

(Jan 31)
2026

Exam portion: Moodle - first 4 chapters of textbook.
Exam time 9-11?? Why not 10-12?

$$P(x|w_1)P(w_1) > P(x|w_2)P(w_2) \Rightarrow w_1$$

Here, the decision is made based on a priori probabilities and class conditional probabilities.

6 Feb

* Parametric methods:

- Maximum Likelihood estimation (MLE):

Bayes' Theorem: $P(\theta|D) \propto P(D|\theta)P(\theta)$
posterior likelihood joint?

$P(x|\theta)$ ← likelihood ∵

$$l(\theta|x) = P(x|\theta) = \prod_{t=1}^N P(x_t|\theta).$$

$X = \{x_t\}_{t=1}^N$ ← Draw some (N) data from a known distribution: $P(x|\theta)$.

$$x_t \sim P(x|\theta); \quad \theta = [\mu] \quad \theta_1 = [\mu_1] \text{ and } \theta_p = [\mu_p]$$

Log likelihood function:

$$\begin{aligned} L(\theta|x) &= \log l(\theta|x) \\ &= \sum_{t=1}^N \log P(x_t|\theta). \end{aligned}$$

Bernoulli Distribution: $\text{Ber}(p)$; binomial(n, p).

$$P(x) = p^x(1-p)^{1-x}; \quad x \in \{0, 1\}.$$

$$L(p|x) = \log \prod_{t=1}^N P(x_t) (1-p)^{(1-x_t)}$$

$$= \sum_{t=1}^N x_t \log p + (N - \sum_t x_t) \log (1-p)$$

Maximize Estimate $\hat{\theta} = \frac{1}{N} \sum_{i=1}^N x_i$

$$\hat{x} = \frac{\sum x_i}{N}$$

\hat{x} is called estimate

* Bias and Variance:

$$x + \theta$$

Bias of a prob distribution not random.

Estimator of $\theta \leftarrow d(x) \rightarrow \hat{\theta}(x)$.

$$\text{Error} = (d(x) - \theta)^2$$

$$\text{MSE} = E[(d(x) - \theta)^2]$$

Mean Square Error.

Bias of an estimator: $b_\theta(d) = E[d(x)] - \theta$.

What does unbiased estimator mean?

For all θ , $b_\theta = 0$, we call it unbiased estimator.

\rightarrow number of samples

$$E[m] = E\left[\frac{1}{N} \sum x_i\right] = \frac{1}{N} \cdot \frac{1}{N} \sum x_i = \frac{N\mu}{N} = \mu$$

\rightarrow concept used similar to.

Law of large numbers: When sample size N , sample mean \approx actual mean

$$\text{Variance: } \text{Var}(m) = \text{Var}\left[\frac{1}{N} \sum x_i\right] = \frac{\sigma^2}{N}$$

$$\text{Var} = E[x^2] - (E[x])^2$$

Mean Square Error:

$$V(d, \theta) = E[(d(x) - \theta)^2]$$

$$V(d, \theta) = E[(d(x) - \theta)^2].$$

$$= E[(d - E[d])^2] + (E(d) - \theta)^2.$$

* Variance + bias.

$$V(d, \theta) = \text{Var}(d) + b_e(d)^2.$$

$$\text{Error} = \text{Variance} + \text{bias}^2.$$

13 Feb.

* Dimensionality Reduction:

1. Linear Discriminant Analysis LDA:

Assume we have 2 classes:

Goal: To reduce the dimension while preserving as much as the class discriminatory information as possible.

D-dimensional data: $\{x_1, x_2, \dots, x_n\}$.

$N_1 \rightarrow w_1 (\text{Class}_1)$.

$N_2 \rightarrow w_2 (\text{Class}_2)$.

Get a scalar $y = w^T x$.

Projecting the data points on a line.

Of all possible lines, we select a line that maximizes the separability of classes; infinite lines are possible, so find the one line among them.

To get a projection vector, define measure of separation.

