

3/10/2023

Tuesday

## Lecture No. 2

→ Two Views of AI.

1. AI agents
2. AI tools.

→ An intelligent agent:

- Perception
- Robotics → performing Actions
- Languages
- Knowledge → Procedural Knowledge  
→ Declarative Knowledge
- Reasoning → Decision Making.
- Learning → OverTime Learning



Requires Thousands of  
examples to Train

Note: Few Short Learning  
↳ chat gpt.

→ AI Tools:

- Poverty Prediction
- Saving Energy by Cooling Data Centers.
- Authentication → Deep Face login

### Limitation:

- Bias in machine translation
  - ↳ Due to collection of Data from Society.

Society → Data → Predictions.

- fairness in criminal risk assessment  
↳ Algorithm will decide either to give bail or not on base of previous data.

Exp: Racism - Black & whites

Summary:

- AI Agents: Achieving human level intelligence, still very far.
- AI Tools: Need to think carefully about real world consequences.

How do we solve AI Tasks?

Paradigm

Modeling

(real world to mathematical object)

Inference

↓ (is a module that answers our questions based on model)

Model Learning

Model ~~without~~ with parameters +

Data

↓ (Learning

Model with parameters

Generalization: When algorithm predicts beyond the data set given to it.

Exp: Heart Patients.

Given Data Set must be rich  
also Model / Skeleton must be good

→ Reflex - Based Models

- Most common models in machine learning
- Fully feed-forward (No Back Tracking)

→ State - Based Models

- Exp: Chess, Pac-Man
- Robotics: Motion Planning.
- Types
  - Search Problems: (You control Everything)
  - Markov Decision Processes:  
Randomness ← Against Nature  
(e.g Blackjack)
  - Adversarial Games:  
Against Opponent (e.g Chess)

→ Variable-Based Models:

- Constraint Satisfaction Problems  
Hard Constraints e.g. Sudoku, Scheduling
- Bayesian Networks  
(Soft Dependencies)

→ Logic-Based Models.

Ex: Virtual Assistant.

6/10/2023

Friday

## Lecture No. 3

→ Reflex Models:

=> Machine Learning



- Linear Predictors
- Deep Learning

↳ Loss Minimizer  
↳ Stochastic Gradient Descent.

- Classification Problems → Discrete Output.  
e.g. Spam Emails
- Regression Problems.

↳ output continuous.

e.g.: Predicting Temp, house prices.

→ Spam Classification : (App)

Input →  $x = \text{email message}$

From:	_____
Subject:	_____
To:	_____
Date:	_____

Output:  $y \in \{\text{spam, Not spam}\}$

Objective: Obtain a predictor  $f$   
processes and then categorize, either spam or not (Example)

$$x \rightarrow [f] \rightarrow y$$

- We need a measure to classify either my objective is good or not.
  - Objective function could be a loss; how many mistakes were made, so in order to reduce it.
  - Function is basically a model/skeleton of a real world scenario.
- Binary Classification:

$$x \rightarrow [f] \rightarrow y \quad y \in \{+1, -1\}$$

↓  
Output Discrete

→ Regression (House Price Prediction)

$$x \rightarrow [f] \rightarrow y \in \mathbb{R}$$

↓  
Output Continuous.

→ Multi-class Classification:

Output is discrete but classes can be more than one.

$$x \rightarrow [f] \rightarrow y \in \{0, 1, 2, 3\}$$



Ex:

Healthy ← mandatory

covid

Lung Cancer

Dengue

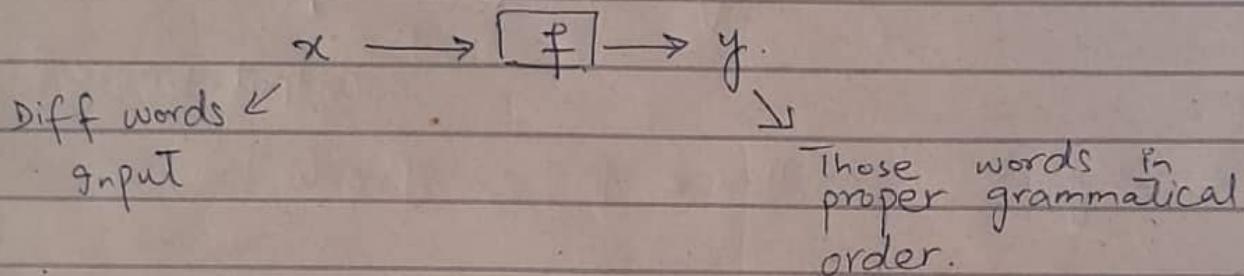
→ Ranking:

$$4, 2, 3, 1, 5 \rightarrow [f] \rightarrow 1, 2, 3, 4, 5$$

Example:

Teams Ranking                      on The  
University Ranking              Basis of certain  
    criteria.

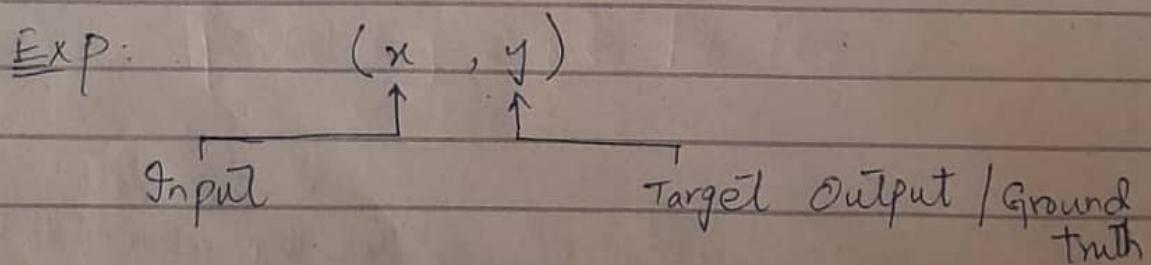
→ Structured Prediction:



• Data

Machine Learning

- Supervised Predictors → (Target output is provided)
- Unsupervised Predictors → (Target output is not provided).

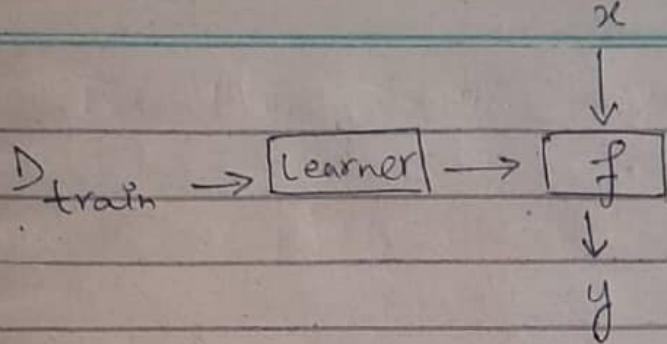


Data Set → List of Examples (Instances / Data point)

$$D_{\text{train}} = \{ ("---", +1), ("----", -1) \}$$

↓  
contains set of examples.

## → Loss Minimization Framework



// Based on data we want to construct function  $f$  - which requires some parameters. So in order to find some parameters for predictor functions we need Learner.

// Data is problem dependent.

// Function  $f$  depends upon features of Data. Function captures these features from Data and to enrich itself. we focus on relevant features.

### • feature Extraction:

Exp:  $x = "abc@hotmail.com"$

predict  $y$ : whether a string  $x$  is an email address

Feature Name	Feature value.
• Length > 10	1 or 0 True      False

• fraction of alphabets      0.85%

• contains @      1

• ends .com      1

$\Rightarrow \phi \rightarrow$  is a function that will extract features.

### Data Set

	length > 10	fraction of Alphabets	contains @	end.com
1.				
2.				
3.				

### Feature Vectors:

$$x = " \underline{\quad} " = \begin{bmatrix} 1 \\ 0.85 \\ 1 \\ 1 \\ 0 \end{bmatrix} \rightarrow \begin{array}{l} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_5(x) \end{array} \in \mathbb{R}^5$$

$$\phi(x) = [\phi_1(x), \phi_2(x), \phi_3(x), \dots, \phi_d(x)]$$

$$\phi(x) \in \mathbb{R}^d$$

### Weight Vector:

for each features  $j$ , we have a real number  $w_j$ . represents contribution of the feature  $j$ .

$$w = \begin{bmatrix} 1.2 \\ -0.8 \\ 3 \\ 1 \\ 0.1 \end{bmatrix} \in \mathbb{R}^d \quad d=5$$

weight of feature.

$\phi(x)$  depends on  $x$ .  
But  $w$  does not directly depend upon  $x$ .

10/10/2023

## Lecture No. 4

Tuesday

→ Reflex Models.

• Linear Predictions:

$$x \rightarrow [f] \rightarrow y$$

$$x \rightarrow \phi \rightarrow \phi(x) \in \mathbb{R}^d$$

F.E  $[\phi_1(x), \phi_2(x), \dots, \phi_d(x)]$

(feature Extractor).

$\phi_i(x) = (\text{feature Name}, \text{feature value})$

↓

is imp than  
feature name as  
it is required for  
computation.

→ Feature Vector:

$$\begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_d(x) \end{bmatrix}$$

→ Weight Vector:

$$w \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

Weights / Parameters

Assume

$$\phi(x) \in [0, 1]$$

$$y \in \begin{cases} +1 & \text{email Add.} \\ -1 & \text{Not email Add.} \end{cases}$$

Labels

⇒ Features are Binary

- If weight  $w_i$  is +ve,  $\phi_i(x) = 1$  (Positive prediction)  
⇒ +ve class label

- If weight  $w_i$  is -ve,  $\phi_i(x) = 0$   
⇒ -ve class label

### Scores

$$\text{Score} = w \cdot \phi(x)$$

$$= \sum_{i=1}^d w_i \cdot \phi_i(x)$$

$$w \quad \phi(x) = \begin{bmatrix} 1.2 \\ 0.6 \\ 3 \\ 2.2 \\ 1.4 \end{bmatrix}, \quad \phi(x) = \begin{bmatrix} 1 \\ 0.85 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{Score} = (1 \cdot 2)(1) + (0 \cdot 6)(0 \cdot 85) + 3(1) + (2 \cdot 2)(1) + (1 \cdot 4)(0)$$

$$\boxed{= 4.5}$$

- +ve Value of Score Represent +ve Classification (+1). (+1)
- Greater the value The more confident we are that it is positive

→ Case of Binary Features:

Let  $w = \begin{bmatrix} 1 \cdot 2 \\ 3 \\ -2 \cdot 2 \\ 1 \cdot 4 \end{bmatrix}$        $\phi(x) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

No contribution ←

+ve feature with +ve weight → +ve classification

+ve feature with -ve weight → -ve classification

→ Linear Predictor:

$$\phi(x) \in \mathbb{R}^d, w \in \mathbb{R}^d$$

• Linear (Binary) Classification.

