

MLT Notes



Prince Kumar

- week 5 to 11(half)
- concept from pyqs
- for last time revesion



Best of luck  

Week-5

Dataset: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^d$

Goal: learn $h: \mathbb{R}^d \rightarrow \mathbb{R}$

$$h(w) = w^T x \quad \forall w \in \mathbb{R}^d$$

Sum of Square Error (SSE)

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

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

Actual labels = y_i

Predicted labels = $w^T x_i$

$$\# f(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 = \|x^T w - y\|_2^2$$

Solⁿ:

$$w^* = (x x^T)^{-1} (x y)$$

$x^T w^*$: projection of labels onto the space spanned by the features.

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

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

① Gradient Descent

② Stochastic gradient Descent

$D = \{x_1, x_2, x_3, \dots, x_n\}$

$$w^{t+1} = w^t - \eta \nabla f(w^t)$$

Step Size

$$\nabla f(w) = 2(x x^T)w - 2(x y)$$

① To sample uniformly the data points —

$$D' = \{x_1, x_T, x\} \quad t = 1, 2, \dots, T$$

$$w^{t+1} = w^t - \eta^T (2(\bar{x} \bar{x}^T)) w^t - 2 \bar{x} \bar{y}$$

②

$$w_{SGD}^T = \frac{1}{T} \sum_{t=1}^T w^t$$

বিদ্যালয় শিক্ষা মন্ত্র, পশ্চিমবঙ্গ সরকার

$(y|x) = w^T x + \epsilon$, $\epsilon \sim \text{Normal}(0, \sigma^2)$

$$(y|x) \sim N(w^T x, \sigma^2)$$

Parameter w

$$L(w) = L(w; \dots) = \prod_{i=1}^n f_{y|x}(y_1, y_2, \dots, y_n, \theta)$$

after solve —

$$w^* = \hat{w}_{ML} = (x x^T)^{-1} x y$$

Conclusion:- Maximum likelihood estimator assuming zero Mean Gaussian Noise is same as linear regression with SQUARED ERROR.

Expected value of \hat{w}_{ML} —

$$\mathbb{E}[\|\hat{w}_{ML} - w\|^2] = \sigma^2 (\text{trace}(x x^T)^{-1})$$

$$\text{Eigen val. of } (x x^T)^{-1} = \left\{ \frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_d} \right\}$$

So, $\mathbb{E}[\|\hat{w}_{ML} - w\|^2] = \sigma^2 \sum_{i=1}^d \frac{1}{\lambda_i}$

So, $\mathbb{E}[\|\hat{w}_{ML} - w\|^2] = \sigma^2 \left(\sum_{i=1}^d \frac{1}{\lambda_i} \right)$

$\hat{w}_{new} = (x x^T + \lambda I)^{-1} (x y)$ mean square error (MSE)

$\text{Trac}(\quad) = \sum_{i=1}^d \frac{1}{\lambda_i + \lambda}$

So, $\mathbb{E}[\|\hat{w}_{new} - w\|^2] = \sigma^2 \sum_{i=1}^d \frac{1}{\lambda_i + \lambda}$

mean square error (MSE)

$\star \quad \text{MSE}(\hat{w}_{new}) < \text{MSE}(\hat{w}_{ML})$

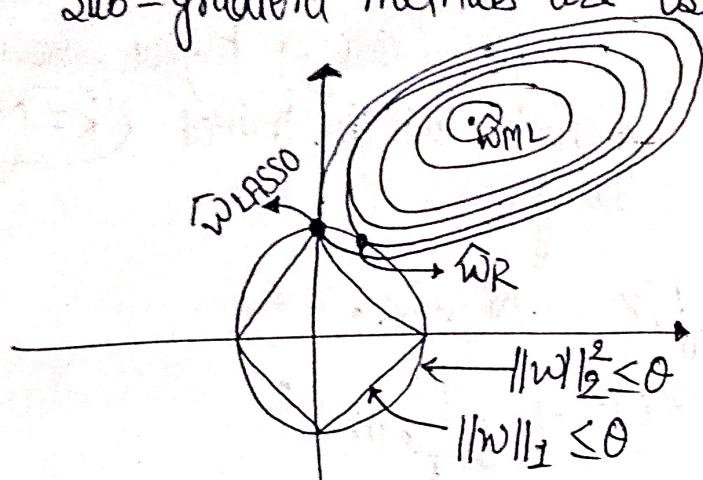
$\lambda > 0, \lambda \in \mathbb{R}$

LASSO Regression

(Least Absolute Shrinkage & Selection Operator)

→ LASSO does not have a closed form soln.

→ Sub-gradient methods are usually used to solve LASSO.



* An alternate way to regularize would then be using $\|w\|_1$ norm instead of $\|w\|_2$ norm

$$\|w\|_1 = \sum_{i=1}^n |w_i|$$

KNN (K-Nearest Neighbors), K is hyper parameter

→ Smaller the k, complicated the decision boundary

→ Solution: Cross-validate for k.

Issues :- choosing a distance function

① Manhattan distance :- $D(x_1, y_1) \& D = (x_2, y_2)$

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

② Euclidean Distance :-

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

:- Prediction is COMPUTATIONALLY EXPENSIVE. It involves sorting the datapoints acc. to its distance from test point.

:- for prediction we have to use the entire dataset. We can't throw the dataset after "learning." Since no MODEL is learnt.

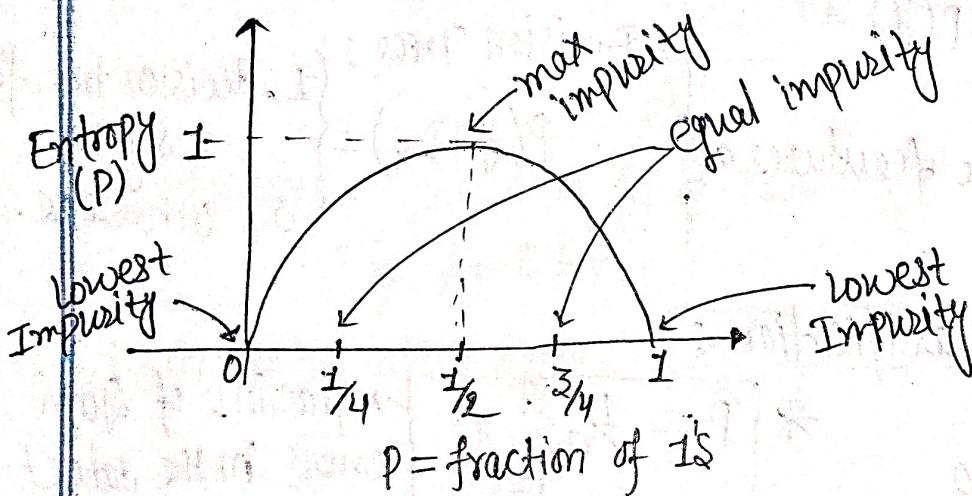
Decision Trees → This will used to make the prediction

* Decision tree Algorithm partitions the feature space (here \mathbb{R}^2) into horizontal & by cutting the line at thresholds. vertical half spaces.