CRAZY



Mom tell me if I could
Send up my heart to you?
So when I die, which I must do,
make it shine down here with you

Mean vector: $\mu_i = \frac{1}{N_i} \sum_{y \in w_i} z$.

$$\tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in w_i} y = \frac{1}{N_i} \sum_{z \in w_i} w^T z = w^T \mu_i.$$

choose the distance b/w projected means as our objective fn.

$$J(w) = |\tilde{\mu}_1 - \tilde{\mu}_2| = w^T (\mu_1 - \mu_2).$$

Variance, $w^T w$, has to be considered here.

Maximize the difference b/w means normalized by ~~a~~ a measure of the within class scatter.

$$\text{variance}_{\tilde{\mu}_i} = \sum_{y \in w_i} (y - \tilde{\mu}_i)^2.$$

Within class scatter $(\tilde{S}_1^2 + \tilde{S}_2^2)$.

of the projected mean samples.

Fisher linear discriminant is defined as a linear fn. $w^T z$ that maximizes the criterion fn.

$$J(w) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{S}_1^2 + \tilde{S}_2^2}.$$

Mean of w_1 & w_2 are close to each other.

We look for a projection where the examples from the same class projected very close to each other and at the same time, the projected mean are farther as possible.

find optimum w?

Sum $S_{B1} + S_{B2}$

S_B

Within class scatter

Between class scatter

Measures of the scatter in feature space x.

$$S_w = \sum_{i=1}^n (x - \mu_i)(x - \mu_i)^T$$

) same concept

Scatter of projection y:

$$S_y = \sum_{i=1}^n (y - \tilde{f}_i)^2$$

but different space.

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

$$= \sum_{i=1}^n w^T (x - \mu_i)(x - \mu_i)^T w$$

$$= w^T S_w w$$

$$\boxed{S_w + S_B = w^T S_w w}$$

Mean: $(\tilde{f}_1 - \tilde{f}_2)^2 = (w^T \mu_1 - w^T \mu_2)^2$

$$= w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w$$

$$= w^T S_B w$$

$$f(w) = \frac{w^T S_w w}{w^T S_B w}$$

$$w^T S_B w$$

$$\frac{df}{dw} = (f(w)) = 0 \Rightarrow \frac{d}{dw} \left[\frac{w^T S_w w}{w^T S_B w} \right] = 0$$

$$\Rightarrow [w^T S_w w] \frac{d}{dw} [w^T S_B w] - [w^T S_B w] \frac{d}{dw} [w^T S_w w] = 0$$

$$[w^T S_w w]^2 S_B w - [w^T S_B w]^2 S_w w = 0$$

25 Feb

Self taught

Self taught / /

Aditya

Dividing by S_{WW} , we get:

$$\begin{bmatrix} w^T S_B w \\ w^T S_W w \end{bmatrix} = \begin{bmatrix} w^T S_B w \\ w^T S_W w \end{bmatrix} = 0.$$

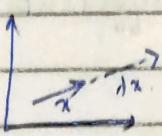
$$S_B w - \frac{\begin{bmatrix} w^T S_B w \\ w^T S_W w \end{bmatrix}}{w^T S_W w} S_W w = 0.$$

$$S_B w - J w = 0.$$

$$S^{-1} S_B w - J w = 0.$$

$$\rightarrow \boxed{S^{-1} S_B w = J w.}$$

Eigen value problem:
 $\lambda_2 = \lambda_1$.



Operator A acts on the vector z to scale z along l .
No change in direction.

Optimized w : w^* .

$$w^* = \text{argmax}_{w^T S_W w} \frac{w^T S_B w}{w^T S_W w} = S^{-1} (\mu_1 - \mu_2).$$

25 Feb.

* Principle Component Analysis:

Linear combination of features that has max. variance.

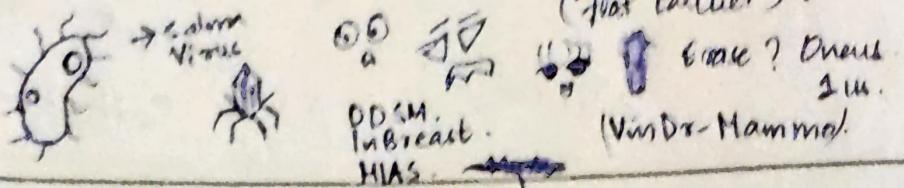
Scaling data - normalization.
Mean $\rightarrow 0$.

