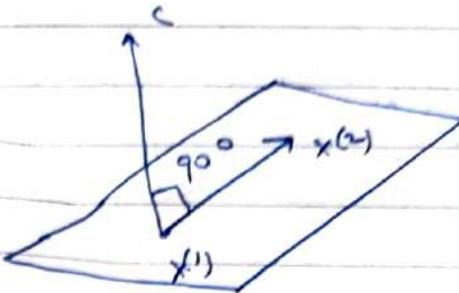
$$\rightarrow \theta^T x^{(1)} + \theta_0 = \theta^T x^{(2)} + \theta_0$$

$$\Rightarrow \theta^T (x^{(1)} - x^{(2)}) = 0$$

H prove θ is normal vector to the surface

$\frac{1}{\|\theta\|}$ is
 $\|\theta\|$.

(For every it's normal direction $\|\theta\|$)



$$g_\theta(x^{(P)}) = 0$$

$$g_\theta(x^{(0)}) = 0$$

③ Pick a point $x^{(P)}$ on H.

④ $x^{(P)}$ is the closest point orthogonal projection to $x^{(0)}$ (shortest distance)

$x^{(0)} - x^{(P)}$ is the normal direction, so for some

scalar $\beta > 0$

$$x^{(0)} - x^{(P)} = \beta \frac{\theta}{\|\theta\|_L}$$

→ As $x^{(P)}$ lies on H i.e. $g_\theta(x^{(P)}) = \theta^T x^{(P)} + \theta_0 = 0$

$$g_\theta(x^{(0)}) = \theta^T x^{(0)} + \theta_0$$

$$= \theta^T \left[x^{(0)} + \beta \frac{\theta}{\|\theta\|_L} \right] + \theta_0$$

$$= \theta^T x^{(0)} + \delta_0 + \beta \frac{\theta^T \theta}{\|\theta\|_2}$$

$$g_\theta(x^{(0)}) = g(x^{(P)}) + \beta \frac{\|\theta\|^2}{\|\theta\|_2} = \alpha + \beta \|\theta\|_2$$

$$\beta = \frac{g_\theta(x^{(0)})}{\|\theta\|_2} \quad \text{proved.}$$

So the closest point $x^{(P)}$ is

$$x^{(P)} = x^{(0)} - \frac{\beta \theta}{\|\theta\|_2}$$

$$x^{(P)} = x^{(0)} - \frac{g_\theta(x^{(0)})}{\|\theta\|_2} \cdot \frac{\theta}{\|\theta\|_2}$$

$\xleftarrow{\text{distance}}$ $\xleftarrow{\text{normal vector}}$

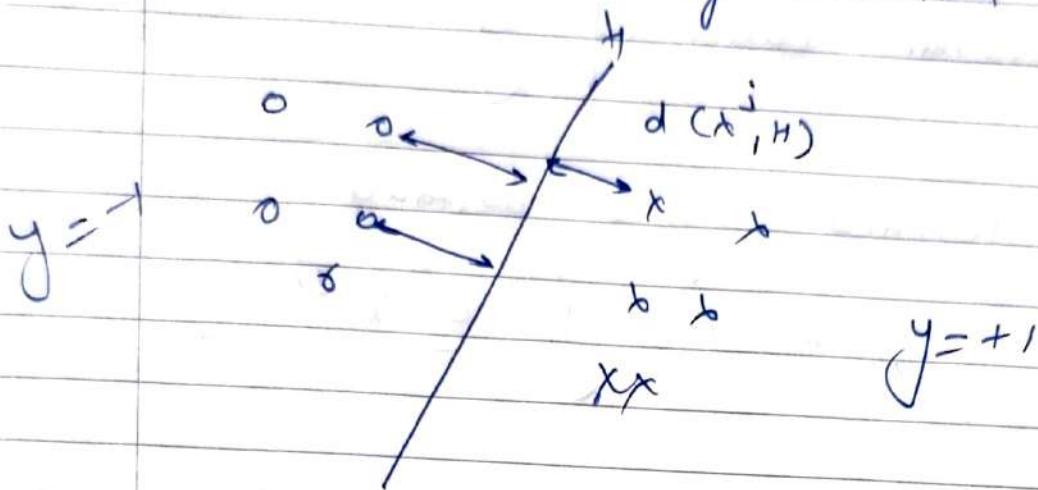
Un signed Distance

⑥ we define the distance by a date point $x^{(j)}$ & a separating hyperplane θ

$$d(x^{(j)}, \theta) = \frac{|g_\theta(x^{(j)})|}{\|\theta\|_2} = \frac{|\theta^T x^{(j)} + \delta_0|}{\|\theta\|_2}$$

