

Machine Learning

(18/02)

Machine learning is a subfield of AI where algorithms learn patterns from data and make predictions.

The three main types of ML are:

① supervised learning (Eg: Linear regression, decision trees)

→ This involves training the model on labelled data, where input-output pairs are explicitly provided

- The algo learns to map inputs to outputs and hence can predict outcomes for new unseen data

Eg: Image recognition, identifying whether it's a dog or cat based on labeled data

Spam detection, classifying emails as spam or not spam

② unsupervised learning (Eg: K-means clustering, PCA)

→ In this, the algorithm is trained on unlabeled data and identifies patterns, clusters within the data without explicit guidance.

Eg: clustering - like grouping customers based on purchasing patterns.

Recommendation system - suggesting products based on user behaviour and similarities.

* semi-supervised learning (Eg: self-trained models)

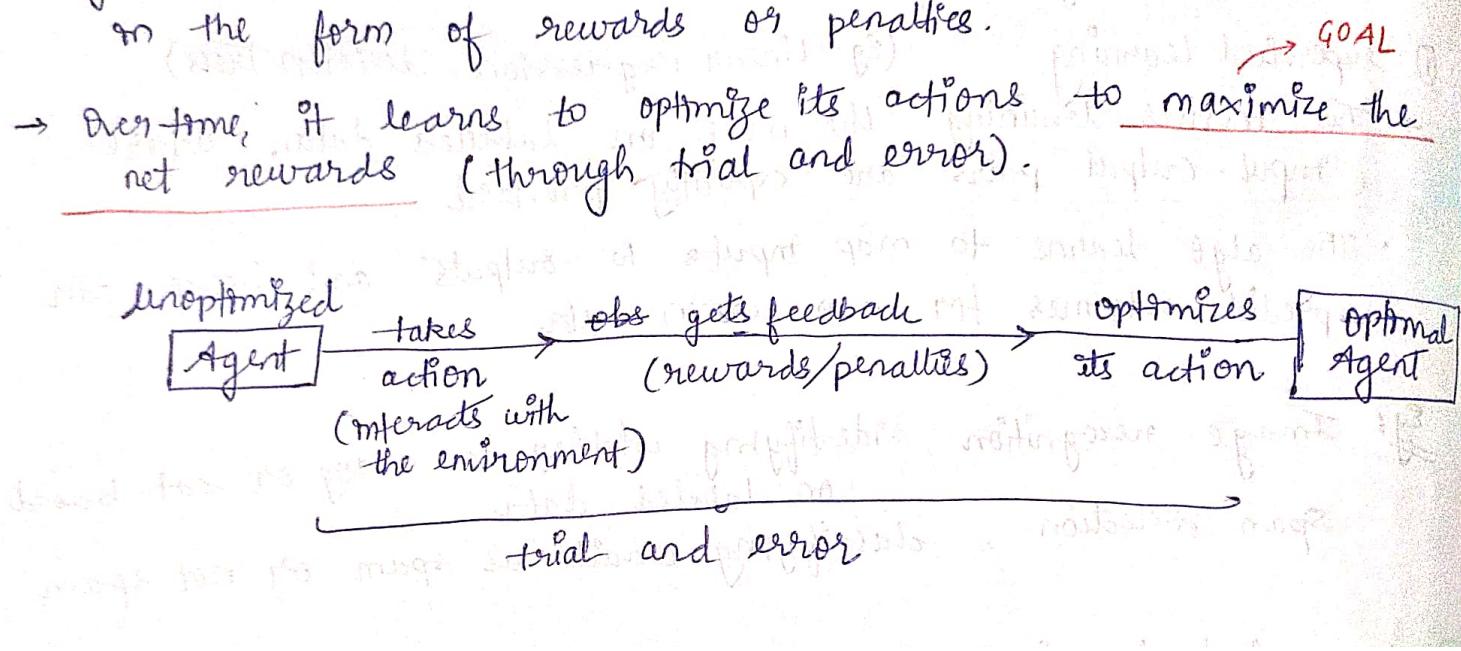
→ this combines a small amount of labeled data with a large amount of unlabeled data.

• Bridges the gap between unsupervised & supervised learning by leveraging both types of data.

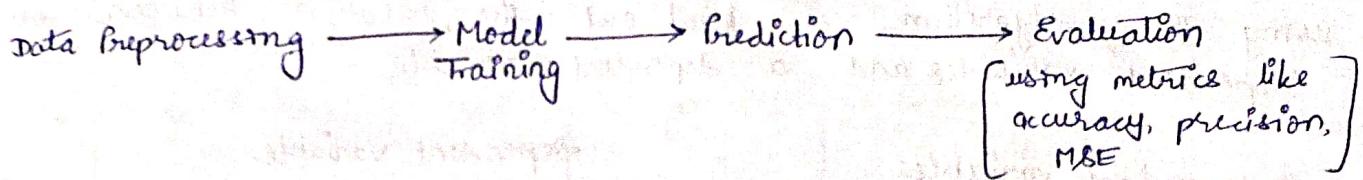
Eg: speech recognition - training models for speech-to-text conversion using limited labeled dat audio clips alongside vast amount of unlabeled audio.

③ Reinforcement learning (RL)

- In this, a agent learns to make decisions by interacting with the environment to achieve a specific goal.
- Agent takes actions, observes the outcomes, and gets feedback in the form of rewards or penalties.
 - Over time, it learns to optimize its actions to maximize the net rewards (through trial and error).



~~Supervised Learning~~



Types

① Regression

→ Regression is used when the target variable is continuous, meaning it can take any value within a range.

- input data & output data, both need to be continuous.
 ↓
 continuous or categorical continuous
 not

e.g. Predicting house prices
Forecasting stock prices.

- * If you need a number,
→ regression
- * If you need a category,
→ classification

Common regression algos :

- i) Linear, Multiple
- ii) Polynomial
- iii) Lasso
- iv) Ridge

② Classification

→ Classification is used when the target variable is categorical, meaning it belongs to a predefined class or category.

e.g. Spam/not spam
Digit Recognition

Common classification algos :

- i) Logistic Regression
- ii) SVM
- iii) Decision Tree
- iv) Random Forest
- v) Neural Networks

Linear Regression

- using this algorithm, we find out the relation between an independent variable and a dependent variable

Independent variables
(Regressors or predictors)

Dependent Variable
(Response)

Simple Linear Regression

- using this algo, we find the linear relation b/w the linearly dependent variable and the independent variable

$$y = mx + c$$

How to Solve

1. Least square Method

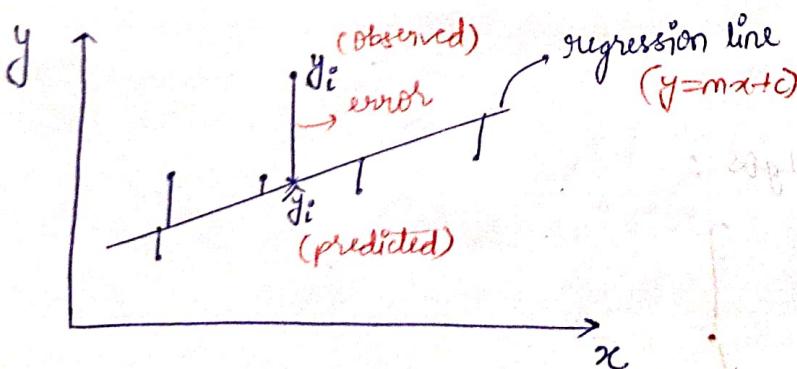
- used to find the eqn of best fitting curve to the set of data points by minimizing the sum of squared differences b/w the observed and predicted values.

