

df.sortvalues(by=['col1'])
'col2'])
02/09/2020

AI and ML

Linear Regression

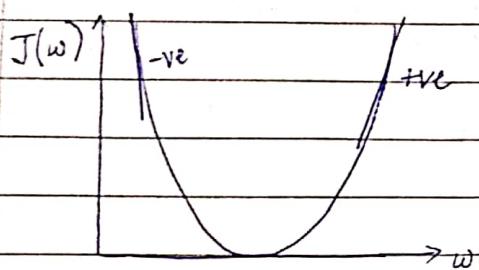
Representation of Hypothesis -

$$h(x) = w_0 + w_1 x_1$$

$$= w_0 x_0 + w_1 x_1$$

$$h(x) = w_0 + w_1 x_1 + w_2 x_2 \dots$$

$$= \sum_{j=0}^n w_j x_j \quad \text{where } x_0 = 1$$



gini
 $1 - (\text{probability of } y_{\text{yes}})^2 - (\text{probability of } y_{\text{no}})^2$

weighted avg =

$$\frac{144(0.375)}{144+159} + \frac{159(0.625)}{144+159}$$

entropy $\rightarrow - \sum p(t) \log(p(t))$

$$- \frac{2}{6} \log \frac{2}{6} - \frac{4}{6} \log \frac{4}{6}$$

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} (J(w))$$

$$\begin{cases} +ve \\ -ve \end{cases}$$

$$IG = \text{entropy}(P) -$$

Benefit of Square - Single Global Minima.
Loss (Convex Function)

(weighted sum of entropy of children nodes)

ML Life Cycle

Data \rightarrow Preprocessing Data



Split Data (Train, test, validate)



Decide problem type (Regression, Classification)



Decide loss type



Decide Hypothesis (Model Structure)



Learning Algo (Gradient Descent/Ascent)

Training $w_j = w_j - \alpha \sum_{i=0}^n (h(x_i) - y_i) \cdot x_i$

$h_w(x) \leftarrow \text{Testing set}$ satisfied?
No

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h(x)^{(i)} - y^{(i)})^2 \leftarrow \text{squared loss}$$

$$h(x) = w_0 x_0 + w_1 x_1 + \dots + w_m x_m.$$

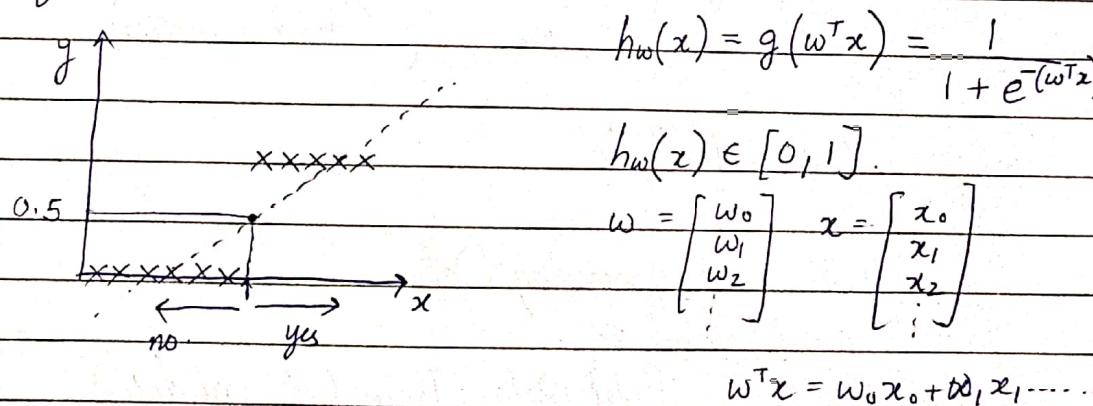
$$w_j = w_j - \alpha \frac{\partial J(w)}{\partial w_j}$$

$$J(w) = \frac{1}{2} (h_w(x)^{(i)} - y^{(i)})^2$$

$$\begin{aligned} \frac{\partial J(w)}{\partial w_j} &= 2 \times \frac{1}{2} (h_w(x)^{(i)} - y^{(i)}) \cdot \frac{\partial}{\partial w_j} (w_0 x_0 + w_1 x_1 + \dots) \\ &= (h_w(x) - y) \cdot x_j \end{aligned}$$

$$w_j = w_j - \alpha (h_w(x) - y) x_j$$

Classification -



likelihood.

$$P(y=1|x; w) = h_w(x)$$

$$P(y=0|x; w) = 1 - h_w(x)$$

$$y = \{0, 1\} \quad h_w(x) \in [0, 1]$$

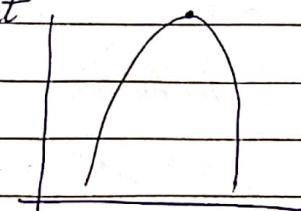
only one term shall be used at one time other terms

$$S(w) = \prod_{i=0}^m P(y^{(i)} | x^{(i)}, w) = \prod_{i=0}^m ((h_w(x^{(i)})^{y^{(i)}} + (1 - h_w(x^{(i)}))^{(1-y^{(i)})})$$

$$\log S(w) = \log \left[\prod_{i=0}^m (h_w(x^{(i)})^{y^{(i)}} + (1 - h_w(x^{(i)}))^{(1-y^{(i)})}) \right]$$

$$= \sum_{i=0}^n y^{(i)} \log h_w(x^{(i)}) + (1-y^{(i)}) \log (1-h_w(x^{(i)}))$$

Gradient Ascent



$$w_j = w_j + \alpha \frac{\partial J(w)}{\partial w}$$

$$+ \alpha \sum_{i=0}^m (y^{(i)} - h_w(x^{(i)})) x_j^{(i)}$$

Softmax - Multiple classes

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix} = w^T x = \begin{bmatrix} 10 \\ 8 \\ 3 \\ -1 \end{bmatrix} \rightarrow \begin{bmatrix} e^{10} \\ e^8 \\ e^3 \\ e^{-1} \end{bmatrix} \rightarrow \begin{bmatrix} e^{10} \\ \text{mean} \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} ? \\ ? \\ ? \\ ? \end{bmatrix}$$

Binary Classification $\rightarrow \begin{bmatrix} C_1 \\ C_2 \end{bmatrix}$ 0s and 1s

values b/w

0 and 1

add up to 1.

$y = f(x_1, x_2)$

Threshold > 0.5

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad J(y) = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} \end{bmatrix}$$

$$u, v = (x_1, x_2)$$

$$J(u, v) = \begin{bmatrix} \frac{\partial u}{\partial x_1} & \frac{\partial u}{\partial x_2} \\ \frac{\partial v}{\partial x_1} & \frac{\partial v}{\partial x_2} \end{bmatrix}$$

Hessian Matrix.

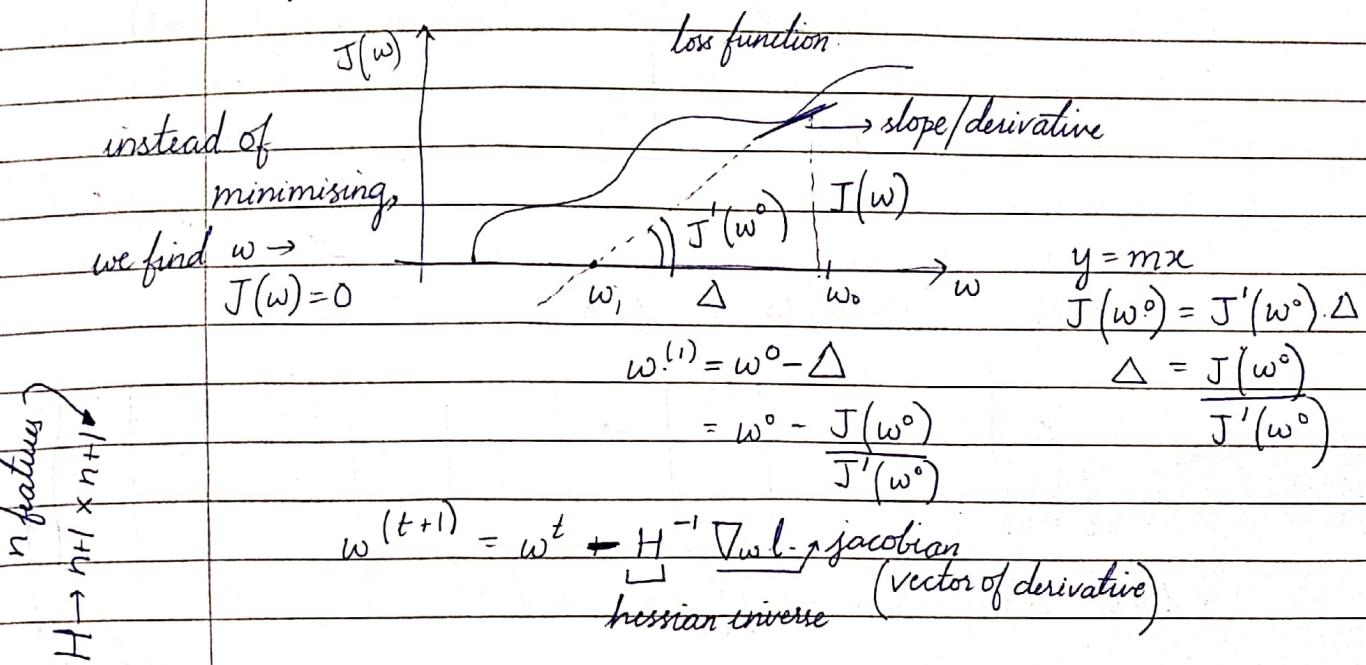
$$\frac{\partial^2 f}{\partial w^2} \cdot \frac{\partial y}{\partial (x_1, x_2)} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial^2 y}{\partial x_1 \partial x_1} & \frac{\partial^2 y}{\partial x_1 \partial x_2} \\ \frac{\partial^2 y}{\partial x_2 \partial x_1} & \frac{\partial^2 y}{\partial x_2 \partial x_2} \end{bmatrix}$$

\rightarrow Hessian Matrix.