$$= \frac{|\theta^T x^{(j)} + \theta_0|}{\|\theta\|_2}$$

$d(x^{(j)}, H)$ is unsigned distance



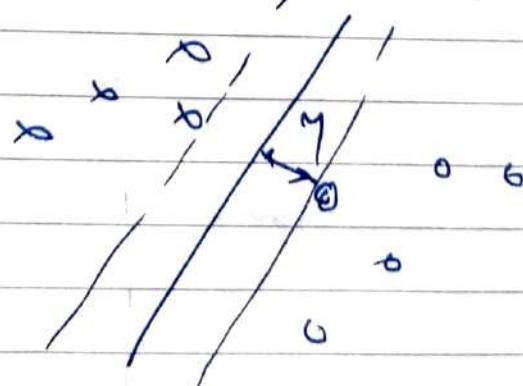
Margin \Leftrightarrow smallest distance among all the unsigned distances

$$\gamma = \min_{j=1, 2, \dots, m} d(x^{(j)}, H)$$

④ γ is a function of θ & θ_0 . (H is function of θ)

⑤ $\gamma > 0$ & is zero when $x^{(j)}$ is on the boundary

abs value



Signed distance

- * $d(x^{(j)}, H)$ is unsigned distance
- * so y does not tell whether the point $x^{(j)}$ is correctly classified or not
- * Assume labels of $y^{(j)} = \{-1, 1\}$

- * we defined signed distance

$$d_{\text{signed}}(x^{(j)}, H) = \frac{y^{(j)} \theta^T x^{(j)} + \theta_0}{\|\theta\|_2}$$

$$= \begin{cases} \geq 0 & \text{correctly classify } x^{(j)} \\ < 0 & \text{incorrectly classify } x^{(j)} \end{cases}$$

Recall perceptron loss

$$L = \max(0, -y^{(j)} (\theta^T x^{(j)} + \theta_0))$$

Unsigned

$$d > 0$$

* d is just a distance

Signed

$$d, < 0$$

* it is a distance plus whether on the correct side or not

max Margin objective

(y)

* Assume data set is linearly separable

It means

$$y^{(j)} \left(\frac{\theta^T x^{(j)} + \theta_0}{\|\theta\|_2} \right) \geq y \quad j = 1, 2, \dots, n$$

Condition ensure
all data is correctly
classified atleast
greater than the
margin y

→ Training samples are
correctly classified

→ All samples are atleast y from the boundary

* Max margin classifier is

$$\begin{matrix} \text{maximize}_{\theta, \theta_0} & y \\ \text{subjected to} & \end{matrix}$$

$$y^{(j)} \left(\frac{\theta^T x^{(j)} + \theta_0}{\|\theta\|_2} \right) \geq y$$

This constraint problem optimization problem
is support vector machine problem.

Problem

$$\begin{matrix} \text{maximize}_{\theta, \theta_0} & y \\ \text{subjected to} & \end{matrix}$$

$$\text{subject to } \gamma^{(j)} \left(\frac{\theta^T x^{(j)} + \theta_0}{\|\theta\|_2} \right) \geq \gamma$$

- * This optimization problem is not easy
 - * γ depends on (θ, θ_0)
 - * The term $\frac{1}{\|\theta\|_2}$ is non-linear
- $\sqrt{\theta_1^2 + \theta_2^2 + \dots + \theta_n^2}$

Convert it into better optimization prob,

scaling

① let $x^{(m)}$ be point closest to the boundary.

