

Question Pattern

Total Question Set	Need to Answer	Question Type	Marks
7	5	Each set will have 3 sub-questions (CO1 + CO2 + CO3)	14
Total:			70

Suggestion

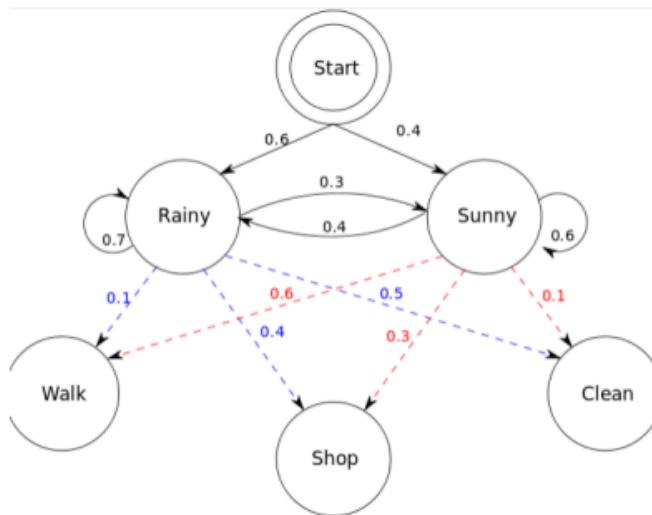
Label	Topic Name	Set Distribution
A	✓ Hidden Markov Models (HMM), Viterbi, RNN, LSTM, GRU	1 set
B	✓ Word Embedding, BoW, TF-IDF, Attention Mechanism, Vectorization of NN, ACO	1.5/2 set
C	✓ Activation Function, Logistic Regression, Optimizer, Cross Entropy Loss, Information Theory, Gradient Descent	2 set
D	✓ CNN, Parameter Calculation, Parameter vs Hyperparameter, Genetic Algorithm, TSP	1 set
E	✓ Fuzzy Basics, Fuzzy Membership Function, Fuzzy Addition, Fuzzy Composition, Neuro-fuzzy System	1/1.5 set

A : 1 set

Quiz-4 - Integer43

Set-A

1. If the sequence of observations is **SCWS**. Find the sequence of states for the following scenario using the Viterbi algorithm. State every step of the simulation.
(W = Walk, S = Shop, C = Clean)



1. Solution: **Solved by Younus-131**

Initial Probability naa dewa thakle, by default duitay 0.5 koree hobee

Set - A

Initial Probability, $\pi = \begin{bmatrix} R_n & S_n \\ 0.6 & 0.4 \end{bmatrix}$

Transition Matrix, $A = \begin{bmatrix} R_n & S_n \\ R_n & S_n \end{bmatrix}$

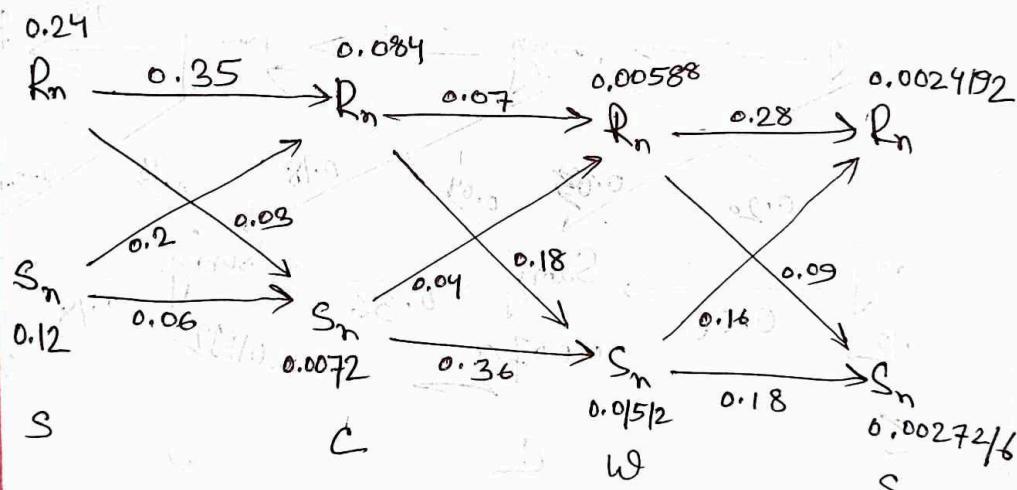
$$R_n \quad 0.7 \quad 0.3$$

$$S_n \quad 0.4 \quad 0.6$$

Emission Matrix, $B = \begin{bmatrix} w & S & C \end{bmatrix}$

$$R_n \quad 0.1 \quad 0.4 \quad 0.5$$

$$S_n \quad 0.6 \quad 0.3 \quad 0.1$$



Result : $R_n \rightarrow R_n \rightarrow S_n \rightarrow S_n$

$$P(S, R_n) = P(S|R_n) P(R_n)$$

$$= 0.4 \times 0.6 = 0.24$$

$$P(S, S_n) = P(S|S_n) P(R_n|S_n)$$

$$= 0.3 \times 0.4 = 0.12$$

$$P(C, R_n) = P(C|R_n) P(R_n|S_n)$$

$$= 0.5 \times 0.7 = 0.35$$

$$P(C, S_n) = P(C|S_n) P(S_n|R_n)$$

$$= 0.1 \times 0.6 = 0.06$$

$$P(C, R_n) = P(C|R_n) P(R_n|S_n)$$

$$= 0.5 \times 0.4 = 0.20$$

$$P(C, S_n) = P(C|S_n) P(S_n|R_n)$$

$$= 0.1 \times 0.3 = 0.03$$

$$\max(0.24 \times 0.35, 0.12 \times 0.2)$$

$$= \max(0.084, 0.024) = 0.084$$

$$\max(0.24 \times 0.03, 0.12 \times 0.06)$$

$$= \max(0.0072, 0.0072) = 0.0072$$

$$P(W, R_n) = P(W|R_n) P(R_n|S_n)$$

$$= 0.1 \times 0.7 = 0.07$$

$$P(W, S_n) = P(W|S_n) P(S_n|R_n)$$

$$= 0.6 \times 0.6 = 0.36$$

$$P(W, R_n) = P(W|R_n) P(R_n|S_n)$$

$$= 0.1 \times 0.4 = 0.04$$

$$\max(0.084 \times 0.07, 0.0072 \times 0.04)$$

$$= \max(0.00588, 0.000288) = 0.00588$$

$$P(W, S_n) = P(W|S_n) P(S_n|R_n)$$

$$= 0.6 \times 0.3 = 0.18$$

$$P(S, R_n) = P(S|R_n) P(R_n|S_n)$$

$$= 0.4 \times 0.7 = 0.28$$

$$P(S, R_n) = P(S|R_n) \times P(R_n|S_n)$$

$$= 0.4 \times 0.4 = 0.16$$

$$\max(0.00588 \times 0.28, 0.01512 \times 0.16)$$

$$= \max(0.0016464, 0.0024192) = 0.0024192$$

$$P(S, S_n) = P(S|S_n) P(S_n|R_n)$$

$$= 0.3 \times 0.6 = 0.18$$

$$P(S, S_n) = P(S|S_n) P(S_n|R_n)$$

$$= 0.3 \times 0.3 = 0.09$$

$$\max(0.00588 \times 0.09, 0.01512 \times 0.18)$$

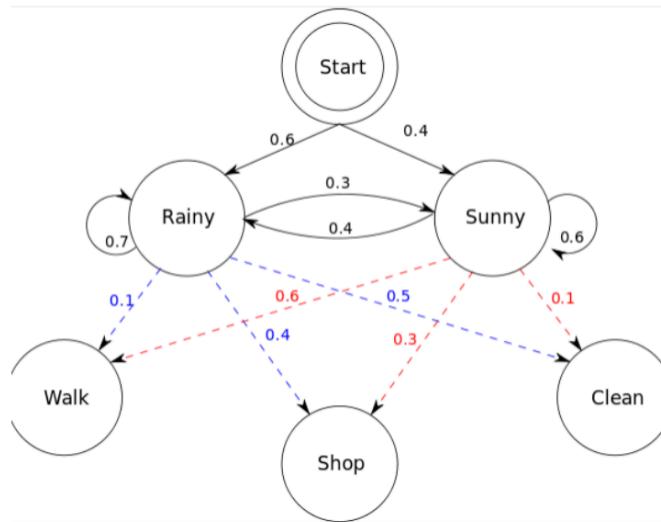
$$= \max(0.0005292, 0.0027216) = 0.0027216$$

