

Ensemble learning

09/09/2025

Combining weak learners to get a strong learner

Bagging - You can replace

Independent weak learners

b - no. of learners

n - no. of samples

$$P(\text{not being chosen}) = \left(1 - \frac{1}{n}\right)^n \approx 0.37$$

n samples being picked n times
samples get replaced so we
always have n samples

Out of Bag - Those sample not being picked in training
↳ Acts as a validation set

If we don't add diversity, we will get same trees.

Subsampling - Subset of sample.

Bagging or Bootstrapping Algorithms

h_b = base learner (independent)

D_b = subset

For classification - majority voting

For regression - average

$$\text{Error} = (\text{Bias})^2 + \text{Variance} + \epsilon$$

Bias is already low in decision tree. It has more chances of overfitting ~~and~~ rather than underfitting.

So we reduce variance to reduce error

$$P(\text{ensemble voting error}) = \sum_{k=\lceil \frac{B}{2} \rceil}^B \binom{B}{k} p^k (1-p)^{B-k}$$

$B \uparrow$ Trees \uparrow Variance \downarrow
Error \downarrow

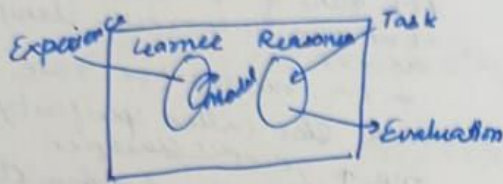
Linear regression - simple model \rightarrow less variance - Bagging not helpful

Machine Learning (DSL-501)

Theory (60%)

Labs (40%)

Learner System



Learner

Reasoner

Mid-20%

End-30%

Assignment-10%

midst (Best 2/3)

Project

Phase I

10% 5% Present action

QL Journal

CSE

A/A+ main

Phase-II

5% Progress

Phase-III

20% Code

Project Report Present

(10 pages) - action

Problem Statement

Related Work

Reference

Pattern Recogn ML-Christopher Bishop

An Intro to Statistical Learning-

Garreth James, Daniel Witten

Hands On ML with Scikit Learn

by Aurélien Géron 3rd Ed

Machine Learning - Tom Mitchell

Deep Learning - Ian Goodfellow

Essential Mathematics for

Machine Learning

Features

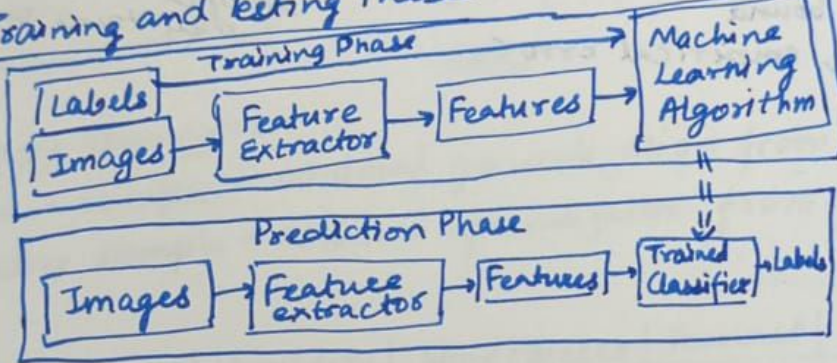
→ Categorical (eg. "A", "AB" ...)

→ ~~Order~~ Ordinal (eg. "Large", "Medium", "Small")

→ Integer valued

→ Real Valued

Training and Testing Phase



Example(x,y) : Instance x with label y

Reinforcement Learning

The learning system, called an agent in this context, can observe the environment, select and perform actions and get rewards in return (or penalties)

Hyperparameters Tuning

- Manual Search
- Grid Search
- Random Search

$$\begin{aligned}
 \text{Expected Error} & \rightarrow EPE = \sum_{i=1}^N (f(x) - \hat{y}_i) \rightarrow \text{prediction function} \\
 & \quad \text{Prediction} \quad \text{True prediction function} \\
 & \quad \text{Bias} \quad \text{Variance} \\
 & \quad \text{Under fit} \quad \text{Overfit} \\
 & \quad \text{Bias - Variance Tradeoff} \\
 & \quad \text{Low Bias} \quad \text{High Bias} \\
 & \quad \text{High Variance} \quad \text{Low Variance} \\
 & \quad \text{Bias} = E[(f(x) - \hat{y}(x))^2] \\
 & \quad \text{Variance} = E[(\hat{y}(x) - E[\hat{y}(x)])^2] \\
 & \quad \text{irreducible error} = \sigma^2 \\
 & \quad \text{Var} = E[f(x)^2] - (E[f(x)])^2 \\
 & \quad \text{precision} = \frac{TP}{TP + FP} \\
 & \quad \text{Recall} = \frac{TP}{TP + FN} \\
 & \quad \text{F1 score} = \frac{2}{(\frac{1}{\text{Prec}} + \frac{1}{\text{Recall}})}
 \end{aligned}$$

