# Supervised Learning: Linear Models
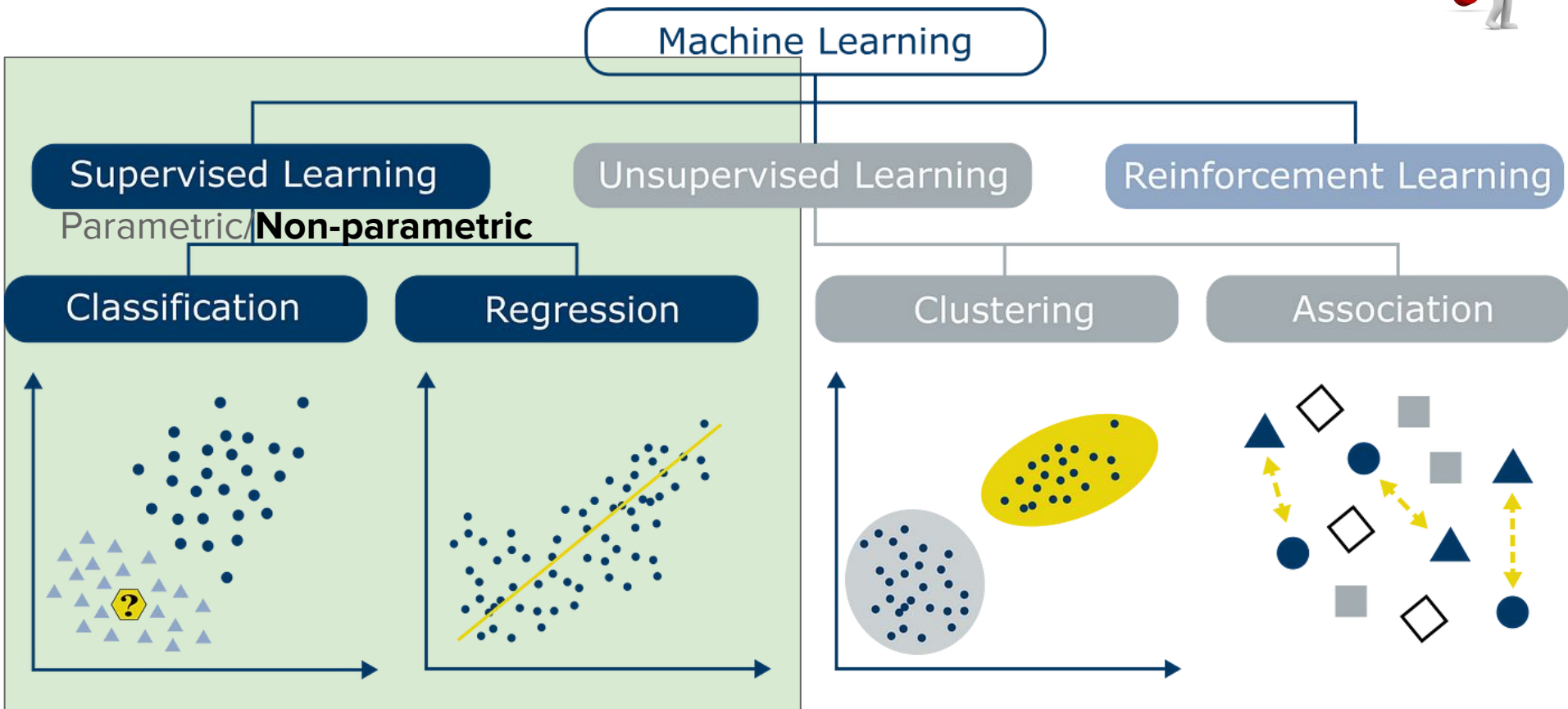
Mahesh Mohan M R
Centre of Excellence in AI
Indian Institute of Technology Kharagpur