* D.T Algorithm partitions the Space into intervals by cutting the line at thresholds. {for \mathbb{R}^1 }

* in $\mathbb{R}^d \rightarrow 3D$ Boxes (rectangular prism)

Pure Data :- If the all labels corresponding to the features is same nature (+ve or -ve)



$$\text{Entropy}(D) = - [P \log_2 P + (1-P) \log_2 (1-P)]$$

where, $P = \frac{\text{fraction of 1's}}{\text{total Data Points}}$

& $(1-P)$ is fraction of another class.

D: $f_k \leq 0$ Entropy (D)

yes

D_{Yes}

D_{No}

$$G = \text{Entropy}(D) \cdot [P_{\text{Yes}} E(D_{\text{Yes}}) + P_{\text{No}} E(D_{\text{No}})]$$

Entropy (D_{Yes})

Entropy (D_{No})

$$\text{Information Gain} = \text{Entropy}(D) - [Y \text{Entropy}(D_{\text{Yes}}) + (1-Y) \text{Entropy}(D_{\text{No}})]$$

(I.G)

$$\text{where, } Y = \frac{|D_{\text{Yes}}|}{|D|}$$

Types of modelling

Generative model

- Joint Distribution: $P(x, y)$

features labels

$$= P(y/x) \cdot P(x)$$

$$= P(x/y) \cdot P(y)$$

- Naive Baye's
- Interested in how the features are generated

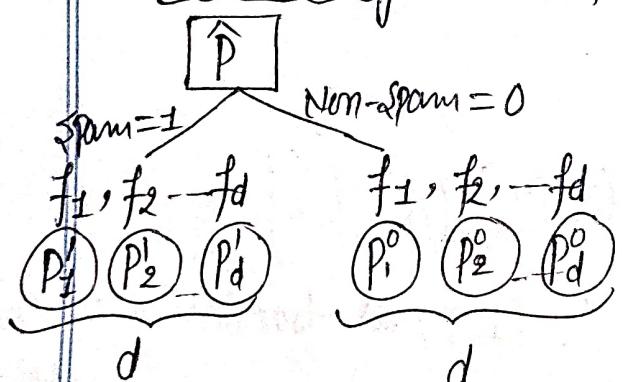
Discriminate models

Conditional dis.: $P(y/x)$

KNN: $P(y=1/x) = \begin{cases} 1, & \text{if majority of} \\ & \text{neighbours say 1} \\ 0, & \text{otherwise} \end{cases}$

Decision Trees: $P(y=1/x) = \begin{cases} 1, & \text{decision tree for} \\ & x \text{ say 1} \\ 0, & \text{otherwise.} \end{cases}$

Naive Baye's classification :-



* $\hat{P} = \frac{1}{n} \sum_{i=1}^n y_i$ → fraction of spam emails in the dataset

No. of paraMetr → $2d+1$

Bayes Rule :-

$$P(y^{\text{test}}=1/x^{\text{test}}) = \frac{P(x^{\text{test}}/y^{\text{test}}=1) \cdot P(y^{\text{test}}=1)}{P(x^{\text{test}})}$$

$$P(y^{\text{test}}=0/x^{\text{test}}) = \frac{P(x^{\text{test}}/y^{\text{test}}=0) \cdot P(y^{\text{test}}=0)}{P(x^{\text{test}})}$$

$$\text{if } \left(\prod_{j=1}^d (\hat{P}_j^1)^{f_j} (1-\hat{P}_j^1)^{1-f_j} \right) \cdot \hat{P} > \left(\prod_{j=1}^d (\hat{P}_j^0)^{f_j} (1-\hat{P}_j^0)^{1-f_j} \right) \cdot (1-\hat{P})$$

then $\hat{y}_{\text{test}} = 1$, else $\hat{y}_{\text{test}} = 0$

Decision fn of Naive Bayes

Class: 0 & 1 predict $y_{\text{test}} = 1$, if $P(y_{\text{test}} = 1/x_{\text{test}}) > P(y_{\text{test}} = 0/x_{\text{test}})$

Given x_{test}

$$\frac{P(x_{\text{test}}/y_{\text{test}} = 1) P(y_{\text{test}} = 1)}{P(x_{\text{test}}/y_{\text{test}} = 0) P(y_{\text{test}} = 0)} > 1$$

Hence, decision fn is of the

form predict $y_{\text{test}} = 1$ if

$$w^T x_{\text{test}} + b \geq 0$$

$\in \mathbb{R}^d$

$$\text{where, } w_i = \log \left(\frac{\hat{P}_i^1 (1-\hat{P}_i^1)}{\hat{P}_i^0 (1-\hat{P}_i^0)} \right)$$

* Conclusion :- Decision fn of Naive Bayes is LINEAR.

