

# SMAI

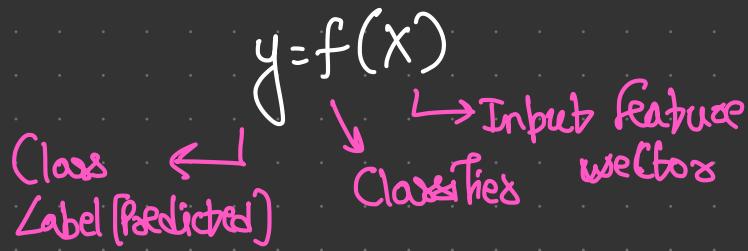
## Quiz

### Notes

## Recognitions:

- Identification of a pattern as member of known category

## Problem formulations



## Pattern:

Opposite of chaos, vaguely defined

### Pattern class:

- Collection of similar objects
- Class defined by samples of class

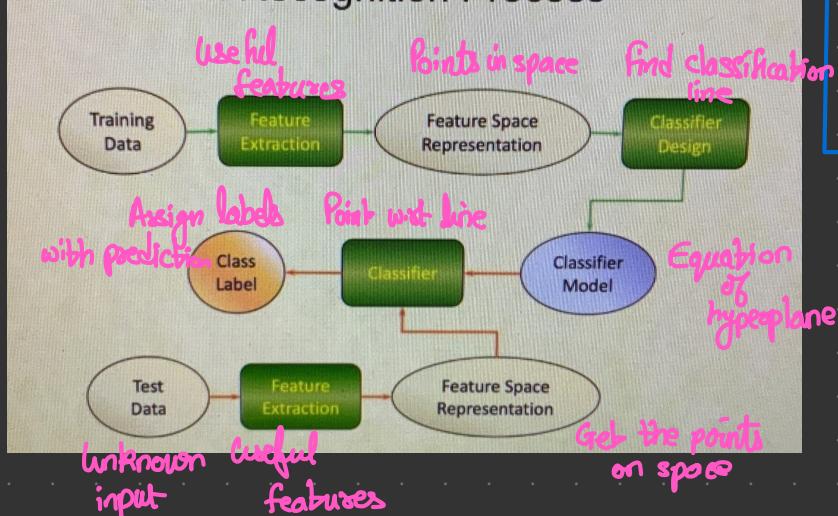
Intra-class variability  
Among groups

Interc-class variability  
Within groups  
Examples:  
Trees

## Pattern Classification

- Remove irrelevant info
- Feature extraction - Extracting useful info. from raw data

## The Pattern Recognition Process



## Pattern Recognitions:

- Learns to classify
- Cognition: forming new classes
- Recognition: Known classes

## Approaches to recognitions:

Patterns represented in feature space

- ↳ Learn separating boundary b/w classes
- ↳ Learn stat. model for pattern generation in feature space

## Nearest Neighbors Classifier:

- Assign label of sample <sup>nearest</sup> to test sample

→

## Voronoi Tesselation:

- Boundary b/w every point
- Then ignore those surrounded by same
- Decision boundary is piecewise linear



### Summary

- Nearest Neighbor is a simple yet effective classifier
- Error rate (with increasing number of training samples) is at-most twice that of the ideal classifier
- Easy to implement, Training is trivial
- Easy to interpret
- Time complexity is high during testing
  - Can be improved in different ways

## k-Nearest Neighbor

- Assign test sample label that is most frequent amongst k-Nearest neighbors
- k-parameters that user selectable
- Output label can change with change of k

## Sample Pruning:

- Remove samples that don't affect decision boundary
- To identify samples to be pruned:
  - ↳ One per class
  - ↳ Class means

↑  
Nearest mean  
classifies  
(Equivalent to  
1-Nearest Neighbor)

### Summary

- Time complexity of Nearest Neighbor is high during testing
  - Can be improved in different ways
- Very popular classifier due to
  - Simplicity
  - Explainability
  - Robustness (fewer assumptions)
  - Ability to do continuous learning
  - Guarantees on Error Rate