Newton's Method.

Logistic Regression \rightarrow Want to maximize $J(w)$.



Newton's method - large size (can't compute efficiently)
as original input $n \times n$.

$$R^{n+1 \times n+1}$$

Advantage - Quadratic Convergence

Decision Trees.

Difference B/w Batch Gradient

Descent and Stochastic Gradient

Descent.

↓
in each step

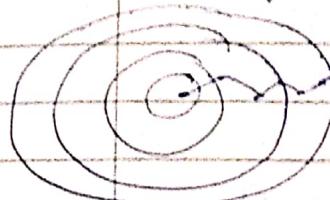
$$l(n) = 0.01$$

$$n+1 = 0.0001$$

$$n+2 = 0.00000001$$

learning rate is hyperparameter

Cost function
(for single training example)
Loss function
over the entire training set (or many batches)
for mini-batch gradient descent)



$$\hat{O}_j = O_j - \alpha \cdot \frac{\partial}{\partial O_j} J(O)$$

$$\min_{m \in M} \sum_{i=1}^m (\hat{y}_i - y_i) x_i'$$

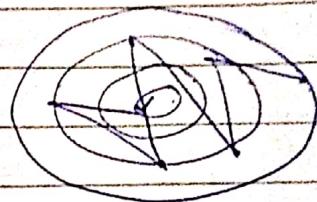
→ Vanilla (Batch) G.D. — (slower but smooth.)

for i in range(m):

$$O_j = O_j - \alpha \cdot \underbrace{(y^{(i)} - \hat{y}^{(i)})}_{\text{one training example}} \cdot x_j' \rightarrow \text{faster but uneven moment}$$

→ noisier random
more random

iteration of one training example at a time



Mini-

→ Batch gradient descent:

Mini-Batch size 1000.

for i in range($0, m, 1000$)

$$O_j = O_j - \alpha \sum_{i=0+P}^{P+1000} (h_w(x(i)) - y(i)) x_j'$$

instead of one sample in SD → we take n data points in each iteration

Decision Trees :- To predict the class value or the target output by learning decision rules inferred from the prior data (training data).

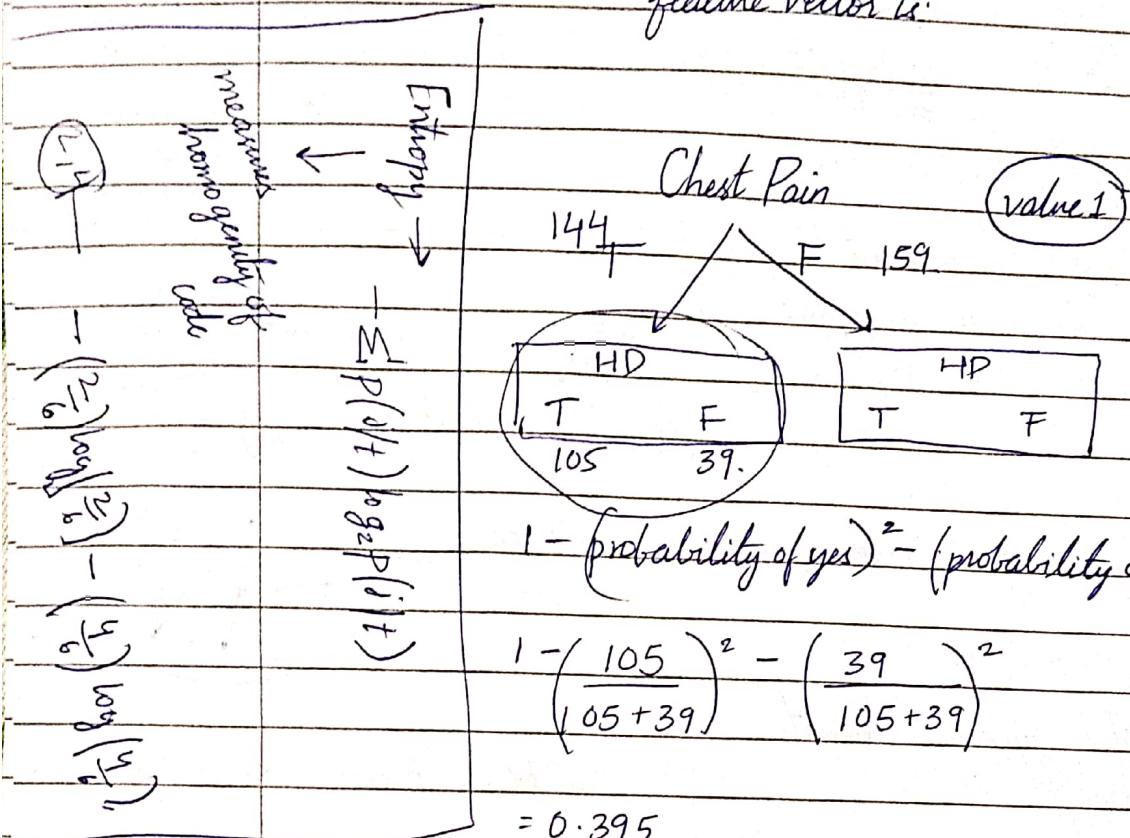
✓ prioritize which to consider first

For predicting a class label - we start from the root of the tree

Tree can have single root

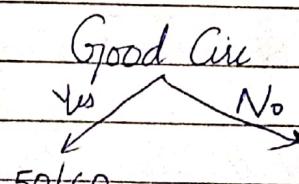
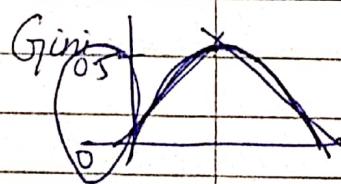
Impurity γ becomes our criteria if we need to quantify the impurity
 then we say → it's a good node

Information Gain - How important a given attribute of the feature vector is:



$$\text{weighted avg} \rightarrow \frac{144}{144+159} (0.395) + \frac{159}{144+159} (0.336)$$

5, 5

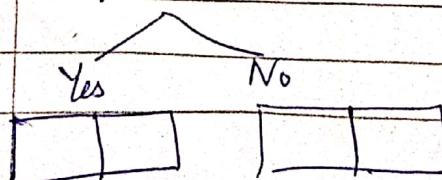


Weight	Heart Disease
220	y
180	y
225	y
190	n
155	n

least value

Weight-sort()
 take avg of this and next val

weight ≤ 167.5



Entropy - Completely pure set \Rightarrow

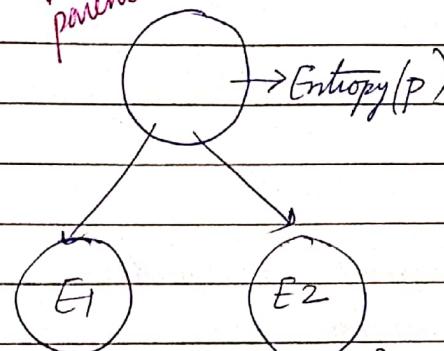
" in " \Rightarrow 1 b/w 0 and 1

1 /

Information Gain :-

Entropy(P) - (weighted sum of entropy of children nodes)

↓
parent



Algorithms
↓

CART \rightarrow Entropy, IG

ID3 \rightarrow Gini Score

(used for splitting)
(in decision trees)

Entropy(P) -
(E_1, E_2) \rightarrow weighted sum = IG.

Regression Metrics

Explained Variance (Similarity Measure)

$$1 - \frac{\text{var}(y - \hat{y})}{\text{var}(y)}$$

$$\text{MAE} = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} |y_i - \hat{y}_i| \quad L_1 \text{ norm}$$

$$\text{MSLE}(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (\log_e(1 + y_i) - \log_e(1 + \hat{y}_i))^2$$

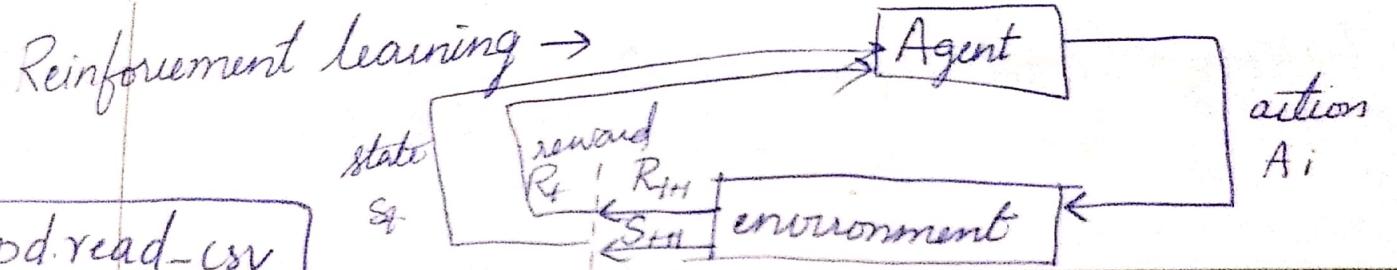
makes down data to a smaller scale

Median Absolute Error (y, \hat{y}) = median($(y_1, \hat{y}_1), \dots, (y_n, \hat{y}_n)$)

R^2 score:

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

$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$
mean.



```
data = pd.read_csv
```

Labs Revision

Agent-Environment Interface

- AI tools → <sup>to solve
problems</sup> Unsupervised Learning — y labels not given - tries to make sense out of data by extracting features
- AI Agents → <sup>like
detecting</sup> Supervised Learning — y labels given and patterns on its own

↓	have intelligence	Search problems
→		Markov decision processes
		Adversarial games

Reflex.	States	Variables	Logic "High Level Intelligence"
"low level intelligence"			↓ Constraint satisfaction problems Bayesian networks.