Gaussian Naive Bayes :-

$$f(x/y) = \frac{1}{\sigma_y \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_y}{\sigma_y^2} \right)^2}$$

Simplest Assumption

$$P(y_i = 1/x) = \begin{cases} 1 & \text{if } w^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Linear Separability Assumption -

$$\exists w \in \mathbb{R}^d \text{ s.t. } \text{Sign}(w^T x_i) = y_i \quad \forall i \in [n]$$

PERCEPTRON :- Input: $\{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$

Iteration

Until Convergence

$$w^0 = 0 \in \mathbb{R}^d$$

- Pick (x_i, y_i) pair from the dataset
- If $\text{Sign}(w^T x_i) = y_i$
do nothing

ELSE

$$w^{t+1} = w^t + x_i y_i \quad \leftarrow \text{Update Rule}$$

end

Update Rule \rightarrow Update Rule pushes w in the "right" direction for x_i if mistake

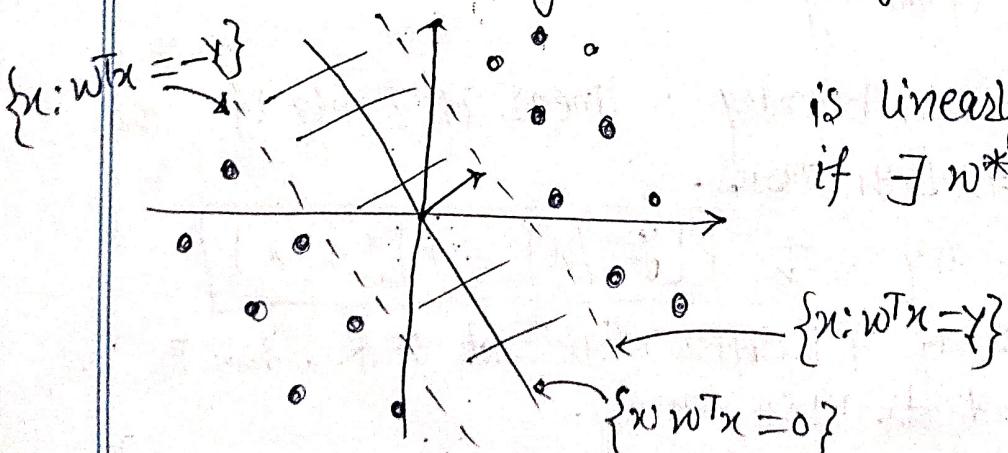
$$w^{t+1} = w^t + x_i y_i$$

Assumption

① Linear Separability with γ -Margin :- A Dataset

$$\{(x_1, y_1), \dots, (x_n, y_n)\}$$

is linearly separable with γ margin
if $\exists w^* \in \mathbb{R}^d$ s.t. $(w^* x_i) y_i \geq \gamma \quad \forall i$ for some $\gamma \geq 0$



② Radius Assumption

$\forall i \in D \quad \|x_i\|_2 \leq R \text{ for some } R > 0$

③ without loss of generality, assume $\|w^*\| = 1$

$$Y = w^*^T x, \quad Y' = \frac{Y}{\|w^*\|}$$

Analysis of "Mistakes" of perception —

$$\rightarrow \|w^l + 1\|^2 \leq lR^2 \quad (1)$$

Mistakenly

$$\rightarrow (w^l + 1)^T w^* \geq l Y \quad (2)$$

$$\rightarrow \|w^l + 1\|^2 \geq l^2 Y^2 \quad (3)$$

$$w^l + 1 = w^l + x \cdot y$$

for any x, y —

$$(x^T y)^2 \leq \|x\|^2 \|y\|^2$$

(Cauchy-Schwarz)

$$\Rightarrow l^2 Y^2 \leq \|w^l + 1\|^2 \leq lR^2$$

from (3)

from (1)

mistakes is bounded
becoz $[Y > 0]$

$$\Rightarrow l \leq \frac{R^2}{Y^2} \quad \begin{matrix} \leftarrow \text{Radius Margin} \\ \text{Bound} \end{matrix}$$

$l \rightarrow \# \text{ mistakes}$

Perception Convergence!

Larger the score ($z = w^T x$), more the probability of being +1.

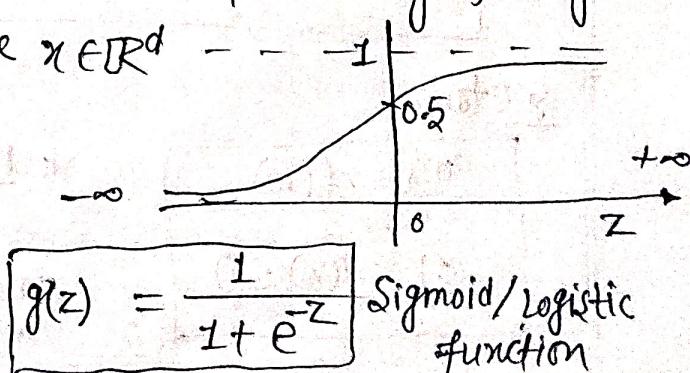
score $z \in [-\infty, +\infty]$ feature $x \in \mathbb{R}^d$

link function

$g(z) = 0.5 \text{ if } z = 0$

$g(z) \rightarrow 1 \text{ as } z \rightarrow \infty$

$g(z) \rightarrow 0 \text{ as } z \rightarrow -\infty$



Model: Logistic Regression

বিদ্যালয় শিক্ষা দপ্তর, পশ্চিমবঙ্গ সরকার

$$P(y=1/x) = \frac{1}{1+e^{-w^T x}} = g(w^T x)$$

Gradient Update Rule :-

$$\begin{aligned} w_{t+1} &= w_t + \eta_t \nabla \log L(w_t) \\ &= w_t + \eta_t \left(\sum_{i=1}^m x_i \left(y_i - \underbrace{\frac{1}{1+e^{-w^T x_i}}}_{g(w_t^T x_i)} \right) \right) \end{aligned}$$

$\in \mathbb{R}^d$ $\{0, 1\}$

Kernel version

- Can argue $w^* = \sum_{i=1}^m \alpha_i x_i$ formal theorem is called the
Regeator theorem

Regularized Version

$$\min_w \sum_{i=1}^m \left[\log \left(1 + e^{-w^T x_i} \right) + w^T x_i (1 - y_i) + \underbrace{\frac{\lambda}{2} \|w\|^2}_{\text{Regularization}} \right]$$

Cross validate hyperparam

Goal: To come up with a formulation that maximizes "margin".