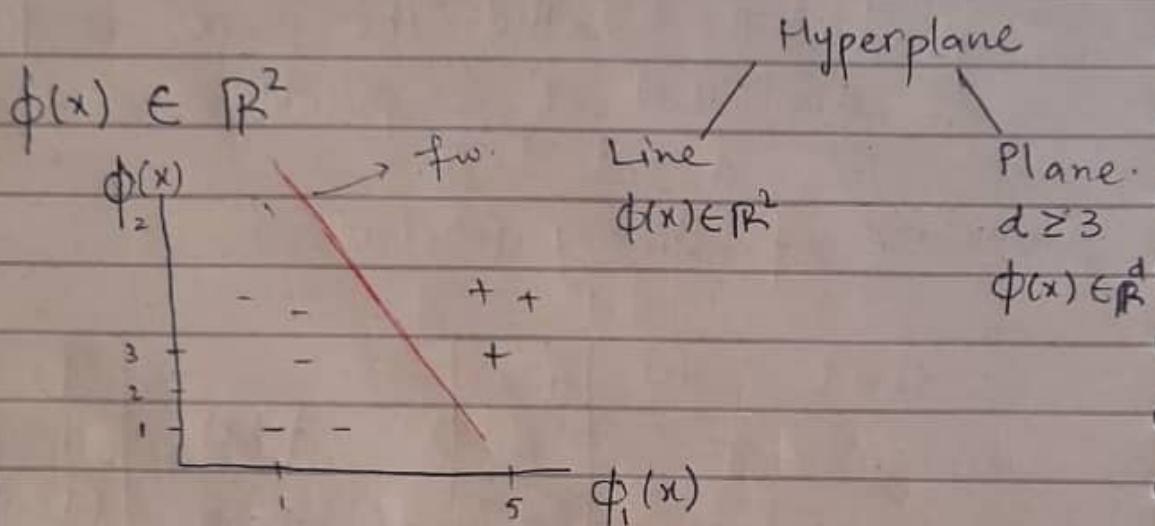
$f_w(x) = \text{sign}(w \cdot \phi(x))$

parameter  $w$       string input

$$= \begin{cases} +1 & \text{if } w \cdot \phi(x) > 0 \\ -1 & \text{if } w \cdot \phi(x) < 0 \\ ? & \text{otherwise, (un decisive)} \end{cases}$$

## Geometric Interpretation:

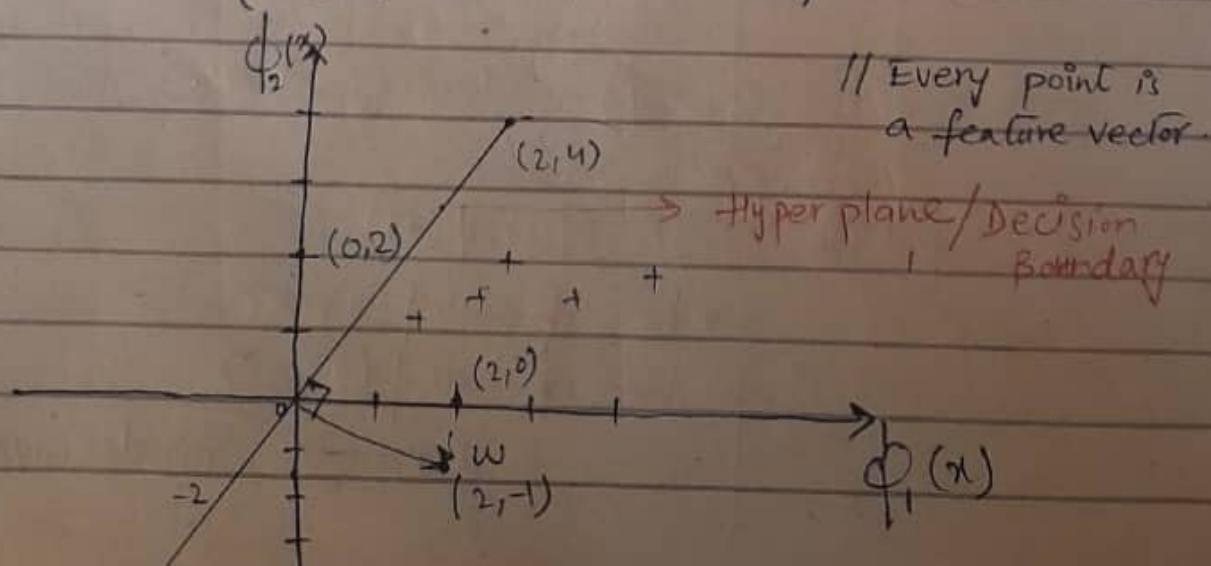
- $f_w$  is a binary classifier
- It defines a hyperplane with normal vector  $w$ .



$$w = [2, -1]$$

$$\phi(x) \in \{[2, 0], [0, 2], [2, 4]\}$$

$$D_{\text{train}} = \{(x_1, y_1), (x_2, y_2), \dots\}$$

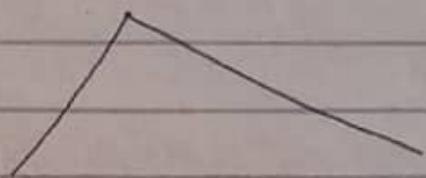


$$w \cdot \phi(x) = 4+0 = 4$$

$$w \cdot \phi(x) = 0-2 = -2$$

$$w \cdot \phi(x) = 4-4 = 0$$

→ Loss Minimization Frame Work.



Optimization  
Problem

Model { what we want to  
compute

Optimization Algorithm  
(How to compute?  
(Learning)).

ground truth  
Loss ( $x, y, w$ )

$$y' = f_w(x)$$

Loss is function  
in which we give  
value, target value

Classification → miss classification matters (loss)  
Regression → diff matters (loss).  
exp: house price, Age.

17/10/2023

Tuesday

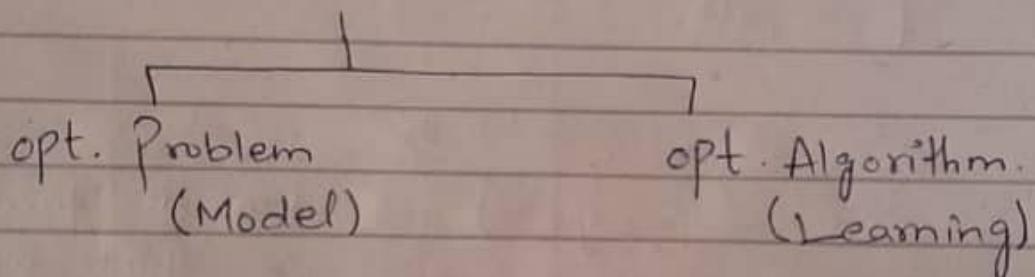
## Lecture No. 5

→ Reflex Models: (cont.)

$$f \leftarrow \phi(x) \\ w$$

$$\text{Score} = w \cdot \phi(x)$$

→ Loss Minimization Frame Work.



• Loss Function:      ground truth  
                          ↓  
 $\text{Loss}(x, y, w) = ?$

• Margin:

$$\begin{aligned} &= \text{Score} * y \\ &= w \cdot \phi(x) * y \quad \text{Margin} \\ &\quad \begin{array}{ll} +ve & +ve \} \rightarrow +ve (+1) \rightarrow \text{classification.} \\ -ve & -ve \} \rightarrow +ve (-1) \end{array} \end{aligned}$$

$$\begin{array}{ll} -ve & +ve \} \rightarrow -ve \\ +ve & -ve \} \rightarrow -ve \end{array}$$

$\Rightarrow$  if margin  $> 0$  Then our prediction is correct.

$\Rightarrow$  if margin  $< 0$  Then our prediction is wrong.

Example :

$$w = [2, -1]$$

$$x = \phi(x) = [2, 0], y = -1$$

$$\begin{aligned}y' &= f_w(x) = \text{sign}(w \cdot \phi(x)) \\&= \text{sign}(4 - 0) \\&= +\text{ve} \Rightarrow +1\end{aligned}$$

$$\begin{aligned}\Rightarrow \text{Margin} &= w \cdot \phi(x) * y \\&= 4 * (-1) = -4.\end{aligned}$$

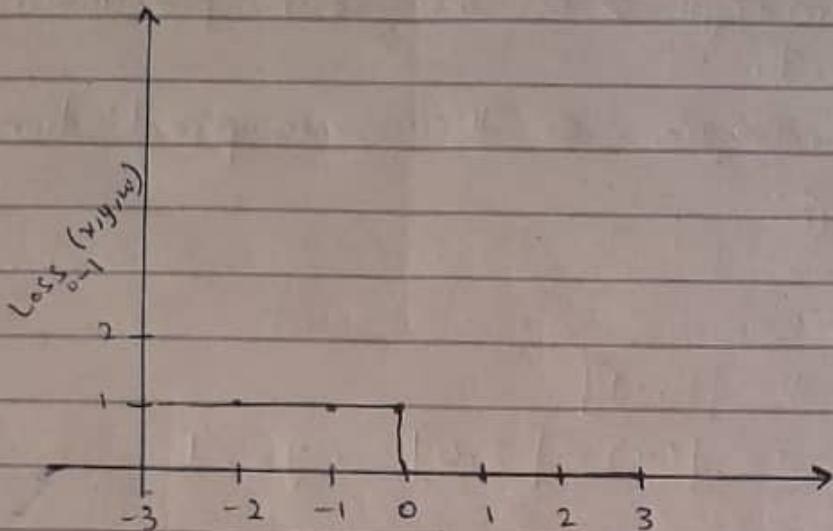
• Loss Function : (Classification Problems).