(descriptive analysis) (predictive analysis) (prescriptive analysis)

Unsupervised Supervised Reinforcement

(Hindsight) (Insight) (Foresight)

(clustering) ← ↓ (regression) →
value, difficulty ramification

Astronomical data analysis	Tumor size: Malignant / Benign	Game Playing Robot in a maze
Market Segment ation		Balance Poll in hand. Credit assignment problem

Organic Computing
Clusters

→ Regression Metrics:-

(1) Error ← dissimilarity ↓

(2) Similarity ← accuracy ↑

→ Mean Absolute Error



$$MAE(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} |y_i - \hat{y}_i|.$$

does a poor job when scale of the data is high.

→ Mean Squared Error

$$MSE(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} (y_i - \hat{y}_i)^2$$

Squaring the value, sum won't be zero

" emphasizes larger differences

→ Median Absolute Error

$$MedAE(y, \hat{y}) = \text{median}(|y_1 - \hat{y}_1|, \dots, |y_n - \hat{y}_n|)$$

→ does not support multiclass output

Median AE prevents outliers contributing more error to model evaluation

Type of prediction
log scales it down

→ Mean Squared Logarithmic Error.

$$MSLE(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} (\log(1+y_i) - \log(1+\hat{y}_i))^2$$

when targets have exponential growth → population counts, average sales of commodity over years

Method penalizes an underpredicted estimate greater than over predicted estimate.

→ Max Err.

$$\text{MaxErr}(y, \hat{y}) = \max(|y_i - \hat{y}_i|)$$

aims to find max error made by model for single sample → finding best model that covers all samples

Similarity Measures

→ Explained Variance Score (y, \hat{y}) Best possible score → 1.0
 $= 1 - \frac{\text{Var}\{y - \hat{y}\}}{\text{Var}\{y\}}$ lower values are worse

Estimates deviation between prediction and actual values

→ R^2 score, the coeff of determination

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

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

↓ calculate unadjusted $\sum_{i=1}^n (y_i - \bar{y})^2$
 R^2 without correcting for bias in sample variance of y

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

Error/Dissimilarity Measure (less the better)

Classification Metrics

	actual	predicted	
True pos	1	1	11 - 1
False neg	0	1	01 - 0
False pos	1	0	10 - 0
True neg	0	0	00 - 1

Accuracy/Similarity Measure (more the better)

$$A\bar{B} + \bar{A}B - \text{EXOR}$$

$$A\bar{B} + \bar{A}\bar{B} - \text{EXNOR}$$

→ Accuracy = $\frac{TP + TN}{TP + TN + FP + FN}$ } all true values
total:

→ Precision = $\frac{TP}{TP + FP}$ } fraction of predicted that are actually true

$$\text{Recall} \quad (\text{Sensitivity}) = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

↗ want
 ↗ so then recall max
 ↗ all want

Total days - 30 (one month)

actual rain - 10

Model - 9 days (5✓, 4✗)

$$\text{Precision} = 5/9; \text{Recall} = 5/10$$

$$\text{Accuracy} := \frac{(\text{TP} + \text{TN})}{\text{total}} \rightarrow \frac{\text{TP} + \text{FP} + \text{TN} + \text{FN}}{30}$$

	Predicted NO	Predicted YES
Actual NO	TN	FP
Actual YES	FN	TP

$$\text{True pos Rate}:- \frac{\text{TP}}{\text{actual yes}} \rightarrow \frac{\text{FN} + \text{TP}}{10}$$

	Actual NO	Actual YES
Actual NO	TN	FP
Actual YES	FN	TP

$$\text{False neg Rate}:- \frac{\text{FP}}{\text{actual no}} \rightarrow \frac{\text{TN} + \text{FP}}{20}$$

	Actual NO	Actual YES
Actual NO	TN	FP
Actual YES	FN	TP

$$\text{True -ve Rate}:- \frac{\text{TN}}{\text{actual no}} \rightarrow \frac{\text{TN} + \text{FP}}{20}$$

$$\text{Precision}:- \frac{\text{TP}}{\text{predicted yes}} \rightarrow \frac{\text{FP} + \text{TP}}{9}$$

$$\text{Misclassified Rate}:- \frac{\text{FP} + \text{FN}}{\text{Total}} \rightarrow \frac{\text{TP} + \text{FP} + \text{TN} + \text{FN}}{30}$$

$$\rightarrow F_1 = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = \frac{2 * (\text{precision} * \text{recall})}{\text{precision} + \text{recall}}$$

↗ 1. 50% focus
 ↗ 2. 50% focus

$$\rightarrow F_B \text{ score} = \frac{(1 + B^2) * (\text{precision} * \text{recall})}{(B^2 * \text{precision} + \text{recall})} \quad B=0 P \uparrow$$

$$B=\infty R \uparrow$$

• $F_{0.5}$ ($\beta = 0.5$) :- Precision \uparrow Recall \downarrow

• F_1 ($\beta = 1.0$) :- Precision 50% Recall 50%

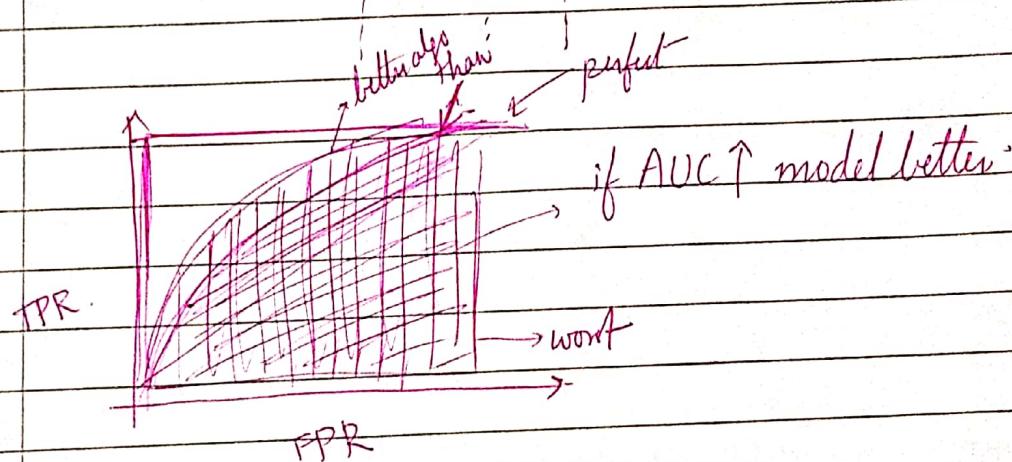
• F_2 ($\beta = 2.0$) :- Precision \downarrow Recall \uparrow

ROC Curve and ROC AUC Score.
helps in understanding balance b/w TPR & FPR

- thresholds \rightarrow all unique prediction probabilities in 'descending order'
- $FPR = \frac{FP}{(FP + TN)}$, for each threshold
- $TPR = \frac{TP}{(TP + FN)}$.

Higher the AUC, better the model is at predicting 0s as 0s and 1s as 1s.

eg:	$thr_1 = 0.31 \downarrow$	fpr	tpr
	$thr_2 = 0.41 \downarrow$	fpr	tpr
	$thr_3 = 0.5 \downarrow$	fpr	tpr
		$thr_1 \ x_1 \ y_1$	$thr_1 \ x_1 \ y_1$
		$thr_2 \ x_2 \ y_2$	$thr_2 \ x_2 \ y_2$
		$thr_3 \ x_3 \ y_3$	$thr_3 \ x_3 \ y_3$



Overfitting and underfitting

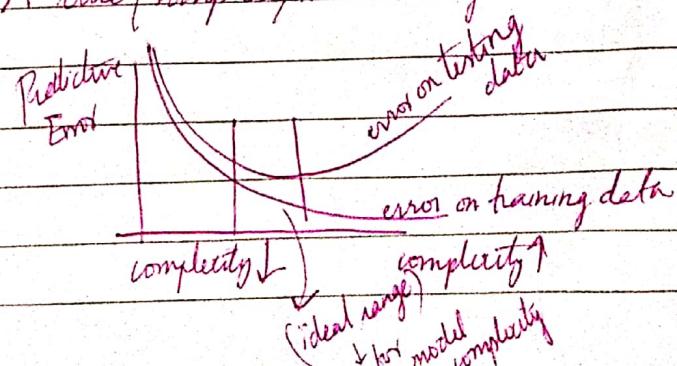
$$x_1 \propto y \quad x_2 \propto \frac{1}{y} \quad x_3 \text{ } \overset{\text{no}}{\text{f}} \text{ impact on } y \\ w_1 \quad w_2 \downarrow \quad w_3 = 0$$

always tries to get the impact on 'y'

Training Error	Testing Error
Training Loss	Testing Loss
Underfit	Overfit
Bias	Variance

\rightarrow A complex model $y = \text{high order polynomial in } X$

\rightarrow A true (simple) model $y = aX + b + \text{noise}$



	training acc	test acc
60	58 - (HB LV)	
60	40 - (HB HK)	
98	70 - (LB HV)	
98	96 - (LB LV)	

✓ Data Distribution to Train and Test should be uniform

like random shuffling

DR

underfit | good split

overfit
train
Murray
test

Multi Split

Ways to Control Logistic Regression

→ Adjust Step Size

→ Adjust Iterations / Stopping

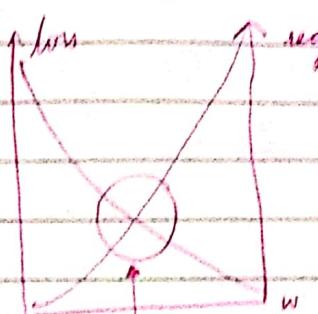
→ Regularization

curves for Gradient Descent
hard constraints
soft constraints
internal penalties
for iterations

over constraint test and stat
(early stopping)

$$L_{\text{reg}}(S) = \sum_{i=1}^n \text{Loss}(\hat{y}_i, y_i) + \lambda \sum_{j=1}^d |w_j| \rightarrow \text{L}_1 \text{ norm}$$

↑ weight
↓ bias



$$J(w) = \frac{1}{2m} \sum_{i=0}^m (\text{hw} x^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^d |w_j| \leftarrow \text{L}_1 \text{ norm}$$

$$J(w) = \frac{1}{2m} \sum_{i=0}^m (\text{hw} x^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=0}^d |w_j|^2 \leftarrow \text{L}_2 \text{ norm}$$