Manifold? - Generalization of Euclidean space.

- Non linear space.

- Eg: Swiss roll?





27 Feb

- * Non linear Dimensionality Reduction :
- Locally Linear Embedding (LLE)

Making sure the relationship (H) among the data points are preserved as far as possible.

LLE is a mapping which preserves the neighbourhood.

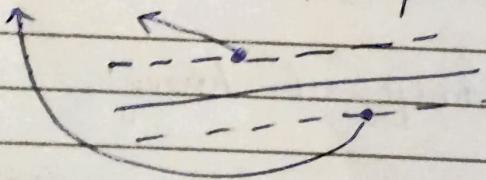
$$\text{Reconstruction: } \hat{\mathbf{x}}(\mathbf{y}) = \sum_i |\mathbf{y} - \mathbf{z}_i \mathbf{w}_i|^2 (\mathbf{z}_i).$$

Note: In higher dimensions, the notion of distance / similarity is different.

4 March

- * SVM - Support Vector Machine:
- Supervised.
- Mainly for classification, can also be used for regression.
- For n dimensional data, $(n-1)$ dimensional hyperplane.

Support vectors: vectors / data points on the margin.



Hard margin: Margin defined clearly.

Soft Margin Classifier:

Instead of hard margin, we allow misclassification.



7 floor

— / —

- Non-linear:

for n dimensions, might need to map to $(n+1)$ dimensions to get a linearly separable space.

Kernel tricks:

- RBF:

$$K(x, x') = e^{-\gamma \|x - x'\|^2}.$$

- can handle infinite dimensions.

γ - gamma - decides how big/small the margin is.

Optimization? finding the optimal hyperplane.

$$z = \phi(x)$$

where $z_j = \phi_j(x)$, $j=1, \dots, k$.

discriminant: $g(z) = w^T z$.

$$\begin{aligned} g(x) &= w^T \phi(x) \\ &= \sum_{j=1}^k w_j \phi_j(x). \end{aligned}$$

Replace the inner product of basis functions,
 $\phi(x^t)^T \phi(x^s)$, by kernel function $k(x^t, x^s)$

$$g(x) = w^T \phi(x) = \sum a^t j^t \phi(x^t)^T \phi(x) = \sum_t a^t j^t k(x^t, x).$$

* Exploration & Exploitation:

6 March.

* Exploration & Exploitation:

Exploring everything and exploiting the optimal one.

* CART: Classification & Regression Trees:

	GPA	Interaction	Knowledge	Softskills	Job Offer
1.	≥ 9	Yes	VGood	Good	Yes
2.	≥ 8	No	Good	Moderate	Yes
3.	≥ 9	No	Avg	Poor	No
4.	< 8	No	Avg.	Good	No
5.	≥ 8	Yes	Good	Moderate	Yes
6.	≥ 9	Yes	Good	Moderate	Yes
7.	< 8	Yes	Good	Poor	No
8.	≥ 9	No	VGood	Good	Yes
9.	≥ 8	Yes	Good	Good	Yes
10.	≥ 8	Yes	Avg.	Good	Yes

$$\text{Gini Index}(T) = 1 - \sum_{i=1}^m p_i^2$$

where m is the number of labels (Yes, No, ≥ 2).

$$\text{Gini Index}(T) = 1 - \left[\left(\frac{7}{10}\right)^2 + \left(\frac{3}{10}\right)^2 \right] = 1 - \frac{49+9}{100}$$

$$= 1 - \frac{58}{100} = \frac{100-58}{100} = \frac{42}{100} = 0.42$$

0.42 for whole dataset.

Gini Index for each attribute now:

Gini Index (GPA) \Rightarrow	≥ 9	Job Yes	Job No.
	≥ 9	3	1
	≥ 8	1	0
	< 8	0	2

Subsets : $\{ \{ \geq 9, \geq 8 \}, \{ < 8 \} \}$

A.

$\{ \{ \geq 8, < 8 \}, \{ \geq 9 \} \}$

B.

$\{ \{ \geq 9, < 8 \}, \{ \geq 8 \} \}$

C.