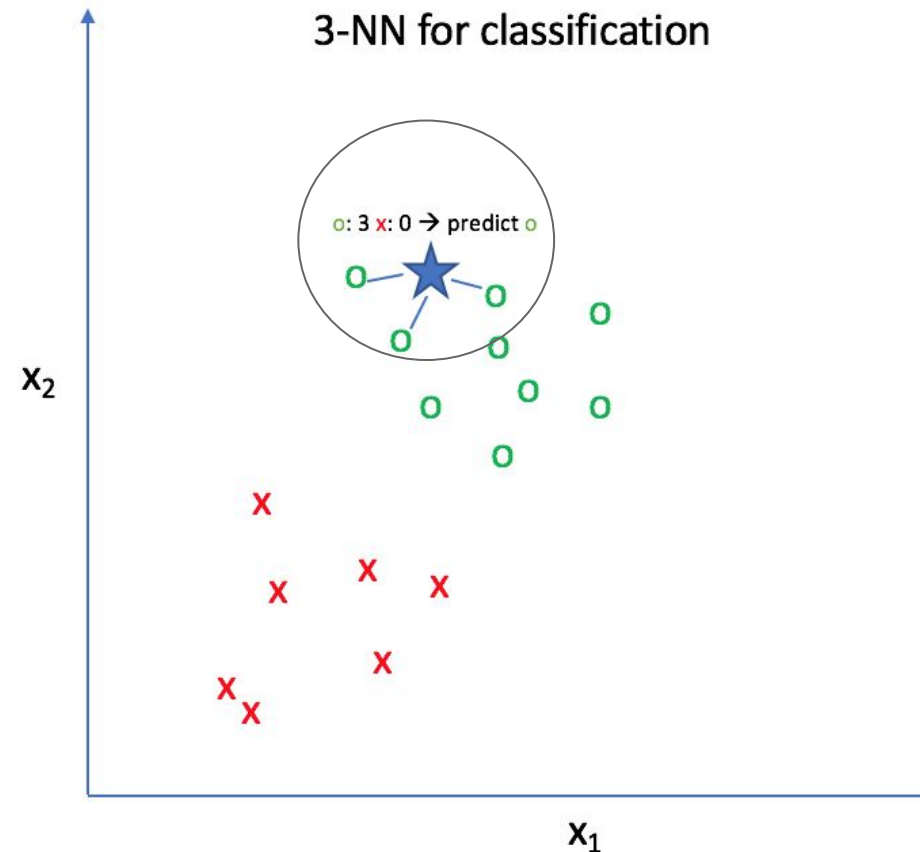# Recap

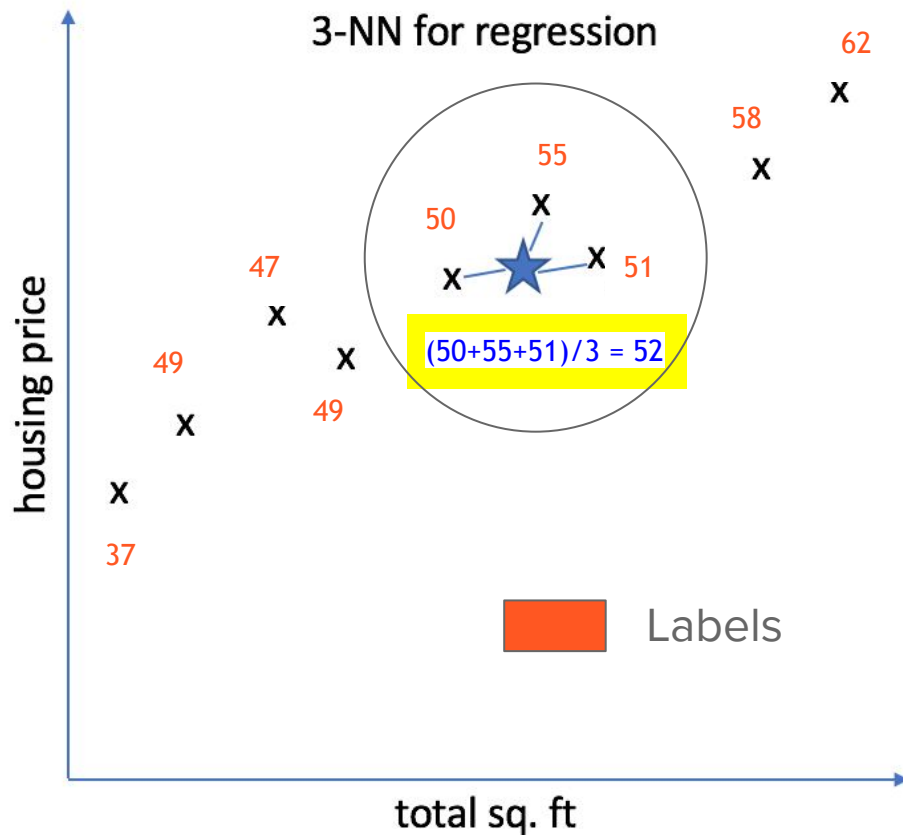# Summary of "Intro to Machine Learning"