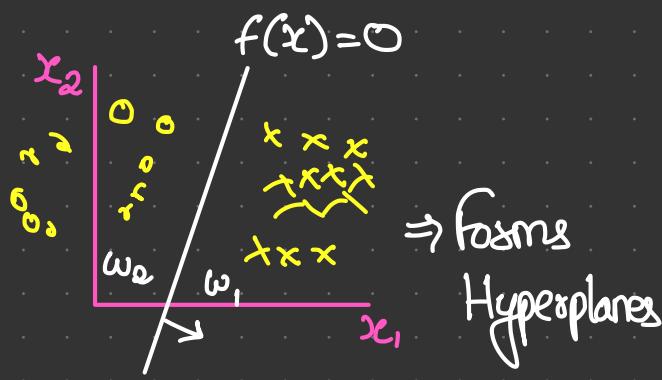
## Lecture 3

### Linear decision boundary:

→ Straight line separating 2 classes

$$y = mx + c$$

↓ Vectors



→ Generalised eqn:

$$w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = 0$$

$$\Rightarrow f(x) = 0$$

-ve                      +ve  
class  $w_2$               class  $w_1$

→ For 2-class,

$$f(x) = w_0 + w_1 x_1 + w_2 x_2 = w \cdot x \quad (\approx) \quad w^T x$$

$$w = [w_0 \ w_1 \ w_2]^T \quad x = [1 \ x_1 \ x_2]^T$$

Decision boundary:  $w^T x = 0$

Learn  $d+1$  values of  $w_i$

Goal:

Learn  $w_i$ 's to draw decision rule

Summary:

- Linear decision boundary are popular in ML due to:
  - Simplicity of model
  - Theoretical result
  - Simplicity of computation

### Gradient Descent

### Perceptron Learning Algorithm

To learn decision boundary:

→ Randomly initialise  $w_i$ :

→ Iteratively modify  $w_i$  so that each sample is correctly classified

To update automatically

→ Define Loss function  $J(w)$  which will be minimum when training samples are correctly classified

$$\nabla J = \frac{\partial J}{\partial w} \quad J$$

Gradient Vector

### Gradient Descent Algorithm:

1. Randomly initialise elements of  $w$ :  $w^0$

2. Compute gradient of  $J$  at  $w$ :  $\frac{\partial J}{\partial w}$

$$3. w^{t+1} = w^t - \eta \frac{\partial J}{\partial w}$$

4. Repeat 2 & 3 steps until convergence ( $\eta \frac{\partial J}{\partial w}$  becomes small)

## Loss functions for GD:

1. No. of misclassified samples  
↳ Not differentiable

2.  $w^T x_i > 0$  for +ve class  
 $< 0$  for -ve class

$y_i(w^T x_i) > 0$  for all samples

$$J(x) = -y_i w^T x_i$$

$$\nabla J = \frac{\partial J}{\partial w} = -y_i x_i$$

Batch perceptron rule:

$$J_p(x) = \sum_{x_i \in X} -y_i w^T x_i$$

↳ Set of misclassified samples

Perceptron update rule

To compute  $J$ :

→ Per training sample  $\Rightarrow$  Single sample or Stochastic GD

→ Summed up over batch  $\Rightarrow$  Batch GD

→ In between  $\Rightarrow$  Mini-Batch GD

Perceptron Loss fn

$$\nabla J_p = \sum_{x_i \in X} -y_i x_i \Rightarrow w^{t+1} = w^t + \eta \sum_{x_i \in X} y_i x_i$$

### Summary

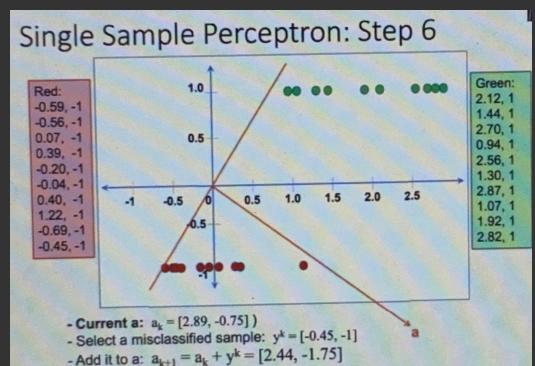
- Linear Decision boundaries are easily expressed as dot product with the co-efficient vector
- Gradient Descent allows us to learn the decision boundary from a set of labelled training samples
- Several variants of GD based on the loss function used
  - Perceptron learning algorithm is one such

## Lecture 4

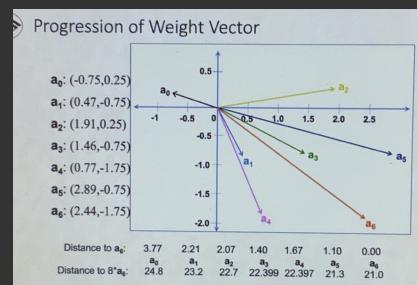
(Revise perceptron loss fn)