$\lambda \uparrow$ (very high) $w \downarrow$, underfit

$\lambda \downarrow$ overfit

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (\text{hw} x^{(i)} - y^{(i)}) x_j - \lambda w_j \right]$$

$$w_j = (1 - \alpha \lambda) w_j - \alpha \sum_{i=1}^m (\text{hw} x^{(i)} - y^{(i)}) x_j$$

regularization factor

Ways to Control Decision Tree Learning

- Increase min To Split
- Increase min Gain To Split
- Limit total number of Nodes.
- Penalize complexity

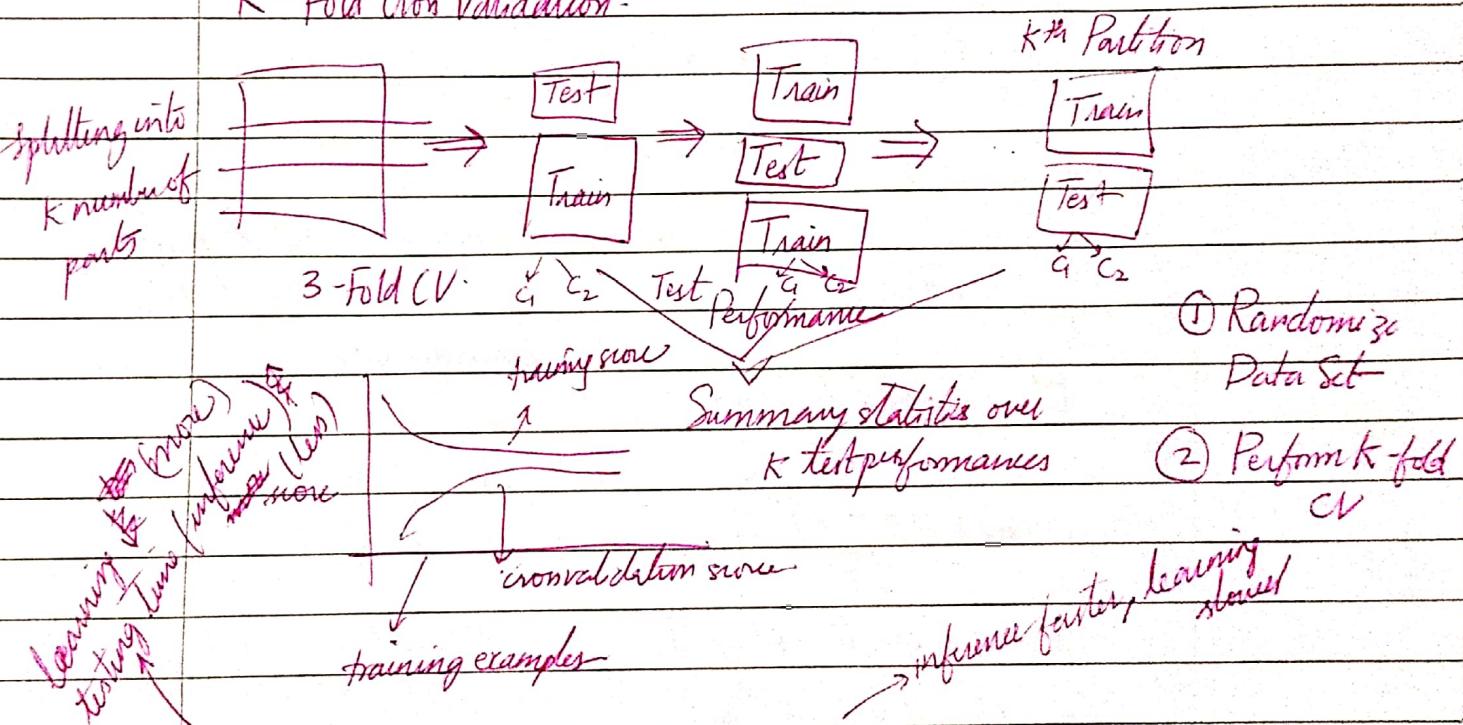
$$\text{Loss}(S) = \sum_{\substack{i=1 \\ \text{node } \uparrow}}^n \text{Loss}(y_i, \hat{y}_i) + \gamma \log_2(\# \text{Nodes})$$

node ↓

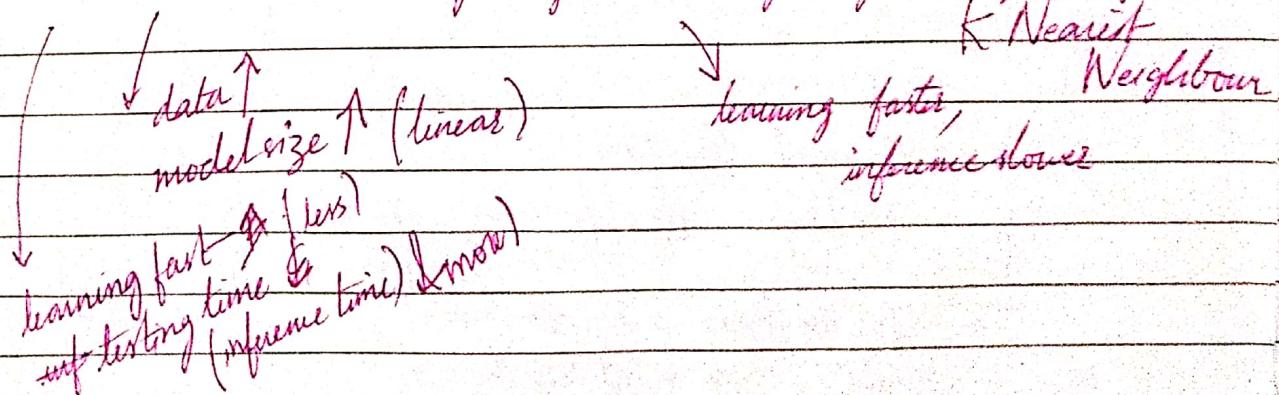
node ↓

Comparing Classifiers → test accuracy

K-Fold Cross Validation



- ⇒ Parametric Learning Algo - Doesn't need I/P training data for inference
- ⇒ Non-Parametric Learning Algo - Locally Weighted Regression



Parametric Learning Algo

Non-Parametric Learning Algo (Training data required for inference.)

Reap Linear Regression.

fit w to minimize

$$\rightarrow J(w) = \sum_{i=0}^m (h_w(x^{(i)}) - y^{(i)})^2$$

Locally Weighted Regression.

fit w to minimize.

$$\sum_{i=0}^n \beta^{(i)} (y^{(i)} - h_w(x^{(i)}))^2. \quad x \rightarrow \text{point where we want to make prediction.}$$

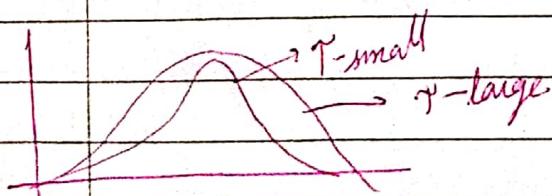
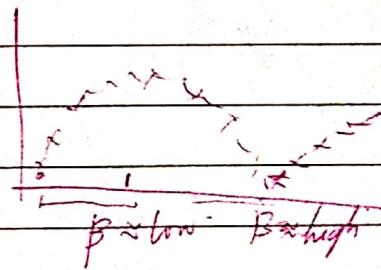
$x^i \rightarrow \text{training data.}$

$\beta' \rightarrow \text{weighted function.}$

$$\beta^{(i)} = \exp\left(-\frac{(x^i - x)^2}{2\gamma^2}\right)$$

$|x^i - x|$ is small, $\beta^{(i)} \approx 1$

$|x^i - x|$ is large $\beta^{(i)} \approx 0$.



γ - small (overfit) \rightarrow narrow focus

γ - large (underfit) \rightarrow broad focus

example

Training Phase \rightarrow save the model.

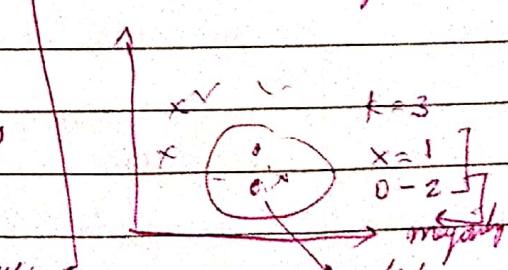
Prediction phase \rightarrow get test data x_t

find k training examples $(x_1, y_1), \dots, (x_k, y_k)$

closer to x_t value

classification $\rightarrow (y_1, y_2, \dots, y_k)$ - majority prediction

K-Nearest Neighbors



K-odd

regression \rightarrow

avg (y_1, y_2, \dots, y_k)

prediction

max/min

Voronoi Diagram

Distance (Euclidean Distance)

$$x_i = \{x_{i1}, x_{i2}, \dots, x_{im}\}$$

$$x_j = \{x_{j1}, x_{j2}, \dots, x_{jm}\}$$

$$D = \sqrt{\sum_{l=1}^m (x_{il} - x_{jl})^2}$$

spherical → in 3D

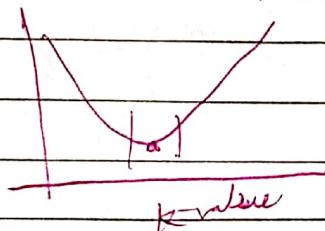
Assumption in traditional KNN

- giving equal weight to all attributes
- scale of attribute ranges are same
- equal variance
- classes are spherical
- attributes have very less noise
- all attributes have same importance

How to overcome?

- Use large K value
- Weighted distance function

validation error

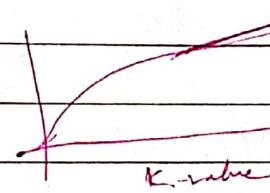


Impact of K

K is less → more discrete, not smooth boundary
 K is more → more smooth boundary

error / (less) boundary

training error



→ Distance weighted KNN.