# Non-parametric Method: K Nearest Neighbo ❓

### 3-NN for classification

o: 3 x: 0 → predict o

$x_2$

$x_1$

### 3-NN for regression

62

58

55

50

47

49

51

49

37

(50+55+51)/3 = 52
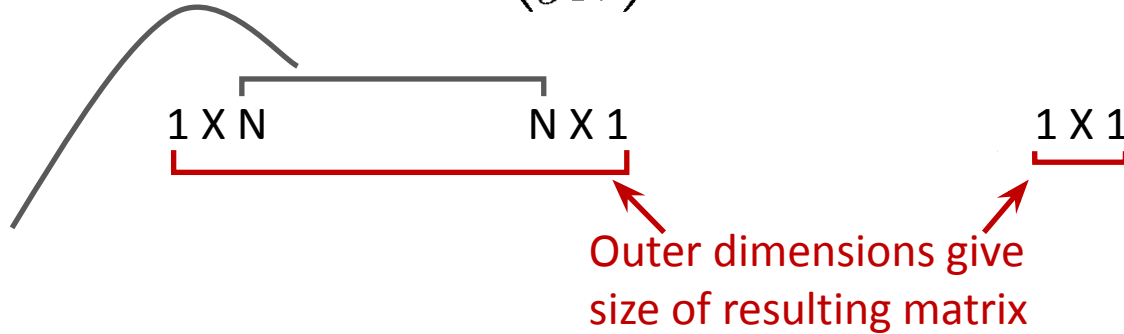
housing price

total sq. ft

Labels

# Multiplication:
## Dot product (inner product)

$$\vec{x} \cdot \vec{y} =$$

$$(x_1 \quad x_2 \quad \cdots \quad x_N) \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = x_1 y_1 + x_2 y_2 + \cdots + x_N y_N$$

1 X N      N X 1         1 X 1

Outer dimensions give
size of resulting matrix

$$\begin{pmatrix} r_1 & r_2 & \cdots & r_N \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{pmatrix} = r_1 w_1 + r_2 w_2 + \cdots + r_N w_N$$

Input neurons'
Firing rates

Synaptic weights

Output neuron's
firing rate

$r_1$

$r_2$

$r_i$

$r_n$

$w_1$

$w_2$

$w_i$

$w_n$

$$\vec{r} \cdot \vec{w} = |\vec{r}||\vec{w}| \cos(\theta)$$

Input neurons' Firing rates

$r_1$

$r_2$

$r_i$

$r_n$

$w_1$

$w_2$

$w_i$

$w_n$

Synaptic weights

Output neuron's firing rate

- <u>Insight</u>: for a given input (L2) magnitude, the response is maximized when the input is parallel to the weight vector
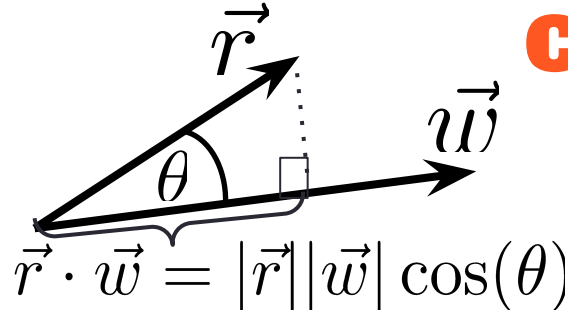- Receptive fields also can be thought of this way

