

⊗ Model Selection:

- model should be generalized
 - for all type data
 - train
 - test
 - validation \Rightarrow unseen data
- model must do well on both train as well as 'unseen' data.

- ⊗ - good on train set but
- not good on unseen data } over-fit model

- ⊗ - does not work well on both train and unseen data
 \Rightarrow under-fit model

- ⊗ we need to test our model with unseen data for validation.
 - divide the data set into two part

- train set
- test set / validation set

most of the time same

\Rightarrow but its better to use three part

- train set \Rightarrow create model
- validation set \Rightarrow select model
- test set \Rightarrow evaluate model, test accuracy

Slide-2 ***

Multi Layer Perceptron (MLP)

MLP Notation:

⇒ Bias notation:

$b_{ij} \Rightarrow$ $i \Rightarrow$ layer number, starts from 1
 $j \Rightarrow$ index of neuron in layer 'i'

⇒ Output notation:

$O_{ij} \left\{ \begin{array}{l} i \Rightarrow \text{layer number} \\ j \Rightarrow \text{index of neuron in layer 'i'} \end{array} \right.$

⇒ Weight notation:

$W_{ij}^k \left\{ \begin{array}{l} k \Rightarrow \text{target layer number} \\ \quad \quad \quad k-1 \text{ to } k \\ i \Rightarrow \text{index of neuron in layer 'k-1'} \\ j \Rightarrow \text{index of neuron in layer 'k'} \end{array} \right.$

Number of trainable parameters:

Σ for each layer \Rightarrow weight = number of nodes in previous layer * current layer
bias = number of nodes in current layer

Input layer:

$X_{ij} \left\{ \begin{array}{l} i \Rightarrow \text{row number} \\ j \Rightarrow \text{number of feature} \end{array} \right.$

L-10/18.02.2025/

⊗ For MLP:

combined probability = \sum probability of each perceptron
 $= z \Rightarrow$ most of the cases greater than 1.

$= \sigma(z) \Rightarrow$ apply sigmoid function
 \Rightarrow without weight & bias

$$z = p_1(y) + p_2(y)$$

$$= 0.8 + 0.7$$

$$= 1.5 > 1$$

$$\sigma(z) = 0.82 \leq 1$$

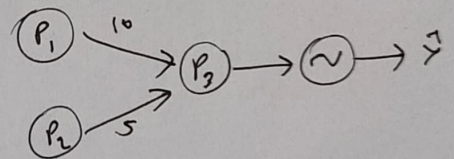
\Rightarrow with weight & bias!

$$z = 10(0.8) + 5(0.7) + 3$$

\swarrow weight \swarrow weight \searrow bias

$$\sigma(z) = \dots$$

Slide-17, 18



⊗ Loss Function

⊗ Loss function!

- for a single sample

$$L = (\hat{y}_i - y_i)^2$$

⊗ Cost function!

- for entire datasets

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

* Loss function!

- mathematical function that quantifies how well a model's predictions match the actual target values.

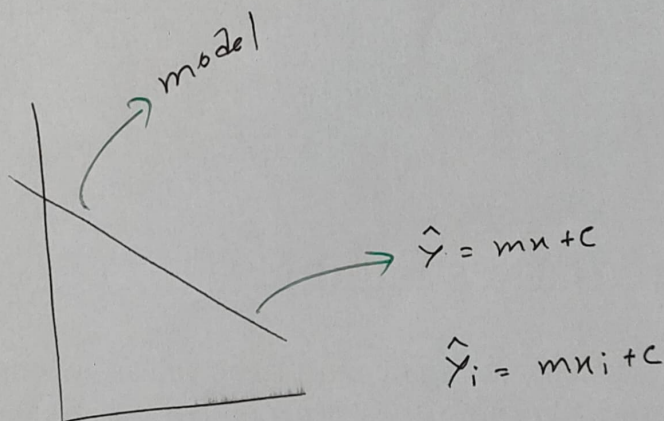
* Advantage & disadvantage of different loss function:

Slide - 4 - 10 ***

Slide - 11 ***

- * keep focus on the error unit & data unit.

*



$$L = (y_i - \hat{y}_i)^2$$

* For project, use:

$$\Rightarrow L(m, c) = \underbrace{(y_i - mx_i - c)^2}_{\text{learnable parameter}}$$