$$x=2, f(x)=5$$

$$D_1: \hat{y}_1(x) = 4.5$$

$$D_2: \hat{y}_2(x) = 5.0$$

$$D_3: \hat{y}_3(x) = 5.5$$

$$\text{Variance} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \sum_{j=1}^M (y^j(x_i) - E[\hat{y}(x_i)])^2$$

\checkmark # of instances in Test
 \rightarrow # of Models
 Here $N=1$

$$\text{Bias} = 0$$

Hypothesis Space -

VC Dimension - how effectively select

Generalisation bound

$$\text{True error} \geq \text{empirical error}$$

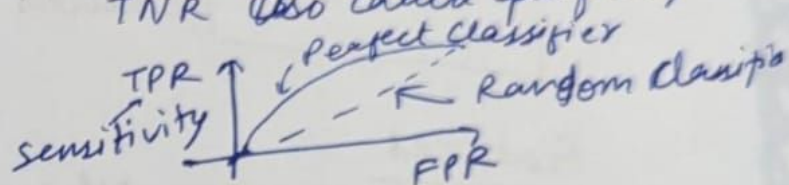
Roc Curve

True +ve rate (Recall)
vs False +ve rate

FPR ratio of -ve instances
that are incorrectly classified
as +ve

$$= 1 - \text{True negative rate}$$

TNR Also called specificity



Cross validation

K Fold \rightarrow Each fold serves as validation set

Leave one out - Each observation

Stratified - Ensures each fold maintains same proportion of class

Concept Learning

Find-S

Candidate Elimination Algo

$$D = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\} \quad \eta \in \{1, 1000\}$$

$$m \in \{\text{true, false}\}$$

$$\Downarrow$$

$$(x_1, c(x_1)), (x_2, c(x_2)), \dots$$

$$\Downarrow$$

$$h(x) = c(x) \quad \forall x \in X$$

Find-S

Not using -ve examples

$$S_0 = \{\phi, \phi, \phi\}$$

$$S_1 = \{\text{red, round, small}\}$$

$$S_2 = \{?, \text{round}, ?\}$$

If values are mismatch in the example then generalise(?)

~~c(x)~~

f₁, f₂, f₃
red, round, small +
yellow, long, medium -
green, round, medium +

Candidate Elimination

Positive sample - minimal generalization from specific boundary

Negative sample - minimal specialization from general boundary

Sky	Air temp	Humidity	Wind	Water	Forecast	PlaySport
			Strong	Warm	Same	Yes
Sunny	Warm	Normal	"	"	"	Yes
Sunny	Warm	High	"	"	Change	No
Rainy	Cold	"	"	cool	"	Yes
Sunny	Warm	"	"	"	"	Yes

$$S_0 = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$$

$$G_0 = \langle ?, ?, ?, ?, ?, ? \rangle$$

$$S_1 = \langle \text{Sunny, warm, normal, strong, warm, same} \rangle$$

$$S_2 = \langle \text{sunny, warm, ?, strong, warm, same} \rangle$$

$$G_1 = \langle ?, ?, ?, ?, ?, ? \rangle$$

$$G_2 = \langle ?, ?, ?, ?, ?, ? \rangle$$

$$G_3 = \langle \text{sunny, ?, ?, ?, ?, ?} \rangle$$

$$G_4 = \langle \text{sunny, warm, ?, strong, ?, ?} \rangle$$

$$S_3 = S_2$$

$$S_4 = \langle \text{sunny, warm, ?, strong, ?, ?} \rangle$$

$$V = 3$$

version space

- $\langle \text{sunny, warm, ?, ?, ?, ?} \rangle$
- $\langle \text{sunny, ?, ?, strong, ?, ?} \rangle$
- $\langle \text{?, warm, ?, strong, ?, ?} \rangle$