$\lambda = \frac{1}{2} \|w\|^2$

$$\boxed{\begin{array}{l} \max_{w, y} \\ \text{s.t. } (w^T x_i) y_i \geq 1 \quad \forall i \\ \|w\| = 1 \end{array}}$$

$\# \quad \max_w \frac{2}{\|w\|^2} = \text{width}(w)$

$$\text{s.t. } (w^T x_i) y_i \geq 1 \quad \forall i$$

Detour

$$\boxed{\begin{array}{l} \min_w f(w) \\ \text{s.t. } g(w) \leq 0 \end{array}}$$

$$* L(w, \alpha) = f(w) + \alpha g(w)$$

$$\max_{\alpha \geq 0} f(w) + \alpha g(w) = \begin{cases} \infty & g(w) > 0 \\ f(w) & g(w) \leq 0 \end{cases}$$

বিদ্যালয় শিক্ষা মন্ত্র, পঞ্জিকবঙ্গ সরকার
 → Substitute back value of w^* in the objective

$$w^* = X Y \alpha$$

$$X = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_n \\ | & | & | & \dots & | \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}_{d \times n} \begin{bmatrix} y_1 & 0 \\ 0 & \dots & y_n \end{bmatrix}_{n \times n} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

$$\frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - (w^T x_i) y_i)$$

on simplification

$$w^* = X Y \alpha$$

$$= \alpha^T I - \frac{1}{2} (X Y \alpha)^T (X Y \alpha) \quad \text{where, } I = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times n}$$

PRIMAL

$$\begin{aligned} \min_{w \in \mathbb{R}^d} & \frac{1}{2} \|w\|^2 \\ \text{s.t. } & (w^T x_i) y_i \geq 1 \end{aligned}$$

DUAL problem

$$\max_{\alpha \geq 0} \alpha^T I - \frac{1}{2} \alpha^T Y^T X^T X Y \alpha \quad \begin{matrix} \in \mathbb{R}^{n \times n} \\ \text{Kernel } K \end{matrix}$$

What have we gained?

- Dual Variable dimension in \mathbb{R}^n , while primal problem dimension in \mathbb{R}^d
- Dual Constraints are "easier"
- More Importantly dual problem on $X^T X$ & So can be "KERNELISED"

Revisiting the Lagrangian —

$$\underbrace{\min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right]}_{w^* \text{ is the primal soln}} \equiv \max_{\alpha \geq 0} \underbrace{\left[\min_w f(w) + \alpha g(w) \right]}_{\text{Dual}}$$

α^* is the dual soln

$$\max_{\alpha \geq 0} f(w^*) + \alpha^* g(w^*) = \min_w f(w) + \alpha^* g(w)$$

After solve :-

$$\alpha^* g(w^*) \geq 0 \quad \text{--- (1)}$$

But we already know $\alpha^* > 0$ & $g(w^*) \leq 0$

So,

$$\alpha^* g(w^*) = 0$$

Complementary Slackness

for multiple constraint -

$$\alpha_i^* g_i(w^*) = 0 \quad \forall i$$

$$\text{If } \alpha_i^* > 0 \xrightarrow{\text{C.S}} 1 - (w^{*T} x_i) y_i = 0$$

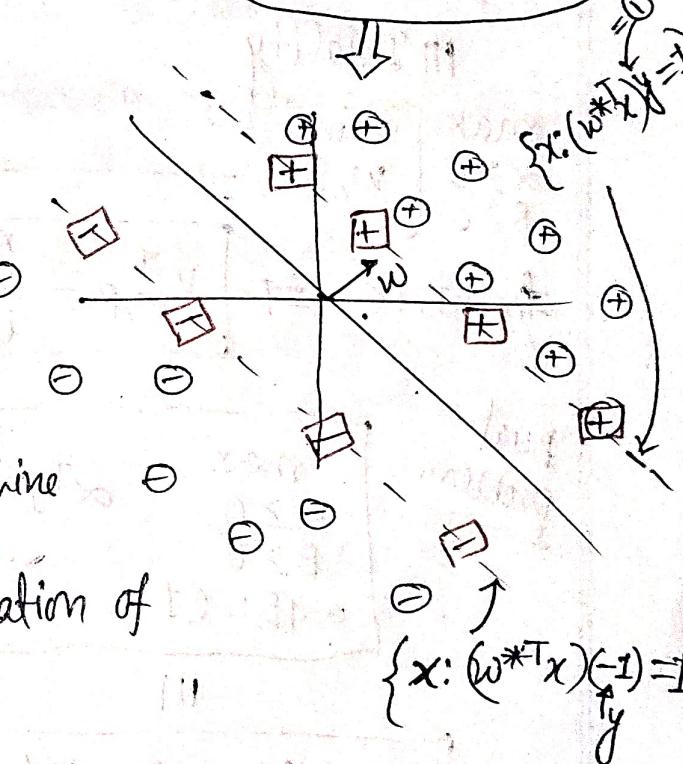
$$(w^{*T} x_i) y_i = 1$$

→ Only the points that are on the "Supporting" hyperplane can contribute to w^*

→ These special points are called "Support Vectors".

→ ALGORITHM → Support Vector Machine (SVM)

→ w^* is a sparse linear combination of the data points



Given x_{test} : $w^*(\phi(x_{\text{test}})) = \sum_{i=1}^n \alpha_i^* y_i k(x_i, x_{\text{test}})$

Soft Margin formulation

→ used: when data contains outliers or, not perfectly linearly separable

Primal

$$\begin{aligned} & \min_{w, \epsilon} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \\ & \text{s.t. } (w^T x_i) y_i + \epsilon_i \geq 1 \quad \forall i \quad \& \epsilon_i \geq 0 \end{aligned}$$

Hyperparameter $C \geq 0$

Small C : Margin wide
(High bias & low var)
Error Jyada Allow
Karta hai

বিদ্যালয় শিক্ষা দপ্তর, পশ্চিমবঙ্গ সরকার

$C=0 \Rightarrow$ Bribes don't cost
 $\rightarrow w=0 \in \mathbb{R}^d$ is the soln & $C=\infty \Rightarrow$ linear separable case
 $\downarrow \rightarrow$ Hard margin

$L(w, \varepsilon, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \varepsilon_i + \sum_{i=1}^n \alpha_i (1 - w^T x_i y_i - \varepsilon_i)$

$+ \sum_{i=1}^n \beta_i (-\varepsilon_i)$

DUAL Problem

$$\min_{w, \varepsilon} \max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \left[\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \varepsilon_i + \sum_{i=1}^n \alpha_i (1 - w^T x_i y_i - \varepsilon_i) + \sum_{i=1}^n \beta_i (-\varepsilon_i) \right]$$

III Duality

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \min_{w, \varepsilon}$$

fix $w, \alpha \Rightarrow w^*_{\alpha, \beta} = \sum_{i=1}^n \alpha_i x_i y_i$ — ① & $\alpha_i + \beta_i = C$ — ②

Dual problem

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0 \\ \alpha + \beta = C}} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T Y^T X^T X Y \alpha$$

no β term

III

$$\max_{0 \leq \alpha \leq C} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T Y^T X^T X Y \alpha$$

Kernelizable

If $C=0$
 $\Rightarrow \alpha^* = 0 \in \mathbb{R}^n$
 $\Rightarrow w^* = 0$

& $C \neq 0 \Rightarrow$ Hard margin

Complementary Slackness {about Soft Margin}

let (w^*, ε^*) be the primal optimal soln

& (α^*, β^*) be the dual optimal soln.