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

$y_i \rightarrow$ observed
 $\hat{y}_i \rightarrow$ predicted

$x \rightarrow$ independent variable on x-coordinate

$y \rightarrow$ dependent on y-coordinate



$$\therefore \bar{y} = m\bar{x} + c$$

$$\Rightarrow c = \bar{y} - m\bar{x}$$

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

$$\Rightarrow m = \frac{n(\sum xy) - \sum x \cdot \sum y}{n \sum x^2 - (\sum x)^2}$$

Q. Find the line of best fit for data points using LSM:

$$(x, y) = (1, 3), (2, 4), (4, 8), (6, 10), (8, 15)$$

x	y	xy	x^2	y^2	$1/x^2$
1	3	3	1	2.6336	0.134
2	4	8	4	4.3106	0.0964
4	8	32	16	7.6646	0.1125
6	10	60	36	11.0186	0.0375
8	15	120	64	14.3726	0.0936
21	40	223	121		

$$m^2 \frac{5(223) - 21(40)}{5(121) - (21)^2} = \frac{275}{164} = 1.677$$

URBAN
EDGE
COLORFUL YOU

$$C = 8 - 1.677(4.2)$$

$$1) C = 8 - 7.0434 = 0.9566$$

$$\text{eqn}, y = 1.677x + 0.9566$$

Assumption of the LSM

* i. to apply LSM to predict x

1. linear relationship b/w the variables

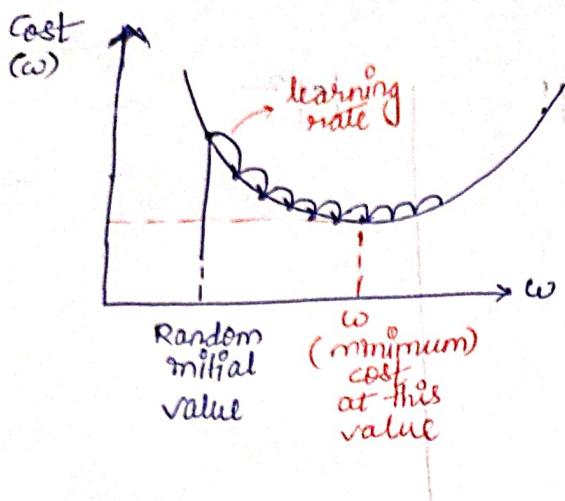
2. the observations are independent of each other.

3. variance of residual is constant with a mean of 0

4. Error are distributed normally.

2. Gradient descent method

* This is used to find the global minima of a differentiable function



Here,

Aim is to find the values of 'm' and 'c' that defines the relationship betⁿ x and y, correctly with the help of loss function, and minimize the cost function.

→ steps involved in linear regression with GD.

- ① Initialize the weight and bias randomly or start with '0'.
- ② Make prediction with the initial weight and bias.
- ③ Calculate loss function
- ④ With the help of differentiation, calculate how loss function changes with respect to weight and bias.
- ⑤ Update weight and bias term so as to minimize the loss function (cost funⁿ).

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

MSE

$$Y = mx + c$$

$$\begin{aligned}
 \frac{\partial (\text{cost } f^m)}{\partial m} &= \frac{\partial}{\partial m} \left(\frac{1}{n} \sum (y_i - \hat{y}_i)^2 \right) \\
 &= \frac{\partial}{\partial m} \left(\frac{1}{n} \sum (y_i - (mx + c))^2 \right) \\
 &= \frac{1}{n} \frac{\partial \sum}{\partial m} (y_i^2 + (mx + c)^2 - 2y_i(mx + c)) \\
 &= \frac{1}{n} \cancel{\frac{\partial}{\partial mx}} \sum (2x(mx + c) - 2y_i x) \\
 &= \frac{2}{n} \sum [y_i - (mx + c)] (-x)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial (\text{cost } f^m)}{\partial c} &= \frac{\partial}{\partial c} \left[\frac{1}{n} \sum (y_i - (mx + c))^2 \right] \\
 &= \frac{1}{n} \cdot 2 \sum (y_i - \cancel{(mx + c)}) (-1)
 \end{aligned}$$

New, update $m + c$ using the following formula

$$\begin{aligned}
 m &= m - \alpha \frac{\partial (\text{loss } f^m)}{\partial m} && \left\{ \begin{array}{l} \alpha \rightarrow \text{learning rate} \\ 0 \leq \alpha < 1 \end{array} \right. \\
 c &= c - \alpha \frac{\partial (\text{loss } f^m)}{\partial c}
 \end{aligned}$$

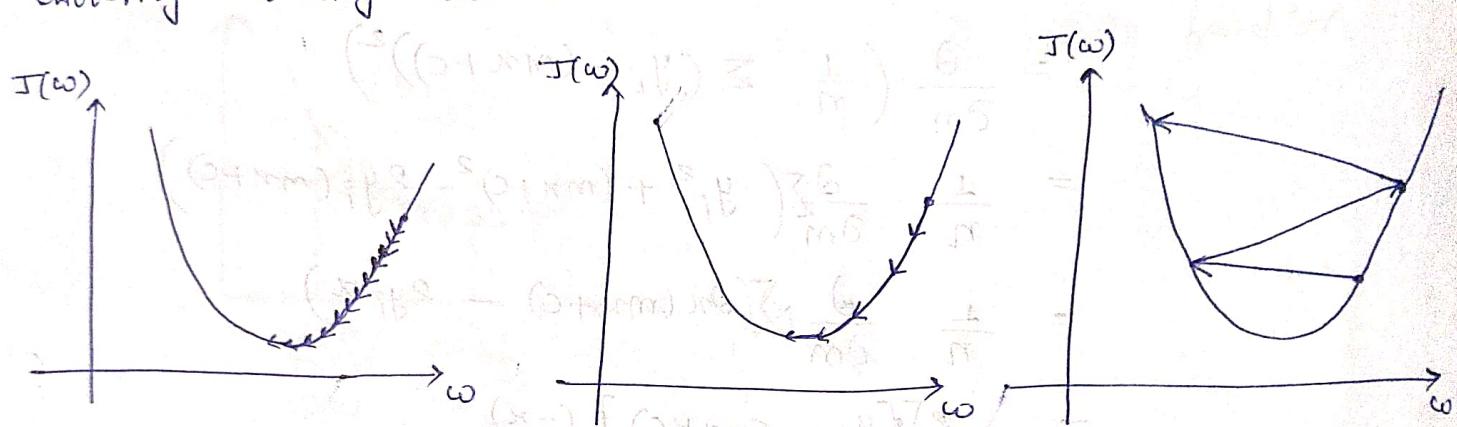
$$\therefore m = m - \alpha \left[\frac{-2}{n} \sum (y_i - (mx + c)) x \right]$$

$$\Rightarrow m = m - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) x \right]$$

$$\therefore c = c - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) \right]$$

Learning Rate ★ Rate of convergence of gradient descent depends upon the learning rate.

choosing learning rate is a matter of trial and error.



★ Too Low

[requires many iterations]

★ Just Right

[swifly reaches the min point]

★ Too High

[causes drastic update which leads to divergent behaviour]

★ α -fixed \rightarrow the algorithm might get trapped in local minima

Q. Apply G.D on following obs and observe MSE value
of 3 iterations

$$\alpha = 0.001 \text{ with } m = 10, b = 300$$

age(x) salary(y)

30	800
32	950
25	600
43	1050
50	1200
29	740
46	1100

initially $m=c=0 \quad m=10 \quad b=300$
 $\therefore \hat{y} = 10x + 300$.

$$\frac{\partial \text{cost}}{\partial m} = -\frac{2n}{n} (y_i - \hat{y}_i)$$

~~cl-a~~

$$\text{cost} = \text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$m = m - \alpha \left[-\frac{2}{n} \sum (y_i - \hat{y}_i)x \right]$$

$$c = c - \alpha \left[-\frac{2}{n} \sum (y_i - \hat{y}_i) \right]$$

$$\sum (y_i - \hat{y}_i) = [800 - 600] + [950 - 620] + [600 - 55] \\ [1050 - 730] + [1200 - 800] + [740 - 590] + [1100 - 760] \\ = 1290$$

$$MSE = \frac{1}{7} \sum (y_i - \hat{y}_i)^2 = \cancel{45} \rightarrow 28.57 \quad 78842.8$$

$$m_{\text{new}} = 10 - 0.0001 \left[-\frac{2}{7} (71560) \right] \\ = 10 + 2.04457 \\ = 12.04457$$

$$b_{\text{new}} = 300 - 0.0001 \left[-\frac{2}{7} (1290) \right] \\ = 299.998 \quad 300.05114$$

$$\hat{y}_i = 12.04457x + 300.05114 \quad (\text{2nd epoch})$$

$$\sum (y_i - \hat{y}_i) \\ = [800 - 661.388] + [950 - 685.477] + [600 - 601.165] \\ + [1050 - 814.967] + [1200 - 902.279] + [740 - 649.39] \\ + [1100 - 854.101] \\ = 1268.283$$

$$MSE = \frac{1}{7} (300349.72) = 42907.040$$

$$m_{\text{new}} = 12.04457 - 0.0001 \left[-\frac{2}{7} (51397.93) \right] \\ = 13.513$$

$$b_{\text{new}} = 300.05114 - 0.0001 \left[-\frac{2}{7} (1268.283) \right] = 300.087$$

Multiple Linear Regression

In multiple linear regression, we model relationship b/w one dependent variable and different multiple independent variables.