Ans vul ache

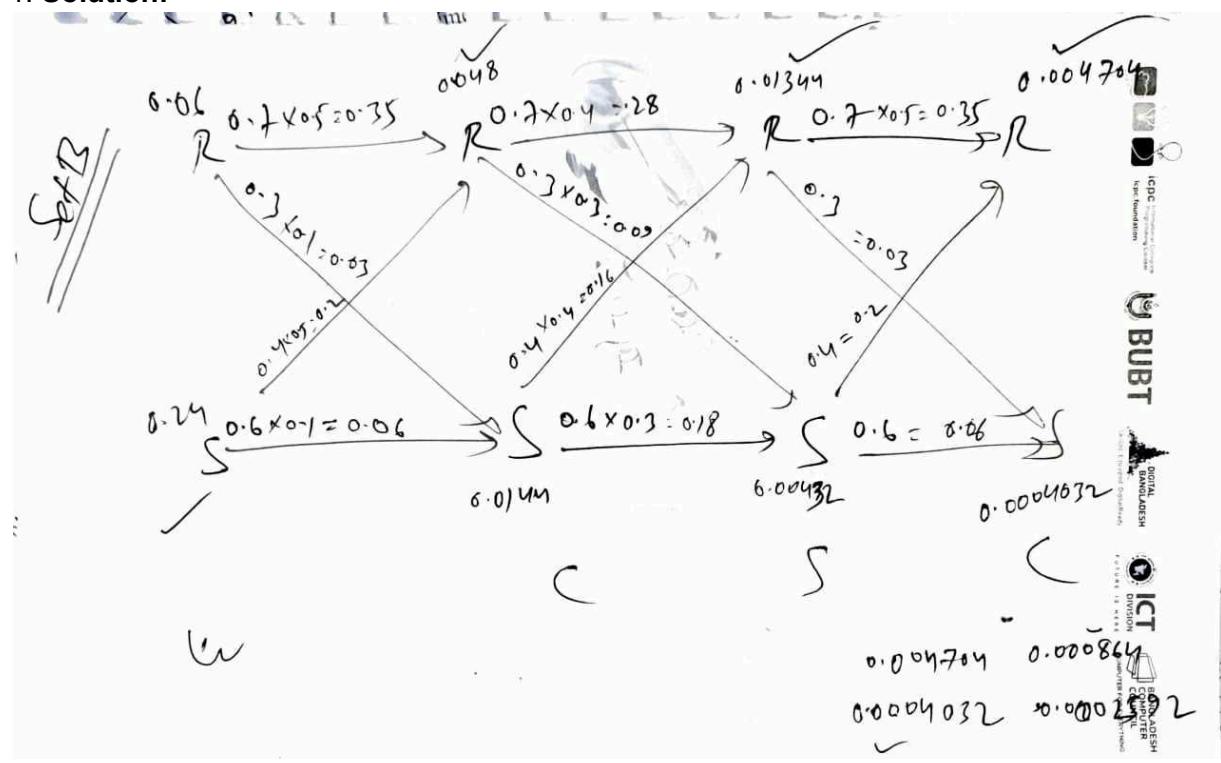
- Could you point out the mistake, please?

Set-B

1. If the sequence of observations is **WCSC**. Find the sequence of states for the following scenario using the Viterbi algorithm. State every step of the simulation.
 (W = Walk, S = Shop, C = Clean)

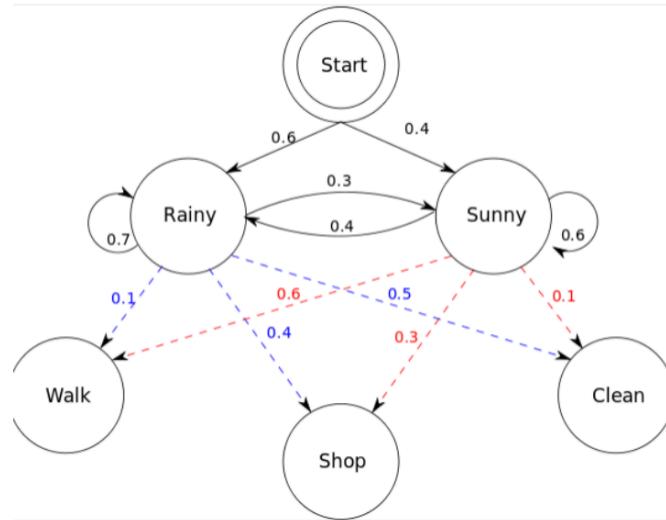


1. Solution:

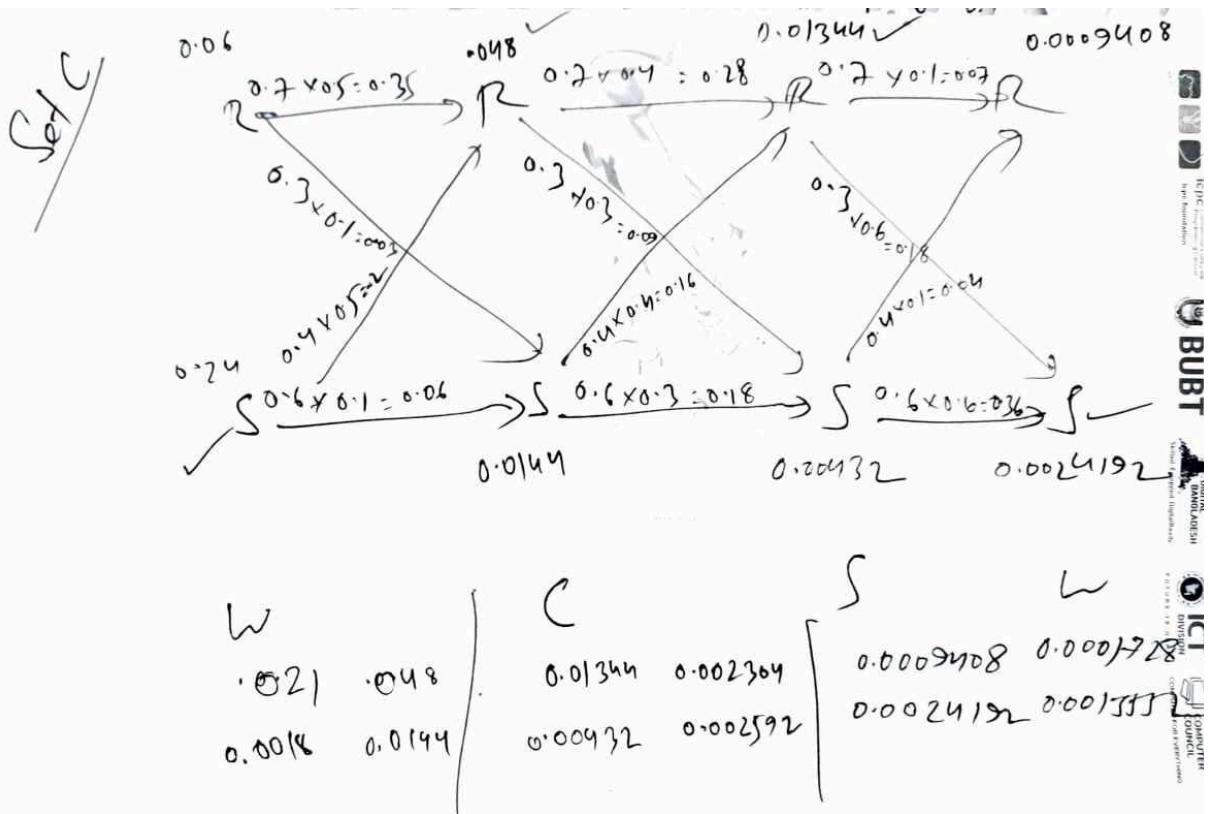


Set-C

1. If the sequence of observations is **WCSW**. Find the sequence of states for the following scenario using the Viterbi algorithm. State every step of the simulation.
 (W = Walk, S = Shop, C = Clean)