C.S

$$\textcircled{1} \quad \forall i \quad \alpha_i^* (1 - w^T x_i y_i - \varepsilon_i^*) = 0$$

$$\textcircled{2} \quad \forall i \quad \beta_i^* (\varepsilon_i^*) = 0$$

Various Cases

$$\textcircled{1} \quad \alpha_i^* = 0 \Rightarrow \beta_i^* = 0 \xrightarrow{\text{C.S } \textcircled{2}} \boxed{\varepsilon_i^* = 0} \xrightarrow{\text{C.S } \textcircled{1}} \boxed{w^T x_i y_i \geq 1}$$

$$\textcircled{2} \quad \alpha_i^* \in (0, C) \Rightarrow \beta_i^* \in (0, C) \xrightarrow{\text{C.S } \textcircled{2}} \boxed{\varepsilon_i^* = 0} \xrightarrow{\text{C.S } \textcircled{1}} \boxed{w^T x_i y_i = 1}$$

$$\textcircled{3} \quad \alpha_i^* = C \Rightarrow \beta_i^* = 0 \xrightarrow{\text{C.S } \textcircled{2}} \boxed{\varepsilon_i^* \geq 0} \xrightarrow{\text{C.S } \textcircled{1}} \boxed{w^T x_i y_i \leq 1}$$

Points where either x_i is incorrectly classified by w^* or
Correctly classified but with margin ≤ 1 .

let's see this from the primal point of view.

$$\text{Case ①} \quad w^T x_i y_i < 1 \Leftrightarrow \varepsilon_i^* \geq 1 - w^T x_i y_i$$

$$\Rightarrow \varepsilon_i^* > 0 \xrightarrow{\text{E.S } \textcircled{2}} \boxed{\beta_i^* = 0} \Rightarrow \boxed{\alpha_i^* = C}$$

$$w^T x_i y_i + \varepsilon_i^* \geq 1$$

$$\text{Case ②} \quad w^T x_i y_i = 1 \Rightarrow \varepsilon_i^* \geq 0 \Rightarrow \boxed{\alpha_i^* \in [0, C]}$$

$$\text{Case ③} \quad w^T x_i y_i \geq 1 \Rightarrow \boxed{\alpha_i^* = 0}$$

বিদ্যালয় শিক্ষা দপ্তর, পশ্চিমবঙ্গ সরকার
Goal: Meta classifiers or Ensemble classifiers.

Weak learner → Strong learner

"The classifiers whose performance is slightly better than the random classifier"

Overfitting - fit noise (High variance)

Underfitting - missing out on stuff thinking it is noise.
(High bias)

$$\rightarrow \text{error} = \text{bias} + \text{Var}$$

Bagging → "Averaging Reduces Variance"

বিদ্যালয় শিক্ষা দপ্তর, পশ্চিমবঙ্গ সরকার

Extra Concept & Must Remembering point -

* \hat{w}_{MLE} : $\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2$

* \hat{w}_L : $\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2 + \gamma \|w\|$

* \hat{w}_R : $\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2 + \gamma / \|w\|^2$

MLE for Bernoulli parameter (\hat{p}) = $\frac{w}{N}$

- $w \rightarrow$ No. of Success
- $N \rightarrow$ total no. of Sample

* $x = \begin{bmatrix} x_1 & x_2 & \dots \end{bmatrix} \leftarrow \text{feature}$

$\uparrow \quad \uparrow$

$x_1 \quad x_2$

$\Rightarrow \text{Solve } x^T = \begin{bmatrix} -x_1 & \dots \\ -x_2 & \dots \\ \vdots & \vdots \end{bmatrix}$

Kernel Regression :-

$$w^* \phi(x_{test}) = \sum_{i=1}^n \underbrace{\alpha_i^* k(x_i, x_{test})}_{\text{Kernel fn}} \quad \text{where, } \alpha^* = K^{-1} y$$

- { ① $(1 + x_i x_j)^2$: Polynomial kernel
- ② $e^{-\|x_i - x_j\|^2 / 2\sigma^2}$: Gaussian kernel

* error = $\sum (h(x_i) - y_i)^2$

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

$CV_{k\text{-fold}} = \frac{1}{k} \sum_{i=1}^k (\text{MSE})$

* MSE (Mean Square Error) = $\frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$ & MAE (Mean Absolute Error) = $\frac{1}{n} \|y_i - w^T x_i\|$

School Education Department, Government of West Bengal

$$\underbrace{*(y - x^T w^*)^T}_{\text{Error vec}} \underbrace{(x^T w^*)}_{\text{vector prediction}} = 0 \quad \left. \begin{array}{l} \text{Error vector is } \perp \text{ (orthogonal) to} \\ \text{Predictions Vector.} \end{array} \right\}$$

Gini Impurity for Binary classification —

$$G_I(P) = 1 - P^2 - (1-P)^2 \rightarrow \frac{1}{2} \geq P \text{ So, } P: 0 \rightarrow \frac{1}{2} \quad G_I(P) \uparrow$$