Ways to solve :

① Gradient Descent Method (Open form equation)

For two independent variables i.e. x_1, x_2 mean squared error

$$y = \omega_1 x_1 + \omega_2 x_2 + b$$

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

$$\omega_1 = \omega_1 - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) x_{1i} \right]$$

$$\omega_2 = \omega_2 - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) x_{2i} \right]$$

$$b = b - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) \right]$$

$$y = (\sum x_i \omega_i) + b$$

Q): Consider the following dataset. Find out linear relation between the response var 'y' and two predictor variables x_1 & x_2 using GD.

y	x_1	x_2	\hat{y}	
140	60	22	147.5	
155	62	25	149.5	
159	67	24	166.1	
179	70	20	179.9	
192	71	15	188.9	
200	72	14	193.1	
212	75	14	202.1	
215	78	11	214.7	

$\omega_1 = 3$
 $\omega_2 = -1.2$
~~alpha~~
 $b = -6.1$
 $\alpha = 0.001$
 $epoch = 2$

$$\Rightarrow y = \omega_1 x_1 + \omega_2 x_2 + b$$

$$y = 3x_1 - 1.2x_2 - 6.1$$

We need to update ω_1 , ω_2 & b .

$$\omega_1 \text{ (new)} = \omega_1 - \alpha \left[\frac{-2}{n} \sum (y_i - \hat{y}_i) x_{1i} \right]$$

$$\begin{aligned} \sum (y_i - \hat{y}_i) &= (140 - 147.5) + (155 - 149.5) + (159 - 166.1) + (179 - 179.9) \\ &\quad + (192 - 188.9) + (200 - 193.1) + (212 - 202.1) \\ &\quad + (215 - 214.7) \\ &= (-7.5) + (5.5) + (-7.1) + (-0.9) + (3.1) \\ &\quad + (6.9) + (9.9) \\ &\quad + (0.3) \\ &= 10.2 \end{aligned}$$

$$\begin{aligned} \sum (y_i - \hat{y}_i) x_{1i} &= (-7.5 \times \cancel{147.5}) + (5.5 \times \cancel{149.9}) + (-7.1 \times \cancel{166.1}) \\ &\quad + (-0.9 \times 70) + (3.1 \times 71) + (6.9 \times 72) + (9.9 \times 75) \\ &\quad + (0.3 \times 78) \\ &= 835.1 \end{aligned}$$

$$\sum (y_i - \hat{y}_i) x_{2i} = 69.1$$

$$\omega_1 = 3 - 0.001 \left[\frac{-2}{8} \times 835.1 \right] = 3.208$$

~~$$\omega_2 = 1.2 - (0.001 \times (\frac{-2}{8} \times 69.1))$$~~

$$\omega_2 = 1.2 - (0.001) \times \left(\frac{-2}{8} \times 69.1 \right) = -1.182$$

$$b = -6.1 + 0.001 \left(\frac{-2}{8} \times 10.2 \right) = -6.09$$

$$\therefore \boxed{y = 3.208x_1 - 1.182x_2 - 6.09}$$

Multiple Linear Regression using Least Square Method

MLR can be expressed as $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$

$$y = X\beta$$

$(y = X\beta + \epsilon)$ error will always be there

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nk} \end{bmatrix}_{(N \times k)} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix} \quad (1)$$

$(N \gg k)$

If we consider β_0 then eq(1) can be expressed as:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{Nk} \end{bmatrix}_{(N \times (k+1))} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

$(N \gg k+1)$

URBAN
EDGE
COLORFUL YOU

To estimate the parameters of β , we need to solve n equations and it will be difficult to solve this way because

$N \gg k+1$

consider the sum of sq of residuals, in order to estimate β we have to minimize

$$\sum_{i=1}^N e_i^2$$

$$e_i^2 = e'e = [e_1 \ e_2 \ \dots \ e_N] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix}$$

This can be written as a optimization problem or equation i.e.

$$\min_{\beta} e'e$$

$$\frac{\partial e'e}{\partial \beta} = 0.$$

$$[\because e = y - X\beta]$$

$$\Rightarrow \frac{\partial}{\partial \beta} [(y - X\beta)'(y - X\beta)]$$

$$\Rightarrow \frac{\partial}{\partial \beta} (y'y - y'X\beta - (X\beta)'y + (X\beta)'X\beta) = 0$$

$$\Rightarrow \frac{\partial}{\partial \beta} (y'y - (y'X\beta)' - X'\beta'y + (X\beta)'X\beta) = 0$$

$$\Rightarrow \frac{\partial}{\partial \beta} (y'y - y'X'\beta - X'\beta'y + (X\beta)'X\beta) = 0$$

$$\Rightarrow \frac{\partial}{\partial \beta} (y'y - 2\beta'X'y + \beta'X'X\beta) = 0$$

Q: Using LSM find out best regression eqⁿ of the following dataset.

x_1	x_2	y
2	1	1.7
5	3	2.6
6	4	2.8
8	9	2.3
10	11	2.7
11	14	2.4

$$\text{here, } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\begin{bmatrix} 1.7 \\ 2.6 \\ 2.8 \\ 2.3 \\ 2.7 \\ 2.4 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 5 & 3 \\ 6 & 4 \\ 8 & 9 \\ 10 & 11 \\ 11 & 14 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \end{bmatrix}$$

$$\text{We know that, } \beta = (x^T x)^{-1} x^T y$$

$$x^T x = \begin{bmatrix} 2 & 5 & 6 & 8 & 10 & 11 \\ 1 & 3 & 4 & 9 & 11 & 14 \end{bmatrix}_{2 \times 6} \begin{bmatrix} 2 & 1 \\ 5 & 3 \\ 6 & 4 \\ 8 & 9 \\ 10 & 11 \\ 11 & 14 \end{bmatrix}_{6 \times 2} = \begin{bmatrix} 350 & -377 \\ -377 & 424 \end{bmatrix}_{2 \times 2}$$

$$[x^T x]^{-1} = \frac{\text{adj}^o(x^T x)}{|x^T x|} = \begin{bmatrix} 424/6271 & -377/6271 \\ -377/6271 & 350/6271 \end{bmatrix}$$

$$(x^T y) = \begin{bmatrix} 2 & 5 & 6 & 8 & 10 & 11 \\ 1 & 3 & 4 & 9 & 11 & 14 \end{bmatrix}_{2 \times 6} \begin{bmatrix} 1.1 \\ 2.6 \\ 2.8 \\ 2.3 \\ 2.7 \\ 2.4 \end{bmatrix}_{6 \times 1}$$

$$= \begin{bmatrix} 105 \\ 104.7 \end{bmatrix}$$

Now $(x^T x)^{-1} (x^T y) = \begin{bmatrix} 424 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

$$\beta = \begin{bmatrix} -0.242 \\ 0.462 \end{bmatrix}$$

$$y = -0.242 x_1 + 0.462 x_2$$

~~Types of gradient descent~~

① Batch GD

→ The entire dataset is used to calculate the gradient and update the parameters.

~~fast~~ → ~~fixed~~ to

② Stochastic GD

→ A single sample is used to update the parameters at a time.

- converges quickly than batch GD
- can escape from local maxima

③ Mini Batch GD

→ It is a fine balance between stochastic and batch GD.
→ In this ~~as~~ the whole dataset is divided into batches and then batch run one after another and update the parameter.

* Normalisation

$$x_{\text{new}} = \left(\frac{x - x_{\min}}{x_{\max} - x_{\min}} \right)$$

→ bringing in (a,b) ranges

$$x_{\text{new}} = a + \left(\frac{x - x_{\min}}{x_{\max} - x_{\min}} \right) (b-a)$$

→ We do feature scaling with the help of these two techniques to have all the features on the same scale for better results (less variance)

* Standardisation

(Z-score normalisation) → * affected by outliers

$$x_{\text{new}} = \frac{x - \bar{x}}{\sigma}$$

- * when we want to ensure zero mean and unit std
- * much less affected by outliers

$$y = f(x) = \hat{f}(x, \theta) + \epsilon$$

↓
error

Irreducible error

Reducible error

Bias & Variance

* Training Error

$$\text{Bias} = E(\hat{Y}) - Y$$

$$\text{Bias} = E[\hat{f}(x, \theta)] - Y$$

Bias: It is defined as the diff. between expectation of estimated output (\hat{Y}) to true output (Y)

- Bias is introduced when a model is approximating a real-world problem (which is actually complex) with a simplified model
- Results in underfitting — fails to capture the underlying patterns in the data

Variance: Happens when a model is too complex.
 → overfitting: the model fits the training data very closely but may fail with the new data

$$\text{var } \hat{f}(x, \theta) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

The tradeoff

Simple model — High bias, low variance
 → Underfit (misses patterns)

Complex model — High variance, low bias
 → Overfit (too sensitive)

* We need to find the balance.