### Single Sample Perception Training

- Steps:
- 1) Augment  $[x, 1]$
  - 2) Assign  $[tx, t]$  based on label
  - 3) Randomly initialise weight ( $a_k$ )
    - ↳ Select misclassified sample ( $y^k$ )
    - ↳ Add  $a_{k+1} = a_k + y^k$
    - ↳ Finally  $a_R = [a_1, a_2]$   
 $a_1 x - a_2 = 0$  a solution



NOTE: Weight vector taken together when multiplied with particular constant is arranged in descending order



# Relaxation Algorithm:

$$J_p(\omega) = \sum_{x_i \in X} (\omega^T x_i)^2 \quad [\text{Augmented & sign multiplied}]$$

- Smoother
- Often too smooth to converge to boundary
- Dominated by large vectors
- Solutions: Norm. & add margin

$$J_\delta(\omega) = \frac{1}{2} \sum_{x_i \in X} \frac{(\omega^T x_i - b)^2}{\|x_i\|^2} \Rightarrow \nabla J_\delta(\omega) = \sum_{x_i \in X} \frac{\omega^T x_i - b}{\|x_i\|^2} x_i$$

$X = \text{Set of samples where } \omega^T x_i < b$

## Summary

- Linear Decision boundaries are easily expressed as dot product with the co-efficient vector
- Gradient Descent allows us to learn the decision boundary from a set of labelled training samples
- Several variants of GD based on the loss function used
  - Perceptron learning algorithm is one such
  - Relaxation Methods
- GD Methods are generic and applicable to other problems as well
  - We will derive more of them later

# Linear Regression

→ find hyperplane that best approximates  $(y_i)$  in terms of  $(x_i)$

→ Define loss function that's minimised when fit better

↳ Use GD

↳ Distance of samples from line  
↳ Squared distance from line

# Polynomial Regression

$$f(x, \omega) = \sum_{j=0}^M w_j x^j$$

$$J_p(\omega) = \frac{1}{2} \sum_i (f(x_i, \omega) - y_i)^2$$

$$\nabla J_p[\cdot] = \sum_i (f(x_i, \omega) - y_i) x_i$$

# Linear Regression with GD

$$J_{l_i}(\omega) = \frac{1}{2} \sum_i (\omega^T x_i - y_i)^2$$

$$\nabla J_{l_i} = \sum_i (\omega^T x_i - y_i) x_i$$

Use GD to minimise loss fn,  
update rule:

$$\omega^{t+1} = \omega^t + \eta \sum_i (y_i - \omega^T x_i) x_i$$

→ Closed form solutions exist

↳ Can do curve fitting as long as it's linear in  $\omega$ .

## → Basic fn Models

$$f(x_i, \omega) = \sum_{j=0}^M w_j \psi_j(x)$$

*Basis functions*

## Logistic Regression

$$\rightarrow \sigma(a) = \frac{1}{1+e^{-a}} \rightarrow \text{Sigmoid fn}$$

$$\rightarrow \frac{\partial \sigma}{\partial a} = \sigma(1-\sigma)$$

→ Use logistic sigmoid of linear fn to approx.  
two class labels: [0, 1]

$$J(\omega) = J_{l_0}(\omega) = \frac{1}{2} \sum_i (\sigma(\omega^T x_i) - y_i)^2$$

$$\nabla J_{l_0} = \sum_i (\sigma(\omega^T x_i) - y_i) \sigma($$

→ Not affected by outliers  
(compared to GD)

*Do this later*

## Summary

- Linear Regression is a generic tool to fit several types of curves to a set of data points
- Multi-linear regression refers to approximating multiple output values for the same input
- Logistic regression uses the regression formulation for 2-class classification
- Will discuss non-linear regression problems later

## MSE Procedures

→ Constrain  $\omega^T x_i > 0$  to  $\omega^T x_i = b$   
for +ve b

→ n eq. with n training samples

$$X_{n \times d} \hat{w}_{d \times 1} = b_{n \times 1} \Rightarrow Xw = b$$

Multiply by  $X^{-1}$  both sides,

$w = X^{-1}b$

→ Need not be invertible

Let  $J_s(\omega)$  be diff b/w  $X\omega$  &  $b$ :

$$J_s(\omega) = \|X\omega - b\|^2 = \sum_{i=1}^n (\omega^T x_i - b_i)^2 \Rightarrow \nabla J_s = \sum_{i=1}^n 2(\omega^T x_i - b_i) x_i$$