0-1 Loss

$$\text{Loss}_{0-1}(x, y, w) = 1 [f_w(x) + y]$$

$$= 1 [(w \cdot \phi(x)) * y \leq 0]$$

Loss is 1 when prediction is wrong i.e Margin is less than equal to 0. otherwise loss is 0.



$$\text{Margin} = (w - \phi(x)) * y.$$

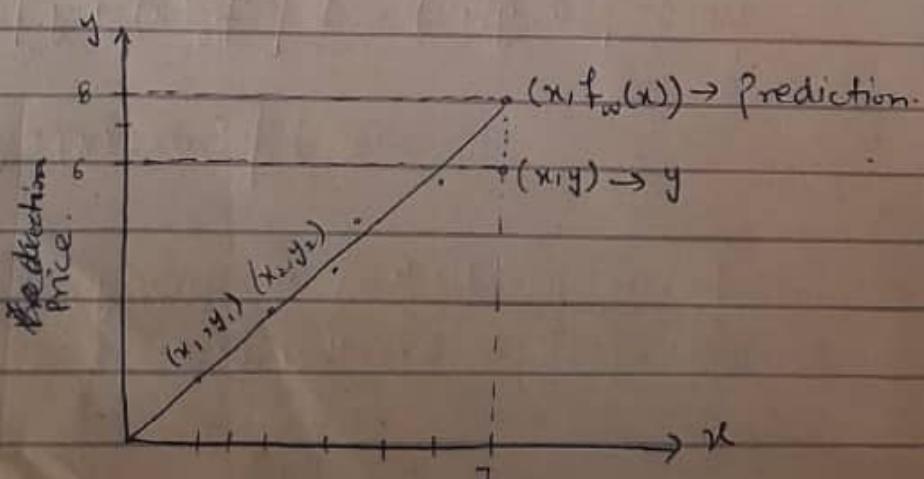
→ Linear Regression:

$\hat{y}' = y$  can't be achieved  
bcz  $y \in \mathbb{R}$ .

⇒ So 0-1 loss can't be used in  
Regression Problems.

⇒ Can be used in Classification Problems.

$\hat{y}' = y$  can be achieved.  $y \in \text{Discrete}$ .



Size of house

$\Rightarrow$  Residual:

Diff b/w our prediction and ground truth.

$$\text{Residual} = (\omega \cdot \phi(x) - y) = (f_\omega(x) - y)$$

Regression Loss Function:

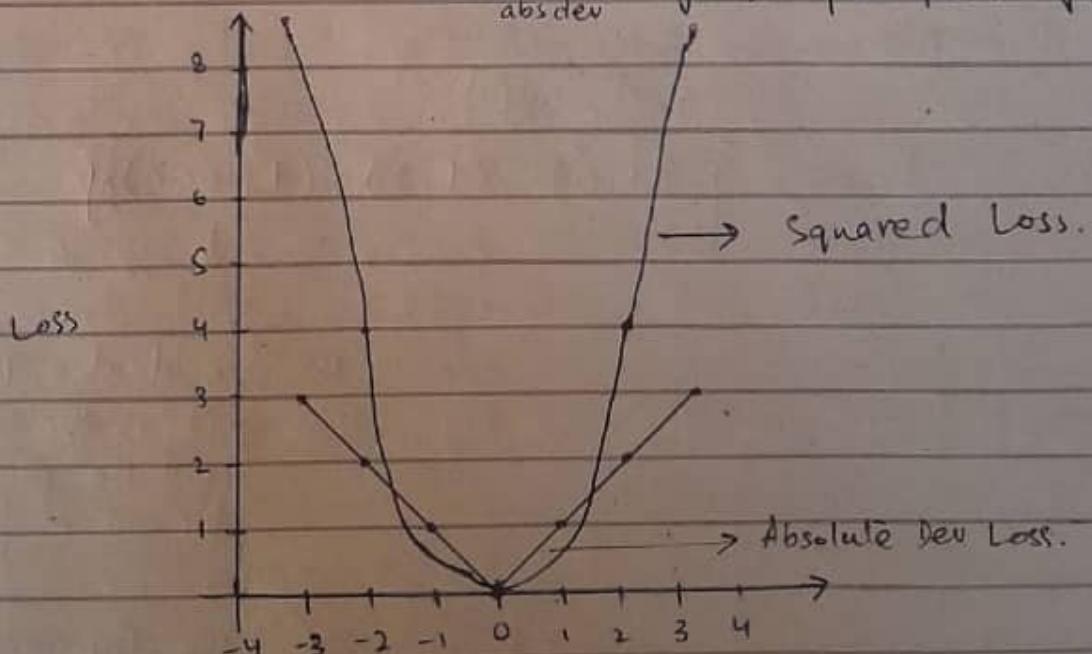
$$f_\omega(x) = \omega \cdot \phi(x).$$

- Squared Loss:

$$\text{Loss}_{\text{sq}}(x, y, \omega) = (\omega \cdot \phi(x) - y)^2$$

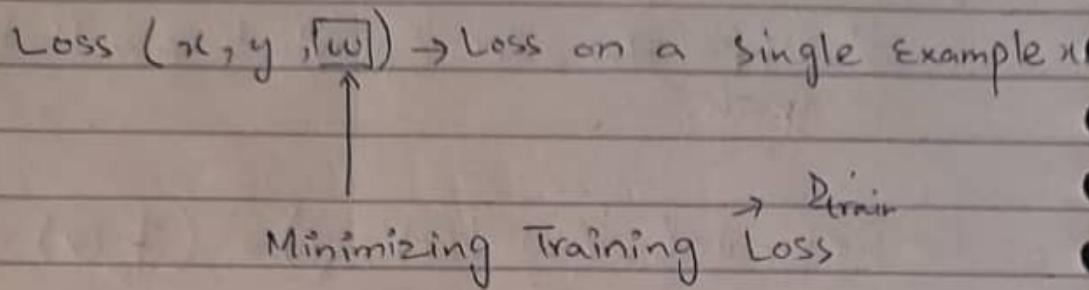
- Absolute Deviation Loss:

$$\text{Loss}_{\text{absdev}}(x, y, \omega) = |\omega \cdot \phi(x) - y|$$



$$\text{Residual} = \omega \cdot \phi(x) - y.$$

## $\Rightarrow$ Loss Minimization Framework.



Train Loss ( $w$ )

$$= \frac{1}{|D_{\text{train}}|} \sum_{(x, y) \in D_{\text{train}}} \text{Loss}(x, y, w)$$

No. of Training Examples.

Minimizing Train Loss ( $w$ ) By  
modifying  $w \in \mathbb{R}^d$

Example:

Input Output.

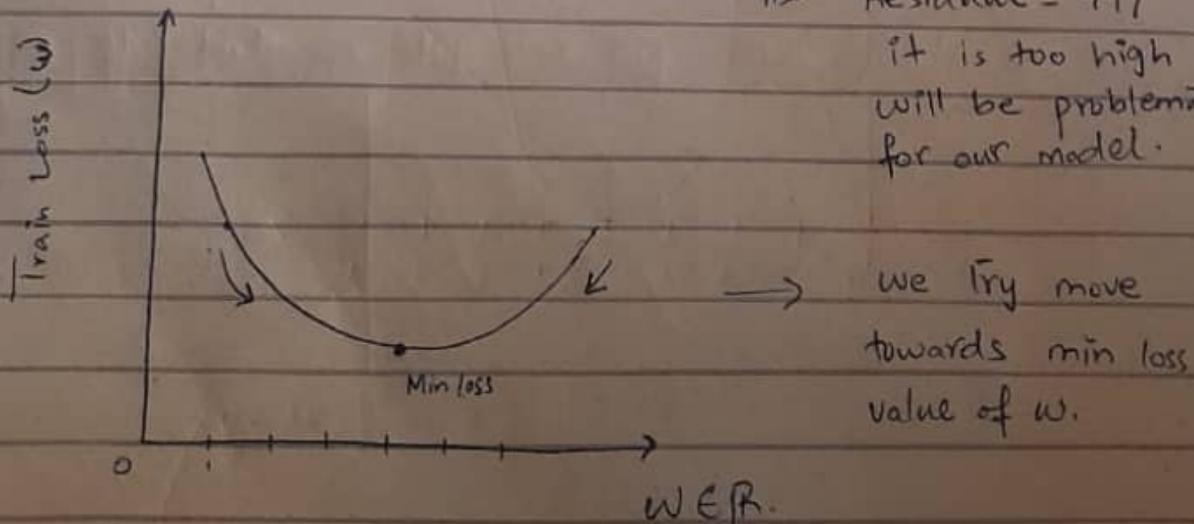
$$D_{\text{train}} = \{(1, 0), (1, 2), \underline{(1, 1000)}\}$$

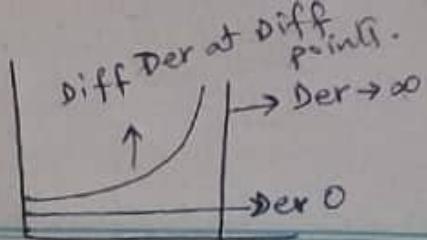
↓

Outlier

Assume Residual = 997

it is too high  
will be problematic  
for our model.





• Gradient / Derivative.

Rate of Change.

$\nabla$  Train loss ( $w$ )

tells us the direction that increases the loss the most.

→ Gradient Descent : (opt. Algorithm).

Objective: Minimal Train Loss ( $w$ ).