- Low variance means $f(x)$ based on different training set doesn't change much.
- High variance means $\hat{f}(x)$ has high variance (small change in the training data set can result large change in $\hat{f}(x)$.)

$$\therefore \text{Bias} = [E[\hat{y}] - y]$$

$$\text{Variance} = E\left[\hat{f}(x, \theta) - E[\hat{f}(x, \theta)]\right]^2$$

$$\Rightarrow \text{MSE} = E[(y - \hat{y})^2] = [(\text{Bias})^2] + E[\hat{f}(x, \theta) - E[\hat{f}(x, \theta)]]^2$$

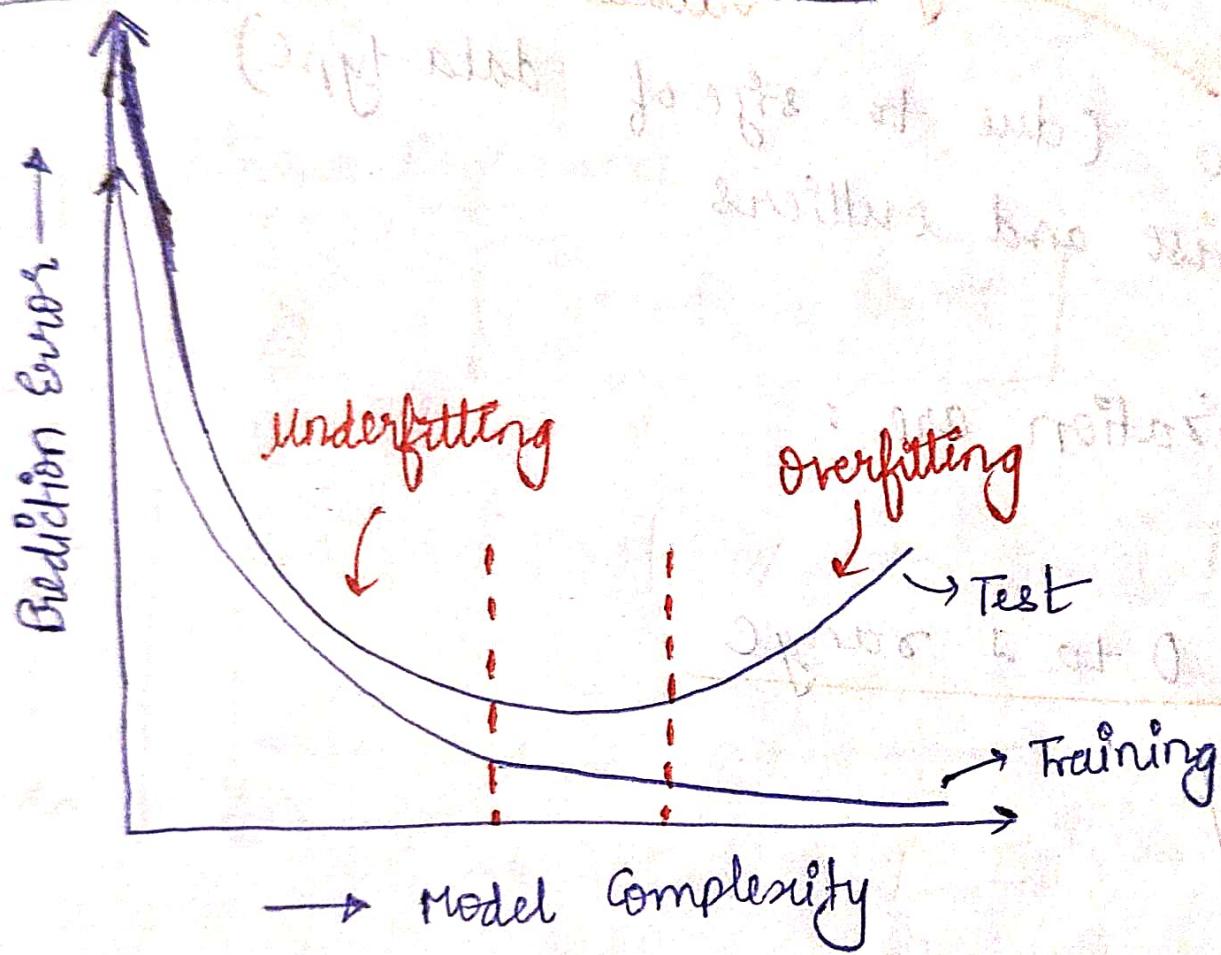
~~MSE = Bias² + Variance~~

~~MSE = Bias² + Variance~~

Bias $\uparrow \Rightarrow$ Underfitting

Variance $\uparrow \Rightarrow$ Overfitting

→ Supervised Model Evaluation



~~Regularisation~~

→ Regularisation is a technique used to reduce overfitting in a model and improve predictive performance.

① Ridge Regularisation (Least Absolute Selection and Shrinkage Operator) LASSO (L₁ Regularisation)

- In this, we add the sum of absolute sum of the coefficients to the MSE (or RSS)

$$\text{loss} = \text{MSE} + \lambda \sum_{i=1}^n |w_i|$$

② Ridge Regularisation (L₂ Regularisation)

- In this, we add sum of square of the coefficients to the MSE (or RSS)

$$\text{loss} = \text{MSE} + \lambda \sum_{i=1}^n w_i^2$$

* In ridge regression, the coefficient shrinks towards 0 whereas in lasso regression, it shrinks to 0.

so when feature selection is done, lasso will case min, the feature is eliminated and in case of ridge, the impact of that feature is reduced as the value of coefficient is shrunked.

Logistic Regression

→ supervised L. Algo used for classification.

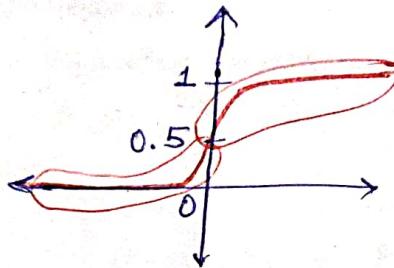
The goal of logistic regression is to predict whether the probability ~~is~~ that an instance belongs to a class or not.

Used for binary classification.

Uses sigmoid fn that produces a probability value

$$\text{class} = \begin{cases} 0 & ; P(z) < 0.5 \\ 1 & ; P(z) \geq 0.5 \end{cases}$$

* Sigmoid Function



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

$$z = \beta_0 + \beta_1 x_1$$

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \beta_0$$

- Q]: Consider the dataset of 5 observations of no. of hours studied by a student.
- Predict the class level of student who studied for 30 hrs, 34 hrs.
 - Calculate how many hours a student should study so that the student will pass with 95% probability.

$$\beta_1 = 2$$

$$\beta_0 = -64$$

Hrs of studies	Pass/Fail (Class)
28	0
14	0
32	1
27	1
38	1

$$\text{[ii]} P = 0.95$$

$$\frac{1}{1+e^{-z}} = 0.95$$

$$\Rightarrow z = 2.94$$

$$\Rightarrow x = 33.47$$

(i) Here, $z = \beta_0 + \beta_1 x$

for 30 $z = -64 + 2(30)$ for 34 $z = -64 + 2(34)$

$$= -4$$

$$= 4$$

$$P = \frac{1}{1+e^{-4}} = 0.98 > 0.5$$

Fail Pass

NAIVE-BAYES CLASSIFIER

→ It is supervised learning algorithm used for classification.

• It is named as 'NAIVE' 'BAYES'

because it assumes that one feature does not affect another feature (feature independence)

based on Bayes Theorem

Assumptions :

- 1) Feature independence
- 2) Features are equally important
- 3) No missing data

Bayes Formula

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)} \quad ; \quad P(X) \neq 0$$

here, $P(Y|X)$ is called posterior probability

$P(X|Y) \rightarrow$ likelihood or conditional probability

$P(Y) \rightarrow$ Prior Probability

$Y \rightarrow$ class variable & $X \rightarrow$ Feature vector of size n

$$P(Y|X=x_1, x_2, x_3, \dots, x_n) = \frac{\overbrace{P(x_1|y) \cdot P(x_2|y) \cdots P(x_n|y) \cdot P(y)}}{P(x_1) \cdot P(x_2) \cdot P(x_3) \cdots P(x_n)}$$

$$P(Y|X=x_1, x_2, \dots, x_n) = \frac{P(Y) \cdot \prod_{i=1}^n P(x_i|y)}{\prod P(x_i)}$$

- Example: consider given dataset & find out the test sample today (x) - {Sunny, Hot, Normal, False} The class level using naive bayes classification.