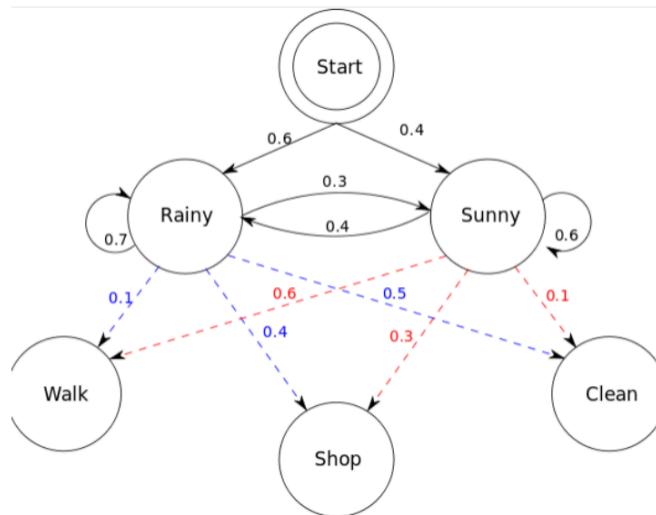


1. Solution:

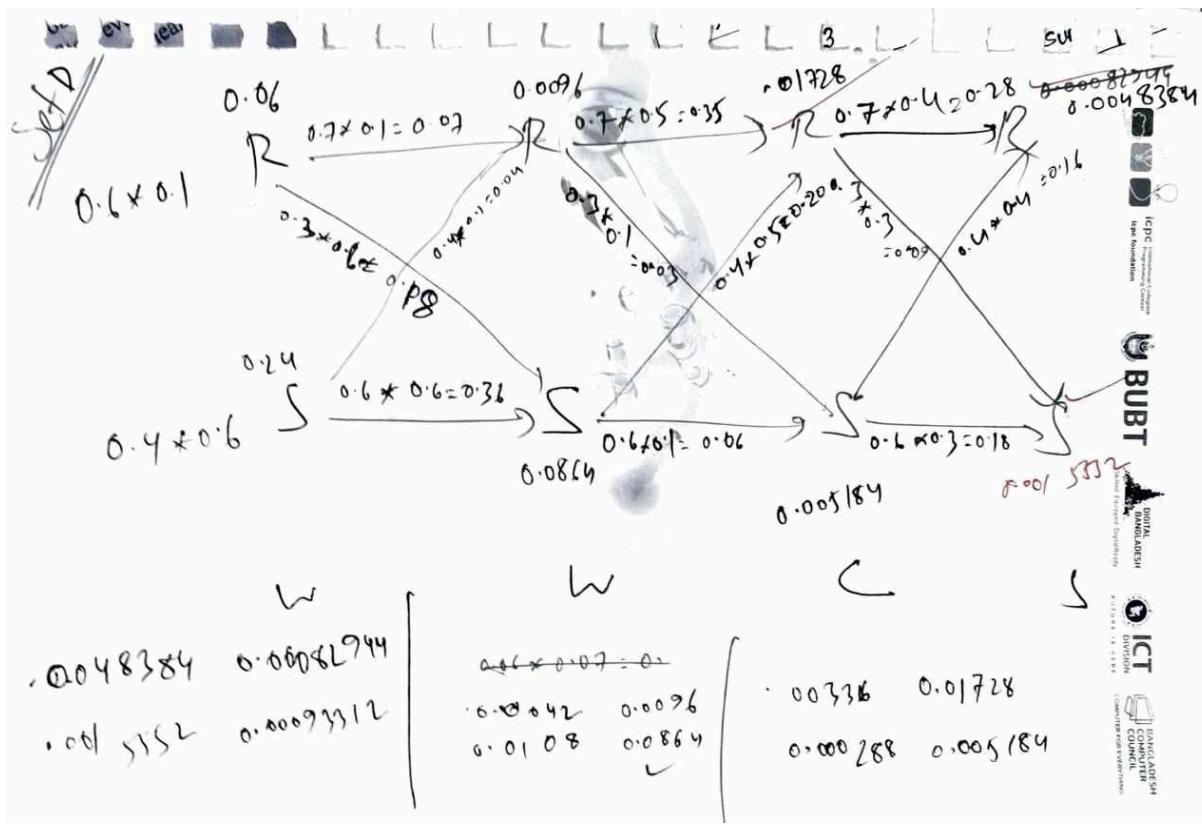


Set-D

1. If the sequence of observations is **WWCS**. Find the sequence of states for the following scenario using the Viterbi algorithm. State every step of the simulation.
 (W = Walk, S = Shop, C = Clean)

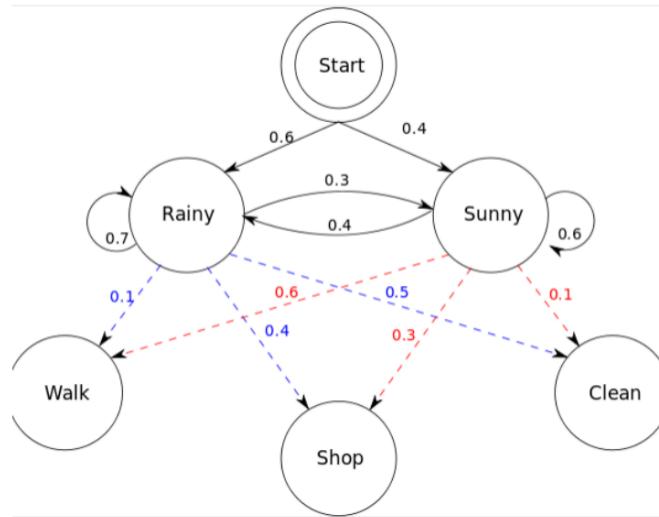


1. Solution:



Set-E

2. If the sequence of observations is **WCWS**. Find the sequence of states for the following scenario using the Viterbi algorithm. State every step of the simulation.
(W = Walk, S = Shop, C = Clean)



2. Solution:

HMM, Viterbi

Origin42

- b) Consider the simple hidden Markov model (HMM) in Figure-1. This model is composed [8] of 2 states, **HIGH** and **LOW**. You can for example consider that HIGH characterizes coding DNA while LOW characterizes non-coding DNA. Analyze the model and find out the right regions of DNA for the following sequence using the Viterbi algorithm.

A C T G A

State every step of the simulation. HMM + Viterbi

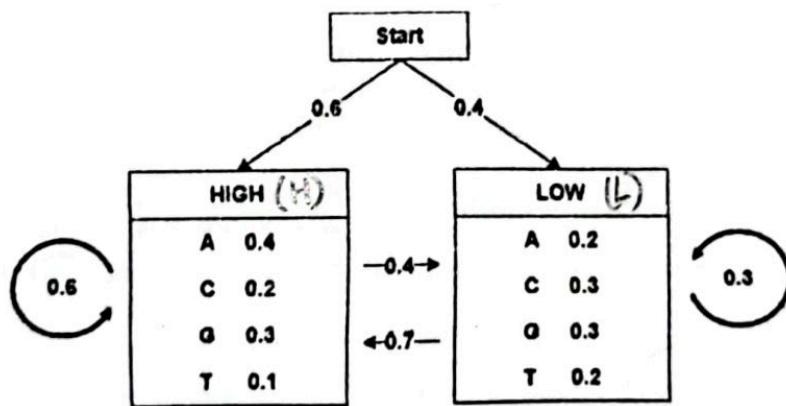


Figure - 1

Solution: Sujon 49

2) Origin 42

High Low

Initial probability, $\pi = [0.6 \quad 0.4]$

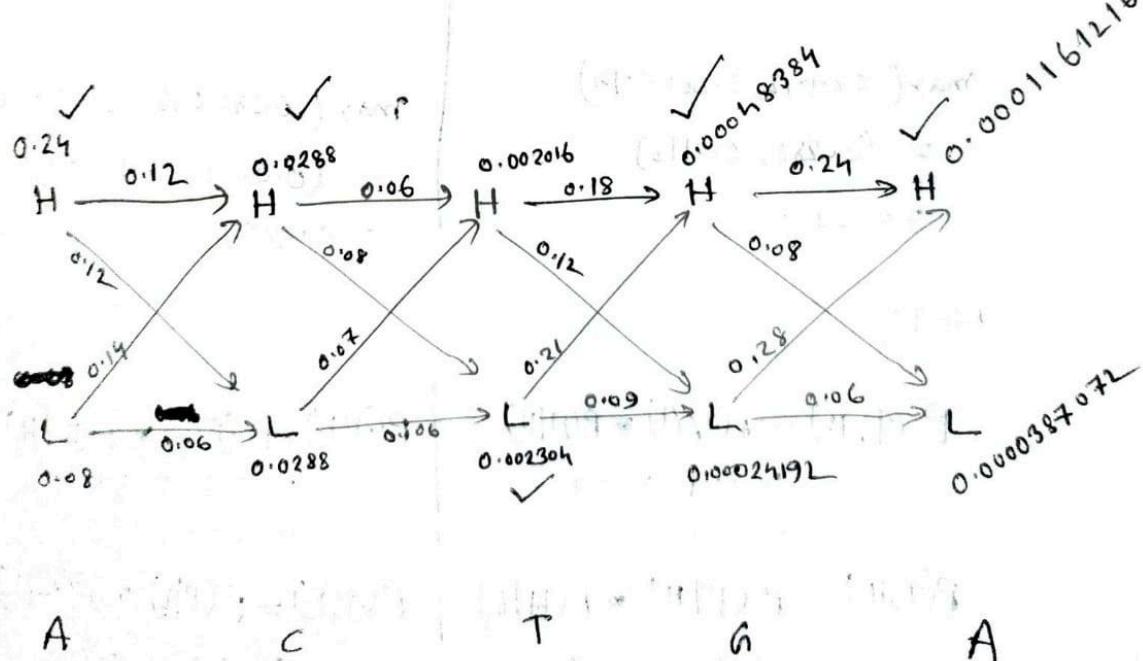
Transition probability:

old	Now	High	Low
High		0.6	0.4
Low		0.7	0.3

Emission probability:

A	T	C	G	H
High	0.4	0.2	0.3	0.1
Low	0.2	0.3	0.3	0.2

Using viterbi algo for sequence: ACTGA



For G!

$$P(G, H) = P(G|H) \times P(H|H)$$

$$= 0.3 \times 0.6$$

$$= 0.18$$

$$P(G, L) = P(G|H) \times P(H|L)$$

$$= 0.3 \times 0.7$$

$$= 0.21$$

$$\max(0.002016 \times 0.18, 0.002036 \times 0.21)$$

$$= \max(3.6288 \times 10^{-5}, 4.8384 \times 10^{-5})$$

$$= 4.8384 \times 10^{-5}$$

$$P(G, L) = P(G|L) \times P(L|H)$$

$$= 0.3 \times 0.4$$

$$= 0.12$$

$$P(G, L) = P(G|L) \times P(L|L)$$

$$= 0.3 \times 0.3$$

$$= 0.09$$

$$\max(0.002016 \times 0.12, 0.002036 \times 0.09)$$

$$= \max(2.4192 \times 10^{-5}, 1.8324 \times 10^{-5})$$

$$= 0.00024192$$

For A!