eliminate inconsistent hypothesis

Decision Tree

How to split at Internal Nodes

- A pure group means that the majority of the inputs have the same label/output

Entropy and Information Gain

S - set of labelled inputs from C classes, p_c as fraction of class c inputs

Entropy $\rightarrow H(S) = - \sum_{c \in C} p_c \log p_c$

split S to S_1 & S_2

Reduction in entropy after split is called information gain

$$IG = H(S) - \frac{|S_1|}{|S|} H(S_1) - \frac{|S_2|}{|S|} H(S_2)$$

Low entropy - High Information Gain - Good split

$$H(S) = 0.94$$

~~IG~~

$$S_{\text{sunny}} = [2+, 0-] \quad S_{\text{rain}} = [3+, 2-], \quad S_{\text{overcast}} = [4+, 0-]$$

$$H(S_{\text{sunny}}) = \left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right) = +0.971$$

~~0.971~~

$$H(S_{\text{rain}}) = 0.971$$

$$H(S_{\text{overcast}}) = 0$$

$$IG(S, \text{outlook}) = 0.94 - \frac{2 \times 5 \times 0.971}{14} = 0.246$$

$$S_{\text{high}} = [3+, 4-] \quad S_{\text{normal}} = [6+, 1-]$$

$$H(S_{\text{high}}) = \frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} = 0.985$$

$$H(S_{\text{normal}}) = \frac{6}{7} \log_2 \frac{6}{7} = 0.19$$

$$IG(S, \text{humidity}) = 0.94 - \frac{7}{14} (0.985) - \frac{7}{14} \times 0.19$$

k-Nearest Neighbour Learning

18/08/2025

Instance based learning \rightarrow 'Lazy / Memory based'

\rightarrow Model is not explicitly trained on a training dataset

Nearest neighbour defined in terms of distance

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n (a_r(x_i) - a_r(x_j))^2}$$

Here i and j are instances, a_r is r th feature

Voronoi Diagram

Draw perpendicular bisectors of each pair of points.

1) Training Algorithm

For each training algorithm $\langle \pi, f(x) \rangle$, add the example to list of training example

2) Classification Algorithm

Given a query instance x_q to be classified

Let x_1, x_2, \dots, x_k denote the k instances from training examples that are nearest to x_q

$$\text{Return } f(x_q) \leftarrow \arg \max_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

$$\delta(a, b) = \begin{cases} 1 & \text{if } a=b \\ 0 & \text{else} \end{cases}$$

For Regression ($f: \mathbb{R}^n \rightarrow \mathbb{R}$)

$$f(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

Here we are assigning the class to unseen/new data points based on True/False (0/1), that is not good.

Distance weighted NN

For classification

$$f(x_q) \leftarrow \arg \max_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i)) \text{ where } w_i = \frac{1}{d(x_q, x_i)^2} \text{ and } \delta(a, b) = \begin{cases} 1 & \text{if } a=b \\ 0 & \text{else} \end{cases}$$

For regression

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}, \text{ where } w_i = \frac{1}{d(x_q, x_i)^2}$$

Linear Regression

21/08/2025

$$y = \theta_1 x_1 + \theta_0 \leftarrow \text{Inductive Bias}$$

$\uparrow x_0 = 1$

$$\hat{y} = h_\theta(x) = \sum_{i=0}^n \theta_i x_i$$

$$x \in \mathbb{R}^d$$

$$\theta \in \mathbb{R}^{d+1}$$

$$J(\theta) = r = y - \hat{y} \quad \checkmark \text{ cost function}$$

$$\text{residual} = \frac{1}{2} \sum_{i=0}^n (\hat{y}_i - y_i)^2$$

$$= \frac{1}{2} \sum_{i=0}^n (h_\theta(x_i) - y_i)^2$$

$$= \frac{1}{2} \sum_{i=0}^n (\theta_0 x_i + \theta_1 x_i - y_i)^2$$

minimise residual based on the value of θ

$$\arg(\min_\theta)$$

$$\arg \min_\theta J(\theta)$$

To get to minima - 3 approach

① Gradient Descent Approach

$$\theta^{(0)} = \text{Init}$$

for all $j \in \{1, \dots, d+1\}$

$$\theta_j^{(1)} = \theta_j^{(0)} - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

$$\theta_j^{(2)} = \theta_j^{(1)} - \alpha \frac{\partial J(\theta^{(1)})}{\partial \theta_j}$$

$$\theta^{t+1} = \theta^t - \alpha \nabla_\theta J(\theta^t)$$

$$= \theta^t - \alpha \nabla_\theta \left[\frac{1}{2} \sum_{i=1}^n (h_\theta(x_i) - y_i)^2 \right]$$

$$\text{diff wrt } \theta = \theta^t - \alpha \nabla_\theta \left[\frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 \right]$$

$$\theta^{t+1} = \theta^t - \alpha \left[\sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right]$$