# Neurons

Presynaptic cell

Postsynaptic cell

Synapses

Presynaptic terminal

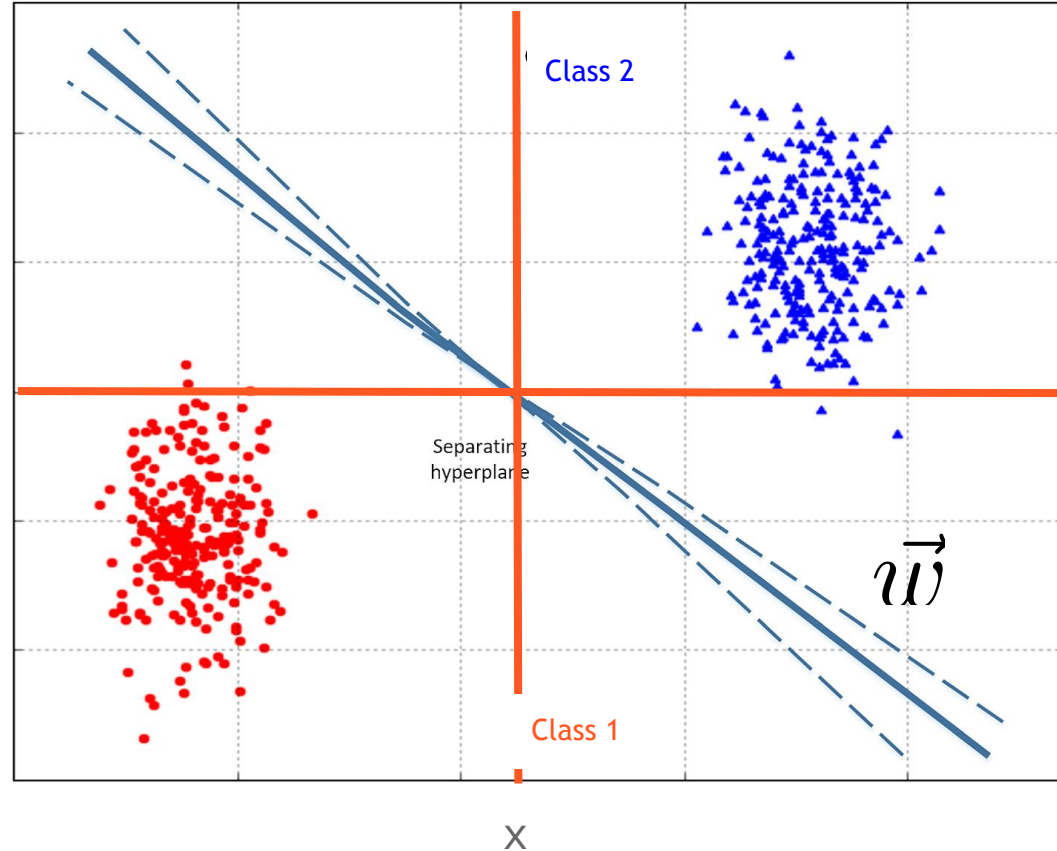Synapses refer to **the points of contact between neurons where information is passed from one neuron to the next**.
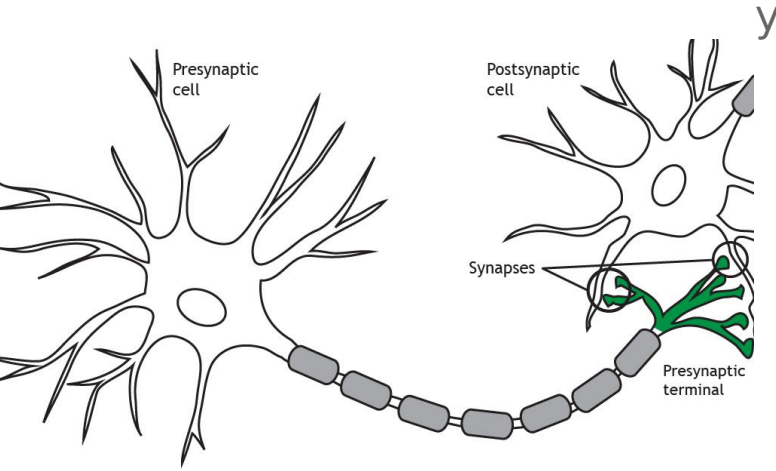
# Linear Models