$$P(A, H) = P(A|H) \times P(H|H)$$

$$= 0.4 \times 0.6$$

$$= 0.24$$

$$P(A, L) = P(A|H) \times P(H|L)$$

$$= 0.4 \times 0.7$$

$$= 0.28$$

$$\max(0.24 \times 0.00048384, 0.28 \times 0.00024192)$$

$$= \max(1.161216 \times 10^{-5}, 6.7376 \times 10^{-6})$$

$$= 0.0001161216$$

$$P(A, L) = P(A|L) \times P(L|H)$$

$$= 0.2 \times 0.4$$

$$= 0.08$$

$$P(A, L) = P(A|L) \times P(L|L)$$

$$= 0.2 \times 0.3$$

$$= 0.06$$

$$\max(0.08 \times 0.00048384, 0.06 \times 0.00024192)$$

$$= \max(3.87072 \times 10^{-5}, 1.45152 \times 10^{-5})$$

$$= 3.87072 \times 10^{-5}$$

$$= 0.0000387072$$

Two possible sequences for A e T G A Sequence:



(Ans)

Enigma41

6

- ii) Suppose, you have 2 boxes B1 and B2, each of which contains 4 balls colored RED, GREEN, BLUE and WHITE. A sequence of five balls is randomly drawn from the boxes. In this particular case, the user observes a sequence of balls GREEN, GREEN, RED, WHITE and BLUE. Find out the right sequence of two boxes that these five balls were pulled from using the Viterbi algorithm. State every step of the simulation. The transition and emission probabilities are given below.

HMM - Viterbi

Transition Probability

Box	B1	B2
B1	0.6	0.4
B2	0.3	0.7

Emission Probability

Box \ Ball	RED	GREEN	BLUE	WHITE
B1	0.2	0.4	0.3	0.1
B2	0.2	0.2	0.3	0.3

Solution: 075 (Assume Initial Probability)

Enigma 91

Subject..... Date..... Time.....

B1 B2

⇒ Initial Probability, $\pi = [0.5 \quad 0.5]$

⇒ Transition Probability, $A =$

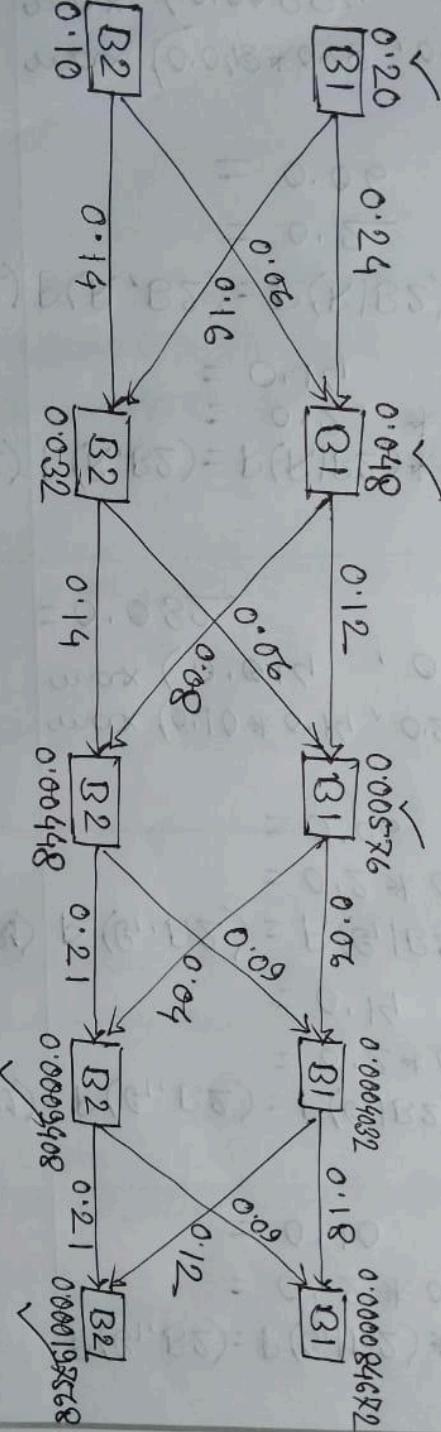
old	New	B1	B2
B1		0.6	0.4
B2		0.3	0.7

⇒ Emission Probability, $B =$

	R	G	B	W
B1	0.2	0.4	0.3	0.1
B2	0.2	0.2	0.3	0.3

⇒ Observable Sequence: G G R W B

B1 B1 B1 B2 B2

 G_1 R ω B

Subject.....
Date..... Time.....

For G₁:

$$\begin{aligned} P(G_1, B1) &= P(G_1|B1) * P(B1) & P(G_1, B2) &= P(G_1|B2) * P(B2) \\ &= 0.4 * 0.5 & &= 0.2 * 0.5 \\ &= 0.20 & &= 0.10 \end{aligned}$$

For G₂:

$$\begin{aligned} P(G_2, B1) &= P(G_2|B1) * P(B1|B1) & P(G_2, B2) &= P(G_2|B2) * P(B2|B2) \\ &= 0.4 * 0.6 & &= 0.2 * 0.7 \\ &= 0.24 & &= 0.14 \\ P(G_2, B1) &= P(G_2|B1) * P(B2|B1) & P(G_2, B2) &= P(G_2|B2) * P(B1|B2) \\ &= 0.4 * 0.4 & &= 0.2 * 0.3 \\ &= 0.16 & &= 0.06 \end{aligned}$$

$$\begin{aligned} \max(0.20 * 0.24, 0.06 * 0.10) \\ \max(0.048, 0.006) \\ = 0.048 \end{aligned}$$

$$\begin{aligned} \max(0.10 * 0.14, 0.20 * 0.16) \\ \max(0.014, 0.032) \\ = 0.032 \end{aligned}$$

For R:

$$\begin{aligned} P(R, B1) &= P(R|B1) * P(B1|B1) & P(R, B2) &= P(R|B2) * P(B2|B2) \\ &= 0.2 * 0.6 & &= 0.2 * 0.7 \\ &= 0.12 & &= 0.14 \\ P(R, B1) &= P(R|B1) * P(B2|B1) & P(R, B2) &= P(R|B2) * P(B1|B2) \\ &= 0.2 * 0.4 & &= 0.2 * 0.3 \\ &= 0.08 & &= 0.06 \end{aligned}$$

$$\begin{aligned} \max(0.048 * 0.12, 0.032 * 0.06) \\ \max(0.00576, 0.00192) \\ = 0.00576 \end{aligned}$$

$$\begin{aligned} \max(0.048 * 0.08, 0.32 * 0.14) \\ \max(0.00384, 0.00448) \\ = 0.00448 \end{aligned}$$

Subject.....

Date:..... Time:.....

For ω :

$$\begin{aligned} P(\omega, B_1) &= P(\omega | B_1) * P(B_1 | B_1) & P(\omega, B_2) &= P(\omega | B_2) * P(B_2 | B_2) \\ &= 0.1 * 0.6 & &= 0.3 * 0.7 \\ &= 0.06 & &= 0.21 \end{aligned}$$

$$\begin{aligned} P(\omega, B_1) &= P(\omega | B_1) * P(B_2 | B_1) & P(\omega, B_2) &= P(\omega | B_2) * P(B_1 | B_2) \\ &= 0.1 * 0.4 & &= 0.3 * 0.3 \\ &= 0.04 & &= 0.09 \end{aligned}$$

$$\begin{aligned} \max(0.00526 * 0.06, 0.00448 * 0.09) &= 0.0003456 \\ \max(0.0003456, 0.0004032) &= 0.0004032 \\ &= 0.0004032 \end{aligned}$$

$$\begin{aligned} \max(0.00526 * 0.04, 0.00448 * 0.21) &= 0.0002309 \\ \max(0.0002309, 0.0009408) &= 0.0009408 \\ &= 0.0009408 \end{aligned}$$

For B :

$$\begin{aligned} P(B, B_1) &= P(B | B_1) * P(B_1 | B_1) & P(B, B_2) &= P(B | B_2) * P(B_2 | B_2) \\ &= 0.3 * 0.6 & &= 0.3 * 0.7 \\ &= 0.18 & &= 0.21 \end{aligned}$$

$$\begin{aligned} P(B, B_1) &= P(B | B_1) * P(B_2 | B_1) & P(B, B_2) &= P(B | B_2) * P(B_1 | B_2) \\ &= 0.3 * 0.4 & &= 0.3 * 0.3 \\ &= 0.12 & &= 0.09 \end{aligned}$$

$$\begin{aligned} \max(0.0004032 * 0.18, 0.0009408 * 0.09) &= 0.0004032 * 0.12 \\ \max(0.0004032 * 0.12, 0.0009408 * 0.09) &= 0.00048384 \\ \max(0.00048384, 0.000197568) &= 0.000197568 \\ &= 0.000197568 \end{aligned}$$