Initialize  $w = [0, 0, \dots, 0]$

For  $t = 1, 2, \dots, T$

We are  
iteratively changing  
the value of  $w$ .

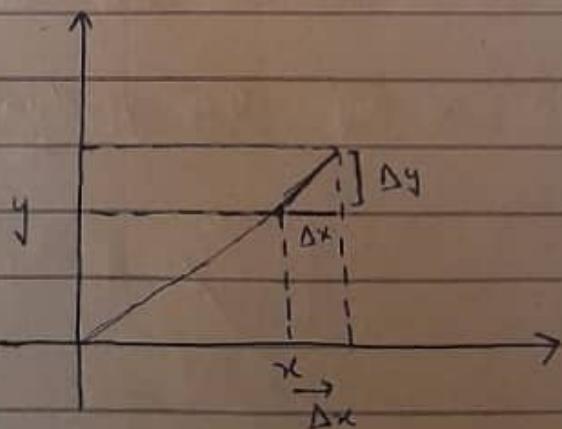
$$w = w - \eta \nabla \text{Train Loss}(w)$$

$\downarrow t \rightarrow$  Step size  $\leq 1$

$T \rightarrow$  no. of iteration

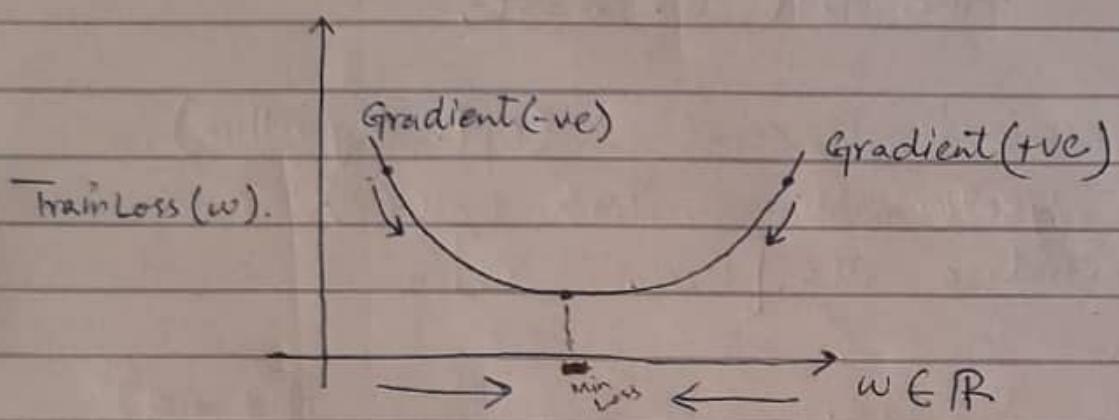
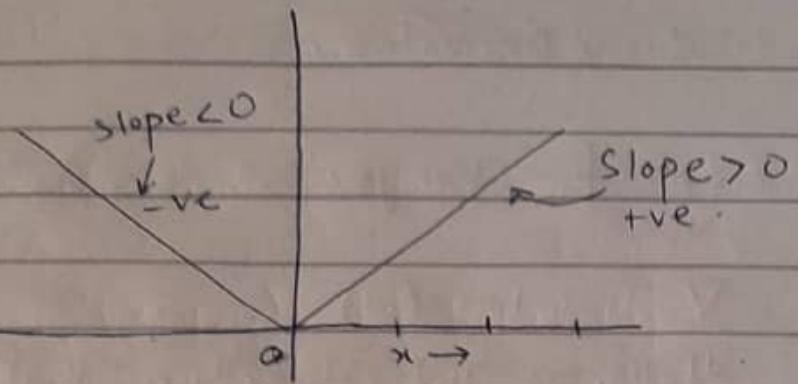
$\frac{\partial}{\partial w} \rightarrow i$

$\eta \rightarrow \text{eta (step size)}$

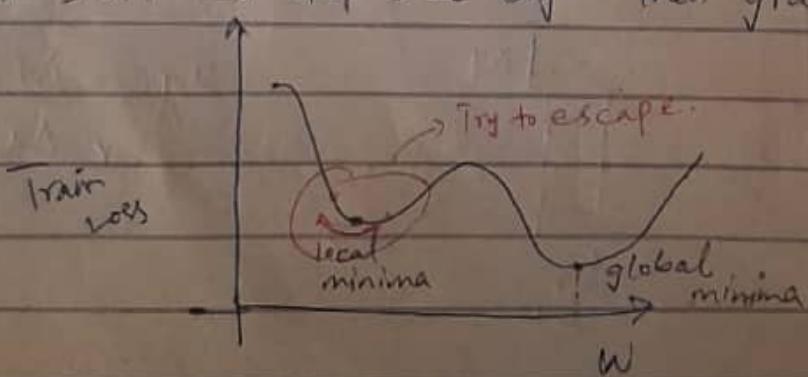


$$\text{slope} f = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1}$$

or



- Small values of step size will make our Algo slow.
- Sometimes we take big steps to escape local minima. After escaping local minima Then we take small steps to reach Global Minima.
- Big steps sometimes cause back & forth movement and we can't achieve our Minima Point.
- At start → step size big Then gradually decrease to small step size



## Lecture No. 6.

→ Gradient Descent : (Cont.)

$$\text{Train Loss}(w) = \frac{1}{|D_{\text{train}}|} \sum_{i=1}^{|D_{\text{train}}|} (w \cdot \phi(x_i) - y_i)^2$$

$$\nabla \text{Train Loss}(w) = \frac{1}{|D_{\text{train}}|} \sum_{i=1}^{|D_{\text{train}}|} (w \cdot \phi(x_i) - y_i) \cdot \phi(x_i)$$

$\because \frac{d}{dx} x^2 = 2x^2 \cdot \frac{d}{dx} (x)$

Here, we took derivative of Train Loss to Get Gradient which gives us direction.

→ Batch Gradient Descent:

↳ Slow.

↳ w will update once for all training data examples

→ Stochastic Gradient Descent:

$$w = [0, 0, 0, \dots]$$

for  $t = 1, \dots, T$

For each  $(x, y) \in D_{\text{train}}$

$$w = w - \eta \nabla_w \text{Loss}(x, y, w)$$

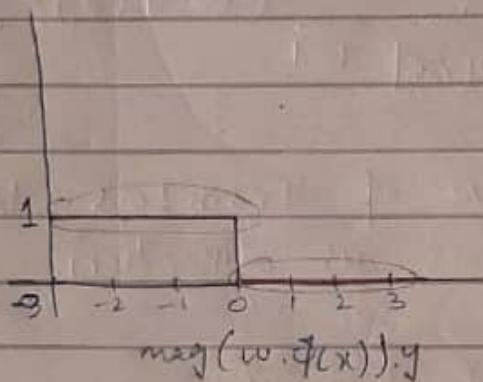
↳ Fast.

↳ Getting Data on the fly / real time.

For Classification Problems.

$$\text{Loss}_{0-1}(x, y, w) = 1 [f_w(x) \neq y]$$

$$= 1 [(\omega \cdot \phi(x))y \leq 0]$$



$$w = w - \eta \nabla \text{loss}(x, y, w)$$

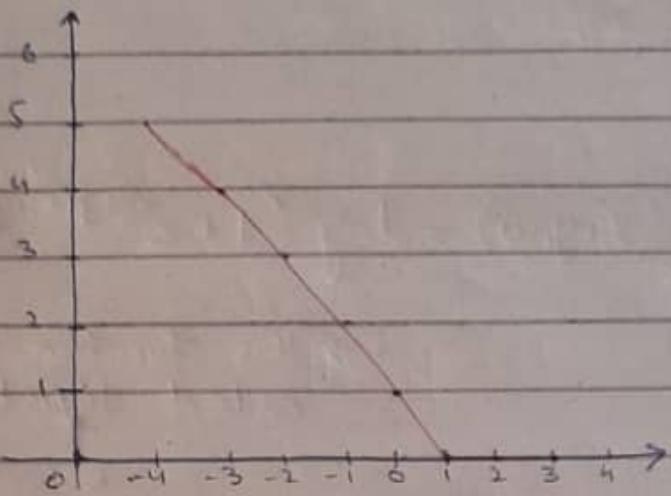
We cannot use  $(0-1)$  Loss for classification Problems in Gradient Descent Algo because  $w$  will not update.

↪ straight line  $\nabla$  is zero. No change.

So we will use another Method.

→ Hinge Loss:

$$\text{Loss}_{\text{Hinge}}(x, y, w) = \max [1 - (\omega \cdot \phi(x))y, 0]$$



So, Hinge loss  
can be used in  
Gradient Descent  
for instead of

$\text{Loss}_{0-1}$

$\text{margin}(\mathbf{w} \cdot \phi(\mathbf{x}))y$ .

$$\mathbf{w} = \mathbf{w} - \eta \nabla_{\mathbf{w}} \text{loss}(\mathbf{x}, y, \mathbf{w})$$

$$\text{Loss}_{\text{Hinge}}(\mathbf{x}, y, \mathbf{w}) = \max [1 - (\mathbf{w} \cdot \phi(\mathbf{x}))y, 0] = \begin{cases} 0 & (\mathbf{w} \cdot \phi(\mathbf{x}))y \geq 1 \\ 1 - (\mathbf{w} \cdot \phi(\mathbf{x}))y & \text{otherwise} \end{cases}$$