To reach out the result we need to calculate

$$P(y|x_1=x_1, x_2=x_2, \dots, x_n=x_n) = \frac{P(x_1|y) \cdot P(x_2|y) \cdot P(x_n|y) \cdot P(y)}{P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_n)}$$

$$\text{OR } P(y|x_1=x_1, x_2=x_2, \dots, x_n=x_n) = P(y) \prod_{i=1}^n P(x_i|y)$$

$$y = \arg \max_y P(y) \prod_{i=1}^n P(x_i|y)$$

	Outlook	Temp	Humidity	Windy	Play Golf
1	Rainy	Hot	High	False	No
2	Rainy	Hot	High	True	No
3	Overcast	Hot	High	F	Yes
4	Sunny	Mild	High	F	Yes
5	Sunny	Cool	Normal	F	Yes
6	Sunny	Cool	Normal	T	No
7	Overcast	Cool	Normal	T	Yes
8	Rainy	Mild	High	F	No
9	Rainy	Cool	Normal	F	Yes
10	Sunny	Mild	Normal	F	Yes
11	Rainy	Mild	Normal	T	Yes
12	Overcast	Mild	High	T	Yes
13	Overcast	Hot	Normal	F	Yes
14	Sunny	Mild	High	T	No

Class	outlook	Temp	Humidity	windy	Prior Prob $P(C)$
Yes	$P(\text{sunny} Y)$ = $3/9$	$P(H Y)$ = $2/9$	$P(N Y)$ = $6/9$	$P(F Y)$ = $6/9$	$P(Y) = 9/14$
No	$P(\text{sunny} N)$ = $2/5$	$P(H N)$ = $2/5$	$P(N N)$ = $1/5$	$P(F N)$ = $2/5$	$P(N) = 5/14$

$$P(C_k | X) = P(X | C_k) \cdot P(C_k)$$

$$\begin{aligned} P(\text{Yes} | \text{today}) &= P(x_1 | C_k) \cdot P(x_2 | C_k) \dots P(x_n | C_k) P(C_k) \\ &= 3/9 \times 2/9 \times 6/9 \times 6/9 \times 9/14 \\ &= 0.021 \end{aligned}$$

$$\begin{aligned} P(\text{No} | \text{today}) &= 2/5 \times 2/5 \times 1/5 \times 2/5 \times 5/14 \\ &= 0.004 \end{aligned}$$

$$y = \operatorname{argmax}_y P(y) * \prod_{i=1}^n P(x_i | y)$$

class label of the test sample today is 'yes'.

Find out the class label of the test sample
 $x = (\text{Rainy}, \text{Hot}, \text{High}, \text{False})$ using naive Bayes

Class	outlook	Temp	Humidity	windy	Prior Prob
Yes	$2/9$	$2/9$	$3/9$	$6/9$	$P(Y) = 9/14$
No	$3/5$	$2/5$	$4/5$	$2/5$	$P(N) = 5/14$

$$P(\text{Yes} | \text{today}) > 2/9 \times 2/9 \times 3/9 \times 6/9 \times 9/14 = 0.007$$

$$\begin{aligned} P(\text{No} | \text{today}) &> 3/5 \times 2/5 \times 4/5 \times 2/5 \times 5/14 \\ &> 0.027 \end{aligned}$$

class is 'no'.

Advantage :

- easy to implement and computationally efficient
- it is effective with a large no. of feature
- it performs well with the limited training data
- it also performs well in the presence of categorical features
- for numerical features, it performs well if it shows normal distribution

Disadvantage :

- it may be influenced by irrelevant attributes so feature selection is required before applying naive Bayes classifier.
- It may assign zero probability to unseen events leading to poor generalisation
- As it assumes all features are independent which may not always hold true in the real world data .

Q. Consider the given dataset. Predict the class level of given test sample $x = \{f_{11}, f_{21}, f_{31}\}$ using naive Bayes classifier based on the following obs

Using NB classifier based on following obs

class label	f_{11}	f_{21}	f_{31}	Total	Summarized table or frequency table
obj1	350	450	0	650	
obj2	400	300	350	400	
obj3	50	100	50	150	

$x = \{f_{11}, f_{21}, f_{31}\}$ class table .

$$P(\text{obj1} | X) = P(X|\text{obj1}) \cdot P(\text{obj1})$$

$$\rightarrow P(f_{11}|\text{obj1}) \cdot P(f_{21}|\text{obj1}) \cdot P(f_{31}|\text{obj1}) \cdot P(\text{obj1})$$

$$\rightarrow \frac{350}{650} \times \frac{450}{650} \times \frac{0}{650} \times \frac{650}{1200}$$

$$P(\text{obj2} | X) = P(X|\text{obj2}) \cdot P(\text{obj2})$$

$$= P(f_{11}|\text{obj2}) \cdot P(f_{21}|\text{obj2}) \cdot P(f_{31}|\text{obj2}) \cdot P(\text{obj2})$$

$$= \frac{400}{400} \times \frac{800}{400} \times \frac{350}{400} \times \frac{400}{1200} = 0.218$$

$$P(\text{obj3} | X) = P(f_{11}|\text{obj3}) \cdot P(f_{21}|\text{obj3}) \cdot P(f_{31}|\text{obj3}) \cdot P(\text{obj3})$$

$$= \frac{50}{150} \times \frac{100}{150} \times \frac{50}{150} \times \frac{150}{1200} = 0.009$$

Gaussian Naive Bayes Classification

$$P(y=c | X) = P(X|y=c) \cdot P(y=c)$$

$$P(X|y=c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(x-\mu_c)^2}{2\sigma_c^2}}$$

μ_c = mean of class c

σ_c = std. deviation of class c

2. Based on given dataset determine class level of the test sample $X = \{x_1=6, x_2=130, x_3=8\}$ using Naive Bayes classifier.

f_1	f_2	f_3	class
6	180	12	C1
5.92	190	11	C1
5.58	170	12	C1
5.78	165	10	C1
5	100	6	C2
5.5	150	8	C2
5.4	130	7	C2
5.72	150	9	C2

$$\mu_{f_1|C_1} = \frac{6 + 5.92 + 5.58 + 5.78}{4} = 5.82$$

$$\mu_{f_1|C_2} = \frac{5 + 5.5 + 5.4 + 5.72}{4} = 5.405$$

$$\sigma_{f_1|C_1} = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} = \sqrt{\frac{(6-5.8)^2 + (5.92-5.8)^2 + (5.58-5.8)^2 + (5.78-5.8)^2}{3}} = 0.185$$

$$\sigma_{f_1|C_2} = \sqrt{\frac{(5-5.405)^2 + (5.5-5.405)^2 + (5.4-5.405)^2 + (5.72-5.405)^2}{3}} = 0.301$$

$$P(\text{expt } p(C_1 | X)) = P(f_1|C_1) \cdot P(f_2|C_1) \cdot P(C_1) \cdot P(f_3|C_1)$$

$$\mu_{f_2|C_1} = \frac{180 + 190 + 170 + 165}{4} = 176.25$$

$$\mu_{f_2|C_2} = \frac{100 + 150 + 130 + 150}{4} = 132.5$$

$$\sigma_{f_2|C_1} = \sqrt{\frac{(180-176.25)^2 + (190-176.25)^2 + (170-176.25)^2 + (165-176.25)^2}{3}} = 11.08$$

$$\sigma_{f_2|C_2} = 23.63$$

$$R_{f_3C_1} = \frac{12+11+12+10}{4} = 11.25$$

$$R_{f_3C_2} = \frac{6+8+7+9}{4} = 7.5$$

$$\sigma_{f_3C_1} = 0.96$$

$$\sigma_{f_3C_2} = 1.3$$

0.14

$$P(f_1|C_1) = \frac{1}{\sqrt{2\pi(0.185)^2}} e^{-\frac{(6-5.82)^2}{2(0.185)^2}} = 0.0243$$

$$P(f_2|C_1) = \frac{1}{\sqrt{2\pi(11.08)^2}} e^{-\frac{(130-126.25)^2}{2(11.08)^2}} = 0.000037$$

$$P(f_3|C_1) = \frac{1}{\sqrt{2\pi(0.96)^2}} e^{-\frac{(8-11.25)^2}{2(0.96)^2}} = 0.000016$$

$$P(Y|x) = P(U) \cdot P(f_1|U) \cdot P(f_2|C_1) \cdot P(f_3|C_1)$$

$$= 0.5 \times 0.143 \times 0.0243 \times 0.000037 = 0.0000063$$

$$= 1.8 \times 10^{-9}$$

$$P(f_1|C_2) = \frac{1}{\sqrt{2\pi(0.3)^2}} e^{-\frac{(6-5.4)^2}{2(0.3)^2}} = 0.0243$$

$$P(f_2|C_2) = \frac{1}{\sqrt{2\pi(23.6)^2}} e^{-\frac{(130-132.5)^2}{2(23.6)^2}} = 0.000009$$

$$P(f_3|C_2) = \frac{1}{\sqrt{2\pi(1.3)^2}} e^{-\frac{(8-7.5)^2}{2(1.3)^2}} = 0.0243$$

$$P(C_2|x) = 0.5 \times 0.143 \times 0.124 \times 0.0004 = 0.0000083$$

$$= 8.3 \times 10^{-9}$$

(2)

KNN Algorithm

- It is a supervised learning algorithm which is used to predict the class of a sample based on the class of K-nearest neighbors.
- used for classification as well as regression.
- It is called Lazy Learner Algorithm because it does not train a model explicitly; it just memorizes the data and makes predictions in real time.
- * choosing the right K-value and distance metric is crucial.