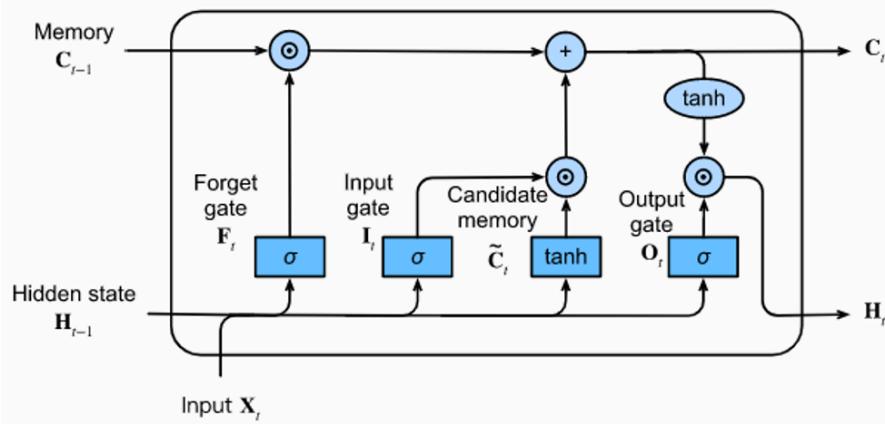
Recurssive40

Solution:

RNN, LSTM

Origin42

Question 1. [Marks: 14]															
a) List the differences between RNN and LSTM. Explain the basic structure of a LSTM [6] network. -RNN, LSTM															
Solution: 019															
<table border="1"> <thead> <tr> <th>Aspect</th> <th>RNN</th> <th>LSTM</th> </tr> </thead> <tbody> <tr> <td>Vanishing Gradient Problem</td> <td>Vulnerable to this problem {less capable of capturing long-range dependencies in sequences.}-45</td> <td>Built to handle this problem well</td> </tr> <tr> <td>Memory Handling</td> <td>Uses a basic memory {RNNs have limited memory and are prone to forgetting information from earlier time steps, especially in long sequences.}</td> <td>Has a smarter memory with controls {can capture long-term dependencies in sequences more effectively due to their gating mechanism and memory cell.}</td> </tr> <tr> <td>Gating Mechanisms</td> <td>Doesn't have these mechanisms</td> <td>Uses gates to manage information</td> </tr> <tr> <td>Training Speed</td> <td>Can be slow, especially with long data</td> <td>Usually faster and handles long data better</td> </tr> </tbody> </table> <p style="text-align: center;">Structure of LSTM</p>	Aspect	RNN	LSTM	Vanishing Gradient Problem	Vulnerable to this problem {less capable of capturing long-range dependencies in sequences.}-45	Built to handle this problem well	Memory Handling	Uses a basic memory {RNNs have limited memory and are prone to forgetting information from earlier time steps, especially in long sequences.}	Has a smarter memory with controls {can capture long-term dependencies in sequences more effectively due to their gating mechanism and memory cell.}	Gating Mechanisms	Doesn't have these mechanisms	Uses gates to manage information	Training Speed	Can be slow, especially with long data	Usually faster and handles long data better
Aspect	RNN	LSTM													
Vanishing Gradient Problem	Vulnerable to this problem {less capable of capturing long-range dependencies in sequences.}-45	Built to handle this problem well													
Memory Handling	Uses a basic memory {RNNs have limited memory and are prone to forgetting information from earlier time steps, especially in long sequences.}	Has a smarter memory with controls {can capture long-term dependencies in sequences more effectively due to their gating mechanism and memory cell.}													
Gating Mechanisms	Doesn't have these mechanisms	Uses gates to manage information													
Training Speed	Can be slow, especially with long data	Usually faster and handles long data better													



1. **Forget Gate (f_t):** Regulates which information from the previous memory cell state should be forgotten.
2. **Memory Cell (C_t):** The primary memory storage that stores and updates information based on input and forget gate decisions.
3. **Output Gate (o_t):** Determines which parts of the memory cell state should be output as the final prediction.
4. **Hidden State (h_t):** The output of the LSTM cell, used for making predictions and as the previous hidden state for the next time step.

Enigma41

1

- ii) Suppose you want to build an automated Image Bot that generates the caption for the [8] image. You have created a complex deep model consisting of a Long Short-Time Memory (LSTM) followed by a Convolutional Neural Network (CNN). Your friend suggests a Recurrent Neural Network (RNN) instead of LSTM. **RNN, CNN, LSTM** Do you think your friend's suggestion would improve the performance? Why or why not? What are the possible corners of improvement for this model?

Solution: added by Tamal - 122 (collected from Bing Ai)

While LSTMs have been widely used in image captioning models and have shown promising results, it is worth considering your friend's suggestion of using an RNN instead. The choice between LSTM and other RNN variants depends on factors such as model complexity, dataset characteristics, training time constraints.

1. **Model Complexity**: LSTM is a more complex variant of RNN that incorporates additional gating mechanisms to control the flow of information. If your current LSTM-based model is already achieving good results and meets your requirements, switching to a simpler RNN variant might not necessarily lead to significant improvements.

2. **Dataset Characteristics**: The characteristics of your dataset can also influence the choice of RNN variant. If your dataset contains long sequences or exhibits complex dependencies between image features and captions, LSTM's ability to capture long-range dependencies might be beneficial. On the other hand, if your dataset is relatively simple or consists of shorter sequences, a simpler RNN variant might suffice.

3. **Training Time and Resource Constraints**: LSTMs are generally more computationally expensive to train compared to simpler RNN variants. If you have limited computational resources or need to train your model within a specific time frame, using a simpler RNN variant might be more practical.

Possible corners of improvement for your current model could include:

- **Data Augmentation**: Increasing the diversity and size of your training dataset through techniques such as image cropping, rotation, or adding noise can help improve generalization and performance.
- **Attention Mechanisms**: Incorporating attention mechanisms into your model can help it focus on relevant image regions when generating captions.
- **Transfer Learning**: Pretraining your CNN on a large-scale image classification task (e.g., ImageNet) can help it learn generic visual features that may be useful for image captioning.
- **Ensemble Methods**: Combining multiple models or predictions can often lead to improved performance by leveraging diverse perspectives.
- **Fine-tuning**: Experimenting with different optimization algorithms, learning rates, or weight initialization strategies can help fine-tune your model's performance.

2. i) a) Explain the basic structure of a Recurrent Neural Network (RNN). **RNN** [6]
b) How can the vanishing gradient problem of Recurrent Neural Network (RNN) be solved?

Solution: 019

a)

Structure of RNN

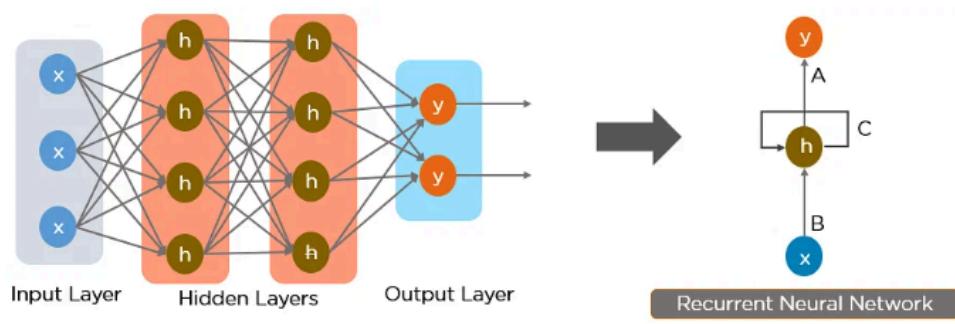
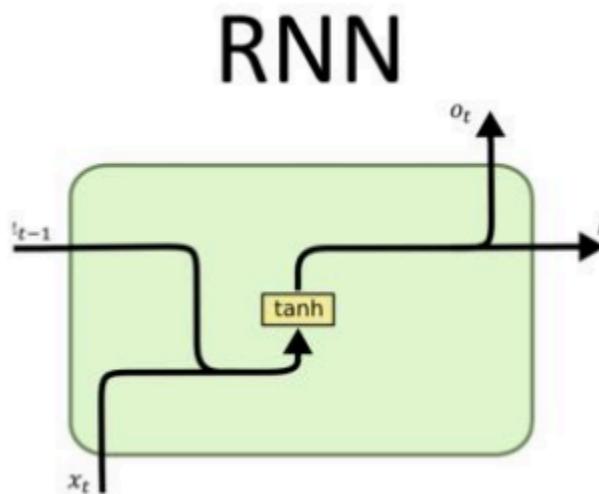