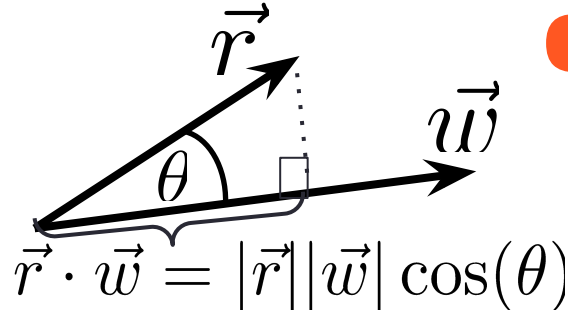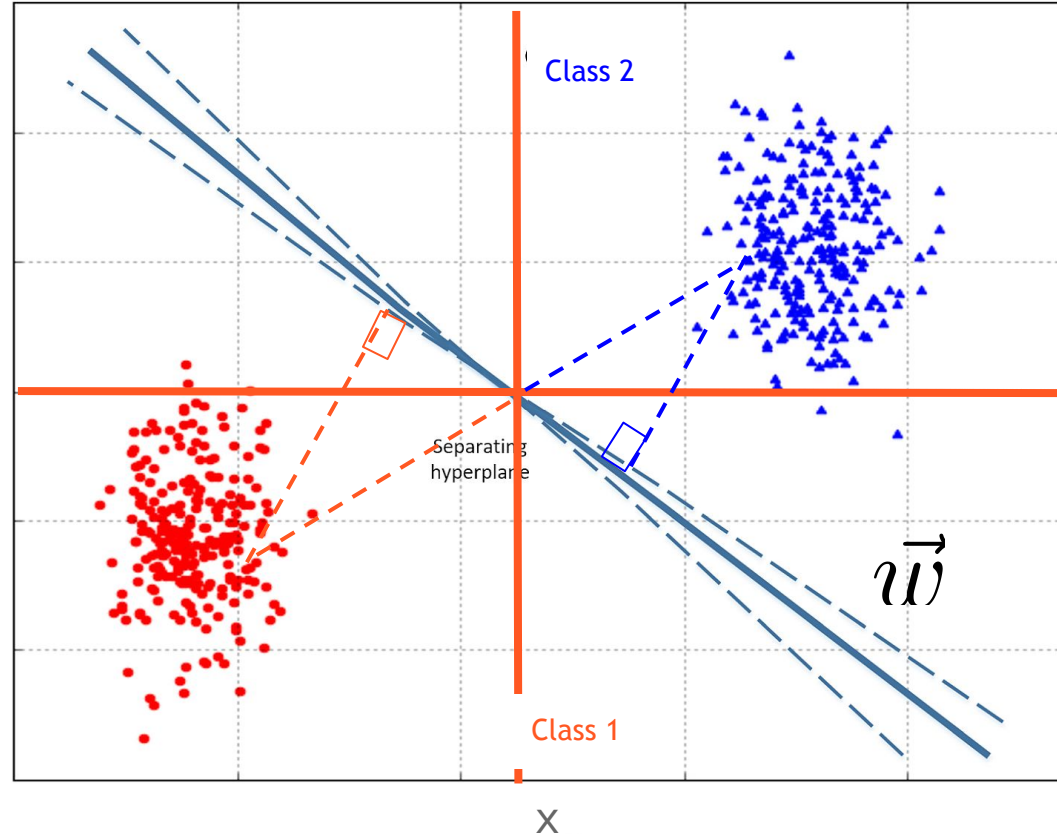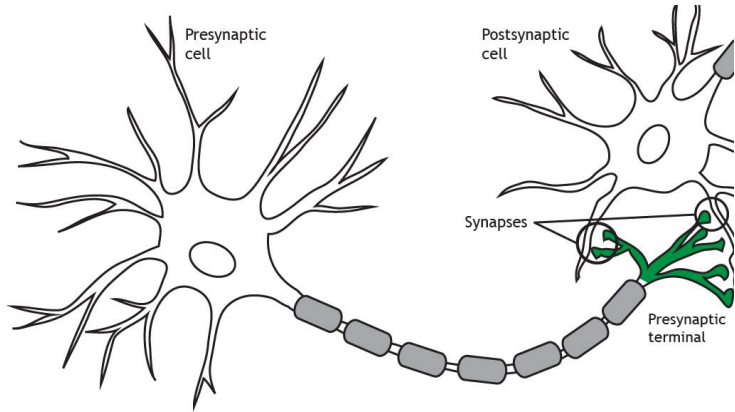# Parametric Method: Linear Models for Classification