Distance Metric

for two points $x(x_1, y_1)$ & $y(y_1, y_2)$

Manhattan distance

$$d = |x_1 - y_1| + |x_2 - y_2|$$

Euclidean distance

$$d = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

Minkowski distance

$$d = \sqrt[p]{(x_1 - y_1)^p + (x_2 - y_2)^p}$$

$\rightarrow p=1$ Manhattan
 $\rightarrow p=2$ Euclidean

steps to perform K-NN

step ① Initialize K

→ many ways to choose K

- odd number
- sq.root of m (idk what m is!?)
- elbow method
- cross validation

step ② calculate distance from each training sample, to the test sample using the dist. formula

step ③ Arrange the distances calculated in ascending order

step ④ Pick the first ^{training} k-samples and get their class.

step ⑤ Predict the class of test sample, based on the majority of class. of k-neighbors.

<u>D:</u>	x_1	x_2	class	of distance (25, 80)
	40	20	C ₁	61.84
	50	50	C ₂	39.05
	60	90	C ₂	36.40
	10	25	C ₁	57.00
	25	80	?	(Test sample)

Let k=3 ✓

$$y_1 = 25 \quad y_2 = 80$$

$$d_1 = \sqrt{(25-40)^2 + (80-20)^2} = 61.84$$

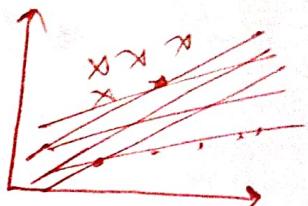
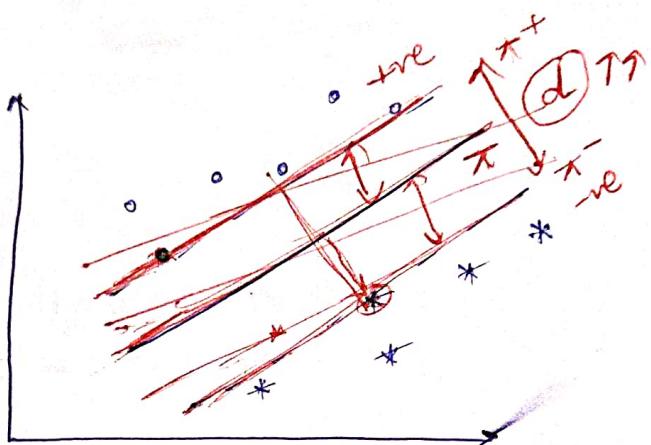
Sorting these distance values,

$$\left. \begin{array}{l} d_3 < d_2 < d_4 < d_1 \\ (C_2) \quad (C_2) \quad (C_1) \quad (C_1) \end{array} \right\}$$

From the 3 nearest neighbors,
we get class of (25, 80)
as C₂.

SVM (Support Vector Machine)

- This is a popular supervised learning algorithm generally used for classification.
- The goal of SVM is to find a hyperplane that can segregate 'n' dimensional space into classes.
hyperplane → best line or decision boundary.



Terminologies in SVM

support vectors : The data points (or vectors) that are closest to the hyperplane & which affect the position of hyperplane are termed as support vectors.

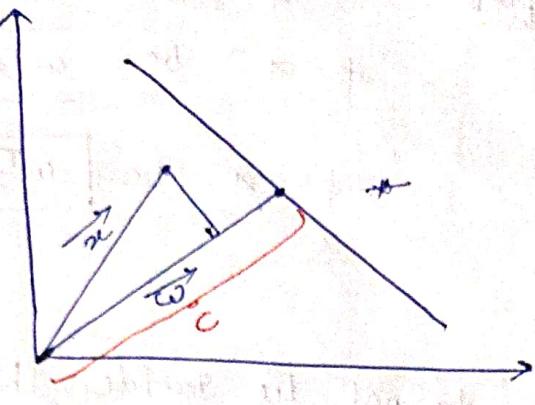
Hyperplane : Multiple decision boundaries can be drawn for segregating the classes; the best boundary is known as hyperplane of SVM.

• Dimension of hyperplane = $(n-1)$ where $n \rightarrow$ no of features in the dataset

Margin : The distance between the support vectors and the hyperplane is called as margin.

→ Goal of SVM is to maximize the margin,
↓
optimal hyperplane

use of dot product in SVM

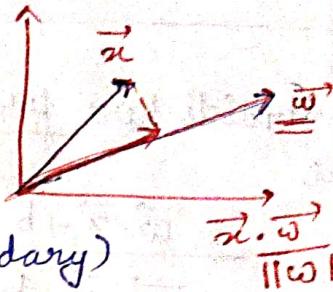


\vec{x} → vector to a random point
 \vec{w} → vector to the hyperplane
 c → distance from origin to the hyperplane

Linear SVM

Mathematically,

$$\left\{ \begin{array}{ll} \vec{x} \cdot \vec{w} = c & (\text{point lies on decision boundary}) \\ \vec{x} \cdot \vec{w} > c & (\text{positive samples}) - \text{right side of hyperplane} \\ \vec{x} \cdot \vec{w} < c & (\text{negative samples}) - \text{left side of hyperplane} \end{array} \right.$$



$$\frac{\vec{x} \cdot \vec{w}}{\|\vec{w}\|}$$

hyperplane

Find value of w and b for which the margin is maximum.

$$\vec{w} \cdot \vec{x} > c \quad (\text{+ve samples})$$

$$\Rightarrow \vec{w} \cdot \vec{x} - c > 0$$

$$\Rightarrow \boxed{\vec{w} \cdot \vec{x} + b > 0} \quad \text{Decision Rule}$$

+ve value for positive samples (+1)

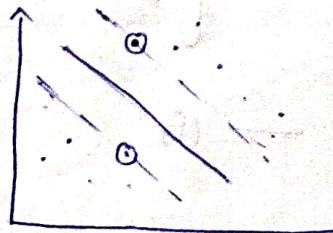
+ve value for negative samples (-1)

$$\vec{w} \cdot \vec{x} - c > 0$$

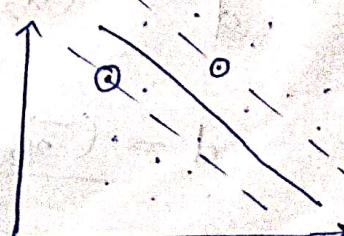
$$\vec{w} \cdot \vec{x} + b < 0$$

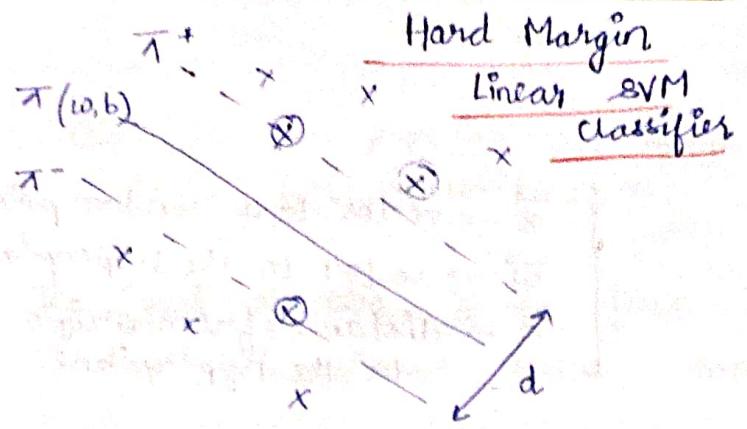
$$\left\{ \begin{array}{l} \text{Decision Rule} \\ \hat{y}_i = \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{x}_i + b \geq 0 \\ -1 & \text{if } \vec{w} \cdot \vec{x}_i + b < 0 \end{cases} \end{array} \right.$$

Hard Margin



soft Margin





Let us assume the eqⁿ of
 if π^+ be $w^T x + b = 1$
 & π^- be $w^T x + b = -1$
 & π be $w^T x + b = 0$

The constraints here for the points to not lie inside this hyperplane are

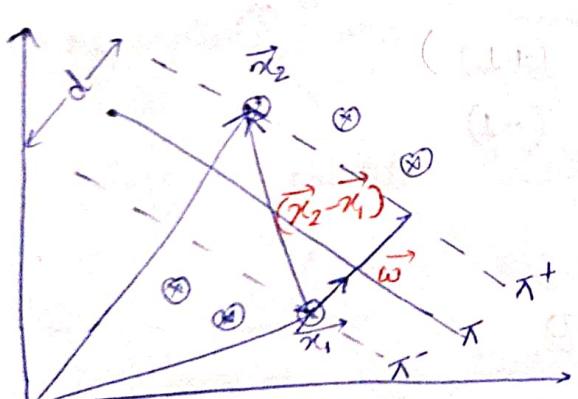
$$\begin{cases} \vec{w} \cdot \vec{x}_i + b \geq 1 \\ \vec{w} \cdot \vec{x}_i + b \leq -1 \end{cases}$$

If we multiply the above constraints with the labels of these points i.e. +1 & -1, we get

$$\begin{cases} \hat{y} (\vec{w} \cdot \vec{x}_i + b) \geq 1 * (+1) \\ \hat{y} * (\vec{w} \cdot \vec{x}_i + b) \geq -1 * (-1) \end{cases} \quad y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

for all vectors

for support vectors, $y_i (\vec{w} \cdot \vec{x}_i + b) = 1$,



To calculate d , we need to find the projection of $(\vec{x}_2 - \vec{x}_1)$ on \vec{w} .