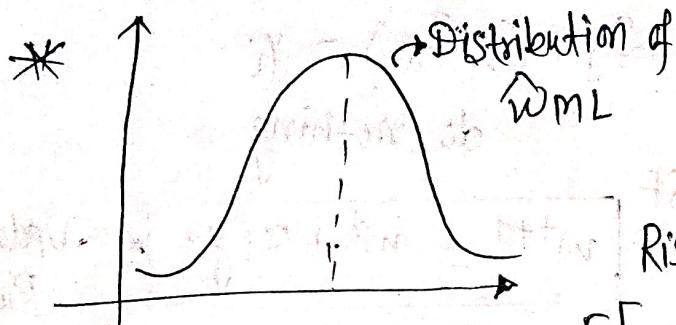
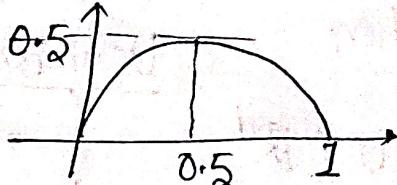
$$P: \frac{1}{2} \rightarrow 1 \quad G_I(P) \downarrow$$

for multiple classes

$$G_I \cdot I = 1 - \sum_{i=1}^m p_i^2$$

$m = \# \text{ of classes}$

$\rightarrow G_I(P)$ attains its max val. at $P = \frac{1}{2}$

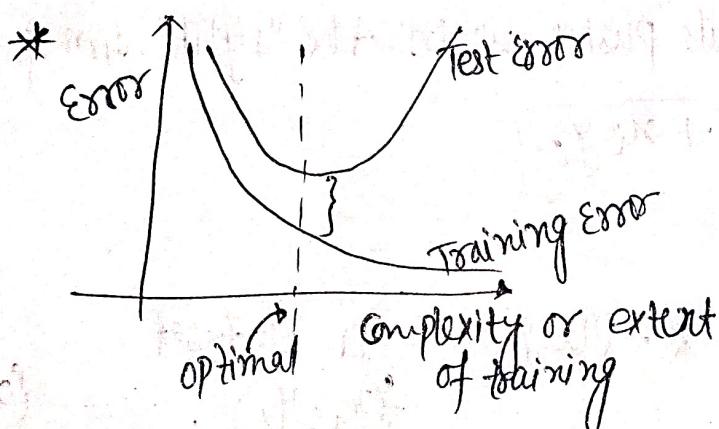


Unbiased estimator of w

$$E[\hat{w}_{ML}] = w$$

Risk = (Bias)² + Variance

$$E[\text{Risk}] = E[\hat{w}_{ML} - w]^2 + E[(\hat{w}_{ML} - w)]^2$$



weight vector will be perpendicular to data separator line.

Theorem :- The decision boundary is linear if & only if the variance are equal for both classes.

The decision boundary is $\Rightarrow P[y=0/x] = P[y=1/x]$

Sub-gradient fact :- if f is differentiable at $a \in \mathbb{R}$, then it has one sub-gradients, ELSE many.

if a matrix is diagonal then

$$\begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix} \rightarrow \text{Inverse of matrix} = \begin{bmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{a} \end{bmatrix}$$

MLE for exponential: $\hat{\lambda}_{MLE} = \frac{n}{\sum x_i}$

with Uniform prior, MAP = MLE (if within bounds)

Exponential distribution (likelihood)

$$X \sim Exp(\lambda)$$

$$L(\lambda | x_1, \dots, x_n) = \lambda^n e^{-\lambda \sum x_i}$$

Prior Dist. $\lambda \sim Uniform(0, n)$

$$P(\lambda) = \begin{cases} \frac{1}{n}, & 0 < \lambda < n \\ 0, & \text{otherwise} \end{cases}$$

posterior dis —

$$P(\lambda | x_1, \dots, x_n) \propto L(\lambda | x_i) P(\lambda)$$

~~$\lambda^n e^{-\lambda \sum x_i}$~~ constant

↓
Apply log & diff &
get MAP estimator

$$\hat{\lambda}_{MAP} = \arg \max_{0 < \lambda < n} [n \ln \lambda - \lambda \sum x_i]$$

PCA Intuition

if Rank of Cov. matrix $\leq n$

then PCA will have — ~~n~~ non zero e.v
-& Remaining e.v = 0

Necessary Condition for Cov. mat —

① Symmetric ($C = C^T$)

② Positive Semi Define (PSD)

- All Eigenval ≥ 0

- Equivalently: all principal minors ≥ 0

if $Uniform(0, \theta)$ [likelihood]

$$\text{Likelihood} = \begin{cases} 0^n, & 0 \geq \max(x_i) \\ \frac{1}{(\theta - a)^n}, & a < \max(x_i) \\ 0, & \text{otherwise} \end{cases}$$

$$\hat{\theta}_{MAP} = \max(x_i)$$

for Uniform Prior

Sigmoid or output -

$$P(y=+1/x) = \sigma(z) = \frac{1}{1+e^{-z}}$$

$$P(y=-1/x) = 1 - \sigma(z)$$

Decision Rule Usually -

$$\text{Predict class} = \begin{cases} 1 & \text{if } \sigma(z) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Where

$$z = w^T x$$

Decision trees are prone to overfitting if the maximum depth is set too deep
 & Too shallow depth → underfitting

Decision tree depend on split selection purity (IG/Gini),
 Order pe nahi.

Decision tree are sensitive [Small perturbations → diff. tree
 in dataset → Sensitive]
 means diff best split
 इसी वजह से decision trees Unstable / high var होते हैं।

Decision trees can handle both numerical & categorical fn
 threshold split, value based split
 (Ex: $x < 5$)

Ridge Regression ka golden Rule -

$$\hat{w}_{\text{Ridge}} = (X X^T + \lambda I)^{-1} X y \quad \text{if } \lambda \uparrow (\text{increase}) \rightarrow \text{Efficient shrink}$$

होता है (O ki taraf)

$\lambda \downarrow$ (decrease) → off. bade ho jate
 thi (ML के pass)

All data Same class then Entropy is 0 & decision tree is pure. & No further split, $IG = 0$.

* Data (50-50 split) \Rightarrow Entropy is maximum.

$IG = \text{Entropy}(\text{Parent}) - \text{Weighted Entropy}(\text{children})$

$\text{Entropy}(\text{parent}) \geq \text{Weighted Entropy}(\text{child})$

Always

So, $IG \geq 0$ এবং IG হলো
ই নেগেটিভ নহীন

Only Lasso can produce exactly zero coefficient.

$\epsilon_i = 0 \rightarrow \text{no bribe}$ $\quad y_i(w^T x_i + b) \geq 1$
 $\epsilon_i > 0 \rightarrow \text{positive bribe}$

If Model \rightarrow Simple (linear) \Rightarrow High bias & Underfitting.