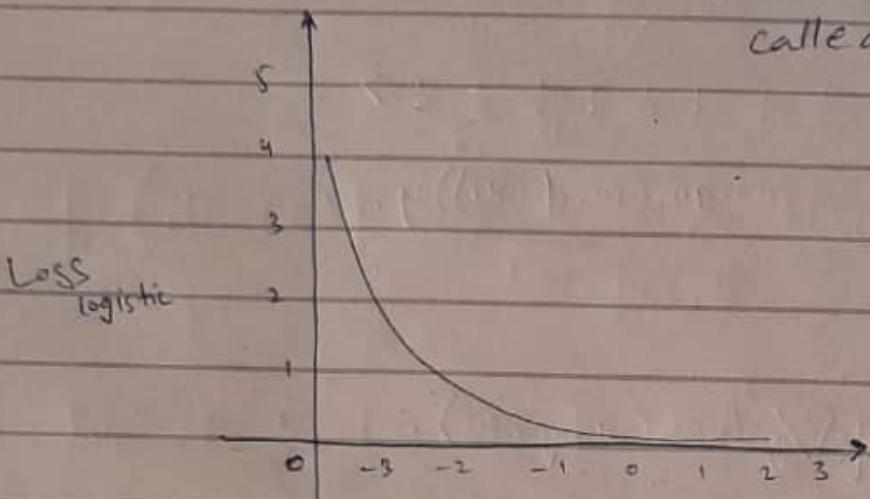
$$= \begin{cases} 0 & (\mathbf{w} \cdot \phi(\mathbf{x}))y \geq 1 \\ -\phi(\mathbf{x})y & \text{otherwise} \end{cases}$$

$$\nabla_{\mathbf{w}} \text{Loss}_{\text{Hinge}}(\mathbf{x}, y, \mathbf{w}) = \begin{cases} 0 & (\mathbf{w} \cdot \phi(\mathbf{x}))y \geq 1 \\ -\phi(\mathbf{x})y & \text{otherwise} \end{cases}$$

→ Logistic Loss:

$$\text{Loss}_{\text{logistic}}(x, y, w) = \log(1 + e^{-(w \cdot \phi(x) \cdot y)})$$

↳ This function is also called logistic Regression.



margin ( $w \cdot \phi(x)$ ) $y$

- Ground truth ( $y$ ) → will remain same.  
So it means Score ( $w \cdot \phi(x)$ ) will change  
because of  $w$ .

~~so loss will not become~~

24/10/2023

Tuesday.

## Lecture No. 7.

→ Python:

// google.colab.com.

// python.org.

→ Operators

→ Variables (we do not declare variables).

→ Range → for i in range(1, 10, 2) → [1, 3, 5, 7, 9]

→ list indexing

↓            ↓              ↓  
Starting    Ending      step size  
(not included)

→ Slicing.

→ Tuple → (1, 2, 3) + (4, 5, 6) = (1, 2, 3, 4, 5, 6)

→ The + Operator

→ concatenation

→ The \* Operator

→ Dictionaries

→ Functions.

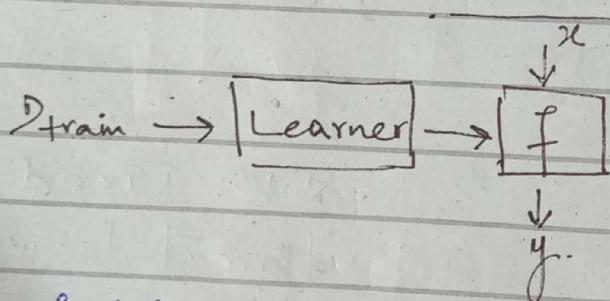
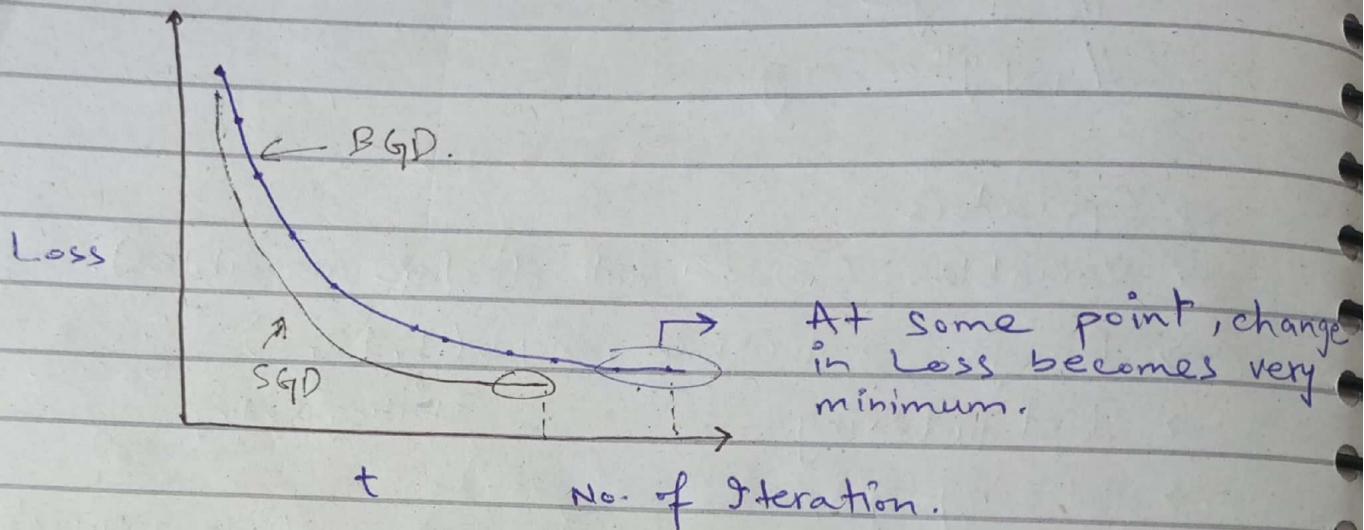
Note: Fixed  
format language.  
Indentations imp

27/10/2023.

Friday.

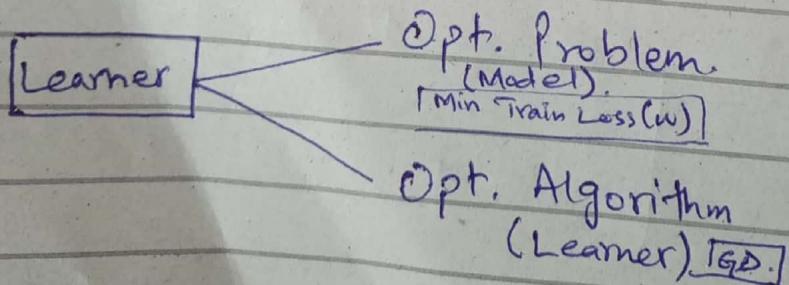
## Lecture No. 8

- Batch Gradient Descent.
- Stochastic Gradient Descent.



$$f_w(x) = w \cdot \phi(x) \quad \leftarrow \text{Regression}$$

$$f_w(x) = \text{sign}(w \cdot \phi(x)) \quad \leftarrow \text{Classification.}$$



## → Regression Example:

Data:

$\begin{bmatrix} \phi(x_1) & \phi(x_2) \\ 1 & 0 \end{bmatrix}$	$y$	Sq. Loss.
$\begin{bmatrix} 1 & 0 \end{bmatrix}$	2	$((w_1 \cdot 1 + w_2 \cdot 0) - 2)^2 = (w_1 - 2)^2$
$\begin{bmatrix} 1 & 0 \end{bmatrix}$	4	$(w_1 - 4)^2$
$\begin{bmatrix} 0 & 1 \end{bmatrix}$	-1	$(w_2 - (-1))^2$
$\begin{bmatrix} 0 & 1 \end{bmatrix}$	x	

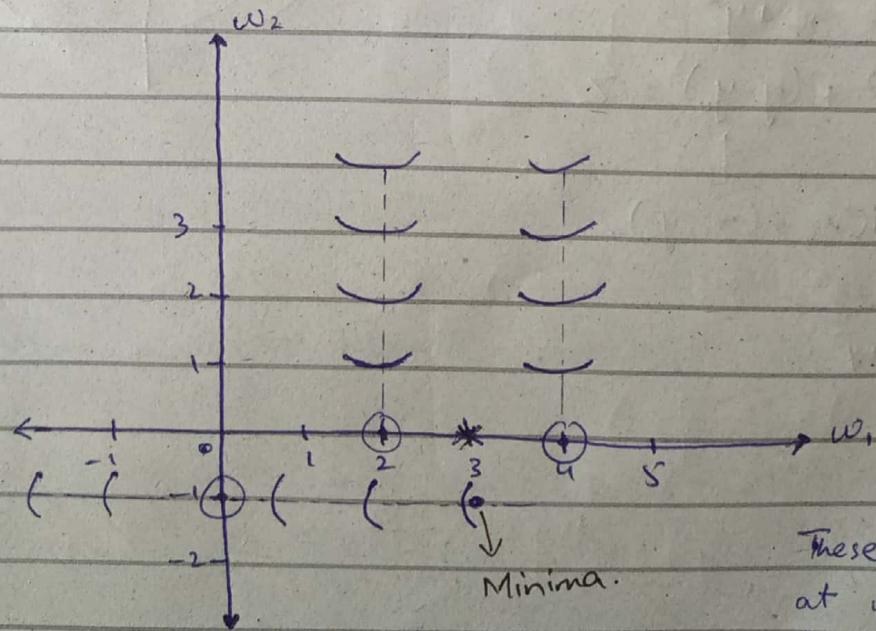
→ Sq. Loss:

$$w = [w_1, w_2]$$

$$\text{Loss}(x, y, w) = (w \cdot \phi(x) - y)^2$$

→ Train Loss:

$$\text{Train Loss}(w) = \frac{1}{3} [(w_1 - 2)^2 + (w_1 - 4)^2 + (w_2 - (-1))^2]$$

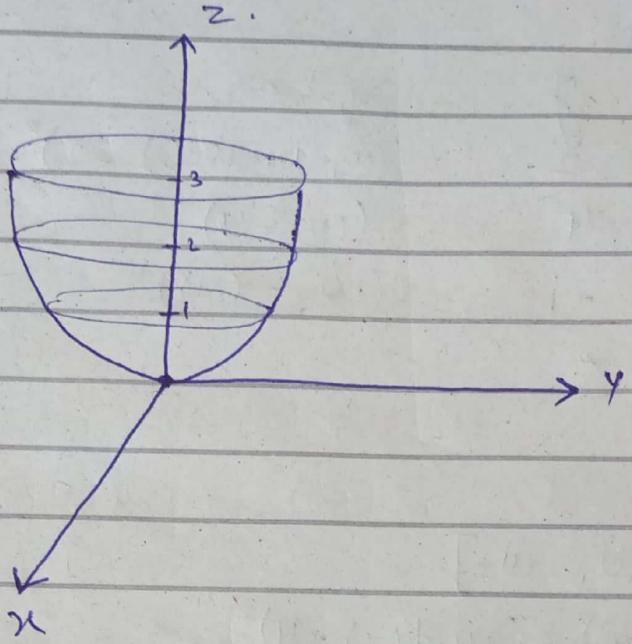


Minima.