② Unsigned distance $\tilde{\gamma} = |\theta^T x^{(m)} + \theta_0|$

$$\gamma = \frac{|\theta^T x^{(m)} + \theta_0|}{\|\theta\|_2} = \frac{\tilde{\gamma}}{\|\theta\|_2}$$

$$\begin{cases} \text{maximize}_{\theta, \theta_0} & \gamma \\ \text{subject to} & \end{cases}$$

$$\gamma^{(j)} \left(\frac{\theta^T x^{(j)} + \theta_0}{\|\theta\|_2} \right) \geq \gamma$$

$$\left. \begin{array}{l} \text{maximize}_{\theta, \phi_0} \quad \tilde{\gamma} \\ \text{subject to} \quad \frac{1}{\|\theta\|_2} \end{array} \right\}$$

$$y^{(j)} \left(\frac{\theta^T y^{(j)} + \phi_0}{\|\theta\|_2} \right) \geq \tilde{\gamma}$$

$\frac{1}{\|\theta\|_2}$ goes to the objective function

$$\text{maximize}_{\theta, \phi_0} \quad \frac{\tilde{\gamma}}{\|\theta\|_2}$$

$$\text{subject to} \quad y^{(j)} \left(\theta^T x^{(j)} + \phi_0 \right) \geq \tilde{\gamma}$$

$j = 1, 2, \dots, n$

→ eliminate $\tilde{\gamma}$
by normalizing

$$\text{maximize}_{\theta, \phi_0}$$

$$\frac{\theta_1 \phi_0}{\tilde{\gamma}}$$

$$\frac{1}{\|\theta\|_2}$$

subject to

$$y^{(j)} \left(\left(\frac{\theta}{\tilde{\gamma}} \right)^T x^{(j)} + \left(\frac{\phi_0}{\tilde{\gamma}} \right) \right) \geq 1$$

$j = 1, 2, \dots, n$

we can define $\theta \leftarrow \frac{\theta}{\gamma}$, $\theta_0 \leftarrow \frac{\theta_0}{\gamma}$

maximize $\frac{1}{2} \|\theta\|^2$
 θ, θ_0

subject to $y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1$
 $\forall j = 1, 2, \dots$

so γ is eliminated

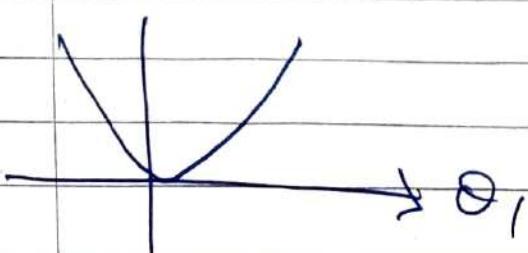
maximize $\frac{1}{2} \|\theta\|^2$
 θ, θ_0

subject to $y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1$
 $\forall j = 1, 2, \dots$

minimize $\frac{1}{2} \|\theta\|^2$ ~~subject~~

subject to $y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1$

$\forall j = 1, 2, \dots$



- * This is a quadratic minimization with linear constraints
- * It is convex
- * Solution to this is called support vector machine

(Hard - sum)

L-H It is a quadratic programming problem

$$\min_u \frac{1}{2} u^T Q u + p^T u$$

$$A u \geq c$$

solt. is $u^* = QP(Q, p, A, c)$. solved using QP solver

$$\left\{ \begin{array}{l} \text{minimize}_{\theta, \theta_0} \frac{1}{2} \|\theta\|^2 \\ \text{subject to } y^{(i)}(\theta_x^T x^{(i)} + \theta_0) \geq 1 \end{array} \right.$$

$$\text{subject to } y^{(i)}(\theta_x^T x^{(i)} + \theta_0) \geq 1$$

$$i = 1, 2, \dots, n$$

How to solve this

Approach :- Using Lagrangian function

We want to solve the following Inequality constrained optimization problem

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\text{subject to } g_i(x) \geq 0, i = 1 \dots m$$

$$h(x) = 0, i = 1 \dots p$$

where $+ h_i$, g_i are differentiable and convex and h_i is affine

$$h(x) = a^T x + b \text{ (affine)}$$

It requires a Lagrangian function

$$L(x, \alpha, \beta) = f(x) - \sum_{i=1}^m d_i g_i(x) - \sum_{j=1}^k \beta_j h_j(x)$$

Here $d_i \in \mathbb{R}^n$, $\beta_j \in \mathbb{R}^n$ are called as 'Lagrange multipliers' or dual variables, x is called as primal variable. Here $d_i \geq 0$ and no constraint on β_j .