\Rightarrow CNN

\Rightarrow GAN

\Rightarrow knowledge distillation

1-1 / 23.02.2025 /

Back-Propagation

⊗

CGPA	IQ	Salary
3.75	80	8
3.00	90	8.3

Salary

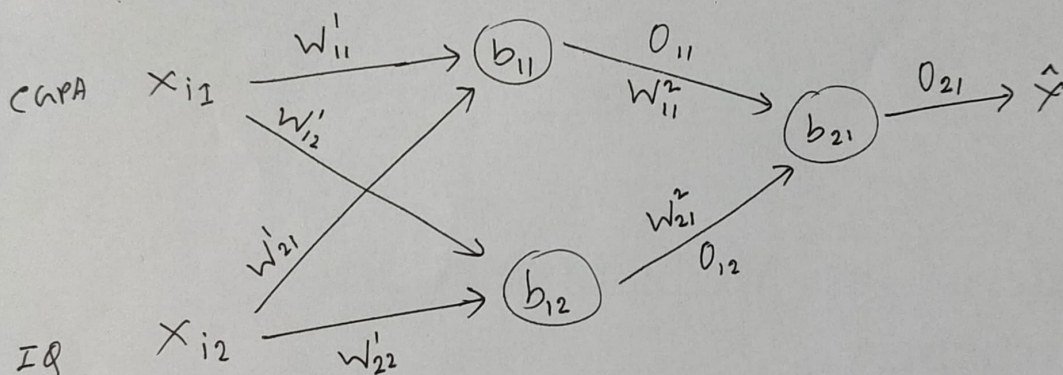
18.96

5.2

⇒ after first iteration, it overfit, change and test the next data

Underfit now, again change and test the next data.

⊗



⇒ activation function = Linear

$$\hat{y} = O_{21} = O_{11} W''_{11} + O_{12} W''_{21} + b_{21}$$

$$L = (y - \hat{y})^2$$

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{dL}{db_{\text{old}}}$$

Partial Derivative

$$\therefore \frac{dL}{d\hat{y}} = \frac{d}{d\hat{y}} (y - \hat{y})^2 = -2(y - \hat{y})$$

⇒ total 9 parameters
find out partial derivative for all of them.

$$O_{11} = x_{i1} W'_{11} + x_{i2} W'_{21} + b_{11}$$

$$O_{12} = x_{i1} W'_{12} + x_{i2} W'_{22} + b_{12}$$

$$\textcircled{i} \quad \frac{dL}{dW_{11}^2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{11}^2} = \frac{dL}{dO_{21}} \cdot \frac{dO_{21}}{dW_{11}^2} = -2(\gamma - \hat{y}) \cdot O_{11}$$

$$\textcircled{ii} \quad \frac{dL}{dW_{21}^2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{21}^2} = -2(\gamma - \hat{y}) \cdot O_{12}$$

$$\textcircled{iii} \quad \frac{dL}{db_{21}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{db_{21}} = -2(\gamma - \hat{y}) \cdot 1$$

$$\textcircled{iv} \quad \frac{dL}{db_{11}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{11}} \cdot \frac{dO_{11}}{db_{11}} = -2(\gamma - \hat{y}) \cdot W_{11}^2 \cdot 1$$

$$\textcircled{v} \quad \frac{dL}{db_{12}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{12}} \cdot \frac{dO_{12}}{db_{12}} = -2(\gamma - \hat{y}) \cdot W_{21}^2 \cdot 1$$

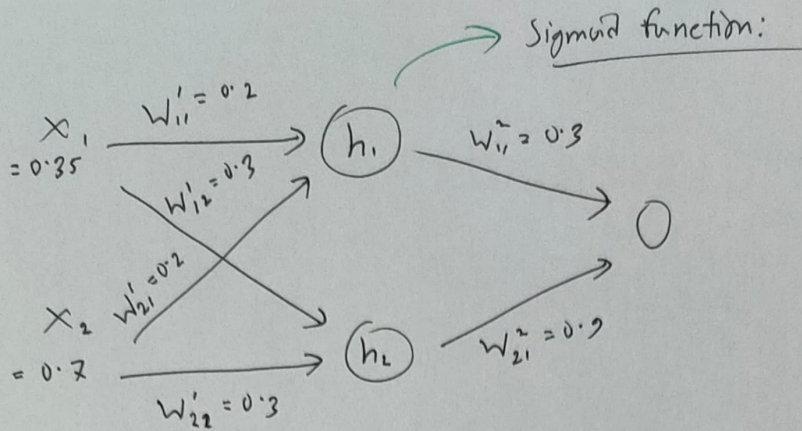
$$\textcircled{vi} \quad \frac{dL}{dW_{11}'} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{11}} \cdot \frac{dO_{11}}{dW_{11}'} = -2(\gamma - \hat{y}) \cdot W_{11}^2 \cdot x_{11}$$

$$\textcircled{vii} \quad \frac{dL}{dW_{12}'} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{12}} \cdot \frac{dO_{12}}{dW_{12}'} = -2(\gamma - \hat{y}) \cdot W_{21}^2 \cdot x_{11}$$

$$\textcircled{viii} \quad \frac{dL}{dW_{21}'} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{11}} \cdot \frac{dO_{11}}{dW_{21}'} = -2(\gamma - \hat{y}) \cdot W_{11}^2 \cdot x_{12}$$

$$\textcircled{ix} \quad \frac{dL}{dW_{22}'} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dO_{12}} \cdot \frac{dO_{12}}{dW_{22}'} = -2(\gamma - \hat{y}) \cdot W_{21}^2 \cdot x_{12}$$