These are weights  
at which loss  
is zero for each  
corresponding example.

↗ Level curves?

→ Level Curves:



$$f(x, y) = x^2 + y^2$$

$$f(x, y) = z.$$

$$z = x^2 + y^2$$

$$0 = x^2 + y^2$$

$$(x - x_0)^2 + (y - y_0)^2 = z$$

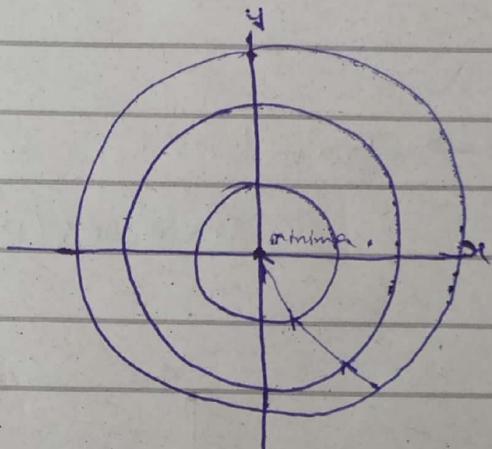
↑

↑

$$\text{Center} = (x_0, y_0),$$

$$1 = x^2 + y^2$$

$$2 = x^2 + y^2.$$



→ Features :  $\phi(x)$ .

- Critical Part of Machine Learning.
- With right features, we can improve the Algo.

$$f_w(x) = w \cdot \phi(x) \rightarrow \text{Score.}$$

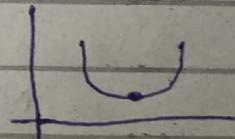
Q: Can Linear Predictors give us a non-Linear Decision Boundary? (Y/N).

Relationship b/w  $w$  and Score → Linear  
Score is <sup>linear in</sup>  $w$  → Linear  
" "  $\phi(x)$  → Linear.

Input is it Linear?

No, Input ( $x$ ) is not linear.

$f_w(x)$ , Loss ← Non-linear



• Feature Template:

It is a group of features all computed in a similar.

Example:

- Length greater than \_\_\_\_\_
- Last three character equals to \_\_\_\_\_
- Contains character \_\_\_\_\_

## • Sparsity in Features: (problem).

abc @ hotmail.com.

ends with <u>aaa</u>	0
" " ab	0
" " ac	0
⋮ ⋮ ⋮	⋮ ⋮ ⋮
" " com	1
⋮ ⋮ ⋮	⋮ ⋮ ⋮
	0 0

One-Hot Representation →

$m \rightarrow$  no. of char

$m^3 \rightarrow$  Combinations

Representations:

- Array  $[0, 0, 0, \dots, 1, 0, 0, \dots]^m$

- Map { ends with .com : 1, Alphabet Alpha = 0.85 }

↓ Save Memory.

• Hypothesis Class: (HC)

$$f_w(x) = w \cdot \phi(x)$$

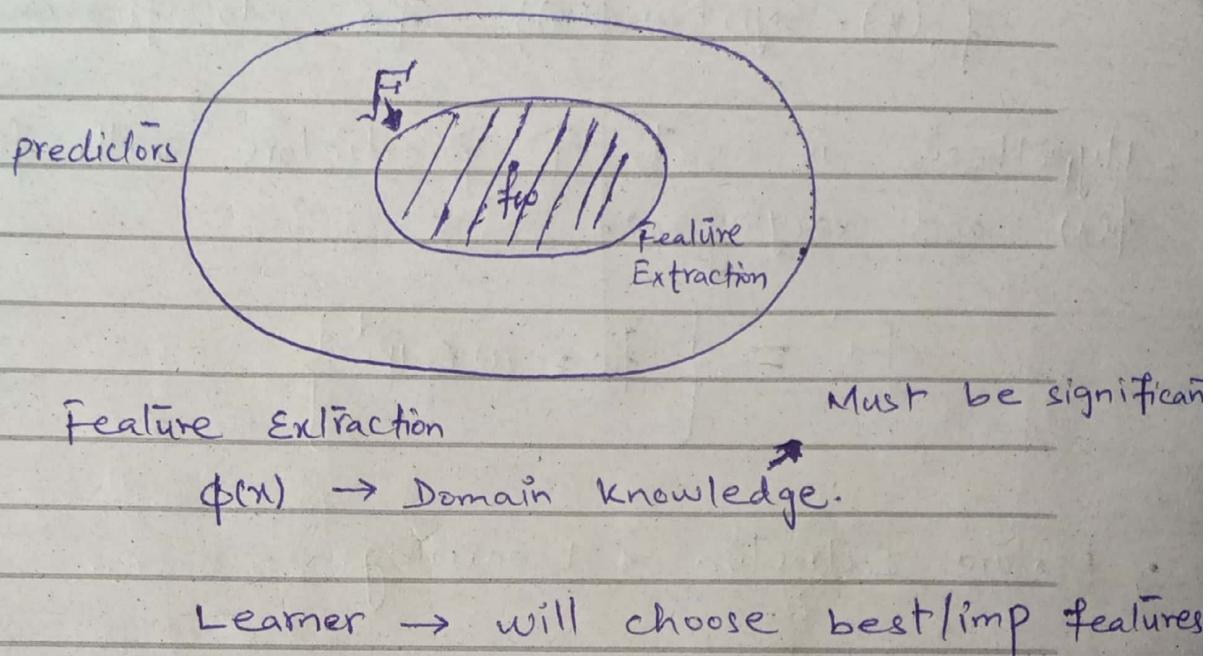
$$f_w(x) = \text{sign}(w \cdot \phi(x))$$

→ We study the impact of Features  
on prediction.

HC :  $\mathcal{F}$  is a set of possible predictors with a fixed  $\phi(x)$  and varying  $w$ .

$$\mathcal{F} = \{ f_w : w \in \mathbb{R}^d \}$$

Exp



F.E  $\rightarrow \mathcal{F}$  is based on Domain Knowledge

Learning  $\rightarrow f_w \in \mathcal{F} \rightarrow$  Based on Data

31/10/2023.

Tuesday

## Lecture No. 9.

→ Hypothesis Class:

Predictors:

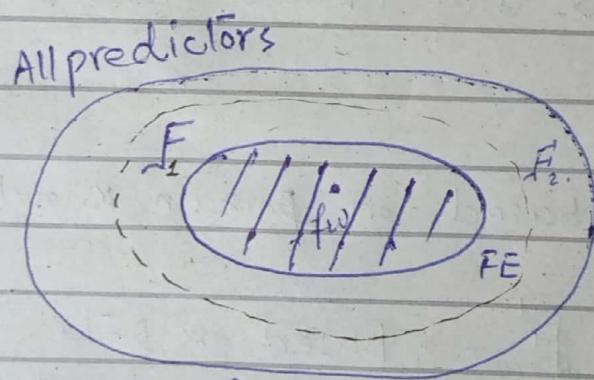
$$f_w(x) = w \cdot \phi(x) \text{ for Regression.}$$

$$f_w(x) = \text{sign}(w \cdot \phi(x)) \text{ for classification.}$$

- Hypothesis is set of predictors with fixed  $\phi(x)$  and varying  $w$ .

$$F = \left\{ f_w : w \in \mathbb{R}^d \right\}$$

- Feature Extraction + Learning:



FE → specify The hypothesis class.

Weights are Learned through data.

- features are very important, the more good features will result in good predictors.
- Learning also depends upon Data.

↳ Data Diversity

Example:

Regression:  $x \in \mathbb{R}$ ,  $y \in \mathbb{R}$

Let's Assume

$$\phi(x) = x.$$

// we are assuming  
input is feature.  
Not extracting here.

$$F_1 = \left\{ w_1 x + w_2 x^2 : w_1 \in \mathbb{R}, w_2 = 0 \right\}$$

↳ Hypothesis class representing Linear Predictors

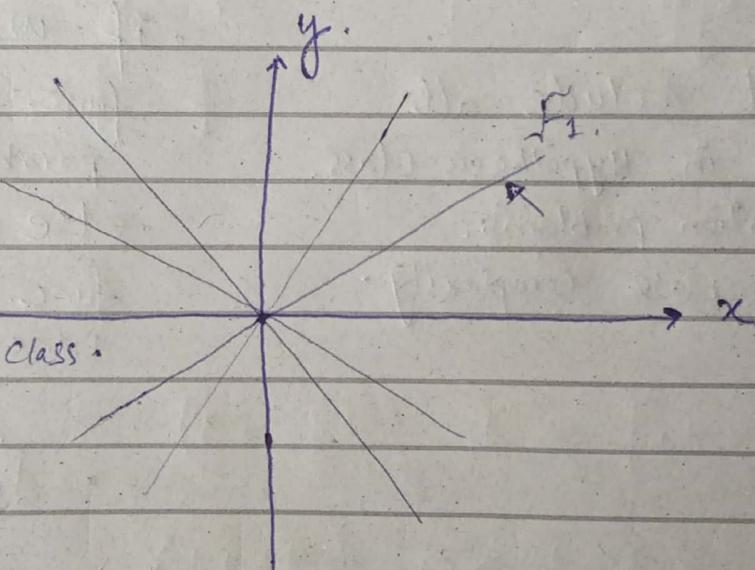
Assumed as zero for linear predictors  
To remove non-linear

$$\therefore y = mx + c$$

$\uparrow$  slope  $\uparrow$   $y$ -  
intercept

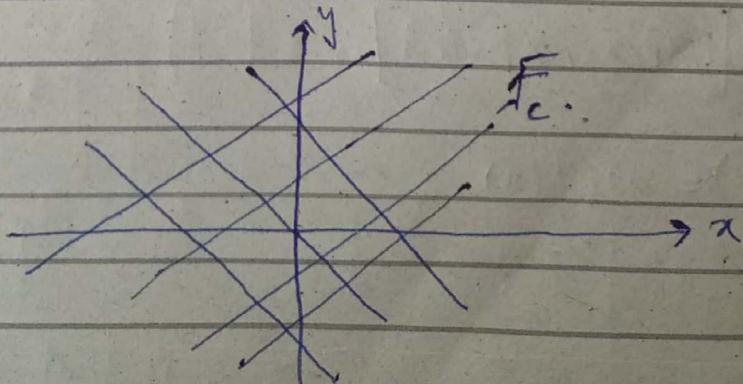
These

Lines are all possible predictors of  $F_1$  Hypothesis Class.