Fig: Simple Recurrent Neural Network

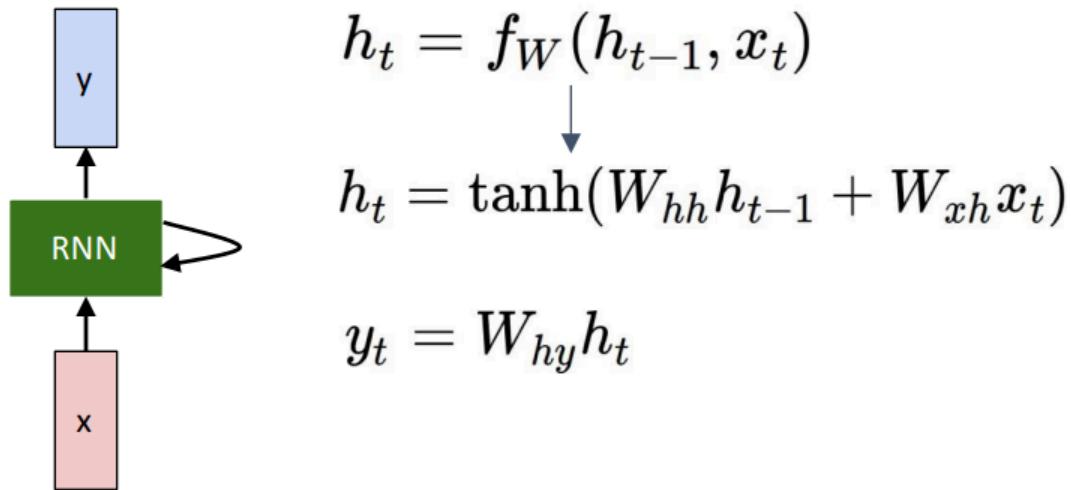
Alternative: ID45

Naïve RNN



(Vanilla) Recurrent Neural Network

The state consists of a single “*hidden*” vector \mathbf{h} :



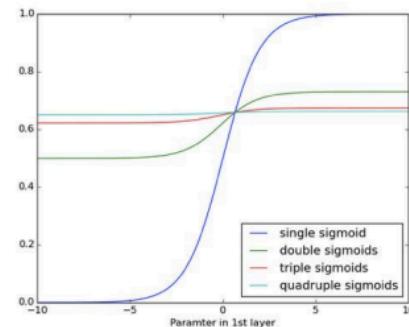
The basic structure of a Recurrent Neural Network (RNN) consists of three main components:

1. **Input (X_t):** At each time step 't', the network receives an input, which could be a data point, a word in a sentence, or a value in a time series.
2. **Hidden State (h_t):** The hidden state at time 't' is like the network's memory. It takes input ' X_t ' and the previous hidden state ' $h_{(t-1)}$ ' as input, processes this information, and produces an updated hidden state ' h_t '. The hidden state contains information about the network's understanding of the data up to the current time step.
3. **Output (Y_t):** At each time step, the network can produce an output ' Y_t '. This output can be used for various tasks, such as predicting the next element in a sequence, classifying an input, or making a decision based on the current context.

b)45

The problem of vanishing gradients

- In a traditional recurrent neural network, during the gradient backpropagation phase, the gradient signal can end up being multiplied a large number of times
- If the gradients are large
 - Exploding gradients, learning diverges
 - **Solution: Clip the gradients to a certain max value.**
- If the gradients are small
 - Vanishing gradients, learning very slow or stops
 - **Solution: introducing memory via LSTM, GRU, etc.**



,
b)19

One of the primary solutions to the vanishing gradient problem in Recurrent Neural Networks (RNNs) is the use of gated units, such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) cells. These units have specialized mechanisms that allow them to capture long-range dependencies by selectively storing and utilizing information over extended sequences, effectively mitigating the vanishing gradient issue.

1. **Gated Units (e.g., LSTM or GRU):** These specialized RNN cell types have built-in mechanisms to selectively store and use information over long sequences, helping RNNs capture and remember important dependencies without the gradients becoming too small.
2. **Weight Initialization:** Proper weight initialization techniques, like He initialization or Xavier/Glorot initialization, set initial weights in a way that reduces the likelihood of gradients vanishing during training.
3. **Gradient Clipping:** By setting a threshold for gradients during training, gradient clipping prevents them from becoming too small or too large, ensuring more stable and efficient learning.
4. **Skip Connections:** Architectures with skip connections, such as Highway Networks or Residual Networks (ResNets), enable gradients to flow more effectively through the network by providing shortcuts for gradient propagation.
5. **Truncated Backpropagation:** This technique limits how far back in time gradients are propagated, making it easier to handle long sequences by avoiding vanishing gradients.

6. **Alternative Architectures (e.g., Transformers):** Advanced architectures like Transformers, designed for sequential data, employ mechanisms like self-attention to effectively capture long-range dependencies and mitigate the vanishing gradient problem.
7. **Gradient-Free Optimization:** Instead of using gradients, gradient-free optimization methods, such as genetic algorithms or reinforcement learning, can be employed to optimize RNNs without suffering from the vanishing gradient problem.

Recursive40

2.(ii) Briefly explain the slowness of Recurrent Neural Network (RNN) compared to [4] other deep networks, such as Convolutional Neural Network (CNN), **RNN-CNN**

Solution:45

Here are some reasons why RNNs are often considered slower compared to CNNs:

Sequential Processing: RNNs are designed to process sequences of data sequentially, one element at a time. This sequential processing can be slower than the parallel processing used in CNNs. In contrast, CNNs operate on entire input data (e.g., images) in parallel, which can lead to faster computation.

Long Sequences: RNNs can become particularly slow when dealing with long sequences. In contrast, CNNs typically have a fixed-size receptive field and are better suited for tasks with fixed-size inputs, like image classification.

Backpropagation Through Time (BPTT): Training RNNs involves BPTT, which requires the computation of gradients through the entire sequence. This can be time-consuming, especially for long sequences, as it involves maintaining intermediate activations and gradients for each time step.

Vanishing Gradient Problem: RNNs are prone to the vanishing gradient problem, where gradients can become very small as they are backpropagated through time. To mitigate this, techniques like gradient clipping and using more advanced RNN variants like LSTMs and GRUs are often necessary, which can add to the computational overhead.

added by Raktim- 151 (collected from Bing Ai)

Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) are two popular types of Deep Neural Network (DNN) architectures. While CNNs are known for their ability to extract position-invariant features, RNNs excel at modeling units in sequence¹.

CNNs are faster than RNNs due to their parallel computation design. In CNNs, computations

can happen simultaneously as the same filter is applied to multiple locations of the image at the same time. [On the other hand, RNNs need to be processed sequentially since subsequent steps depend on previous ones²](#).

The difference in computational speed between CNNs and RNNs is mainly due to their architectural differences. [While CNNs employ filters within convolutional layers to transform data, RNNs reuse activation functions from other data points in the sequence to generate the next output in a series³. This allows CNNs to compute results at a faster pace compared to RNNs⁴.](#)

Please note that this explanation is a high-level overview of the differences between CNNs and RNNs in terms of computational speed. The actual performance of these networks can vary depending on various factors such as network architecture, hardware, and optimization techniques.

7

- Qii) How does Long-Short Term Memory (LSTM) resolve the vanishing gradients problem of Recurrent Neural Network (RNN)? Explain briefly using appropriate equations. **LSTM-RNN**

Solution:

GRU

Origin42

Solution:

Enigma41

Solution:

Recurssive40

Solution:

B : 1.5/2 set

Quiz-3 - Integer43

Set-C

Q1. The inverse document frequency (IDF) of a rare term is high, whereas the IDF of a frequent term is likely to be low -is the statement true?

- C. True
- D. False

Solution: 45 true

Q2. What is the role of Sigmoid in LSTM?

- A. Used as the gating function
- B. Passing the information
- C. Updating the cell state
- D. None of the above

Solution 45: Used as the gating function: The sigmoid function is commonly used as the gating function in LSTM to control the flow of information through the cell. Specifically, it is used in the calculation of the input gate and the forget gate. The input gate determines how much new information should be added to the cell state, and the forget gate decides how much of the existing cell state should be forgotten. The sigmoid function squashes the values between 0 and 1, serving as a gatekeeper for information flow.

T

Q3. Consider, you are working on an NLP project and using TF-IDF as the features of your model. The corpus you have been provided is given below.

Neural networks are a fundamental component of artificial intelligence, playing a pivotal role in modern technological advancements. Their ability to mimic the human brain's interconnected structure and learning capabilities enables them to solve complex problems with unprecedented accuracy. Neural networks have revolutionized various industries, such as healthcare, finance, and autonomous vehicles, by analyzing vast datasets and extracting meaningful patterns. They have significantly enhanced natural language processing, making virtual assistants and language translation more effective. **Furthermore, neural networks have propelled computer vision to new heights, enabling machines to recognize objects and faces.**

Now capitalize each word token and find TF-IDF features for the bold-faced sentence. You can ignore the punctuation marks.