Gini Index (T , CGPA w.r.t $\{ \geq 9, \geq 8 \}$)

$$= 1 - \left[\left(\frac{7}{10} \right)^2 + \left(\frac{3}{10} \right)^2 \right] = 1 - \frac{64}{100} = 0.21875$$

Gini Index (T , CGPA w.r.t $\{ < 8 \}$) = $1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] = 0$.

Combining :

A. Gini Index (T , CGPA w.r.t $\{ \{ \geq 9, \geq 8 \}, \{ < 8 \} \}$)

$$= \frac{|S_1|}{|T|} G_I(S_1) + \frac{|S_2|}{|T|} G_I(S_2)$$

$$= \frac{8}{10} \cdot 0.21875 + \frac{2}{10} \cdot 0 = 0.175$$

Gini Index (T , CGPA w.r.t $\{ \geq 8, < 8 \}$) = $1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right] = 1 - \frac{20}{36} = 0.44$.

Gini Index (T , CGPA w.r.t $\{ \geq 9 \}$) = $1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 1 - \frac{10}{16} = 0.375$.

B. Gini Index (T , CGPA w.r.t $\{ \{ \geq 8, < 8 \}, \{ \geq 9 \} \}$)

$$= \frac{|S_1|}{|T|} G_I(S_1) + \frac{|S_2|}{|T|} G_I(S_2)$$

$$= \frac{6}{10} \cdot 0.44 + \frac{4}{10} \cdot 0.375 = 0.414$$

Gini Index (T , CGPA w.r.t $\{ \geq 9, < 8 \}$) = $1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 1 - \frac{12}{36} = 0.5$.

Gini Index (T , CGPA w.r.t $\{ \geq 8 \}$) = $1 - \left[\left(\frac{9}{4} \right)^2 \right] = 0$.

C. Gini Index (T , CGPA w.r.t $\{ \{ \geq 9, \geq 8 \}, \{ \geq 8 \} \}$)

$$= \frac{|S_1|}{|T|} G_I(S_1) + \frac{|S_2|}{|T|} G_I(S_2)$$

$$= \frac{6}{10} \cdot 0.5 + \frac{4}{10} \cdot 0 = 0.3$$

$$G_I(CGPA) \geq 130.175 \quad \checkmark$$

$$\therefore 0.42$$

$$= 0.3$$

Compute $\Delta Gini$:

$$\begin{aligned}\Delta Gini(CGPA) &= Gini(T) - Gini(T, CGPA) \\ &= 0.42 - 0.175 \\ &= \underline{\underline{0.245}}\end{aligned}$$

Interactive:

	Job Yes	Job No.
Yes	5	1
No	2	2

$$G_I(T, \text{Interactive } \in \{1, 0\}) = 1 - \left[\left(\frac{5}{6}\right)^2 + \left(\frac{1}{6}\right)^2 \right] = 1 - \frac{26}{36} = \underline{\underline{0.28}}$$

$$G_I(T, \text{Interactive } \in \{1, 0\}) = 1 - \left[\left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right] = 1 - \frac{8}{16} = \underline{\underline{0.5}}$$

$$G_I \text{ combined } (T, \text{Interactive } \in \{1, 0\}, \text{Job } \in \{1, 2\}) = \frac{|S_1|}{|T|} G(S_1) + \frac{|S_2|}{|T|} G(S_2)$$

$$= \frac{6}{10} \cdot 0.28 + \frac{4}{10} \cdot 0.5$$

$$= 0.168 + 0.2 = \underline{\underline{0.368}}$$

$$\Delta Gini(\text{Interactive}) = Gini(T) - Gini(T, \text{Interactive})$$

$$= 0.42 - 0.368$$

$$= \underline{\underline{0.052}}$$

May Attribute: Gini Index ΔG_i

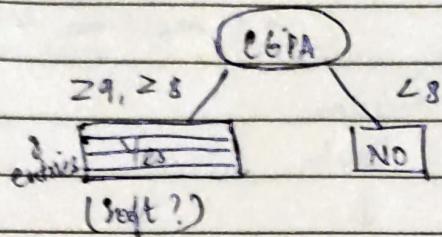
$$\rightarrow CGPA \quad 0.175 \quad 0.245$$