$\because \vec{w}$ is perpendicular to all three lines

$$d = (\vec{x}_2 - \vec{x}_1) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

$$\Rightarrow d = \frac{(\vec{x}_2 \cdot \vec{w} - \vec{x}_1 \cdot \vec{w})}{\|\vec{w}\|} \quad - (1)$$

$\left\{ \frac{\vec{w}}{\|\vec{w}\|} \rightarrow \text{unit vector in direction of } \vec{w} \right\}$

Now, for x_2 :

$$\begin{aligned} y_i (\vec{\omega} \cdot \vec{x}_i + b) &= 1 \\ 1 (\vec{\omega} \cdot \vec{x}_2 + b) &= 1 \\ \Rightarrow \vec{\omega} \cdot \vec{x}_2 &= (1-b) \end{aligned}$$

— (ii)

for x_1 :

$$\begin{aligned} (-1) (\vec{\omega} \cdot \vec{x}_1 + b) &= 1 \\ \Rightarrow -\vec{\omega} \cdot \vec{x}_1 - b &= 1 \\ \Rightarrow \vec{\omega} \cdot \vec{x}_1 &= (-1-b) \end{aligned}$$

— (iii)

Putting values from eqn (ii) & (iii) in eqn (i),

$$d = \frac{(1-b) - (-1-b)}{\|\omega\|} = \frac{2}{\|\omega\|}$$

$$\Rightarrow \boxed{d = \frac{2}{\|\omega\|}}$$

HARD
MARGIN
SVM

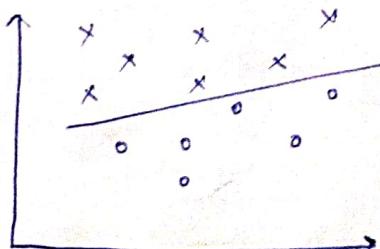
* $\underset{\omega^*, b^*}{\operatorname{argmax}} \frac{2}{\|\omega\|}$

given that $y_i (\vec{\omega} \cdot \vec{x}_i + b) \geq 1$

find such values of

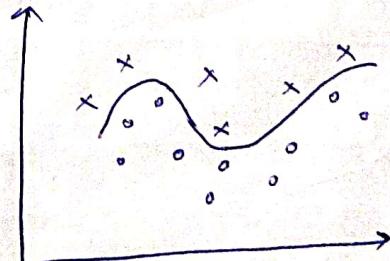
$\omega + b$
for which $\frac{2}{\|\omega\|}$ is max^M

Linear SVM



The two classes can be separated using a straight line.

Non-linear SVM

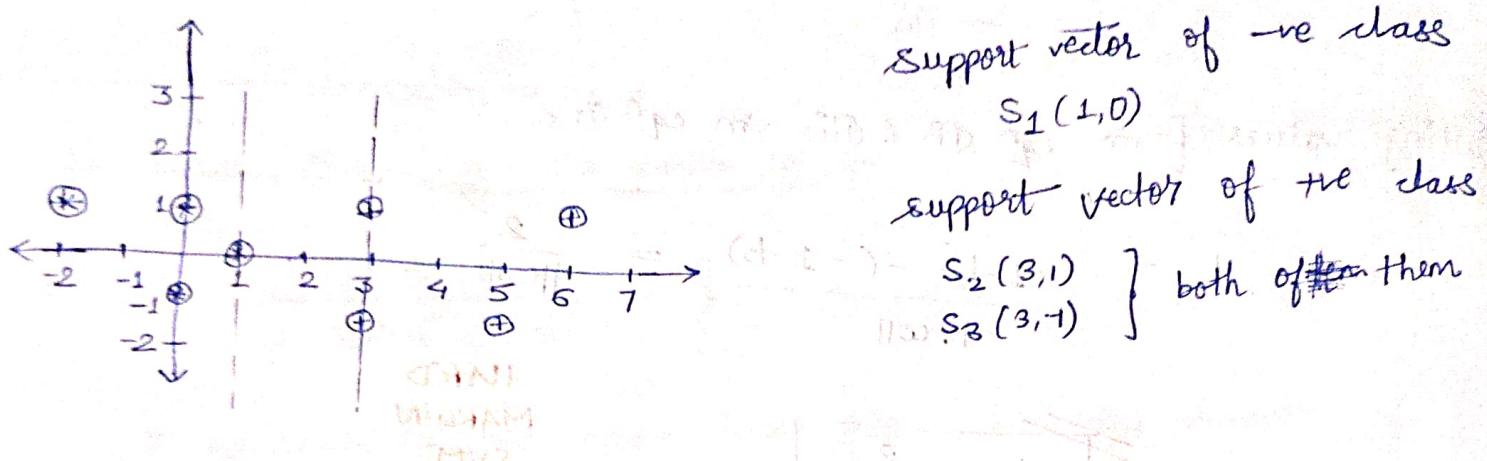


The two classes cannot be separated using a straight line.

Q: Determine the eqn of hyperplane that divides the datapoints into two classes.

+vely labeled $[(3,1), (3,-1), (6,1), (5,-1)]$

-vely labeled $[(1,0), (-2,1), (0,-1), (0,1)]$



⇒ We will get the support vectors from the graph

$$S_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix}, S_2 \begin{pmatrix} 3 \\ 1 \end{pmatrix}, S_3 \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

assigning 1 in all 3 support vectors as bias

$$\bar{S}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \bar{S}_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \bar{S}_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

Let's take these parameters α_1, α_2 & α_3 for S_1, S_2 & S_3 calculate the weight vector, we need to go for the following eqn.

$$\alpha_1 \bar{S}_1 \bar{S}_1 + \alpha_2 \bar{S}_2 \bar{S}_1 + \alpha_3 \bar{S}_3 \bar{S}_1 = -1 \quad \text{--> -vely labelled}$$

$$\alpha_1 \bar{S}_1 \bar{S}_2 + \alpha_2 \bar{S}_2 \bar{S}_2 + \alpha_3 \bar{S}_3 \bar{S}_2 = 1 \quad \text{--> +vely labelled}$$

$$\alpha_1 \bar{S}_1 \bar{S}_3 + \alpha_2 \bar{S}_2 \bar{S}_3 + \alpha_3 \bar{S}_3 \bar{S}_3 = 1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = -1$$

$$\Rightarrow 2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1 \quad \text{--- (i)}$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} = 1$$

$$\Rightarrow 4\alpha_1 + 11\alpha_2 + \alpha_3 = 1 \quad \text{--- (ii)}$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} = 1$$

$$\Rightarrow 4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1 \quad \text{--- (iii)}$$

$$\alpha_1 = -3.5, \alpha_2 = 0.75, \alpha_3 = 0.75$$

$$\bar{\omega} = \sum \alpha_i \bar{s}_i = \alpha_1 \bar{s}_1 + \alpha_2 \bar{s}_2 + \alpha_3 \bar{s}_3$$

$$\Rightarrow \bar{\omega} = -3.5 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