If (x^*, α^*, β^*) is the soln. to the constrained optimization then all following conditions should hold

Karush - Kuhn - Tucker (KKT conditions)

$$\frac{\partial L}{\partial x}(x^*, \alpha^*, \beta^*) = 0$$

① stationary condition

② The primal variable should be stationary

$$(1) g_i(x^*) \geq 0 \quad \& \quad g_i(x^*) = 0 \text{ for all } i, j$$

(*) Primal feasibility

(**) Ensure that constrained constraints are satisfied

$$(11) \alpha_i^* \geq 0 \text{ for all } i \neq j$$

(*) Dual feasibility

(**) Require $\alpha_i^* \geq 0$ but no constraint

$$\beta_{ij}^*$$

$$(14) \alpha_i^* g_i(x^*) + \sum_j \beta_{ij}^* g_j(x^*) = 0 \text{ for all } i \neq j$$

(*) complementary slackness

(**) either $\alpha_i^* = 0$ or $g_i(x^*) = 0$ (or both)

$$\begin{array}{c} \alpha_i^* g_i(x^*) > 0 \\ \downarrow \\ \alpha_i^* > 0 \end{array} \quad \begin{array}{c} \downarrow \\ g_i(x^*) \end{array} \quad \begin{array}{c} \downarrow \\ g_i(x^*) = 0 \end{array}$$

Lagrangian function is given as

$$L(\theta, \theta_0, \lambda) = \frac{1}{2} \|\theta\|_2^2 + \sum_{j=1}^m \lambda_j \left[1 - g_j^{(j)} (\theta^T x + \theta_0) \right]$$

solution \rightarrow happens at

$$\nabla_{(\theta, \theta_0)} L(\theta, \theta_0, \lambda) = 0 \quad \& \quad \nabla_\lambda L(\theta, \theta_0, \lambda) = 0$$

\leftarrow primal variables

(a)

Let re-write the lagrangian function again,

$$L(\theta, \theta_0, \lambda) = \frac{1}{2} \|\theta\|^2 + \sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)]$$

$\Leftarrow \quad \Rightarrow$

complementary condition:-

- ① If $\lambda_j > 0$, then $1 - y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 0$
- ② If $\lambda_j = 0$ then $[1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)] < 0$

useful in picking
up every

so if we take $\nabla_{\lambda} L(\theta, \theta_0, \lambda) = 0$

$$\nabla_{\lambda} L(\theta, \theta_0, \lambda) = 0 + \nabla_{\lambda} \left[\sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)] \right]$$

$\Leftarrow \quad \Rightarrow$

$= 0$

only possible if

$$\sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)]$$

maximum or minimum

So here we point ^{because}, it is linear in λ

$$\lambda_j [1 - y^{(j)} (\theta^T x + \theta_0)] \rightarrow \text{It is bounded minimum to find } \nabla_{\lambda}(\cdot) \text{ we must maximize this argument}$$

\longleftrightarrow

By doing this we are finding the maximum pt. in the dual space

→ we can partition the optimization problem into two steps

① Let λ^* be the maximizer

$$\lambda^* = \underset{\lambda \geq 0}{\operatorname{argmax}} \left[\sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x + \theta_0)] \right]$$

② Then primal problem is (substituting result of ① in ② Page 1)

$$\min_{\theta, \theta_0} L(\omega, \omega_0, \lambda^*) = \min_{\theta, \theta_0} \left\{ \frac{1}{2} \|\theta\|_2^2 + \max_{\lambda \geq 0} \left[\sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x + \theta_0)] \right] \right\}$$

$$\min_{\theta, \theta_0} L(\omega, \omega_0, \lambda^*) = \min_{\theta, \theta_0} \left\{ \max_{\lambda \geq 0} L(\theta, \theta_0, \lambda^*) \right\}$$

This is min-max problem

$$\min_{\theta_0} \max_{\theta, \lambda \geq 0} L(\theta, \theta_0, \lambda)$$

min max problem

Support vector machine! Duality

Two Reasons to go (duality) dual space!

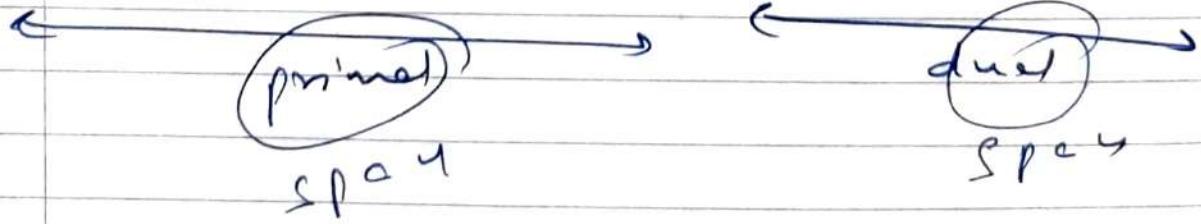
(1) can see why it is called as support vector