Exam question
will be like
this



⇒ Forward Pass:

$$O_{11} = \sigma[0.35 \times 0.2 + 0.7 \times 0.2]$$

$$= \frac{1}{1 + e^{-0.21}} = 0.55$$

$$O_{21} = \sigma[0.35 \times 0.3 + 0.7 \times 0.3]$$

$$= \frac{1}{1 + e^{-0.315}} = 0.579$$

L-12 / 25.02.2025 /

Online Session
28.02.2025
6:30pm



Backpropagation:

initialize weights and bias

for $i=1$ to epochs:

for $j=1$ to row:

select one row (random)

apply forward pass

calculate loss

update weights and bias:

$$W_n = W_0 - \eta \frac{dL}{dW_0}$$

* Gradient Descent:

(i) Batch GD:

- calculate the total error then update the weights and bias
- also known as ~~vanilla~~ vanilla GD.

~~(ii)~~

- time consuming but accurate.

(ii) Stochastic GD:

- weights and biases are updated based on the loss due to single data-point.
- faster but not so accurate.

(iii) Mini batch GD:

- weights and biases are updated based on the loss due to mini batches.

* With two hidden layer diagram with two features:

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dL}{dW_{\text{old}}}$$

$$\therefore W_{11}^1 = W_{11}^1 - \eta \frac{dL}{dW_{11}^1}$$

$$\therefore \frac{dL}{dW_{11}^3} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{11}^3}$$

$$\therefore \frac{dL}{dW_{11}^2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{11}^3} \cdot \frac{dW_{11}^3}{dW_{11}^2}$$

$$\frac{dL}{dW'_{11}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{11}^3} \cdot \frac{dW_{11}^3}{dW_{11}^2} \cdot \frac{dW_{11}^2}{dW'_{11}} +$$

$$\frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW_{21}^3} \cdot \frac{dW_{21}^3}{dW_{12}^2} \cdot \frac{dW_{12}^2}{dW'_{11}}$$

* Vanishing GD Problem:

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\begin{aligned} \frac{d}{dz} \sigma(z) &= \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) \\ &= \frac{(1+e^{-z}) \frac{d}{dz} 1 - 1 \frac{d}{dz} (1+e^{-z})}{(1+e^{-z})^2} \end{aligned}$$

$$= \frac{e^{-z}}{(1+e^{-z})(1+e^{-z})}$$

$$= \sigma(z) \frac{e^{-z}}{1+e^{-z}}$$

$$= \sigma(z)(1 - \sigma(z))$$

$$\begin{aligned} 1 - \sigma(z) &= 1 - \frac{1}{1+e^{-z}} \\ &= \frac{1+e^{-z}-1}{1+e^{-z}} \\ &= \frac{e^{-z}}{1+e^{-z}} \end{aligned}$$

* ~~ReLU~~ ReLU = $\max(0, z)$

$$\frac{d}{dz} = (0, 1)$$

Some more,

⇒ Leaky ReLU (No zero gradient)

⇒ dying ReLU problem

⇒ exploring gradient descent

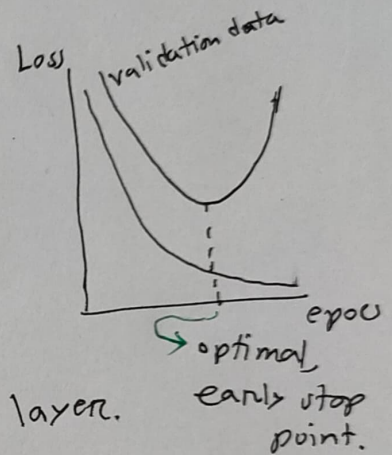
⇒ vanishing gradient descent problem.

1-13/02.03.2025 /

⊕ How to improve performance in a NN?

① Fine tuning the hyper-parameters in a NN.

- learning rate
- batch size \Rightarrow M.B.G.D. $\begin{cases} \rightarrow \text{small (8-32)} \Rightarrow \text{generalize better} \\ \rightarrow \text{large (2-10k)} \end{cases}$
- numbers of epochs \Rightarrow early stop
- number of hidden layers
- activation function
- optimizer
- loss function
- number of neurons per hidden layer.



② By solving problem a NN may encountered.

- vanishing/exploding gradient descent
- not sufficient data \Rightarrow transparent data can be good choice.
- slow training
- Overfits

⊕ Problem with NN:

- vanishing/exploding gradient descent
- weights initialization
- activation function
- Batch normalization

- gradient clipping (exploding ~~gd~~ GD problem, solved)

- not sufficient data,

⇒ data augmentation

⇒ transparent learning

- slow training

- optimizer

- learning rate scheduler

- Overfit:

- L_1/L_2 regularization

⊕ What is diverse data?

Midterm
02.03.2025
<u>upto this</u>