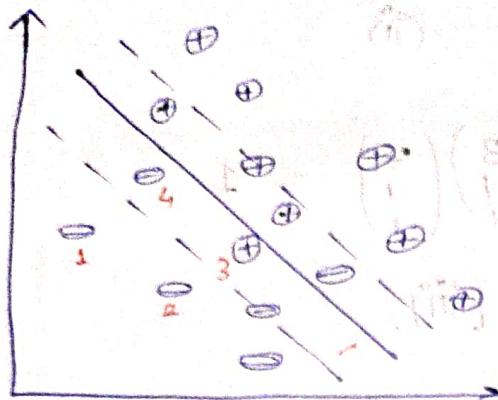
Equating it with hyperplane offset;

$$y = \omega x + b \quad \text{where } \omega = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, b = -2$$

$$\frac{x}{d/\omega_1} + \frac{y}{d/\omega_2} = 1 \Rightarrow \frac{x}{2/1} + \frac{y}{2/0} = 1 \Rightarrow x = 2$$

$$\left\{ d = \frac{2}{\|\omega\|} = \sqrt{\frac{2}{1^2+0^2}} = 2 \right\}$$

~~soft Margin Linear SVM classifier~~



If the datapoints are almost linearly inseparable, that's why hard margin is not possible.

$$\text{From hard margin, } \underset{\omega^*, b^*}{\operatorname{argmax}} (\omega^*, b^*) \frac{2}{\|\omega\|} + \text{const}$$

$$\Rightarrow \underset{\omega^*, b^*}{\operatorname{argmin}} (\omega^*, b^*) \frac{\|\omega\|}{2}$$

Now, we add a hyperparameter to account for classification

$$\underset{\omega^*, b^*}{\operatorname{argmin}} (\omega^*, b^*) \frac{\|\omega\|}{2} + C \sum_{i=1}^n \xi_i \text{ (slack variable)}$$

margin error classification error

we need to minimize both

Here, c is regularisation constant.

A higher value of c implies a stricter penalty for margin violation.

This leads to a smaller margin with fewer misclassifications.

* Hinge loss

Hinge loss is the loss function used in SVM that ensures that the misclassified points contribute to the loss and the correctly classified points do not.

* It penalizes points that are within the margin or misclassified. (uses slack variable to do so)

- $\xi = 0$: point is correctly classified and outside the margin.
- $0 < \xi < 1$: point is inside the margin but still correctly classified (on the correct side) of the decision boundary
- $\xi > 1$: the point is misclassified (and on the wrong side of the dec. boundary)

Advantages

- simple algos, less complex
- very few parameters to be considered
- does not require model training (lazy learner algo)
- works well for small dataset

Max 2 iterations

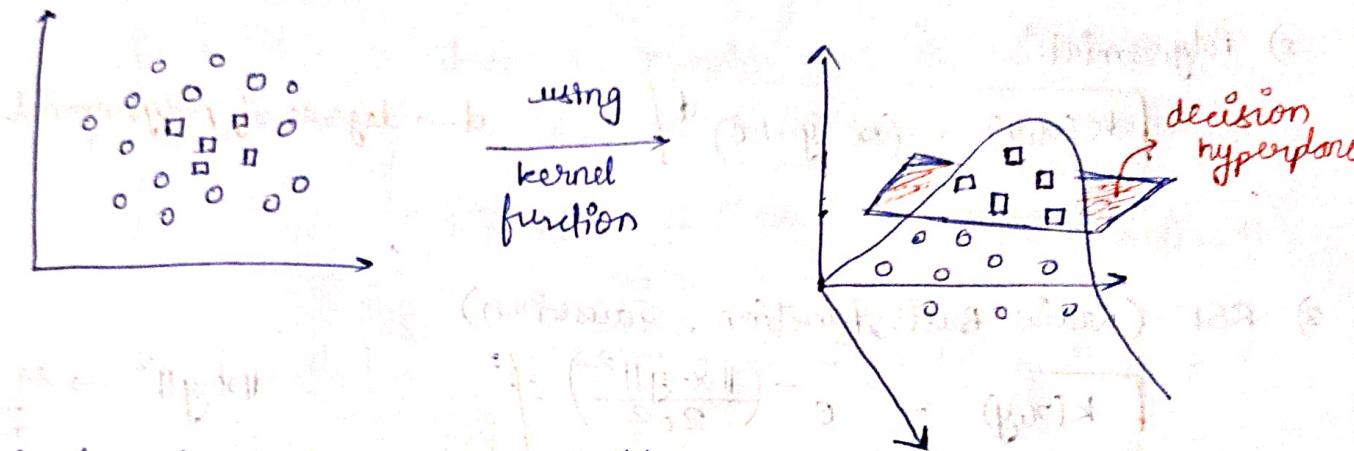
for MLR

Disadvantages

- slow for large datasets
- choosing the right k is tricky

Non-Linear SVM

When \rightarrow the datapoints are linearly not separable, we use non-linear SVM with kernel function to make it separable with a decision boundary, hyperplane.



- The kernel function allows the SVM to solve both linear and non-linear problems, by mapping the data points into a higher dimensional space.
And then uses a decision hyperplane to separate them into classes.

→ different types of kernel functions:

- Linear Kernel F.
- Polynomial KF
- Radial Basis Function (RBF) — Gaussian Kernel
- Bessel Function
- Sigmoid KF

Non-linear

1) linear kernel

$$k(x, y) = x^T y$$

2) Polynomial

$$k(x, y) = (x^T y + c)^d \quad d \rightarrow \text{degree of polynomial}$$

3) RBF (Radial Basis Function, Gaussian)

$$k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

$\|x-y\|^2 \rightarrow$ sq. euclidean dist betn x & y vectors

e.g.: $x(1, 2)$, $y(3, 4)$, $\sigma=1$ \rightarrow parameter controlling the spread

$$k(x, y) = e^{-\frac{8}{2(1)^2}} = e^{-4} \quad \|x-y\| = \sqrt{2^2+2^2} = \sqrt{8}$$

$$\Rightarrow k(x, y) = 0.018$$

4) Sigmoid

$$k(x, y) = \tanh(\alpha x^T y + c)$$

$\alpha \rightarrow$ slope parameter
 $c \rightarrow$ constant

e.g.: $x(1, 2)$, $y(3, 4)$, $c=1$, $\alpha=0.5$

$$k(x, y) = \tanh(0.5(11) + 1) = 0.99$$

Primal Problem

The primal problem is the direct formulation of the SVM optimization problem.

- It aims to find the optimal hyperplane $w^T x + b = 0$

maximize margin minimize misclassifications

Dual Primal Problem

∴ primal problem involves inequality constraints, solving it directly is difficult.

so we convert it into a dual problem using Lagrange's multipliers.

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + w_0)) \quad : \alpha_i \geq 0$$

$\hookrightarrow \text{eq ①}$

To solve, we find out derivatives wrt w & w_0

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \alpha_i^* y_i x_i = 0$$

$$\Rightarrow \boxed{w = \sum_{i=1}^n \alpha_i^* y_i x_i}$$

$$\frac{\partial L}{\partial w_0} = - \sum_{i=1}^n \alpha_i^* y_i = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^n \alpha_i^* y_i = 0}$$

$$\|\omega\|^2 = \omega^T \omega$$

$$\begin{aligned}\|\omega\|^2 &= \left(\sum_{j=1}^n \alpha_j y_j x_j \right)^T \left(\sum_{j=1}^n \alpha_j y_j x_j \right) \\ &= \sum_{j=1}^n \sum_{i=1}^n \alpha_i \alpha_j y_i y_j x_i x_j^T\end{aligned}$$

Putting these values
in eq①

~~$\omega^T \omega$~~

$$\begin{aligned}\omega^T x_i &= \left(\sum_{j=1}^n \alpha_j y_j x_j \right)^T x_i \\ &= \sum_{j=1}^n \alpha_j y_j x_j^T x_i\end{aligned}$$

$$L(\omega, \omega_0, \alpha) = \frac{1}{2} \left(\sum \sum \alpha_i \alpha_j y_i y_j x_i x_j^T \right) + \sum \alpha_i - \sum \sum \alpha_i \alpha_j y_i y_j x_i x_j^T - \cancel{\sum \alpha_i y_i \omega_0}$$

$$L(\omega, \omega_0, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^n \alpha_i \alpha_j y_i y_j x_i x_j^T$$

maximize such that $\alpha_i \geq 0 \quad \forall i$
 $\quad \quad \quad + \sum \alpha_i y_i = 0$

KKT conditions
For a solution to be optimal, the following KKT conditions:

① Stationarity ~~$\frac{\partial}{\partial \omega} L(\omega, \omega_0, \alpha)$~~ $\Rightarrow \frac{\partial L}{\partial \omega} = 0 \Rightarrow \omega = \sum_{i=1}^n \alpha_i y_i x_i$
 $\Rightarrow \frac{\partial L}{\partial \omega_0} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$

② Bound Feasibility $y_i (\omega^T x_i + b) \geq 1 \quad \forall i$

③ Dual Feasibility $\alpha_i \geq 0 \quad \forall i$

④ Universal Complementary Slackness $\alpha_i [y_i (\omega^T x_i + b) - 1] = 0 \quad \forall i$