second will be discussed later

* Our problem is quadratic programming

* strong duality hold for quadratic programming

$$\min_{\theta, \theta_0} \max_{\lambda \geq 0} L(\theta, \theta_0, \lambda) = \max_{\lambda \geq 0} \min_{\theta, \theta_0} L(\theta, \theta_0, \lambda)$$



$$\hat{\theta}, \hat{\theta}_0 = \max_{\lambda \geq 0} \min_{\theta, \theta_0} L(\theta, \theta_0, \lambda)$$

$$\left\{ \begin{array}{l} \text{minimize}_{\theta, \theta_0} \frac{1}{2} \|\theta\|_2^2 \\ \text{subject to} \end{array} \right.$$

$$y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1, \quad \forall j = 1, 2, \dots$$

Lagrangian can be written as

$$L(\theta, \theta_0, \lambda) = \frac{1}{2} \|\theta\|_2^2 + \sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)]$$

In Dual setting, we will minimize w.r.t θ & θ_0 first

$$\nabla_{\theta} L(\theta, \theta_0, \lambda) = \theta - \sum_{j=1}^m \lambda_j y^{(j)} x^{(j)} = 0$$

$$\nabla_{\theta_0} L(\theta, \theta_0, \lambda) = \sum_{j=1}^m \lambda_j y^{(j)} = 0$$

stationary condition i.e KKT condition 1st

interesting observations

$$\theta = \sum_{j=1}^m \lambda_j y^{(j)} x^{(j)}$$

hyperplane is written as a linear combination of Data points.

- Here λ_j is either zero or greater than zero
- (*) Comes from complementary condition
 - (**) Either use data point or not use it

Lagrangian function

$$\begin{aligned} L(\theta, \theta_0, \lambda) &= \frac{1}{2} \|\theta\|_2^2 + \sum_{j=1}^m \lambda_j [1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)] \\ &= \frac{1}{2} \left\| \sum_{j=1}^m \lambda_j y^{(j)} x^{(j)} \right\|_2^2 \\ &\quad + \sum_{j=1}^m \lambda_j [1 - y^{(j)} \left(\left(\sum_{i=1}^m \lambda_i y^{(i)} x^{(i)} \right)^T x^{(j)} + \theta_0 \right)] \end{aligned}$$

$$-\sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} (x^{(i)})^T x^{(j)} - \sum_{j=1}^m \lambda_j y^{(j)} \theta$$

\longleftrightarrow
stationary condition

$$L(\theta, \theta_0, \lambda) = -\frac{1}{2} \sum_{j=1}^m \sum_{i=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} (x^{(i)})^T x^{(j)} + \sum_{j=1}^m \lambda_j$$

Now we do max ~~and minimize~~ as m_1 is already evaluated for dual setting

$$\max_{\theta, \theta_0} \min_{\lambda} L(\theta, \theta_0, \lambda)$$

S. dual problem

maximize $\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \lambda_i y^{(i)} y^{(j)} (\mathbf{x}^{(i) T} \mathbf{x}^{(j)}) + \sum_{i=1}^m \lambda_i$

subject to $\sum_{j=1}^n \lambda_j y^{(j)} = 0$

→ Solution to this will be λ^*

complementary condition: $\lambda_j^* (1 - y^{(j)} (\mathbf{\theta}^* T \mathbf{x}^{(j)} + \delta_0^*)) = 0$

* If $\lambda_j^* > 0$ then $1 - y^{(j)} (\mathbf{\theta}^* T \mathbf{x}^{(j)} + \delta_0^*) = 0 \Rightarrow$

for all the pts. lying on the boundary
 $y^{(j)} (\mathbf{\theta}^* T \mathbf{x}^{(j)} + \delta_0^*) = 1$

the value of $\lambda_j^* > 0$

* If $\lambda_j^* = 0$, then $[1 - y^{(j)} (\mathbf{\theta}^* T \mathbf{x}^{(j)} + \delta_0^*)] < 0$