- Take large value of K

$$\text{Prediction}_{\text{test}} = \frac{\sum_{i=0}^K w_i \cdot \text{class}_i}{\sum_{i=0}^K w_i} \quad w_K = \frac{1}{\text{Dist}(x_k, x_{\text{test}})} \quad \begin{matrix} \text{more focus on closer} \\ \text{further} \end{matrix}$$

→ Weighted Distance function-

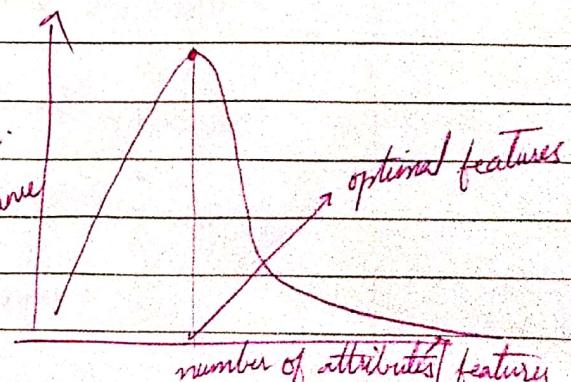
$$D(x_i, x_j) = \sqrt{\sum_{i=1}^m w_i (x_{im} - x_{jm})^2}$$

$w_i \rightarrow$ high (more imp. feature)
 $w_i \rightarrow$ low (less imp. feature)
 $w_i = 0$ (no impact feature)

Feature Engineering.

Curse of Dimensionality.

decreasing performance



- Irrelevant feature
- Redundant feature
- limited training examples
- limited computational resources

Steps

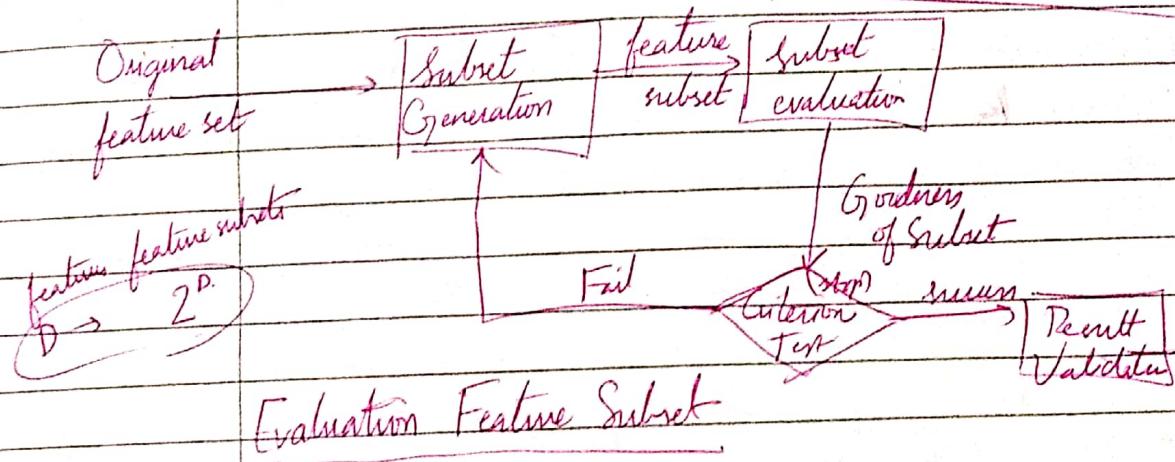
Feature Selection → maintains accuracy or increase

Feature Extraction

(higher dimension to lower dimension)

feature

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{high} \geq \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$



- Supervised Method (Wrapper method)
 - Training using selected feature subset
 - Estimate error on validation dataset
- Unsupervised (Filter method)
 - Look at input only
 - Select the subset that has most information

Filter Method

original features

↓
importance score

↓
selecting the best

↓
learning algorithm

↓
performance

Wrapper Method

Original features

↓
generate a subset

↓
selecting best subset

↓
learning algo

↓
performance

Method for filter method - ① Select uncorrelated feature

② Forward Selection Algo

start with empty feature
try each remaining feature
eliminate word (best forward)
select feature that gives next improvement
stop where there is no improvement
same for all

③ Backward Selection Algo

↓
start with full feature set
try removing features
Drop feature with smallest impact on error

Stopping criteria → Recursive Feature Elimination At $\leq A_m$

• Use Linear Classifier to find W (weight of each feature).

• Steps → Compute W on all features

(w_1, w_2, \dots, w_M) • Remove feature with smallest $|w_i|$ → magnitude!

least w value again compute (remove) recompute W on reduced data (feature)

repeat this till stopping criteria doesn't meet

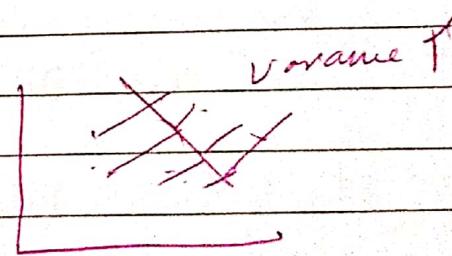
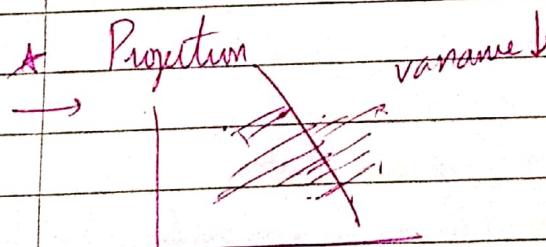
(w_1, w_2, \dots, w_{M-1})

new set → Feature Extraction - $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_{M-1} \end{bmatrix} = f \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

Projection of higher dimension to lower dimension $N > M$.

→ Uncorrelated, can not reduce further

→ high variance



$$z = w^T x$$

$$\dim z < \dim x$$

Variance →

$$\sigma^2$$

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

mean

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

Standard deviation
sq root of variance

Covariance

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

$$\rightarrow \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\cdot \text{Cov}(x, n) = \text{Var}(x)$$

$$\cdot \text{Cov}(x, y) = \text{Cov}(y, x)$$

\rightarrow Covariance Matrix \rightarrow is symmetric

$$\text{Cov}(\Sigma) = \begin{bmatrix} \text{cov}(n_1, n_1) & \text{cov}(n_1, n_2) & \dots & \text{cov}(n_1, n_m) \\ \vdots & \text{cov}(n_2, n_2) & & \vdots \\ \text{cov}(n_m, n_1) & \dots & \dots & \text{cov}(n_m, n_m) \end{bmatrix}$$

$$\text{cov}(\Sigma) = \frac{1}{m} (X - \bar{X})(X - \bar{X})^T; \text{ where } X = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_m \end{bmatrix}$$

↑ vectorized

\rightarrow diagonal elements are variances, i.e. $\text{cov}(n, n) = \text{Var}(n)$

\rightarrow covariance matrix is symmetric

\rightarrow it is a (true) semi definite matrix

- All eigenvalues must be real.
- All eigenvectors corresponding to different eigenvalues are orthogonal.
- All eigenvalues are greater than or equal to zero.
- Covariance matrix can be diagonalized.

$$\text{cov} = P D P^T$$

\rightarrow How to Compute PCs \rightarrow Covariance Matrix S

\rightarrow Compute its eigenvectors: $\{a_i\}_{i=1}^d$

\rightarrow first eigenvector $\{a_1\}_{i=1}^p$ corresponds to p largest eigenvalues from the

$\begin{array}{c} \text{sign} \quad \checkmark \\ \downarrow \quad \downarrow \quad \text{left} \\ \checkmark \quad \checkmark \quad \checkmark \\ v_1, v_2, v_3, \dots, v_p \\ \text{eigen vector} \\ \text{highest covariance} \\ \rightarrow \text{transformation } G \text{ consists of } \{p\} \text{ PCs.} \\ G = [v_1 \ v_2 \ \dots \ v_p] \end{array}$

→ Most weights are in fully connected layers ~~fully connected~~

$$W = USV^T$$

• $W \in R^{m \times n}$, $U \in R^{m \times m}$, $S \in R^{m \times n}$, $V^T \in R^{n \times n}$

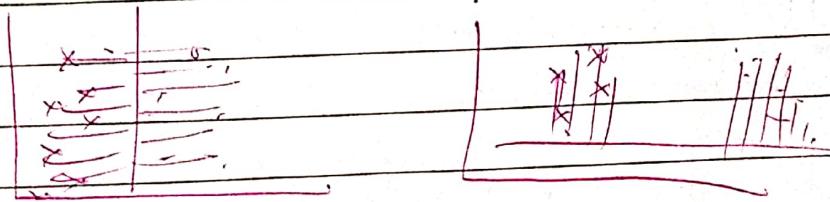
$S \rightarrow$ diagonal, decreasing magnitudes along diagonal.

$$[A] = [U] \begin{bmatrix} S & \\ & I \end{bmatrix} [V^T]$$

• $U \in R^{m \times k}$, $S \in R^{k \times k}$, $V^T \in R^{k \times n}$.

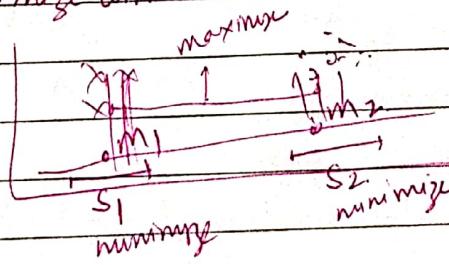
$k \rightarrow$ hyperparameter

→ PCA good for classification?