$$\vec{r} \cdot \vec{w} = |\vec{r}||\vec{w}|\cos(\theta)$$

# Parametric Method: Linear Models for Classification


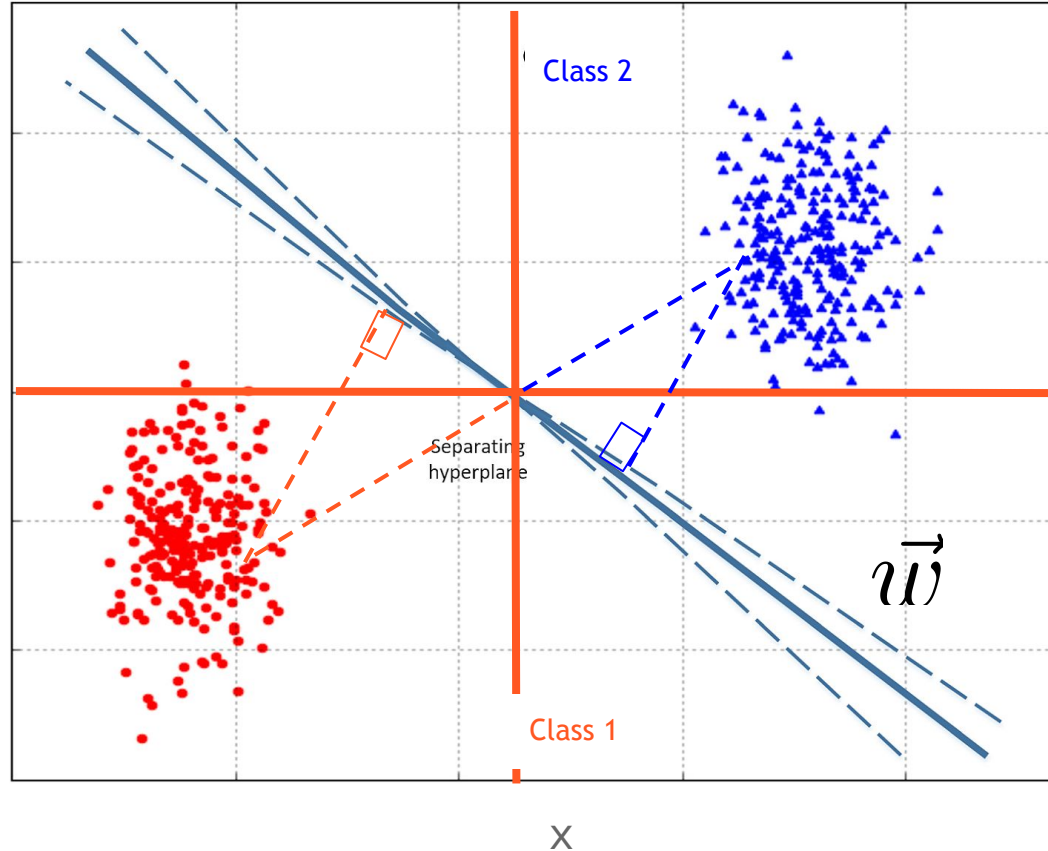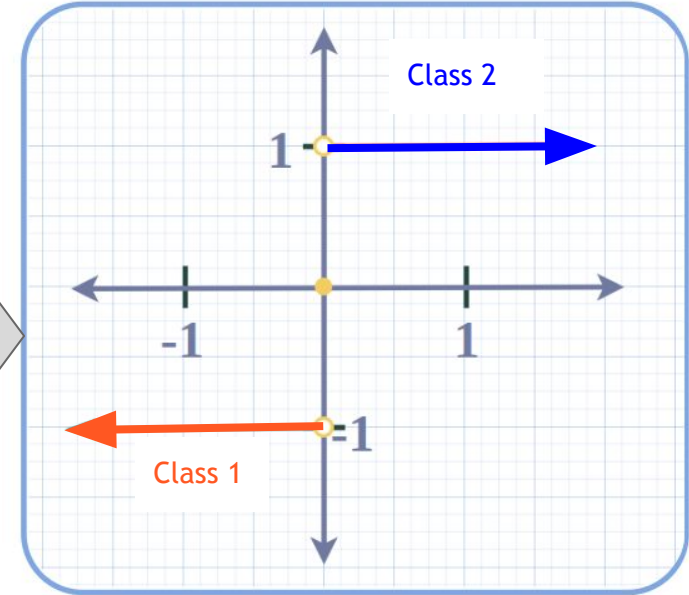
$$\vec{r} \cdot \vec{w} = |\vec{r}||\vec{w}| \cos(\theta)$$
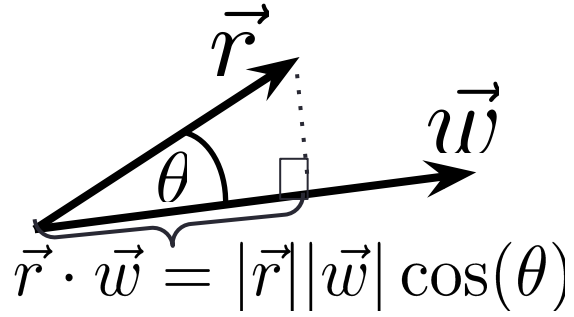
# Parametric Method: Linear Models for Classification