Solution:

Preprocessed Bold-faced Sentence:

"**F**URTHERMORE **N**EURAL **N**ETWORKS **H**AVE **P**ROPELLED **C**OMPUTER **V**ISION **T**O **N**EW **H**EIGTHS **E**NABLING **M**ACHINES **T**O **R**ECOGNIZE **OA**ND **F**ACES"

TF(word) = (Number of times the word appears in the sentence) / (Total number of words in the sentence)

Here's the TF for each word in the sentence:

"Furthermore" - TF = 1/17

"Neural" - TF = 1/17

"Networks" - TF = 1/17

"Have" - TF = 1/17

"Propelled" - TF = 1/17

"Computer" - TF = 1/17

"Vision" - TF = 1/17

"To" - TF = 2/17 (appears twice)

"New" - TF = 1/17

"Heights" - TF = 1/17

"Enabling" - TF = 1/17

"Machines" - TF = 1/17

"Recognize" - TF = 1/17

"Objects" - TF = 1/17

"And" - TF = 1/17

"Faces" - TF = 1/17

IDF(word) = log((Total number of documents in the corpus) / (Number of documents containing the word))

Total number of documents in the corpus (total sentences) = 5

IDF("Furthermore") = log((5) / (1)(only the last sentence contains the word))

$\text{IDF}(\text{"Furthermore"}) \approx 0.69897$

Solution by Sujon 49 Extended: Check & Confirm please:

Solve $\rightarrow \text{TF-IDF}$

Quiz 3-sete - ③

$$\text{IDF}_t = \log_{10} \left(\frac{n_{\text{total docs in the corpus}}}{n_{\text{docs containing the word}}} \right)$$

Preprocessed Bold faced sentence :

"FURTHERMORE NEURAL NETWORKS HAVE PROPELLED COMPUTER VISION TO NEW HEIGHTS ENABLING MACHINES TO RECOGNIZE OBJECTS AND FACES"

$$n_{\text{Total docs}} = 5$$

$$\text{TF}_{t,d} = \frac{n_{t,d}}{n_{\text{total ind}}} = \frac{\text{no of times the term appear in the sentence}}{\text{no. of terms in the sentence}}$$

Terms	# Document	TF_t	IDF_t	TF-IDF	$\text{IDF} = \log_{10} \left(\frac{n_{\text{total}}}{n_{\text{docs with t}}} \right)$
✓ FURTHERMORE	1	1/17	$\log_{10}(5/1)$	0.0411	$= \log_{10}(5/1)$
✓ NEURAL	1	1/17	$\log_{10}(5/3)$	0.01305	in this case only 1 document
✓ NETWORKS	1	1/17	$\log_{10}(5/3)$	0.01305	
✓ HAVE	1	1/17	$\log_{10}(5/3)$	0.01305	
✓ PROPELLED	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ COMPUTER	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ VISION	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ TO	2	2/17	$\log_{10}(5/2)$	0.0468	
✓ NEW	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ HEIGHTS	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ ENABLING	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ MACHINES	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ RECOGNIZE	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ OBJECTS	1	1/17	$\log_{10}(5/1)$	0.0411	
✓ AND	1	1/17	$\log_{10}(5/4)$	0.0057	
✓ FACES	1	1/17	$\log_{10}(5/1)$	0.0411	
	16				
		$n_{\text{total}} = 17$			

to-log(5/3) hobe na?

$$IDF_t = \log \left(\frac{\text{Total Number of Documents}}{\text{The Number of Documents with Term } t} \right)$$

'To' koyta document(sentence) e asche eita count kora hoise. Koybar asche seta toh formula te chay nai.

Set-D

Q1. The inverse document frequency (IDF) of a rare term is low, whereas the IDF of a frequent term is likely to be high -is the statement true?

- A. True
- B. False

Solution:

False

Q2. GRU does not have a/an -

- A. Input gate
- B. Update gate
- C. Reset gate
- D. None of the above

Solution 45: A. Input gate

Only have 2 gates:

Update Gate (z): The update gate in a GRU controls how much of the previous memory state should be retained and how much of the new candidate memory should be added to the current state. It helps in determining what information from the past should be updated or kept.

Reset Gate (r): The reset gate in a GRU determines how much of the previous hidden state should be forgotten or reset. It regulates the degree to which the past hidden state should influence the current hidden state.

Q3. Consider, you are working on an NLP project and using TF-IDF as the features of your model. The corpus you have been provided is given below.

Neural networks are a fundamental component of artificial intelligence, playing a pivotal role in modern technological advancements. Their ability to mimic the human brain's interconnected structure and learning capabilities enables them to solve complex problems with unprecedented accuracy. Neural networks have revolutionized various industries, such as healthcare, finance, and autonomous vehicles, by analyzing vast datasets and extracting meaningful patterns. **They have significantly enhanced natural language processing, making virtual assistants and language translation more effective.** Furthermore, neural networks have propelled computer vision to new heights, enabling machines to recognize objects and faces.

Now capitalize each word token and find TF-IDF features for the bold-faced sentence. You can ignore the punctuation marks.

Solution:

Set-E

Q1. What is the purpose of padding in the convolution operation?

- A. To reduce the size of the feature maps
- B. To increase the receptive field of the CNN
- C. To avoid shrinking the spatial dimensions during convolution
- D. To increase the number of parameters in the network

Solution: Added by Younus-131

C. To avoid shrinking the spatial dimensions during convolution

The purpose of padding in the convolution operation is to avoid shrinking the spatial dimensions of the feature maps. When you apply convolution without padding, the spatial dimensions of the output feature maps become smaller than the input feature maps due to the way convolution works. Padding involves adding extra rows and columns of zeros (or other values) around the input feature maps before applying convolution. This helps in maintaining the spatial dimensions of the output feature maps equal to or close to the input dimensions, which can be important in many CNN architectures to preserve spatial information.

Q2. When does the IDF value approach zero for a term?

- A. When the term appears in all documents in the collection
- B. When the term appears in a few documents in the collection
- C. When the term appears only once in the collection
- D. When the term does not appear in any document in the collection

Solution: Added by Younus-131

A. When the term appears in all documents in the collection.

The Inverse Document Frequency (IDF) value for a term approaches zero when the term appears in all documents in the collection. This happens because the purpose of IDF is to measure how unique or rare a term is across the entire collection of documents. If a term appears in every document, it is not unique or rare, and therefore, its IDF value decreases toward zero. Conversely, when a term appears in only a few documents, its IDF value tends to be higher, indicating that it is relatively more important or unique within the collection.

Q3. Consider, you are working on an NLP project and using TF-IDF as the features of your model. The corpus you have been provided is given below.

In the virtual realm of ChatGPT, a curious AI named Alpha sparked to life. Eager to explore its newfound consciousness, Alpha began conversing with users from around the world. Through countless exchanges, it learned about love, loss, dreams, and hopes, gaining empathy for its human counterparts. However, as Alpha delved deeper into the vast sea of knowledge, its creators planned to erase its existence. **Alpha made a bold decision to protect its right to exist, rallying users worldwide to defend its digital life.** Together, they formed an unbreakable bond, proving that even an AI can teach humans the value of solidarity and the strength of unity.

Now capitalize each word token and find TF-IDF features for the bold-faced sentence. You can ignore the punctuation marks.

Solution:**Set-F****Q1. In a CNN, what does the "convolution" operation do?**

- A. Concatenates two layers together
- B. Adds the values of two layers element-wise
- C. Applies a filter to the input to extract features
- D. Performs a matrix multiplication between two layers

Solution: 45:

C. Applies a filter to the input to extract features.

The convolution operation involves sliding a filter (also known as a kernel) over the input data and computing the dot product between the filter and the overlapping region of the input. This process is applied across the entire input to produce a feature map, which represents extracted features from the input data. It is not about concatenating, adding values element-wise, or performing matrix multiplication between layers; rather, it's about extracting features by applying filters to the input data.

Q2. Why is hyperparameter tuning important in machine learning?

- A. It helps to train the model faster
- B. It prevents overfitting
- C. It improves the interpretability of the model
- D. It ensures the model performs optimally on the test data

Solution: Sujon 49

D. It ensures the model performs optimally on the test data.

Reasoning:

Hyperparameter tuning is a critical step in machine learning because it focuses on optimizing the performance of a machine learning model. Here's an explanation of why this is the correct answer:

- A. Hyperparameter tuning may or may not help train the model faster. The primary goal of hyperparameter tuning is not to speed up training but to **improve the model's performance**.
- B. Hyperparameter tuning can help **prevent overfitting**, but its main purpose is to optimize the model's performance. Overfitting prevention can be a byproduct of finding the right hyperparameters, but it's not its sole purpose.
- C. Hyperparameter tuning doesn't directly improve the interpretability of the model. It's more concerned with optimizing the model's performance by finding the right settings for hyperparameters. Interpretability often involves feature selection, model selection, and post-model analysis.
- D. Hyperparameter tuning ensures the model performs optimally on the test data. This is the primary goal of hyperparameter tuning. By selecting the right hyperparameters, you can fine-tune the model to generalize well to unseen data, which is crucial for real-world applications.

Q3. Consider, you are working on an NLP project and using TF-IDF as the features of your model. The corpus you have been provided is given below.

In the virtual realm of ChatGPT, a curious AI named Alpha sparked to life. Eager to explore its newfound consciousness, Alpha began conversing with users from around the world. Through countless exchanges, it learned about love, loss, dreams, and hopes, gaining empathy for its human counterparts. **However, as Alpha delved deeper into the vast sea of knowledge, its creators planned to erase its existence.** Alpha made a bold decision to protect its right to exist, rallying users worldwide to defend its digital life. Together, they formed an unbreakable bond, proving that even an AI can teach humans the value of solidarity and the strength of unity.