\rightarrow Complex (polynomial) \Rightarrow high variance & overfitting

~~PRIOR~~ $\rightarrow \beta(\alpha, \beta)$ Posterior mean = $\frac{\alpha + k}{\alpha + \beta + n}$ Use prior

Beta(α, β) where, $n \rightarrow$ trials
 $k \rightarrow$ success MLE $[\hat{P}_{MLE}] = \frac{k}{n}$ not use in prior
 $\hat{P}_{MAP} = \frac{\alpha + k - 1}{\alpha + \beta + n - 2}$ Use prior

FADA-BOOST Algorithm:-

Initialize $D_t(i) = \frac{1}{n}$ \leftarrow total no. of sample point
Iteration

for $t = 0, \dots, T$

~~Error(E)~~ \leftarrow ~~0~~
 $E_i = \text{Error}(i) \times D_t(i)$
training Error of f_t

$$\alpha_t = \ln \frac{1 - \epsilon_i}{\epsilon_i} = \frac{1}{2} \ln \left(\frac{1 - \epsilon_i}{\epsilon_i} \right)$$

weight ↑

$$\text{Update Rule: } D_{t+1}(i) = \begin{cases} D_t(i) e^{\alpha_t} & \text{if misclassified, } \epsilon_i \neq y_i \\ D_t(i) e^{-\alpha_t} & \text{if correctly classified} \end{cases}$$

weight ↓

$$D_{t+1}(i) = \frac{D_t(i)}{\sum_j D_t(j)}$$

$$h_T^*(x) = \text{Sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

One can prove:-

$$\text{If } T \geq \frac{1}{2} \frac{1}{\epsilon^2} \ln(2n), \text{ then Training Error} = 0$$

→ How good is my weak learner?

Low threshold $\Rightarrow FN = \text{Very low}$ & $FP = \text{No High}$
 → Spam, fraud, disease detection

$FN \downarrow$
 false Neg.

High threshold $\Rightarrow FN = \text{No High}$ & $FP \rightarrow \text{low}$
 → Emails filtering, Recommendation

$FP \downarrow$
 false positive

$$\# \hat{P}_i^y = P(x_i = 1 | y) \\ = 1 - P(x_i = 0 | y)$$

or

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$\# \text{Hinge loss } (L_i) = \max(0, 1 - y_i f(x_i)) \text{ where, } f(x_i) = w^T x_i + b$$

~~if $1 - y_i f(x_i) \leq 0$ then 0~~

Soft Margin SVM

$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

C : misclassified/margin violation ka
penalty weight
Margin width $\propto \frac{1}{\|w\|}$

Small C :- Error Zayda Allow करता है।

:- Margin wide hota hai

:- High bias, Low var.

Large C :- Error ka hard punish karta hai

→ Margin narrow (छोटा) होता है।

→ Low bias, high var.

Hinge loss - SVM
Logistic loss - logistic regression

modified hinge loss
→ perceptron

Squared loss
→ least square classification

for logistic Regression with threshold T ,

the decision boundary is —

$$* P(Y=1/x) = \sigma(z)$$

$$\sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} = T \Rightarrow w^T x = \ln\left(\frac{T}{1-T}\right) *$$

feature x_1 & x_2 ↗ binary classification problem [Total]

At prediction (\hat{y}) = 1 if $w^T x \geq 0$ for Angle θ —

Else (\hat{y}) = 0 if $w^T x \leq 0$

$$w = \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}$$

Total no. of cluster assignments possible for k -cluster & n -datapts.

$$= k^n$$

~~Decision~~ Logistic regression : $P(Y=1/x) = \sigma(w^T x)$

→ Decision boundary : $w^T x = 0$

→ probability if $w^T x$ pe defend करती है।

→ w की dim में जितना आगे, prob. बढ़ती है।

→ boundary ke opposite side पर, probability small.

School Education Department, Government of West Bengal

- # In AdaBoost, a weak learner must have misclassification rate strictly less than 0.5.
- # Deep decision trees → low bias, high variance → training accurate → test risky

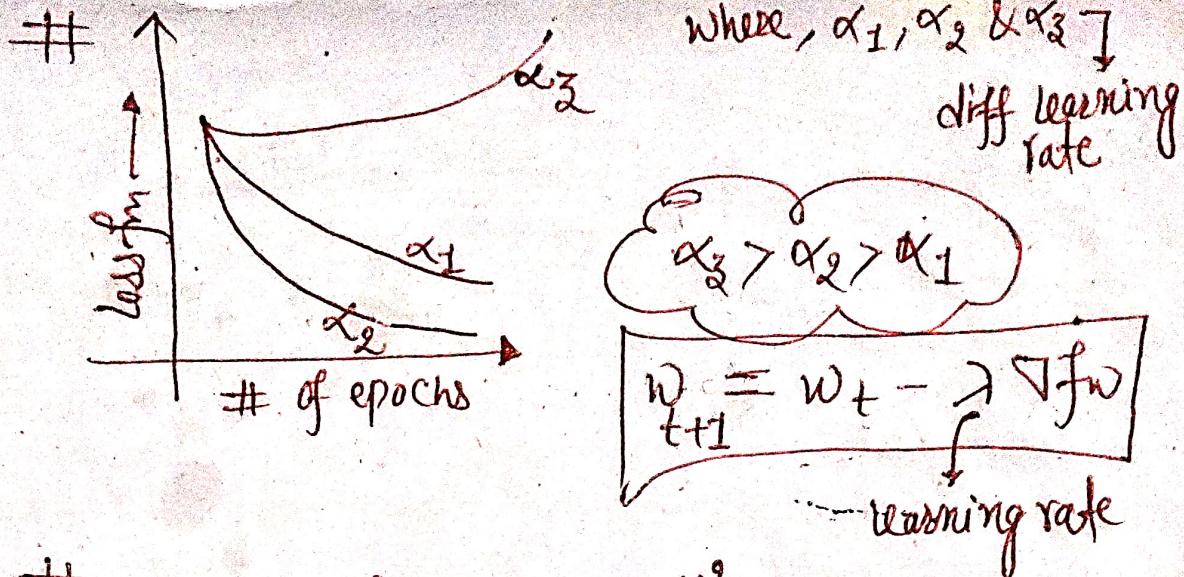
$$\text{If } \underbrace{\text{Normal}(\mu, \sigma^2)}_{\text{PDF } f_X(x)} := \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

if The mean of both the clusters are the same then
these will be no reassignment.