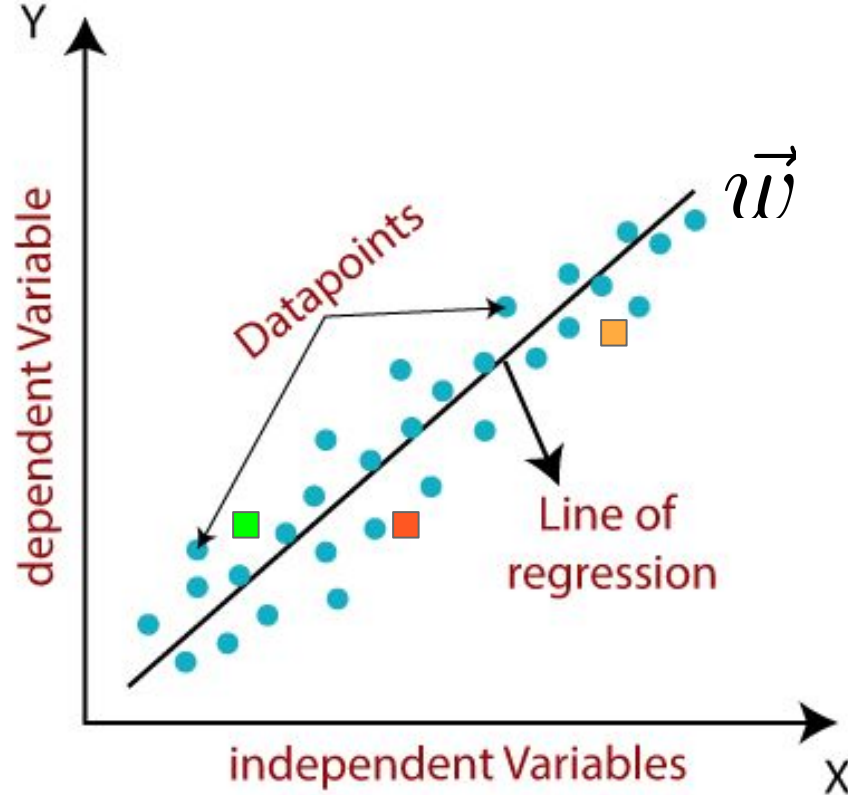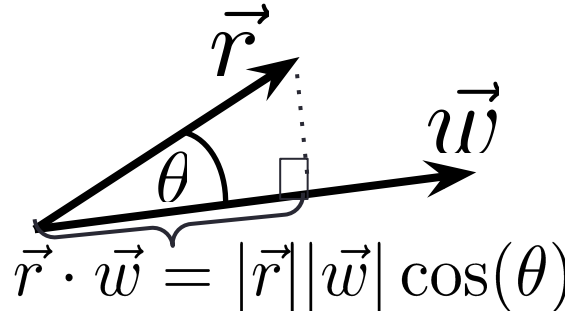# Parametric Method: Linear Models for Regression