→ Linear Discrimination Analysis (LDA)

- Maximize b/w class distance
- Minimize within class distance



maximize dist b/w m_1, m_2
minimize dist within s_1, s_2

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \rightarrow \text{std dev}$$

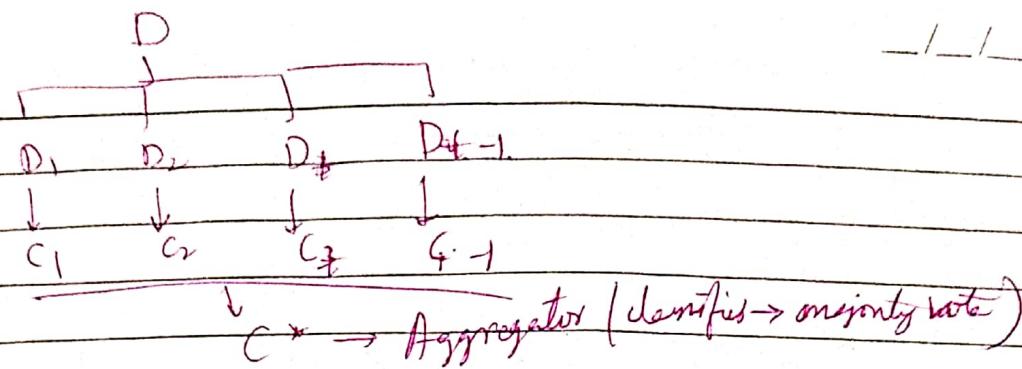
→ Fisher Linear Discriminants

- Means as far away as possible
- Scatter as small as possible

$$J(w) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

→ Ensemble Learning: Strong Learner Weak Learner

single or multiple
combine \rightarrow just better than random guess
accuracy $\rightarrow 70\%$



conditions → classifies outputs for same input should be diverse
 → classifies operate independent or partially independent

→ classifies have unique or local knowledge

Hybridization (when we can't live)
 we can combine → works if there is some 'random factor'

we can combine the same learning algo trained over diff subsets of training data

for certain we use same algo over same data but with different weights over data instances

Randomization of DTs.

→ Data Randomization → diff subsets of data for each tree

→ Random Subspace → each tree grows with random subset of features

→ Algorithm randomization → instead of best split attribute choose randomly

→ Aggregation Methods

- Majority Voting

- Weighted Majority Voting → can weight each classifier { o/p is classified }

- mean, product, minimum, max, weighted mean, avg

+ o/p is numeric

→ stacking → instead of above, we can train another classifier on o/p values of base classifiers

Generate Base learners
 Diverse
 most joint accuracy

Bagging
Random Forest
Boosting
AdaBoost

Bagging = Bootstrap + Aggregating

→ uses Bootstrap resampling

Sampling with replacement

Bagging \rightsquigarrow number of instances

For each of T iterations
Sample n instances from training set
(with replacement)

apply Learning algo to sample
store resulting model

For each of T models:
predict class of instance
using model
return class that is predicted most

A	B	C
D	E	F
G	H	I

E G
H I

A H
E C

F A
G C

model 3.

model 1

model 2

Random Forest: Input: training data set S_n, T, m .

① Choose T - no. of trees

② Choose m - no. of variables to split

③ Grow T trees

④ To classify point X , collect votes
from each tree

one tree

+ random feature split

}

→ many trees by bootstrapped data

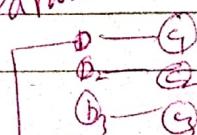
used for testing purposes

Out of the Bag } Samples not in the training data

↓ proportion of these samples wrongly classified

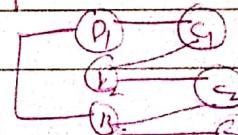
Bagging

↓ Parallel



Boosting

↓ Sequential



'Out of Bag Error'

Boosting

(Classification)

$$D_t = \frac{1}{m} m \text{ data points} - \{x_i\}_{t=1}^T$$

For $t=1$ to T :

Train learner using D_t

Get Weak Classifier h_t

choose $\alpha_t \leftarrow f(\epsilon)_{\text{error}}$

update $D_{t+1} = f(D_t)_{\text{error}}$

prev learner $\sum D = 1$

output $\rightarrow H(x) = \operatorname{sgn} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$

~~$D_{t+1}(i)$~~

$$D_{t+1}(i) = D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

Z_t → normalization factor

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

choose α_t to minimize training error

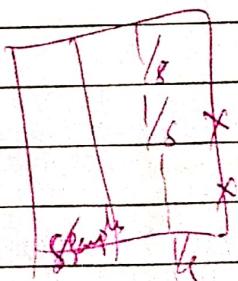
$$\epsilon_t = \frac{1}{2} \ln \left(\frac{1 - E_t}{E_t} \right)$$

$$E_t = \sum_{i=1}^m D_t(i) S(h_t(x_i) \neq y_i)$$

not learned
correctly

$$y_i \quad w^{(n)} \quad x_1 \quad x_2 \quad x_3 \quad \dots \quad x_n \quad y_{\text{most}} = 1$$

$$e^{-2} \quad e^{-2} \quad e^{-2} \quad e^{-2} \quad e^{-2} \quad \dots \quad e^{-2} \quad y^{(n)}$$



$$ex = y_S / 8$$

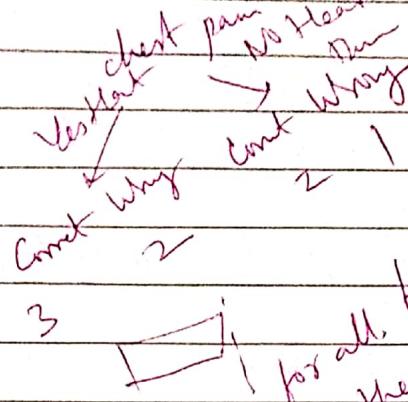
$$\boxed{\sum_{t=1}^m D_t = 1}$$

Random Forest

→ full sized trees

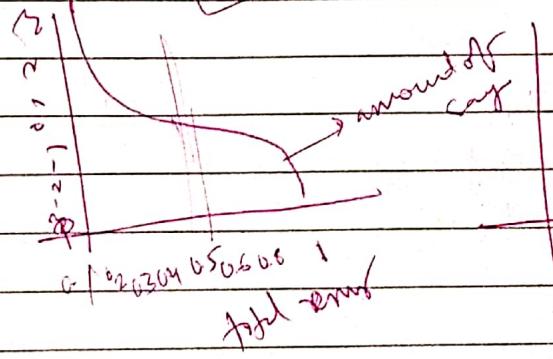
→ equal vote on first classification

Initially sample weights all same like $\frac{1}{n}$



for all β
then get
giving some

$$\text{Amount of Say} = \frac{1}{2} \log \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)$$



New sample = sample
weight \times amount
e.g.

The diagram illustrates the perceptron learning rule. It starts with a table of weights:

	old weight	New weight
$1/r$	0.05	0.05
$1/4$	0.17	0.17
$1/8$	0.07	0.07
$1/r$	0.14	0.14
$1/4$	0.33	0.33
$1/8$	0.07	0.07
$1/r$	0.14	0.14
$1/4$	0.33	0.33
$1/8$	0.07	0.07

Annotations explain the process:

- $H(a) = \text{sign}(w^T x)$ (written vertically on the left)
- $D_{t+1}(i) = D_t(i) e^{-(d_t x)}$ (at the top right)
- $\text{sum all of them and divide each by it}$ (with arrows pointing to the bottom row of weights)
- $0.21 + 0.14 = 0.35$ (with arrows pointing to the bottom row of weights)
- $1/r \approx 0.21 \text{ and } 0.33$ (with arrows pointing to the bottom row of weights)
- $w' = w + \alpha d_t x$ (written at the bottom right)
- $\text{small range} \rightarrow \text{cift in between 0 to 1 then include 0.07 to 0.33}$ (written at the bottom center)
- $8 \text{ data points} \rightarrow \text{generate 8 random numbers b/w 0 and 1}$ (written at the bottom left)
- $\text{Finally we want sign}(w^T x = 0)$ (written vertically on the left)
- $\text{Driver goes on until we have sufficient samples}$ (written at the bottom left)

Gradient Boosting

= Gradient Descent + Boosting.

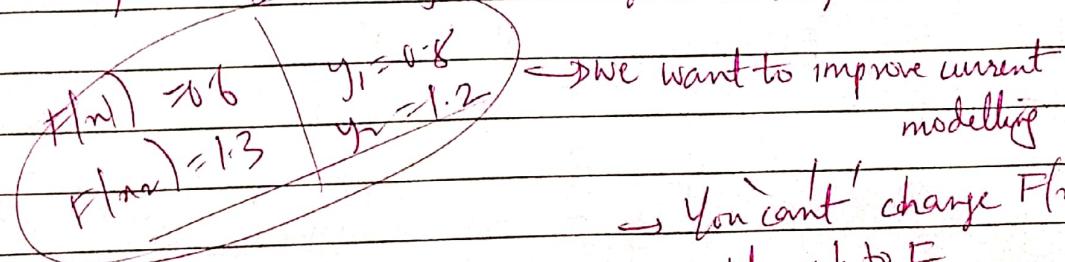
$$h(n) = \arg \left(\sum_{t=1}^T \alpha_t h_t(x) \right) \leftarrow \text{Regression}$$

$$h(n) \text{ or } \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(n) \right) \leftarrow \text{Binary classification}$$

AdaBoost: Shortcoming of Previous Model.

↳ high weighted data.

GB → Shortcomings are identified using Gradient

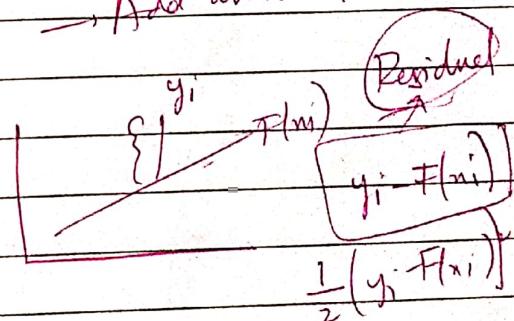


$F(n) + h(n) \rightarrow$ final model

$$F(n_1) + h(n_1) \approx y_1$$

$$F(n_2) + h(n_2) \approx y_2$$

$$h(n_1) = y_1 - F(n_1)$$



we want to approximate $y_i - F(n)$ to achieve $h(n)$.