$$\Rightarrow F_c = \{w_1 x + c : w_1 \in \mathbb{R}, c \in \mathbb{R}\}$$

↳

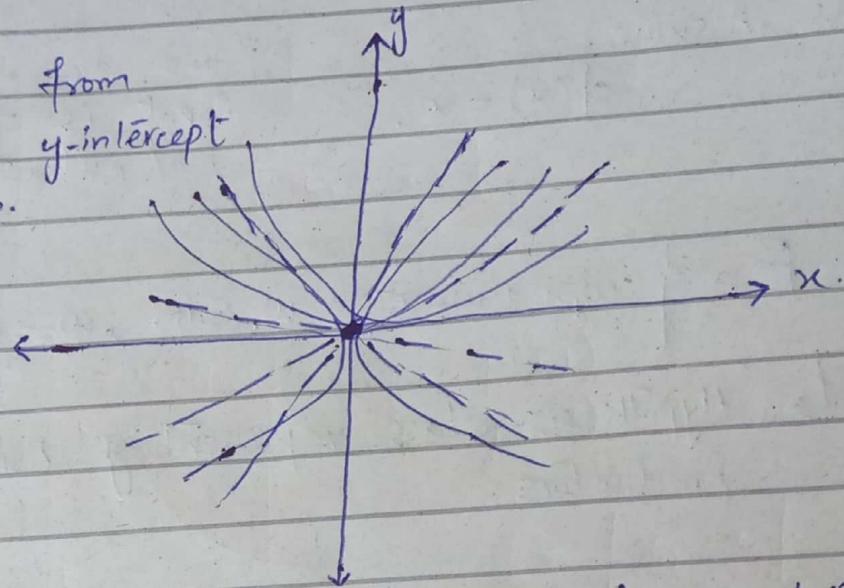


## Assignment:

- Common words.
- Time lapse → Not a good feature.
- Question: Mark it in msg.
- Beta use  
Balanced Data
- Logistic Regression.

$$\mathcal{F}_2 = \left\{ w_0 + w_1 x + w_2 x^2 : w_0, w_1 \in \mathbb{R}, w_2 \in \mathbb{R} \right\}$$

All passing from origin bcz y-intercept (c) is zero.



We can't include all Predictors in Hypothesis Class.  
 → Computation problems.  
 → will increase complexity.

• if  $w_2 \neq 0$ .  
 func. will be parabolic.  
 else if  $w_2 = 0$ .  
 func. will be linear. (straight).

$$\mathcal{F}_1 \subset \mathcal{F}_2$$

•  $\phi(x) = [1, x_1, x_2, x_1^2, x_2^2]$

$$\phi(x) = [1, x_1, x_2, x_1^2, x_2^2]$$

$$f_w(x) = \underbrace{w \cdot \phi(x)}_{\text{Score.}}$$

→ Linear in  $w$ .

→ Linear in  $\phi(x)$

→  $x$  is not Linear

∴  $f_w(x)$  is also non-Linear in  $x$ .

→ Loss function is also non-Linear.

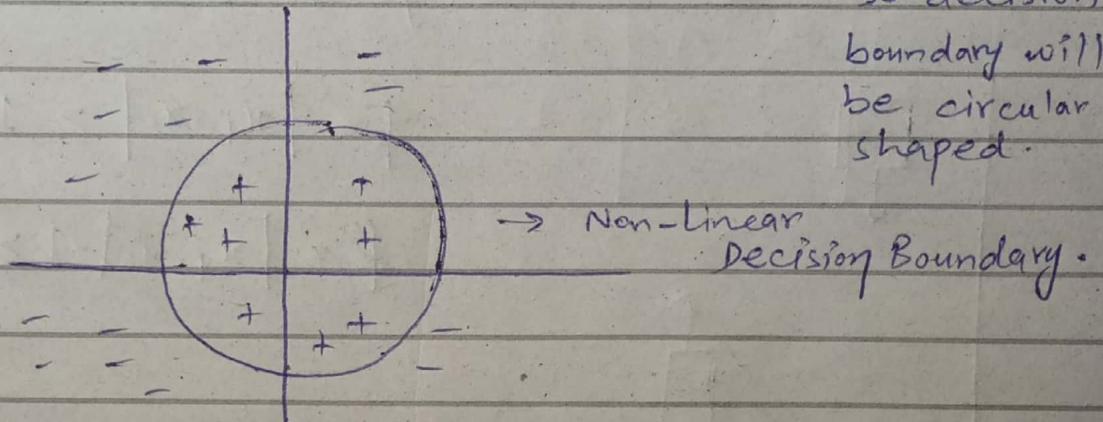
$$w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + 1.$$

Suppose  $w_1 = w_2 = 0$

$w_3 = w_4 = 1$ .

$$y' = f_w(x) = x_1^2 + x_2^2 + 1$$

→ Similar  
to circle  
equation.  
So decision  
boundary will  
be circular  
shaped.



→ NEURAL NETWORKS:

↳ Learns good feature automatically  
from raw / crude data.

ADV:

Crude Data → NN → Learns  
FE

So, focus focus on Learning only.

Example:

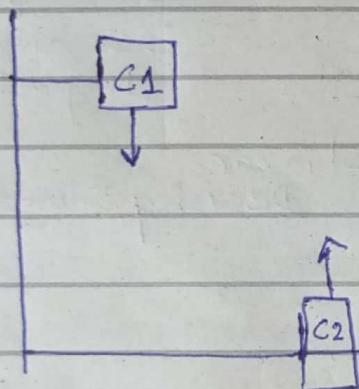
## Predicting Car Collisions

Inp: Position of two incoming cars.

$$x = \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$\downarrow$        $\downarrow$   
car1    car2.

Output: whether safe ( $y=+1$ )  
or collide ( $y=-1$ ).



$$f = |x_1 - x_2| \geq 1$$

$$f' \neq |x_1 - x_2| - 1 \geq 0$$

$$|x_1 - x_2| \geq 1 \quad (y = +1)$$

$$|x_1 - x_2| - 1 \geq 0 \quad (y = +1).$$

$$y' = f_w(x) = \text{sign}(|x_1 - x_2| - 1) \begin{cases} +1 \\ -1 \end{cases}$$

Exp Data Set

$x$	$y$
$[1 \ 3]$	+1
$[3 \ 1]$	+1.
$[1 \ 0.5]$	-1.

1/11/2023.

Wednesday.

## Lecture No. 10

### → Neural Networks:

#### ↳ Artificial Neural Network

Example:

Input: Position of two incoming cars  $x = [x_1, x_2]$

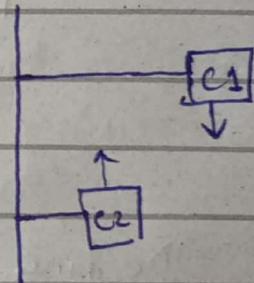
Output: whether safe ( $y = +1$ ) or  
collide ( $y = -1$ )

$$y = \text{sign}(|x_1 - x_2| - 1 \geq 0)$$

• if car1 is far right of car2.

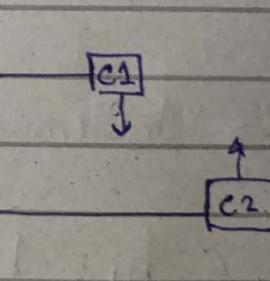
$$h_1 = 1[x_1 - x_2 \geq 1]$$

↳ Indicator Func. / Step Func.



• if car2 is far right of car1.

$$h_2 = 1[x_2 - x_1 \geq 1]$$



• Cars are safe if atleast one is true.

$$y = \text{sign}(h_1 + h_2)$$

## Data Set

$x$	$y$
$[3 \ 1]$	+1
$[1 \ 3]$	+1
$[1 \ 0.5]$	-1

$x$	$h_1$	$h_2$	$y' \leftarrow \text{Predicted}$
$[3 \ 1]$	1	0	+1
$[1 \ 3]$	0	1	+1
$[1 \ 0.5]$	0	0	-1.