Can do GD to find optimal  $\omega$  [Widrow-Hoff]

To get closed form solution for MSE,

$$\nabla J_s = \sum_{i=1}^n 2x_i (\omega^T x_i - b_i)$$

In matrix,

$$\nabla J_s = 2X^T(X\omega - b) \Rightarrow$$

At minimum  $G=0$ ,

$$\underbrace{X^T X}_{\hat{\omega}} \omega = \underbrace{X^T b}_{\hat{b}}$$

$$\Rightarrow \boxed{\omega = (X^T X)^{-1} X^T b}$$

On comparison,

$$(X^T X)^{-1} X^T = X^{-1}$$

$$X^+ = (X^T X)^{-1} X^T \rightarrow \text{pseudo inverse}$$

↳ Square & Non-singular  $X^+ = X^{-1}$

$$\Rightarrow X^+ = \lim_{\epsilon \rightarrow 0} (X^T X + \epsilon I)^{-1} X^T$$

$\omega = X^+ b$  is MSE solution to  $X\omega = b$

NOTE:

→ Solution depends on  $b$  and need not be on training samples  
→ Computationally heavy

### Ho-Kashyap Procedure

- If samples are linearly separable, there exists  $\hat{a}$  and  $\hat{b}$  such that:  $Y\hat{a} = \hat{b}$ , where  $b_i > 0$  for all  $i$ .

$$J_s(a, b) = \|X\omega - b\|^2$$

- As we do not know both  $\hat{a}$  and  $\hat{b}$ , we will search for both.

$$\nabla_w J_s = 2X^T(X\omega - b), \text{ and}$$

$$\nabla_b J_s = -2(X\omega - b).$$

- For any value of  $b$ , we can set  $\omega = X^+ b$ , and if  $e = (X\omega - b)$ ,

$$b_{k+1} = b_k + 2\eta e^+,$$

where  $e^+$  is  $e$  with all its negative components set to 0.

### Summary

- MSE procedures try to minimize the mean-squared error between the predicted value by the linear classifier and a set of constants,  $b$ .
- Depending on the choice of  $b$ , and the optimization procedure, we have different approaches.
- Each one has different numerical properties

### Approach 2: One vs One (Pairwise)

→ Compute 1 line per class pair (A vs B etc)

→ Can lead to 3-way ties but lesser area relatively

### Linear to MultiClass (Lec 6)

#### Approach 1: One vs Rest

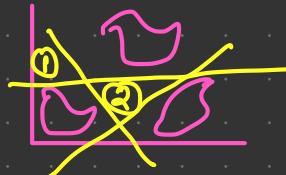
→ Compute 1 line per class (A vs B etc etc)

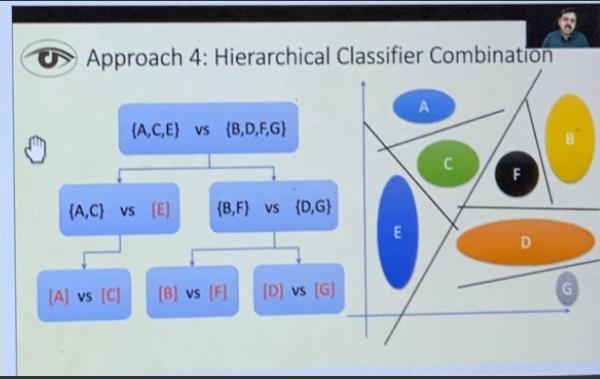
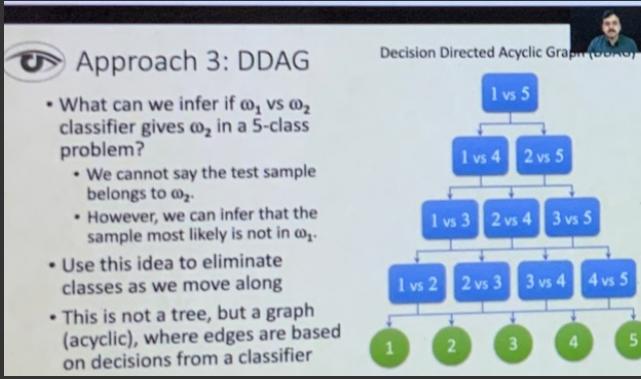
→ Can lead to 2-way, 3-way ties

→ Use weighted voting as solution

Summary: Simple Approaches

- 1-vs-Rest:
  - Build  $n$  classifiers and vote (efficient)
  - Difficult to train as one-vs-rest may not be linearly separable
  - Several possibilities for indecision/tie
- 1-vs-1 (Pairwise):
  - Build  $\binom{n}{2}$  classifiers and vote (compute heavy)
  - Easier to build as pairwise decision boundaries are simpler
  - Fewer possibilities for ties
- Weighted voting can break ties. Will explore ways to avoid it.