for all points correctly classified i.e. lying
 in the interior of the boundary

$$\{y^{(j)} [\mathbf{\theta}^* T \mathbf{x}^{(j)} + \delta_0^*] > 1\}$$

the value of $\lambda_j^* = 0$

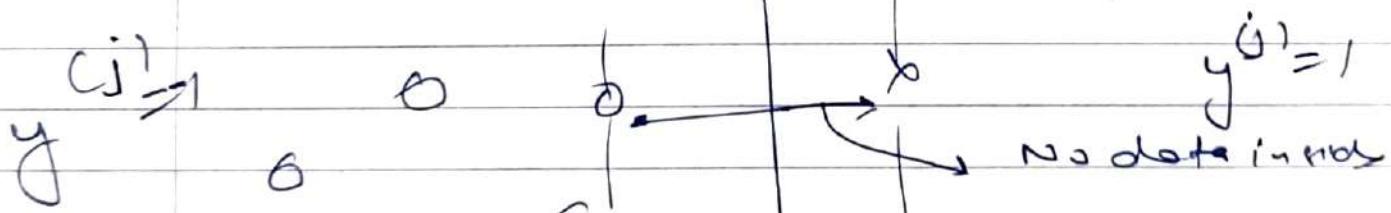
We can define sum as
 $w = [j \mid \lambda_j^* > 0] \Rightarrow x^{(j)} \text{ is on the margin}$

so the optimal wts. $\theta^* \sum_{j \in w} \lambda_j^{(j)} x^{(j)}$

→ The optimal separating hyperplane depends on the support vector points, closest to the decision boundary. It is different from least square solution, which depend on all training data

$$y \geq 0 \quad y^{(j)}(\theta^T x^{(j)} + \theta_0) \geq 1$$

$$y \frac{L-15}{\theta^T x^{(j)} + \theta_0} \geq 1$$



$$\theta^T x^{(j)} + \theta_0 = -1$$

$$a, y^{(j)} = -1$$

$$\theta^T x^{(j)} + \theta_0 = 1$$

$$a, y^{(j)} = 1$$

$$\theta^T x^{(j)} + \theta_0 = 0$$

$$y^{(j)} (\theta^T x^{(j)} + \theta_0) = 1$$

What about θ_0^* ?

- ① Pick any support vector $x^+ \in$ +ve class
 $\& x^- \in$ -ve class

then we have

$$\theta^{*T} x^+ + \theta_0^* = 1 \quad * \quad \theta^{*T} x^- + \theta_0^* = -1$$

$$\Rightarrow \theta_0^* = -\frac{\theta^{*T} (x^+ + x^-)}{2}$$

Support vector machine :- Duality

Reason to go to dual space

- ② Dual enable kernel trick.

$$\begin{aligned} &\text{maximize}_{\lambda} \quad -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} (x^{(i)T})^T x^{(j)} \\ &x \geq 0 \end{aligned}$$

$$\text{subject to } \sum_{j=1}^m \lambda_j y^{(j)} = 0$$

In matrix-vector format

$$\begin{aligned} &\text{max} \quad -\frac{1}{2} \lambda^T P \lambda + \mathbf{1}^T \lambda \\ &\lambda \geq 0 \end{aligned}$$

$$\text{s.t. } \mathbf{y}^T \lambda = 0$$

Here

$$P = \begin{pmatrix} y^{(1)} y^{(1)} (x^{(1)})^T x^{(1)} & \cdots & y^{(1)} y^{(m)} (x^{(1)})^T x^{(m)} \\ y^{(2)} y^{(2)} (x^{(2)})^T x^{(2)} & \cdots & y^{(2)} y^{(m)} (x^{(2)})^T x^{(m)} \\ \vdots & \ddots & \vdots \\ y^{(m)} y^{(m)} (x^{(m)})^T x^{(m)} & \cdots & y^{(m)} y^{(m)} (x^{(m)})^T x^{(m)} \end{pmatrix}$$

Any observation?