→ Learning Strategy:

Define  $\phi(x) = \begin{bmatrix} 1 & x_1 & x_2 \end{bmatrix}$   
Bias

→ Intermediate hidden subproblems.

$$h_1 = 1[V_1 \cdot \phi(x) \geq 0] \quad \text{where } V_1 = [-1 \ 1 \ -1]$$

$$\left[ \begin{array}{l} \text{From } h_1 = 1[x_1 - x_2 \geq 0] \\ \Leftrightarrow -1 + x_1 - x_2 \geq 0 \end{array} \right] \quad \phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$$h_2 = 1[V_2 \cdot \phi(x) \geq 0] \quad \text{where } V_2 = [-1 \ -1 \ 1]$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}.$$

Final Prediction

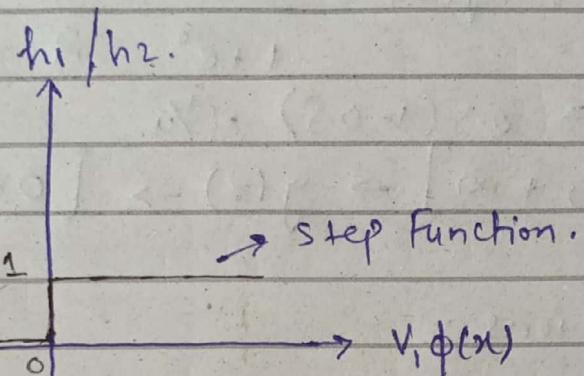
$$f_w = \underbrace{w \cdot \phi(x)}_{v_1, v_2}$$

$$f_{v,w}(x) = \text{sign}(w_1 h_1 + w_2 h_2) \quad w = [1, 1]$$

Gradient Descent Algo

$$w = w - \eta \nabla_w TL$$

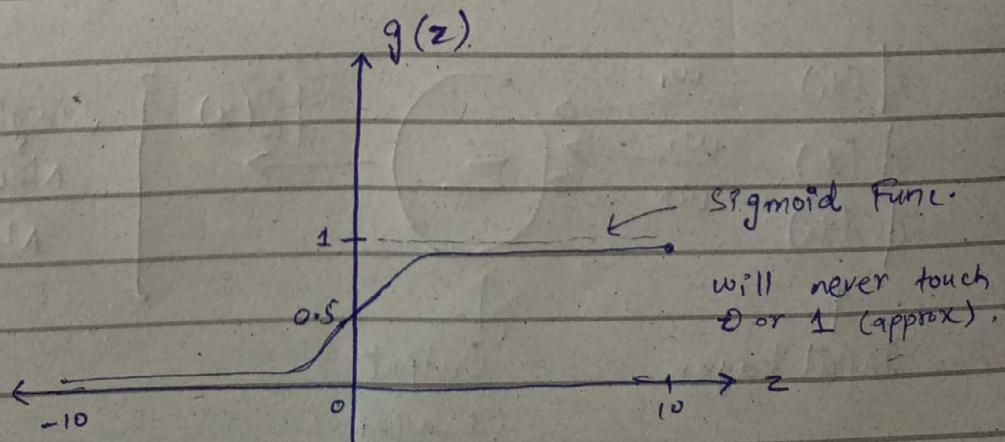
$$v = v - \eta \nabla_v TL$$



Slope is zero which means gradient/derivative is always zero and it will not update weights. So we cannot use it  $\rightarrow$  (gradient descent).

$\rightarrow$  Sigmoid Function:

$$g(z) = \frac{1}{1 + e^{-z}}, \quad g'(z) = g(z)(1 - g(z))$$



$$\text{Let } z = 10 \quad g(z) = \frac{1}{1+e^{-10}} \approx 1$$

$$\Rightarrow g'(z) = 1(1-1) \approx 0$$

$$\text{Let } z = -10 \quad g(z) = \frac{1}{1+e^{10}} \approx 0$$

$$\Rightarrow g'(z) = 0(1-0) \approx 0$$

$$\text{Let } z = 0, \quad g(z) = \frac{1}{1+e^0} = \frac{1}{2} = 0.5$$

$$\Rightarrow g'(z) = 0.5(1-0.5) = 0.25$$

$$[-\infty, +\infty] \rightarrow g(z) \rightarrow [0 \sim 1]$$

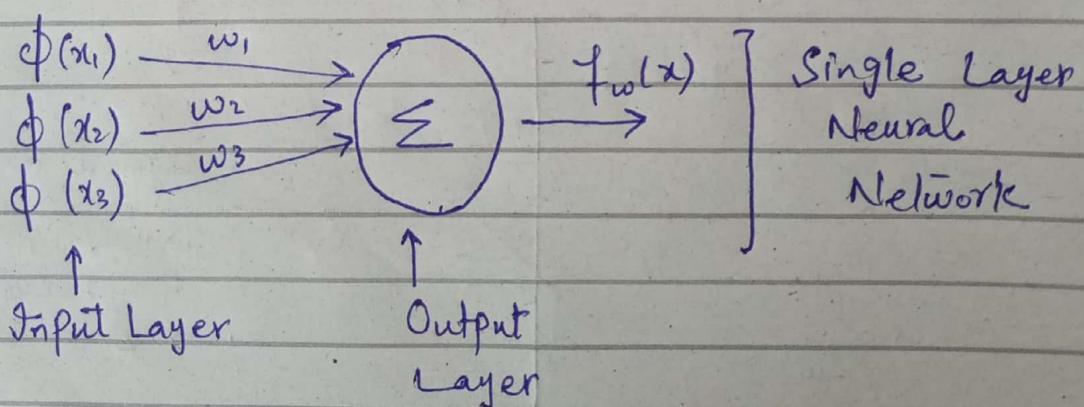
## Sigmoid Functions

$$h_1 = f(v_1 \cdot \phi(x))$$

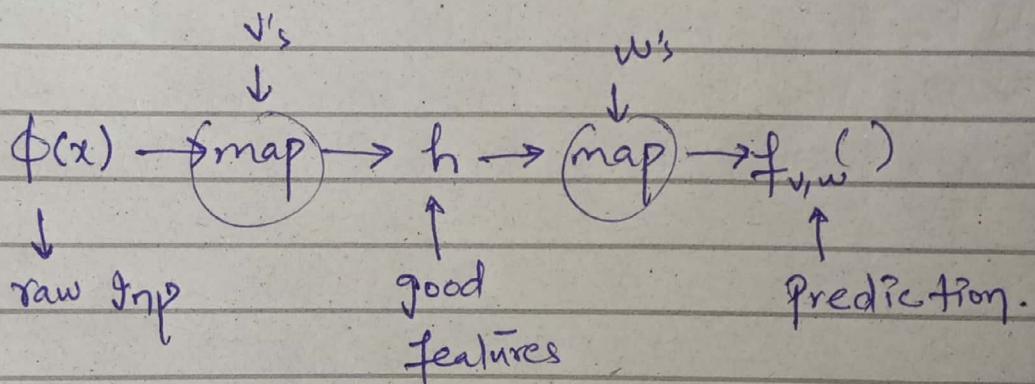
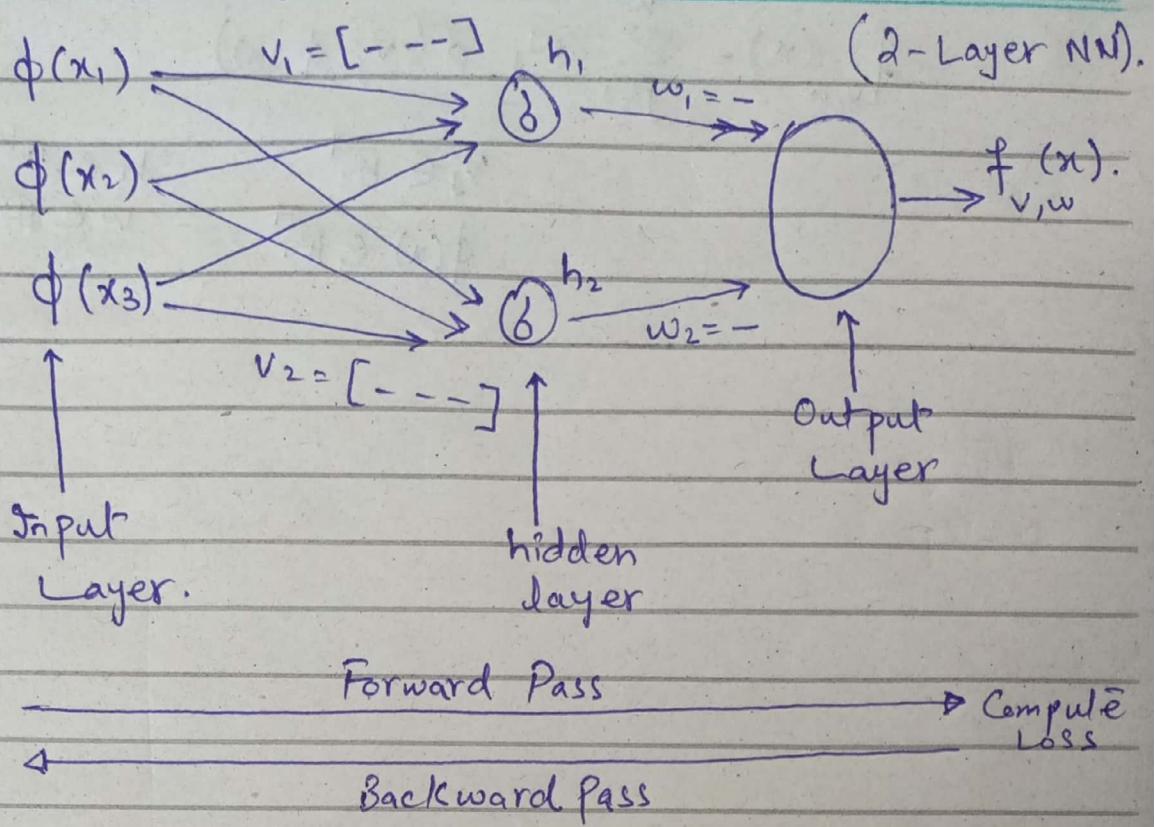
$$h_2 = f(v_2 \cdot \phi(x))$$

We have discussed in Linear Predictors.

$$f_w = w \cdot \phi(x) = \sum_{i=1}^d w_i \phi(x)_i$$



→ Multi-Layer Neural Network.



In Linear Predictors:

$$\frac{1}{|\mathcal{D}_{\text{train}}|} \leq \underset{(x,y) \in \mathcal{D}_{\text{train}}}{\text{Loss}(x, y, w)} \rightarrow = (f_w(x) - y)^2$$

$w \cdot \phi(x)$

So, In Neural Network

$$\text{Loss}(x, y, v, w) = (f_{v,w}(x) - y)^2$$

$$f_{v,w}(x) = \sum_{j=1}^k w_j \cdot \delta(v_j \cdot \phi(x))$$

No. of Hidden units.

$v_j \in \mathbb{R}^d$

$\phi(x) \in \mathbb{R}^d$

$w \in \mathbb{R}^k$