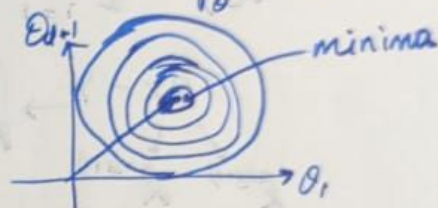
How many iterations?

Repeat Until Convergence

$$\| \theta^{(t)} - \theta^{(t-1)} \|$$

$$| J(\theta^t) - J(\theta^{t-1}) |$$

$$\frac{\partial J(\theta^t)}{\partial \theta_j} \text{ is too small}$$



For each sample updating θ as we calculate

for i in n :

$\theta = \text{Init}$

for j in d :

$\theta_j \leftarrow \nabla_{\theta_j} J(\theta_i)$

Stochastic Grad. Descent

?? In large no. of sample, takes less time to converge

Considering n samples, storing gradients and then updating.

$\theta = \text{Init}$

for i in n :

for j in d :

$\nabla_{\theta_j} J(\theta_i)$

$\theta \leftarrow \nabla_{\theta_j} J(\theta_i)$

Batch Grad. Descent

② Closed form Solⁿ

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x^i - y^i)^2$$

$$\begin{array}{ccc} \theta^T x - y & \rightarrow & \mathbb{R}^n \\ \downarrow & & \downarrow \\ \mathbb{R}^{n \times (d+1)} & & \mathbb{R}^{n \times (d+1)} \\ \mathbb{R} & & \mathbb{R} \end{array}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x - y)^T (\theta^T x - y)$$

$$J(\theta) = 0$$

$$= \frac{1}{2} \sum_{i=1}^n \left[(x\theta)^T x\theta - \underbrace{y^T x\theta - (x\theta)^T y}_{\text{scalar}} + y^T y \right]$$

$$= \frac{1}{2} \sum_{i=1}^n \left[\theta^T (x^T x) \theta - 2 \theta^T (x^T y) + y^T y \right] \quad \left\{ \theta^T A \theta \Rightarrow [A^T A]^T \theta \right.$$

$$= \frac{1}{2} \sum_{i=1}^n \left[2(x^T x) \theta - 2 \theta^T (x^T y) \right]$$

$$(x^T x) \theta - \theta^T (x^T y) = 0$$

$$\left[(x^T x) \theta = \theta^T (x^T y) \right] \leftarrow \text{Normal eq}^n$$

$$\boxed{\theta = (x^T y) (x^T x)^{-1}}$$

Diff between Probability and Likelihood

Probability - ^{when} We know the parameters

Likelihood - When we don't know the parameters

$$y^{(i)} = \theta^T x^{(i)} + \overset{\text{noise}}{\epsilon^{(i)}} \quad \overset{\text{gaussian}}{\epsilon^{(i)} \in N(0, \sigma^2)}$$

$i \in \text{i.i.d.}$

$$\epsilon^{(i)} = y^{(i)} - \theta^T x^{(i)} \sim N(0, \sigma^2)$$

$$\text{If } y^{(i)} - \theta^T x^{(i)} \sim N(0, \sigma^2)$$

$$y^{(i)} \sim N(0, \sigma^2)$$

~~mean~~ would actually not be 0 but depend on $\theta^T x^{(i)}$

$$y^{(i)} \sim \cancel{N(0, \sigma^2)} N(\theta^T x^{(i)}, \sigma^2)$$

$$p(y^{(i)} | x^{(i)}; \theta) \Rightarrow \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\textcircled{3} \text{ Likelihood } L(\theta) \Rightarrow \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma^2}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2\right)$$

$$\log(L(\theta)) = l(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

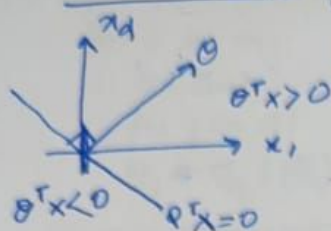
$$l(\theta) = k - \frac{1}{\sigma^2} \sum_{i=1}^n \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2$$

$$\operatorname{argmin}_{\theta} J(\theta) = \operatorname{argmax}_{\theta} l(\theta)$$