Kernel trick

$$\text{maximize}_{\lambda \geq 0} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} (x^{(i)})^T x^{(j)} + \sum_{j=1}^m \lambda_j$$

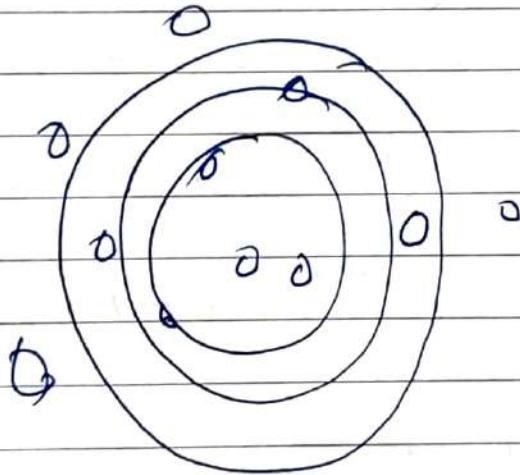
$$\text{subject to } \sum_{j=1}^m \lambda_j y^{(j)} = 0$$

$$\text{maximize}_{\lambda \geq 0} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \underbrace{\frac{g(x^{(i)})^T g(x^{(j)})}{+ \sum_{j=1}^m \lambda_j}}_{K(x^{(i)}, x^{(j)})}$$

$$\text{subject to } \sum_{j=1}^m \lambda_j y^{(j)} = 0$$

$$P = \left[\begin{array}{cccc} y^{(1)} & y^{(1)} & \dots & y^{(1)} \\ y^{(2)} & y^{(2)} & \dots & y^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ y^{(m)} & y^{(m)} & \dots & y^{(m)} \end{array} \right] K \left[\begin{array}{cc} x^{(1)}, x^{(1)} \\ x^{(2)}, x^{(2)} \\ \vdots \\ x^{(m)}, x^{(m)} \end{array} \right]$$

can we write
linear sum
here?



second order
Hard sum
(Kernel Based)

soft sum

Allow error into the margin

Hard sum :- not allow any error into the margin

$\odot \delta$

$$\mathbf{w}^T \mathbf{x}_i + b_0 = +1$$

$$\mathbf{w}^T \mathbf{x}_k + b_0 = 0$$

$$\mathbf{w}^T \mathbf{x}_j + b_0 = -1$$

\odot

\odot

$$\text{for point } G \quad y^{(a)}(\theta^T x^{(a)} + \theta_0) = -1.5$$

$$\text{as } y^{(a)} = 1 \Rightarrow \theta^T x^{(a)} + \theta_0 = -1.5$$

$$\theta^T x^{(a)} + \theta_0 = 1 - (2.5) \quad \epsilon_1$$

$$\text{for point } B \quad y^{(b)}(\theta^T x^{(b)} + \theta_0) = 0.5$$

$$\Rightarrow \theta^T x^{(b)} + \theta_0 = 1 - (0.5) \quad \epsilon_2$$

$$\text{for point } C \quad y^{(c)}(\theta^T x^{(c)} + \theta_0) = -0.5$$

$$(\theta^T x^{(c)} + \theta_0) = 1 - (1.5) \quad \epsilon_3$$

$$\text{for point } D \quad y^{(d)}(\theta^T x^{(d)} + \theta_0) = 1.5$$

$$\text{as } y^{(d)} = -1$$

$$\Rightarrow \theta^T x^{(d)} + \theta_0 = -1.5 = 1 - (2.5)$$

ε -ft margin

→ changing hard margin $y^{(j)}(\theta^T x^{(j)} + \theta_0) \geq 1$ to

Hard margin

ε -ft margin $y^{(j)}(\theta^T x^{(j)} + \theta_0) \geq 1 - \epsilon_j, \epsilon_j > 0$

if $\epsilon_j > 1$

$\sum \epsilon_j$

optimization framework

$$\min_{\theta, \theta_0, \epsilon} \frac{1}{2} \|\theta\|_2^2 + C \|\epsilon\|_1$$

$$\text{s.t. } y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1 - \epsilon_j$$

$$\epsilon_j \geq 0 \quad \text{for } j = 1, 2, \dots, n.$$

If $C \uparrow$, $\epsilon \rightarrow 0 \rightarrow$ tendency to make mistake
on training data \rightarrow overfitting

~~As~~ As $C \downarrow$, $\epsilon \uparrow$ i.e. tendency to allow
mistake to enter the margin \uparrow \rightarrow high bias
(under fit)

Other variant:

$$\min_{\theta, \theta_0, \epsilon} \frac{1}{2} \|\theta\|_2^2 + C \|\epsilon\|_1$$

$$\text{s.t. } y^{(j)} (\theta^T x^{(j)} + \theta_0) \geq 1 - \epsilon_j$$

$$\epsilon_j \geq 0 \quad \text{for } j = 1, 2, \dots, n$$

\rightarrow Encourage ϵ to be sparse

\rightarrow only few samples are allowed to lie
in the margin

connection with Perceptron Algorithm

In soft margin SVM, $\epsilon_j \geq 0$ & $y^{(j)}(\theta^T x^{(j)} + \theta_0) \geq 1 - \epsilon_j$