$$\rightarrow \text{Interactive} \quad 0.368 \quad 0.052$$

$$\rightarrow \text{Knowledge} \quad 0.3054 \quad 0.1196 (?)$$

$$\rightarrow \text{Soft Skills} \quad 0.175 \quad 0.245 (?)$$

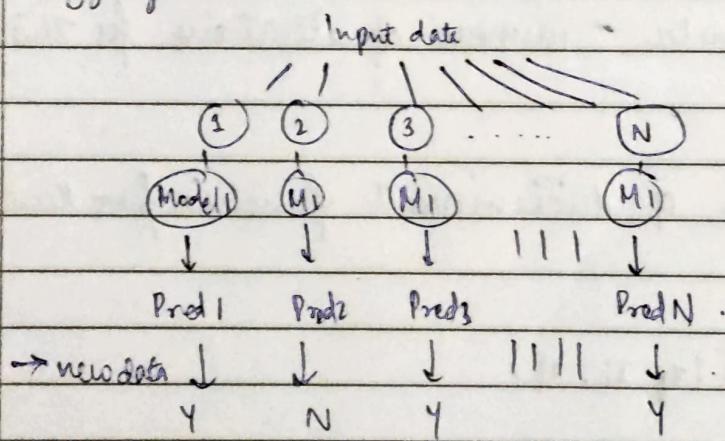
So CGPA vs Soft Skills:



* Ensemble Learning:

- Supervised Machine learning.
- Combines predictions from multiple models.
- Voting (for classification) / Averaging (for regression).
- Bagging (Bootstrap AGGregation).
- Random forests.
- Boosting.
- Stacked Generalization (Blending).

© Bagging:



Voting: #Y , #N \Rightarrow Majority wins \rightarrow classified.

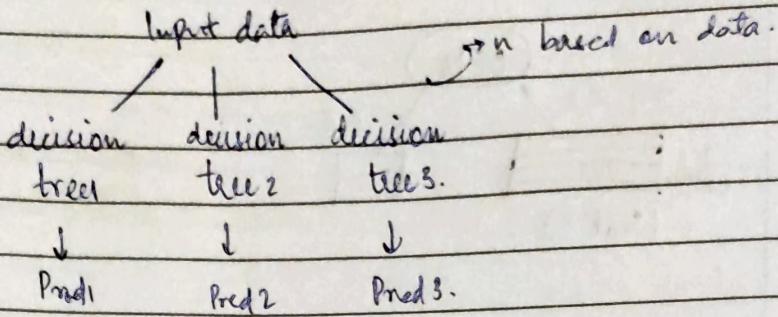
???

Thank you!!!

My mook! too cold...!
It's dying!!

— / /

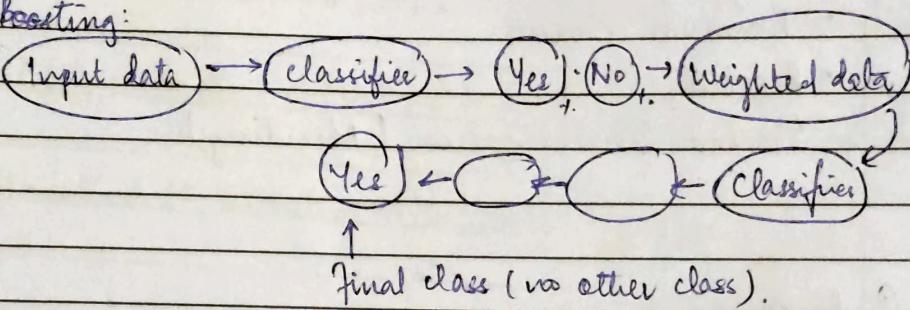
② Random Forests:



Voting /Averaging based on classification / regression -

random forest classifier for final prediction.

③ Boosting:



Hyperparameter = number of iterations for this.

* Derivation of discriminant function from Kernel fn:

$$K(x_i, y)$$

$$k_{ij} = k(x_i, x_j)$$

$$w = \arg \min \sum_{j=1}^J \alpha_j$$

[from pdf]. $\alpha_j \Rightarrow$ Lagrange multipliers.

$$f(x) \in \mathbb{R}, \text{ s.t. } k(m; x) + b,$$

where b = bias term.

\rightarrow hyperparameter.