Classification - using Perceptron (built on Linear regression)

$$y^{(i)} = g(\underbrace{\theta^T x^{(i)}}_{\text{Linear regression}}) \quad y^{(i)} \in \{0, 1\}$$

Normal vector



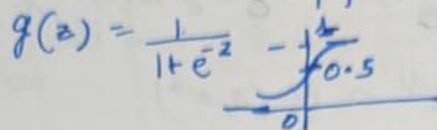
$$D = \{(x_i, y_i)\}_{i=1}^n, y_i \in \{0, 1\}$$

$$y = \theta_0 + \theta_1 x_1$$

$$y = g(\theta^T x)$$



$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{Step function} \quad \begin{cases} 1 & \text{if } g(\cdot) > 0 \\ 0 & \text{if } g(\cdot) < 0 \end{cases}$$



$$\begin{cases} 1 & \text{if } g(z) \geq 0.5 \\ 0 & \text{else} \end{cases}$$

$$h_\theta(x) = g(\theta^T x)$$

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

$$P(y^{(i)} = 1 | x^{(i)}; \theta) = h_\theta(x)$$

$$P(y^{(i)} = 0 | x^{(i)}; \theta) = 1 - h_\theta(x)$$

$$P(y | x; \theta) = (h_\theta(x))^y \times (1 - h_\theta(x))^{(1-y)} \cdot \{(h_\theta(x))^1 \times (1 - h_\theta(x))^{(1-1)}\}$$

$$L(\theta) = -\prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)$$

$$g(z') = g(z) \cdot -g(z)$$

$$l(\theta) = -\sum_{i=1}^n \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$$

Gradient Descent $\rightarrow \nabla_\theta l(\theta)$

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_\theta l(\theta)$$

For single example

$$y \log g(\theta^T x) + (1 - y) \log(1 - g(\theta^T x))$$

$$\nabla_\theta l(\theta) = y \cdot \frac{1}{g(\theta^T x)} g'(\theta^T x) \cdot x + (1 - y) \cdot \frac{1}{1 - g(\theta^T x)} (-1) g'(\theta^T x) \cdot x$$

$$= y \cdot \frac{1}{g(\theta^T x)} g(\theta^T x) (1 - g(\theta^T x)) x + (1 - y) \cdot \frac{1}{1 - g(\theta^T x)} (-1) g(\theta^T x) \cdot (1 - g(\theta^T x)) x$$

$$= y (1 - g(\theta^T x)) x + (1 - y) (-1) g(\theta^T x) x$$

$$= x (y (1 - g(\theta^T x)) + (1 - y) (-1) g(\theta^T x))$$

$$= x (y - y g(\theta^T x) - g(\theta^T x) + y g(\theta^T x))$$

$$\nabla_\theta l(\theta) = x (y - g(\theta^T x))$$

$$\theta_{\text{new}} = \theta_{\text{old}} + \eta [x (y - g(\theta^T x))]$$

SVM - Support Vector Machine

Along with classification also add confidence

γ = min Distance (point, hyperspace)

objective: $\max \gamma$ constraint $(w^T x + b)$

arg max γ
 w, b

$\{(x_i, y_i)\}_{i=1}^n \quad y_i \in \{-1, +1\}$

Max distance with respect to hyperplane (decision surface)

$$\text{distance}(x_i, w^T x_i + b) = \frac{|w^T x_i + b|}{\|w\|}$$

$$\text{dis}(x_i, w^T x_i + b) = \frac{|w^T x_i + b|}{\|w\|}$$

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

for positive example $y_i = +1$

$$w^T x_i + b \geq 1$$

for negative examples $y_i = -1$

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

$$w^T x_i + b \leq -1$$

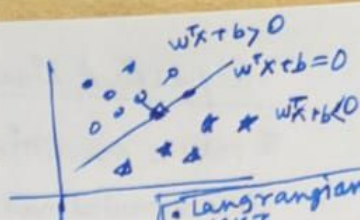
$$\text{dis}(x_i, \text{hyperplane}) = \frac{y_i |w^T x_i + b|}{\|w\|}$$

$$\hat{\gamma}_i = \min_i (\text{dis}(x_i, \text{hyperplane}))$$

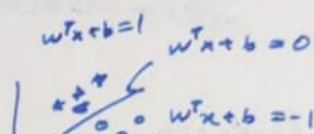
$$(\text{margin}) \rightarrow \gamma = \max(\hat{\gamma}_i)$$