$$\vec{r} \cdot \vec{w} = |\vec{r}||\vec{w}| \cos(\theta)$$

# Parametric Method: Linear Models for Regression

# Linear Models with Bias

## Simple Regression



W.x + b

Body Mass index — Y

Weight — X

# Recap: AI: How to Solve it?

1. Collect Labelled Dataset ✓
2. Design ANN Architecture ✓
3. Define Loss Function ✓
4. Optimize weights

4. Optimize Weight

784 input Neurons



1. Collect Labelled dataset

2. Design Artificial Neural Network

3. Loss Function

# Recap: Closed Form Expression

|  | Maximum | Minimum |
|---|---|---|
| Necessary condition | $\dfrac{dy}{dx} = 0$ | $\dfrac{dy}{dx} = 0$ |
| Sufficient condition | $\dfrac{dy}{dx} = 0$ ; $\dfrac{d^2y}{dx^2} < 0$ | $\dfrac{dy}{dx} = 0$; $\dfrac{d^2y}{dx^2} > 0$ |

$\dfrac{dy}{dx}$ is negative

$\dfrac{dy}{dx}$ is positive

$\dfrac{dy}{dx}$ is zero

Local maximum

A

C

Local minimum

B

# Recap: Gradient Descent

**Definition 1** *Suppose f is a real valued function and a is a point in its domain of definition. The derivative of f at a is defined by*

$$\lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

*provided this limit exists. Derivative of f(x) at a is denoted by f'(a).*



Loss: Denoted by f()

Absolute minimum

w

w

w

w

w w

● iter= 0  ● iter= 1  ● iter= 2  ● iter= 3  ● iter= 4  ● iter= 5

ANN weight: denoted by *a*

What happens to Loss *f(a+h) - f(a)* if the weight update *h* is -*λ* f'(a)?
Note: *λ* >0

Answer: f(a+h) - f(a) is negative, which means f(a+h) ≤ f(a).

Optimizing ANN: Update each ANN weight a as a - *λ* f'(a), where *λ* is the learning rate.

# Linear Regression

Labelled Data $\begin{bmatrix} x_{k1} \\ x_{k2} \\ \ldots \\ x_{kn} \end{bmatrix}$ $\implies$ $y_k$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \ldots & x_{1n} \\ 1 & x_{21} & x_{22} & \ldots & x_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{d1} & x_{d2} & \ldots & x_{dn} \end{bmatrix} \qquad y = \begin{bmatrix} y_1 \\ y_2 \\ \ldots \\ y_d \end{bmatrix}$$

$$w = (X^T X)^{-1} X^T y \qquad \nabla w = (X^T X)w - X^T y$$

$$L = \frac{1}{2} \sum_{n=1}^{N} \left( x_n^T w_0 + w_b - y_n \right)^2$$

where $w_0$ is the linear elements

$w_b$ is the bias.

$$= \frac{1}{2} \| X w - y \|^2$$

$$= \frac{1}{2} \left( w^T X^T X w - 2 w^T X^T y + y^T y \right)$$

* $\frac{\delta}{\delta w} \ w^T a = a$

* $\frac{\delta}{\delta w} \ w^T A w = 2 A w$ $\longrightarrow$ Exercise to be done.

$A$ is similar matrix

i.e; $A = A^T$

$$\frac{\delta L}{\delta w} = (X^T X) w - X^T y \quad \text{——①}$$

Closed form

$$X^T X w - X^T y = 0$$

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

Gradient descent

$$w_{i+1} = w_i - \lambda \frac{\delta L}{\delta w}$$

use ① for $\frac{\delta L}{\delta w}$

18/11/2024