- # Variance of the dataset along a Principal Component
 = eigenvalue corresponding to that component
- $$\lambda_1 = \mathbf{v}_1^T C \mathbf{v}_1$$
- # The Trace of $\mathbf{X}\mathbf{X}^T$ is equal to the trace of $\mathbf{X}^T\mathbf{X}$
- # any kernel of the form $(\mathbf{x}_1^T \mathbf{x}_2 + C)^P$ is valid if $P \in \mathbb{N}$
 $P > 0$
 $C \geq 0$
- # RBF/Gaussians for valid : $e^{-\frac{1}{2}\sigma^2 \|\mathbf{x}_1 - \mathbf{x}_2\|^2}$
- # $X \rightarrow \text{Geo}(P)$ $p(x=k) = (1-P)^{k-1} \cdot P$
 required to obtain
- # Binomial distribution (n, P) where
 $P = P \cdot \binom{n}{x} (1-P)^{n-x}$ \rightarrow Likelihood
- # Beta(a, b) \rightarrow Prior $= P^{a-1} (1-P)^{b-1}$
- $$E = \frac{a}{a+b}$$
- $n = \text{total no. of samples}$
 $p = \text{success denotes}$
 $x \rightarrow \text{no. of success from samples}$
- $$C = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

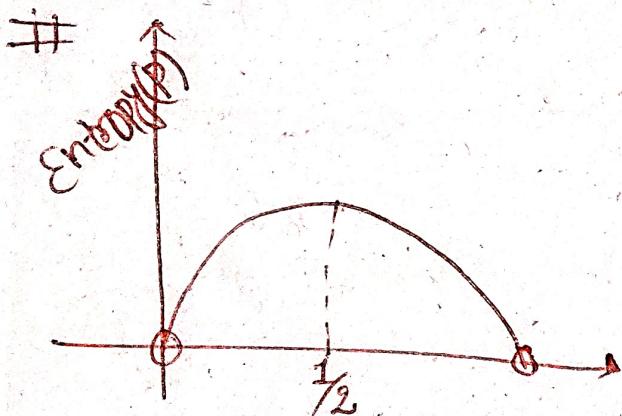
Part I \rightarrow loss function

\downarrow R \rightarrow rates (basic terms to get)



$\|E[\theta_1] - \theta\|^2 > \|E[\theta_2] - \theta\|^2$

θ is parameter
 θ_2 is a better estimator for θ when compared to θ_1 because it is having the lesser bias.



P \rightarrow fraction of 1's in y

$L(w; x, y) = \prod_{i=1}^n f(y_i | x_i)$

$L(w; x, y) = \frac{1}{n} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} (y_i - w)^2}$

$\sigma = \pm$

Gaussian Distribution:-

$$f(x|y_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x - \mu_k)^2}{2\sigma^2}}$$

Posterior \propto prior \times likelihood

$$\propto f(w) \prod_{i=1}^n f(y_i|x_i)$$

$$\propto e^{-\frac{1}{2\sigma^2} \|w\|^2} - \frac{1}{2\sigma^2} \|y - x^T w\|^2$$

◆ Hard-margin vs Soft-margin SVM (Complete Table)

Feature	Hard-margin SVM	Soft-margin SVM
Data assumption	Data perfectly linearly separable	Data not perfectly separable (noise / outliers allowed)
Margin violation	✗ Allowed nahi	✓ Allowed
Misclassification	✗ Allowed nahi	✓ Allowed
Slack variable (ξ_i)	✗ Not used	✓ Used
Constraint	$y_i(w^T x_i + b) \geq 1$	$y_i(w^T x_i + b) \geq 1 - \xi_i$
Objective function	$\min \frac{1}{2} \ w\ ^2$	$\min \frac{1}{2} \ w\ ^2 + C \sum \xi_i$
Regularization parameter C	✗ Not present	✓ Present
Support vectors lie where?	Only on supporting hyperplanes	On hyperplanes, inside margin, or even wrong side
Points inside margin	✗ Impossible	✓ Possible
Misclassified points	✗ Impossible	✓ Possible
Robust to outliers	✗ No	✓ Yes

11. Bias-Variance Tradeoff

Model Type	Bias	Variance	Error Behavior
Underfitted	High	Low	Poor train/test performance
Overfitted	Low	High	Good train, poor test
Well-generalized	Balanced	Balanced	Good train/test

14. Bagging & Boosting

Technique	Parallel?	Base Learners
Bagging	Yes	High variance (e.g., deep trees)
Boosting	No	Weak learners

C. Ensemble Methods

Technique	Key Characteristics	Use Case / Purpose
Bagging	<ul style="list-style-type: none">- Uses deep decision trees (high variance models)- Parallel execution (models are independent)- Bootstrap sampling (with replacement)- All models usually have equal weight	<ul style="list-style-type: none">- Reduces variance- Handles overfitting- Example: Random Forest
Boosting	<ul style="list-style-type: none">- Uses decision stumps / weak learners- Sequential execution (each model learns from previous errors)- Data-point weights change over iterations- Models have different importance	<ul style="list-style-type: none">- Reduces bias- Improves weak models- Examples: AdaBoost, Gradient Boosting

Contrast with Boosting (important 🔥)

Method	Parallel training?	Reason
Bagging	<input checked="" type="checkbox"/> Yes	Models independent
Random Forest	<input checked="" type="checkbox"/> Yes	Bagging + feature randomness
AdaBoost	<input checked="" type="checkbox"/> No	Each model depends on previous

Why? (Bias–Variance intuition)

Model complexity aur test error ka relation monotonic nahi hota.

- **Jab model bahut simple hota hai**
 - **High bias**
 - **Underfitting**
 - **High test error**
- **Complexity badhaane par**
 - **Bias kam hota hai**
 - **Test error pehle kam hota hai**
- **Complexity bahut zyada ho jaaye**
 - **High variance**
 - **Overfitting**
 - **Test error phir badhta hai**

🔑 Sabse pehle clear distinction

- ♦ **Hard-margin SVM**
- **No point margin ke andar hota hai**
- **No misclassification**
- **Constraints:**

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

♦ Soft-margin SVM

- Points margin ke andar bhi ho sakte hain
- Misclassification allowed
- Constraints:

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

💡 Ab tumhare statement par directly aate hain

“kuchh points to support hyperplane ke beech me bhi hota hai”

✓ Ye baat soft-margin SVM me true hai

✗ Lekin hard-margin SVM me false hai