## Linear Discriminant fn's

IF  $g_1(x), g_2(x)$  are linear fn's such that  $g_1(x) > g_2(x)$ , then  $x$  is more likely to belong to  $\omega_1$ , else  $\omega_2$

→ Intersection of fn's becomes linear decision boundary  
 $f(x) = g_1(x) - g_2(x)$

### Summary: Efficient Approaches

#### • DDAG:

- Build  $n(n-1)/2$  classifiers and classify using  $n-1$  of them
- Compute heavy to build, but efficient in inference
- Effective in classification (uses 1-v-1)

#### • Hierarchical

- Build  $n$  classifiers by appropriate grouping of classes and pass through  $\log(n)$
- Far more efficient than 1-vs-1 (or DDAG)
- Similar in accuracy to 1-v-1 and DDAG.
- Difficult to find the right grouping of classes

#### • No Ties

For multiple classes,  
 $w_i = \operatorname{argmax}_i g_i(x)$

→ partitions f-space  
 using piecewise linear  
 decision boundaries

→ converts multiclass  
 to two-class problems

### Summary: Linear Machines

We have several methods to convert linear two-class classifiers into multi-class linear machines.

- Some of the methods are simple and intuitive, others are more efficient and applicable in practice, and the linear discriminant functions provide mathematically elegant solutions for analysis and proofs.
- The methods are applicable irrespective of the learning method used to learn  $w$  (also to support vector machines).
- Note that the discriminant approach also generalizes to non-linear functions.

## Lec 7

## Dimensionality Reduction:

→ Reduce no. of features of feature vector

- Less stuff to measure
- Efficiency
- Simpler boundaries
- Ability to visualise

→ How?

- Select subset of features
- Define new features as linear function of original

## Feature Selection:

→ Select subset ( $V$ ) of original features ( $X$ )

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}$$

↑  
Selected 2<sup>nd</sup> &  
3<sup>rd</sup> feature

# Sequential Search

## Forward

- Start with  $V = \{y\}$ ,  $X = \{x_1, \dots, x_d\}$
- Select  $x_i$  which ↑ accuracy when added to  $V$
- Move  $x_i$  from  $X$  to  $V$
- Repeat 2 & 3 until satisfaction \*
- Accuracy doesn't ↑ anymore
- No. of features in  $V$  is desired

## Backward

- Start with  $V = \{x_1, x_2, \dots, x_d\}$
- $X = \{\}$
- Select  $x_i$  from  $V$  which ↑ acc or reduces ↓ in accuracy when removed from  $V$
- Move  $x_i$  from  $V$  to  $X$
- Repeat - 2 & 3 until satisfaction \*

## Feature Extraction:

- Create new features as lin. combos of original
- If two are highly correlated, can remove one

## feature spaces

- If confined close to a plane in

- 3D space
- Express them with 2 basis vectors (can get most info)
- Basis is orthogonal

↳ Find basis vectors, then shift plane to pass through origin, then shift points to be centred around origin

**Summary: Feature Selection**

- Feature selection is a way for dimensionality reduction
- Helps in reducing the number of features to be measured
- Mathematically: Matrix multiplication
- Sequential selection approaches
  - Greedy methods: Does not guarantee optimal accuracy
  - Explicit search for the best subset need  $2^d$  trials.

**Representing points on the plane**

- Center the points at origin
- The plane passes through origin
- Find the basis vectors
- Project each point to each of the basis vectors

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix} \begin{bmatrix} ht \\ wt \\ lt \end{bmatrix}$$

**Summary: Feature Extraction**

- Find a sub-space in which most points lie (or are close to)
- Find the basis vectors and stack them to form a matrix
- Mathematically: multiplication with the above matrix
- Can reduce dimensionality with minimal reduction in accuracy
- Optimal (MSE) methods exists to find the extraction matrix