(we can fit regression tree h to define $y_1 - F(n_1)$, $y_2 - F(n_2)$)

$h(n) \rightarrow$ unable to satisfy → add more trees

$$h_1(n) \rightarrow \frac{h_2(n)}{y_1 - (F(n_1) + h_1(n))}$$

Review G.D

$$w_i = w_i - \alpha \frac{\partial J(w)}{\partial w_i}$$

$$L(w) = \frac{1}{2} (y - h_w(x))^2$$

we want to minimize

$$J(w) = \sum_{i=1}^m L(w)$$

$$L(y, F(n)) = \frac{1}{2} (y - F(n))^2$$

minimize

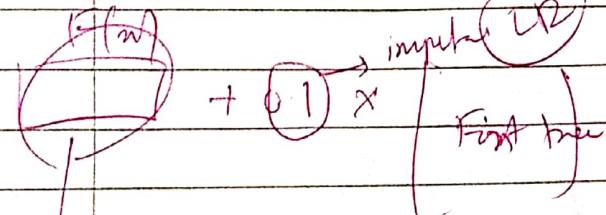
$$J = \sum_{i=1}^m L(y_i, F(n_i))$$

adjust $F(n)$

$$\frac{\partial J}{\partial F(n_i)} = \frac{\partial \frac{1}{2} (y - F(n))^2}{\partial F(n)}$$

$$y_i - F(n_i) = -\frac{\partial J}{\partial F(n_i)} = -F(n) + y$$

Gradient Boost



$$F(n) = F(n) + h(n)$$

$$F(n) = F(n) - \frac{y - F(n)}{\frac{\partial J}{\partial F(n)}} = h(n)$$

based on residual

$$+ 0.1 \times \text{Second tree}$$

Adaboost \rightarrow small tree
(called stump)

Average Weight
of all
weights

avg weight = 1/2

Gradient Boost \rightarrow
we chose starting

trees larger
than stumps
from one leaf and
restricting on no. of
leaves

Flight Below Color weight Period

88 16.1

76

1

81

1

73

1

77

1

57

1

tree made with simple DT
method

restrict leaves = 4

80 11.2

Gender \neq

Flight < 16

color = blue

14.2 - 15.2

14.8

14.5 - 16.8

16.8

14.7 - 15.7

12.9

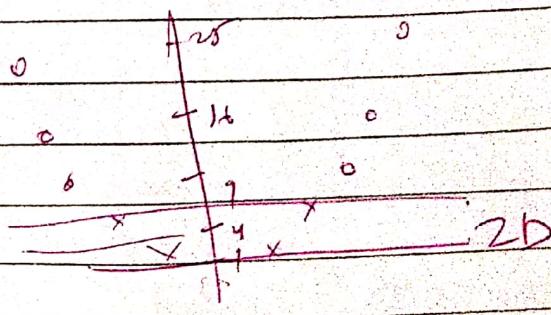
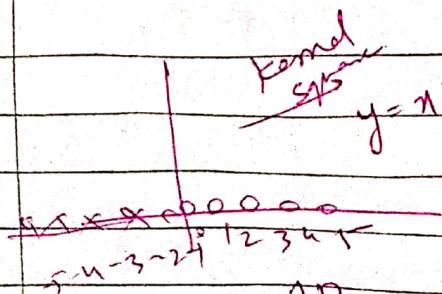
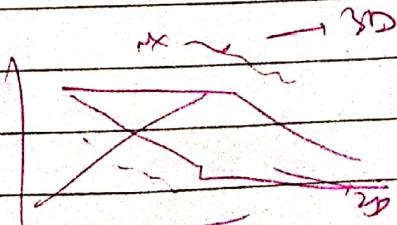
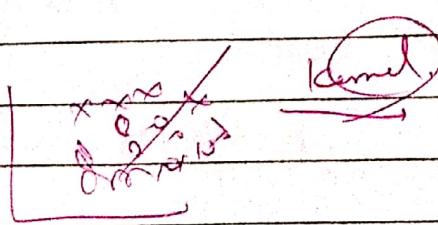
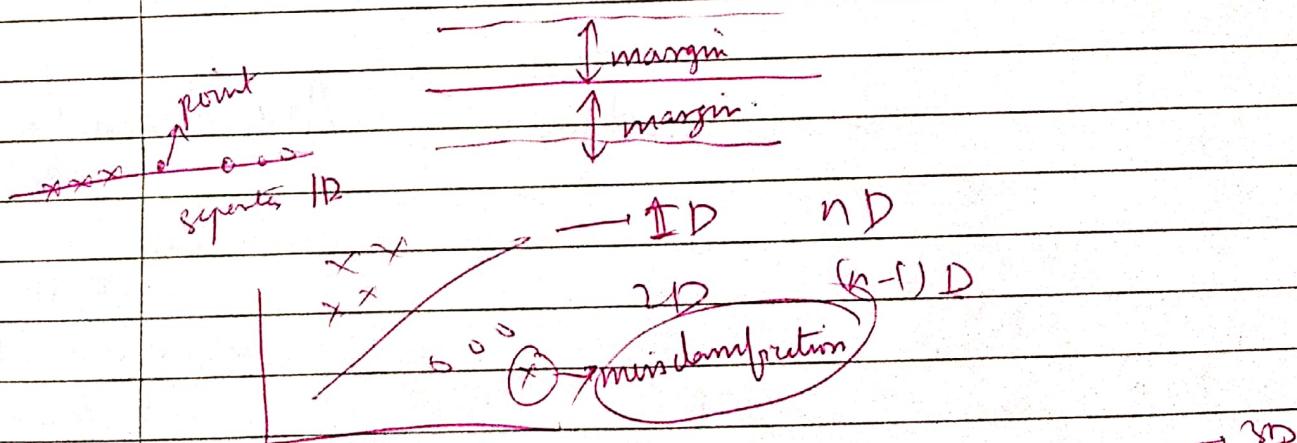
3 - 8

16.8 + 17.2

$f(x) = f_0 + h_1(x)$
 Rerudant values (overall)
 $1.8 \rightarrow 15.1$
 $4.8 \rightarrow 4.3$
 $-15.2 \rightarrow -15.2$
 $1.8 \rightarrow 1.4$
 $5.1 \rightarrow 5.4$
 $-14.2 \rightarrow -12.7$
 small step in right direction
 (minimize residual)

$$f(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x)$$

Support Vector Machine

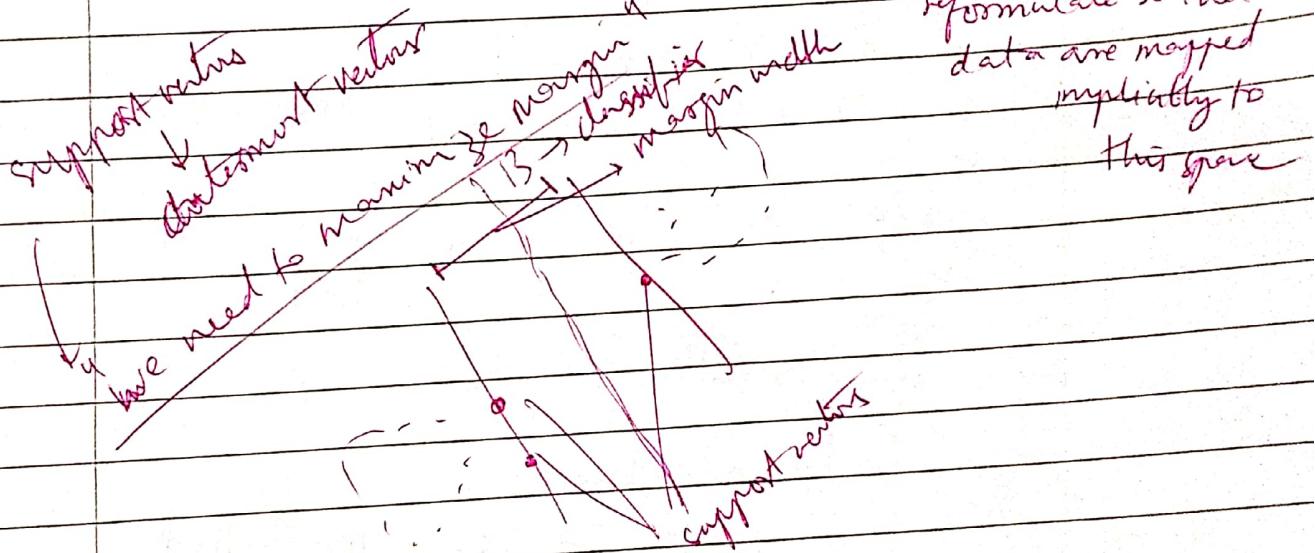


SVM.