Functional margin - $y_i (w^T x_i + b)$

Geometric margin $\frac{y_i (w^T x_i + b)}{\|w\|}$



- Lagrangian
- KKT
- Primal
- Dual



$$y = \max\left(\frac{1}{\|w\|_2}\right) \Rightarrow \min(\|w\|_2) \quad \text{Primal}$$

$$y \Rightarrow \min\left(\frac{1}{2} \|w\|_2^2\right) \text{ s.t. } y_i (w^T x_i + b) \geq 1$$

So we use Lagrangian $L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b))$ with this method

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

$$\min_{w, b} \left(\max_{\alpha \geq 0} (L(w, b, \alpha)) \right)$$

$$\max_{\alpha \geq 0} \left(\min_{w, b} (L(w, b, \alpha)) \right)$$

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

$$\frac{\partial L}{\partial b} = 0 = - \sum_{i=1}^n \alpha_i y_i \quad 08/09/2025$$

$$\max_{\alpha \geq 0} \left(\frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b)) \right)$$

$$\max_{\alpha \geq 0} \left(\frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (w^T x_i + b) \right)$$

$$= \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (w^T x_i + b)$$

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

$$w = \sum_{i=1}^n \alpha_i y_i x_i \quad \max_{\alpha} w(\alpha) \quad Q = y_i y_j (x_i^T x_j)$$

$$\text{KKT} \quad \max_{\alpha} w(\alpha) = \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

- ① Primal Feasibility: $y (w^T x_i + b) \geq 1$
- ② Dual Feasibility: $\alpha_i \geq 0$
- ③ Stationarity: $w = \sum_{i=1}^n \alpha_i y_i x_i ; \sum_{i=1}^n \alpha_i y_i = 0$
- ④ Slackness: $\alpha_i (y_i (w^T x_i + b) - 1) = 0$

$$f(n) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i (x_i^T x) + b \right)$$

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \text{ s.t. } y_i (w^T x_i + b) \geq 1 - \xi_i$$

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

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

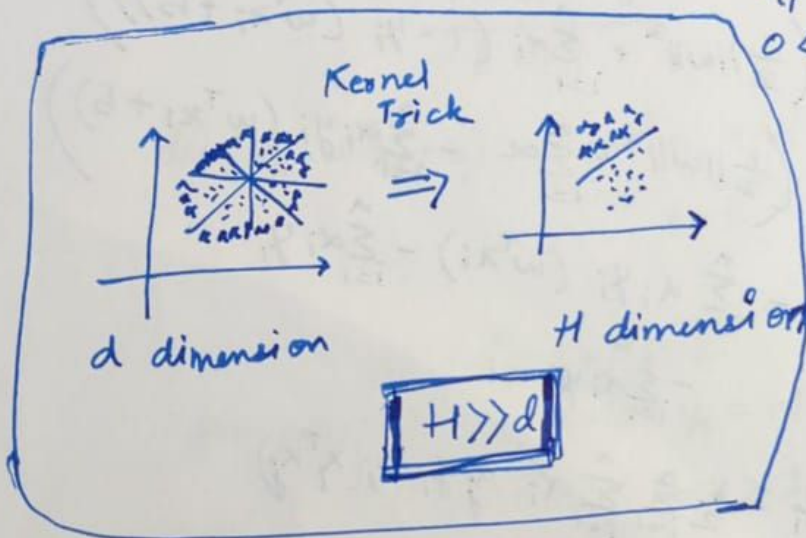
$$\min_{w, b, \xi, \alpha} \frac{1}{2} \|w\|^2 + \sum_i (\xi_i) + \alpha_1 (1 - \xi_i - y_i (w^T x_i + b))$$

$$\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial \xi}$$

$$\max_{\alpha} w(\alpha) = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

$$\text{s.t. } 0 \leq \alpha_i \leq C ; \sum_i \alpha_i y_i = 0$$

\downarrow
 $\alpha_i = 0$ non support vector
 $\alpha_i = C$ inside margin / misclassified
 $0 < \alpha_i < C$ on the margin



Pattern

MSQ

Fill in blank

Some straight forward

Some numerical

Scenario based

Kernel function

Poly:- $K(x_i, x_j) = (x_i^T x_j + c)^q$

eg of more fns: RBF, Gaussian, Sigmoid