$\Rightarrow \epsilon_j \geq 0$ and $\epsilon_j \geq 1 - y^{(j)}(\theta^T x^{(j)} + \theta_0)$

These conditions can be combined together as

$$\epsilon_j \geq \max[0, 1 - y^{(j)}(\theta^T x^{(j)} + \theta_0)]$$

If SVM (soft) with L norm is used

$$L(\theta, \theta_0, \epsilon) = \frac{1}{2} \|\theta\|_2^2 + C \sum_{j=1}^m \epsilon_j$$

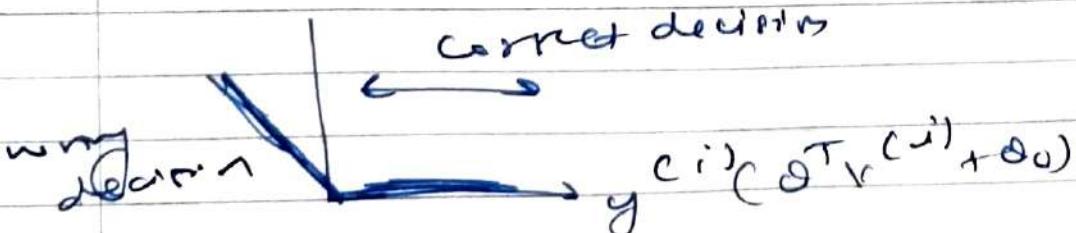
$$\geq \frac{1}{2} \|\theta\|_2^2 + C \sum_{j=1}^m \max[0, 1 - y^{(j)}(\theta^T x^{(j)} + \theta_0)]$$

so the training loss becomes

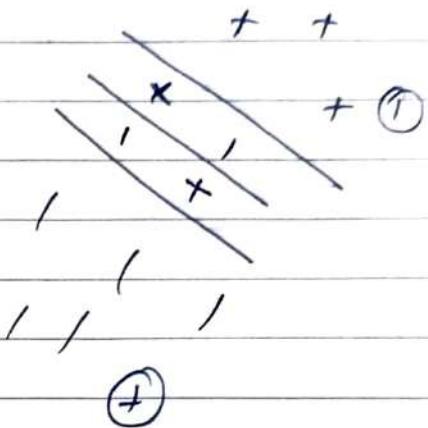
$$L(\theta, \theta_0) = \sum_{j=1}^m \max[0, 1 - y^{(j)}(\theta^T x^{(j)} + \theta_0)] + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\lambda = \frac{C}{2}$$

$$\max[0, 1 - y^{(j)}(\theta^T x^{(j)} + \theta_0)]$$



$$\max(0, 1 - y^{(i)}(\theta^T \mathbf{x}^{(i)} + \theta_0))$$



L₂ loss

$$\sum_{i=1}^m \{y^{(i)} - \theta^T \mathbf{x}^{(i)} - \theta_0\}^2$$

As $y^{(i)} = \pm 1$ in case of classification

$$\sum_{i=1}^m \{1 - y^{(i)}(\theta^T \mathbf{x}^{(i)} + \theta_0)\}^2$$

by taking $y^{(i)}$ outside possibilities
to squaring

Rewriting hinge loss

$$L(\theta, \theta_0) = \sum_{i=1}^m \max[0, 1 - y^{(i)}(\theta^T \mathbf{x}^{(i)} + \theta_0)] + \frac{\lambda}{2} \|\theta\|^2$$

$$\lambda = \frac{1}{C}$$

If $\lambda \rightarrow 0 \Rightarrow C \rightarrow \infty$

So loss becomes

$$L(\theta, \theta_0) = \sum_{j=1}^m \max[0, 1 - y^{(j)} (\theta^T x^{(j)} + \theta_0)]$$

Perception loss.

$$L = \sum_{j=1}^m \max[0, y^{(j)} x^{(j)} \theta]$$

Conclusion: sum (soft) L, norm penalty w.r.t.

$\lambda \rightarrow 0$ is perception (extreme situation)

In general sum allows regularization to be added to perception loss thus regularizing the solution & bring the benefit of norms

~~one~~ solution & bring the book.