① Define optimal hyperplane is (taking into account that it needs to be computed efficiently) maximize margin

② Generalize to non-linearly separable problems:

③ Map data to high dimensional space where it is easier to classify with linear decision surface



Linear Separating Hyperspaces

→ this linear separating plane need not pass through the origin of our feature space, i.e. it does not need to include the zero vector as an entity within the plane. Such hyperplanes are called 'affine'

→ Real dimensional ~~plus~~ plus polynomial feature spaces known mathematically as \mathbb{RP}^n → our linearly separating plane is an affine $p-1$ dimensional space embedded within it

$$y_i \in \{-1, +1\}$$

$$f(x) = \text{sign}(w \cdot x + b)$$

$$wx + b > 0$$

$$wx + b \leq 0$$

$$wx + b = +1$$

$$wx + b = -1$$

Dist from a point (x_0, y_0) to
a line
 $ax + by + c = 0$

$$\frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

$$wx + b = +1$$

$$wx + b = -1$$

$$\left. \begin{array}{l} x_i \cdot w + b \geq +1 \text{ when } y_i = +1 \\ x_i \cdot w + b \leq -1 \text{ when } y_i = -1 \end{array} \right\} \text{combined } y_i(x_i \cdot w) \geq 1$$

$$\frac{w \cdot x + b}{\|w\|} \rightarrow k = \frac{1}{\|w\|} \rightarrow 2k = \frac{2}{\|w\|} \quad \max \frac{2}{\|w\|}$$

So the problem becomes →

$$\max \frac{2}{\|w\|} \quad \text{or} \quad \min \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1, \forall x_i$$

Linear, Hard-Margin SVM Formulation

Find w, b that solve $\min \frac{1}{2} \|w\|^2$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1$$

Quadratic program: quadratic objective, linear (in)equation constraint

- There is also a unique minimizer, i.e. w and b values to provide the minimum.
- No solution if the data are not linearly separable

Lagrange Multipliers

$$\underset{w, b, \alpha}{\text{Minimize}} \quad \phi(w) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (w \cdot x_i + b) + \sum_i \alpha_i$$

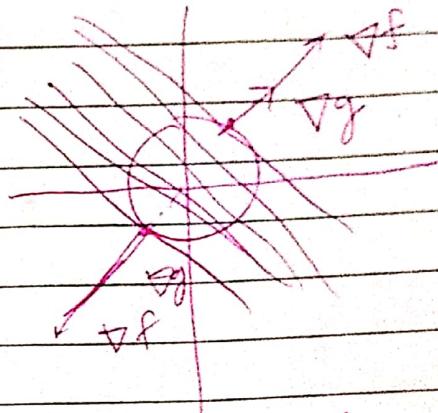
- Convex quadratic programming problem
- Duality theory applies

$$f(x, y) = 3x + 4y \leftarrow \text{Max/Min}$$

s.t. $x^2 + y^2 \leq 1$

$$\nabla f = \lambda \nabla g$$

$$\nabla f - \lambda \nabla g = 0 \text{ max/min}$$



$$f = 3x + 4y \quad g = x^2 + y^2 - 1$$

$$\frac{\partial g}{\partial x} = 2x + 0 + 0$$

$$\nabla f = \lambda \nabla g$$

$$\nabla f - \lambda \nabla g = 0$$

$$\frac{\partial f}{\partial x} - \lambda \frac{\partial g}{\partial x} = 0$$

$$\frac{\partial f}{\partial y} - \lambda \frac{\partial g}{\partial y} = 0$$

$$3 - \lambda 2x = 0$$

$$4 - \lambda 2y = 0$$

$$x = \frac{3}{2\lambda}$$

$$y = \frac{4}{2\lambda}$$

$$\min f(x, y)$$

$$\lambda = -\frac{3}{2}$$

$$\frac{9}{4\lambda^2} + \frac{16}{4\lambda^2} = 1$$

$$4\lambda^2 = 25$$

$$\lambda = \pm \frac{5}{2}$$

$$x = -\frac{3}{5}, y = \frac{4}{5}$$



$$\min \frac{1}{2} \|w\|^2$$

$$\text{s.t } y_i(w \cdot x_i + b) \geq 1, \forall x_i$$

$$L = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y_i(w \cdot x_i + b) - 1]$$

$$\frac{\partial L}{\partial w} = w - \frac{\partial}{\partial w} \left[\sum_i \alpha_i y_i w \cdot x_i + \sum_i \alpha_i y_i b - \sum_i \alpha_i \right]$$

$$\Rightarrow w - \sum_i \alpha_i y_i x_i = 0$$

$$w = \sum_i \alpha_i y_i x_i$$

$$w = \sum_i \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_i \alpha_i y_i = 0$$

$$L = \frac{1}{2} \left(\sum_i \alpha_i y_i x_i \right) \left(\sum_j \alpha_j y_j x_j \right) - \sum_i \alpha_i y_i w x_i +$$

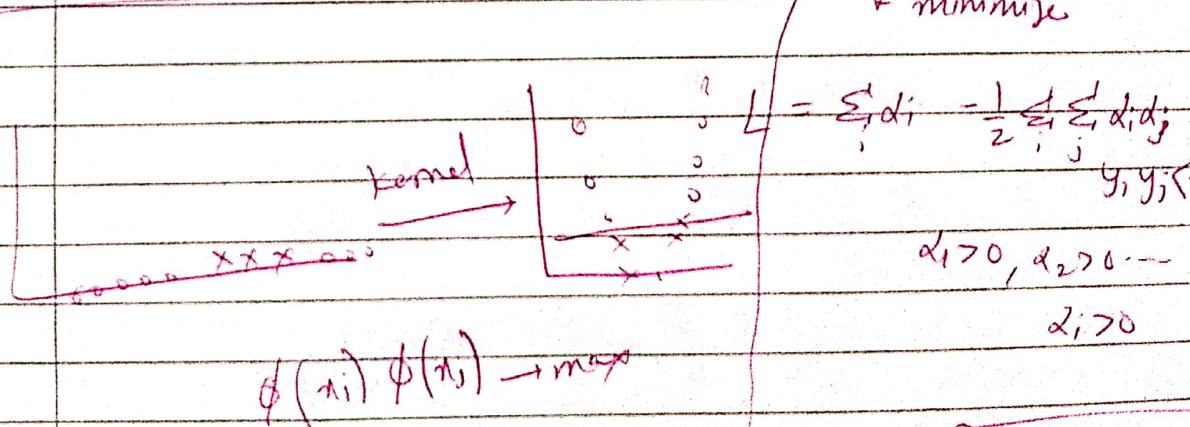
$$\sum_i \alpha_i y_i b +$$

$$= \frac{1}{2} \left(\sum_i \alpha_i y_i x_i \right) \sum_j \alpha_j y_j x_j - \left(\sum_i \alpha_i y_i x_i \right) \left(\sum_j \alpha_j y_j x_j \right) +$$

$$= -\frac{1}{2} \left(\sum_i \alpha_i y_i x_i \right) \left(\sum_j \alpha_j y_j x_j \right) + \sum_i \alpha_i$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j$$

+ minimize



$$\phi(x_i) \phi(x_j) \rightarrow \max$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi(x_i) \phi(x_j)$$

if we pass through kernel

once you pass
through kernel
simply pass out
dot product

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

$$k(\alpha_j, \alpha_j) = \langle \phi(\alpha_1), \phi(\alpha_j) \rangle.$$

\implies Polynomial kernel

$$k(x_i, x_j) = \underbrace{(x_i \cdot x_j - 1)}_{P}$$

\Rightarrow Gaussian kernel \rightarrow also in sklearn library uses

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

→ Find max and min values of $f(x, y) = 8x^2 + y^2$ subject to constraint $4x^2 + y^2 = 9$.

$$f(n) = 81n^2 + y^2 \quad g(n) = 4n^2 - y^2 = 9 \quad L = (81n^2 + y^2) - x(4n^2 - y^2) = 9$$

$$162n^2 - 72y^2 = 0 \quad 2n(81 - 4y^2) = 0$$

$$n=0, y=81/4$$

$$162x - 18x = 0 \quad 2x(81 - 9) = 0$$

$$x=0, \lambda=81/4$$

$$\begin{aligned}
 & \text{The company produces 100 units annually} \\
 f(x+y, x) & \Rightarrow 8x^2y^2 - 200(x+y)^2 = 0 \\
 & 2y(1-x) = 0 \\
 & y=0, x=0 \quad (0, 0) \\
 & y=\pm 3, x=1 \quad (1, 3) \\
 & y=\pm 3, x=-1 \quad (-1, 3) \\
 & y=\pm 3, x=2 \quad (2, 3) \\
 & y=\pm 3, x=-2 \quad (-2, 3)
 \end{aligned}$$

$$\begin{aligned}
 & g(n_1, y_1) \Rightarrow xy + z = 100 \\
 & f(n_1, y_1, z) - yg(n_1, y_1, z) \text{ min value} = f(-3, 21, 0) = 81 \times -3 + 21 \text{ (min)} \\
 & f(n_1, y_1, z) = 81 \times -y + 21 \text{ (max)} \\
 & \frac{729}{y} \\
 & \cancel{\frac{729}{y}} \quad 81 \times +9 \text{ (max)} \\
 & f(n_1) = 20ht \\
 & g(n_1) = 20ht \\
 & \frac{20}{2} = 10 \\
 & \cancel{\frac{20}{2}} = 10 \\
 & \frac{20}{xy} = 10 \\
 & \cancel{\frac{20}{xy}} = 10 \\
 & \text{budget constraint} \quad \frac{20}{xy} - 10 \\
 & \left(\begin{array}{c} \text{eq1} \\ \text{eq2} \\ \text{eq3} \end{array} \right) = \begin{array}{l} \text{max } y \rightarrow \text{not going to} \\ \text{revenue model} \\ \text{so costs} \end{array} \\
 & \text{multiplex}
 \end{aligned}$$

$$h^{(n_1)} y_{3L} = f^{(n_1+1)}$$

$$\frac{yL}{x} =$$

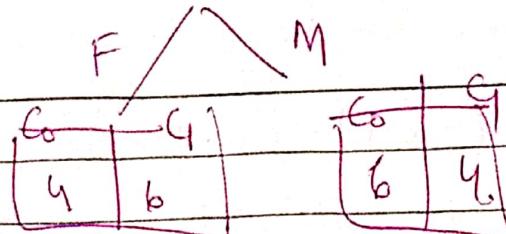
~~6/1
a/(his) - 20
1705
2000
21
23~~

budget constraint

~~eg1~~ = ~~most n/m~~
~~eg2~~ ~~single n/m~~
~~eg3~~ ~~as cols~~ ~~multiple~~
~~eg4~~ ~~rows~~

multiplex
choices
(check out, some problems)

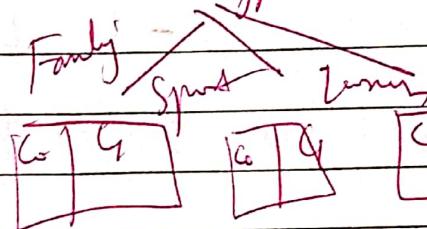
Gender



$$G_1 = 1 - \left(\frac{4}{10}\right)^2 - \left(\frac{6}{10}\right)^2 \quad \text{weighted sum}$$

$$G_2 = 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 \quad \frac{(10) G_1 + G_2 (10)}{20}$$

CarType



Information gain \uparrow Split that
Gini impurity \downarrow split that