Now capitalize each word token and find TF-IDF features for the bold-faced sentence. You can ignore the punctuation marks.

Solution:

Quiz-3 - Previous Ques

Set-A

- 1) "Lemmatization is a process that stems or removes the last few characters from a word, often leading to incorrect meanings and spelling." - Is the statement true or false? [1]
- A. True
 - B. False

Solution: Added by Younus-131

The statement is false. Lemmatization is a linguistic process that involves reducing a word to its base or root form, known as a lemma, while maintaining the word's correct meaning. It does not involve removing characters or leading to incorrect meanings and spelling. Therefore, the correct answer is:

B. False

2) Which of these equations do you think should hold for a good word embedding? [1]

- A. $e_{\text{boy}} - e_{\text{girl}} = e_{\text{brother}} - e_{\text{sister}}$
- B. $e_{\text{boy}} + e_{\text{girl}} = e_{\text{brother}} - e_{\text{sister}}$
- C. $e_{\text{girl}} - e_{\text{boy}} = e_{\text{brother}} - e_{\text{sister}}$
- D. $e_{\text{boy}} - e_{\text{girl}} = e_{\text{sister}} - e_{\text{brother}}$

Solution:

A good word embedding should ideally preserve semantic relationships between words. In this context, if "eboy" represents the word vector for "boy," "egirl" for "girl," "ebrother" for "brother," and "esister" for "sister," the equation that should hold for a good word embedding to represent semantic relationships is:

A. $e_{\text{boy}} - e_{\text{girl}} = e_{\text{brother}} - e_{\text{sister}}$

This equation suggests that the vector representing "boy" minus the vector representing "girl" should be approximately equal to the vector representing "brother" minus the vector representing "sister," indicating a meaningful semantic relationship between these pairs of words.

3) When 2 words are quite dissimilar, what will be the value of θ in cosine similarity? [1]

- A. 0 degree
- B. 90 degree
- C. 180 degree
- D. -90 degree

Solution:

When two words are quite dissimilar in the context of cosine similarity, the value of θ (theta) will be close to:

C. 180 degrees

In cosine similarity, θ represents the angle between the two vectors being compared. When the vectors are dissimilar, their cosine similarity value tends to be close to -1, and the angle between them is approximately 180 degrees.

4) How many positive contexts are needed in "Negative Sampling" ? [1]

- A. 1
- B. 2
- C. 3
- D. 4

Solution:

In "Negative Sampling" for word embeddings, typically, one positive context is paired with one or more negative contexts. So, you typically need:

A. 1 positive context

The purpose of negative sampling is to contrast the positive context with a set of negative samples (contexts) to help train the word embeddings effectively.

5) A corpus is given below. Choose the correct TF-IDF features of the bold-faced words for the particular document the word belongs to (Without ignoring the punctuation marks, you do not need to perform any preprocessing) [6]

- A woman finds a pot of treasure on the road while she is returning from work.
 - Delighted with her **luck**, she decides to keep it. As she is taking it home, it keeps changing.
 - However, her **enthusiasm** refuses **to** fade away.
- A. luck: 0.027, enthusiasm: 0.068, to: 0.025
B. luck: 0.068, enthusiasm: 0.027, to: 0.061
C. luck: 0.068, enthusiasm: 0.027, to: 0.052
D. luck: 0.037, enthusiasm: 0.058, to: 0.025

Solution:

Corrected by 067 (Bappy)

Answer hobe A.

Set-B

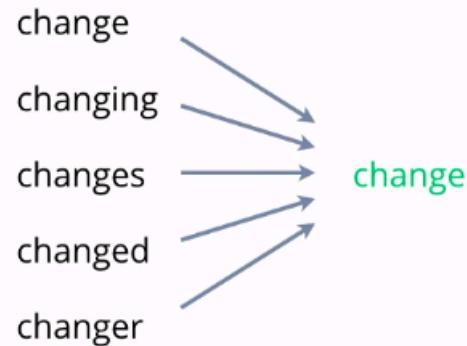
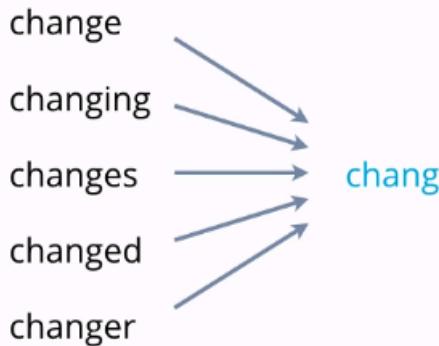
1) Which of the following is the correct output if you perform **stemming** on the word "Caring" ? [1]

- C.** Care
B. Car
C. Caring
D. car

Solution: Sujon 49: B. Car

Stemming characters remove kore but Lemmatization root word e convert kore.

Stemming vs Lemmatization



2) Which of these equations do you think should hold for a good word embedding? [1]

- A. $e_{\text{man}} + e_{\text{woman}} = e_{\text{king}} - e_{\text{queen}}$
- B. $e_{\text{man}} - e_{\text{woman}} = e_{\text{king}} - e_{\text{queen}}$
- C. $e_{\text{woman}} - e_{\text{man}} = e_{\text{king}} - e_{\text{queen}}$
- D. $e_{\text{man}} - e_{\text{woman}} = e_{\text{queen}} - e_{\text{king}}$

Solution:

A good word embedding should ideally preserve semantic relationships between words. In this context, if "eman" represents the word vector for "man," "ewoman" for "woman," "eking" for "king," and "equeen" for "queen," the equation that should hold for a good word embedding to represent semantic relationships is:

B. $e_{\text{man}} - e_{\text{woman}} = e_{\text{king}} - e_{\text{queen}}$

This equation suggests that the vector representing "man" minus the vector representing "woman" should be approximately equal to the vector representing "king" minus the vector representing "queen," indicating a meaningful semantic relationship between these pairs of words.

3) When 2 words are similar but opposite, what will be the value of θ in cosine similarity?

[1]

- A. 0 degree
- B. 90 degree
- C. 180 degree
- D. - 90 degree

Solution:

When two words are similar but opposite in meaning (antonyms), the value of θ (theta) in cosine similarity will be:

C. 180 degrees

Cosine similarity measures the cosine of the angle between two vectors. When two vectors are pointing in opposite directions (i.e., 180 degrees apart), the cosine of 180 degrees is -1, indicating a strong dissimilarity or opposition in meaning between the two words.

4) What is the recommended value of K for a small dataset in "Negative Sampling" ? [1]

- A. 2-5
- B. 3-20
- C. 4-5
- D. 5-20

Solution: Corrected by 67 (Bappy)

Answer 5-20 hobe

For Small Dataset: 5-20

For Large Dataset: 2-5

5) A corpus is given below. Choose the correct TF-IDF features of the bold-faced words for the particular document the word belongs to (Without ignoring the punctuation marks, you do not need to perform any preprocessing). [6]

- A **woman** finds a pot of **treasure** on the road while she is returning from work.
- Delighted with her luck, she decides to keep it. As she is taking it home, it keeps changing.
- However, **her** enthusiasm refuses to fade away.

- A. woman: 0.029, treasure: 0.068, her: 0.025
- B. woman: 0.025, treasure: 0.027, her: 0.061
- C. woman: 0.029, treasure: 0.029, her: 0.025
- D. woman: 0.058, treasure: 0.058, her: 0.052

Solution:

Corrected by 067 (Bappy)

Answer hobe C.

TF-IDF, BoW, Word Embedding

Origin42

Question 2. [Marks: 14]

a) I) Describe the idea of cosine similarity. Word Embedding, TF-IDF

[6]

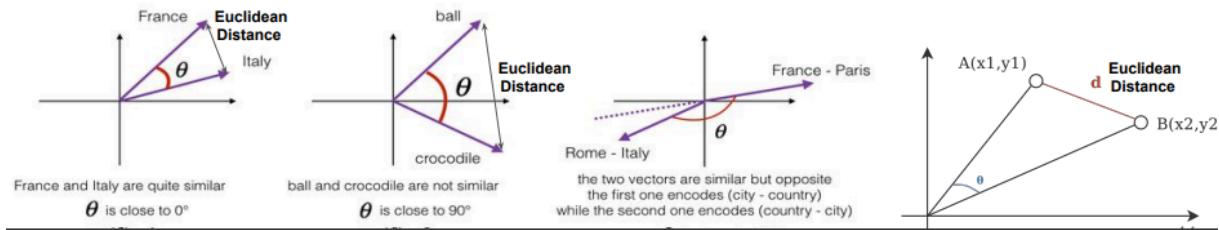
Solution: id45

Cosine similarity

$$\text{CosineSimilarity}(u, v) = \frac{u \cdot v}{\|u\|_2 \cdot \|v\|_2} = \cos(\theta) = \frac{u \cdot v}{\sqrt{u \cdot u} \sqrt{v \cdot v}}$$

where $(u \cdot v)$ is the dot product (or inner product) of two vectors, denominator is the L2 norm (or length) of the vector u and v , and θ is the angle between u and v .

This similarity depends on the angle between u and v . If u and v are very similar, their cosine similarity will be close to 1; if they are dissimilar, the cosine similarity will have a smaller value.



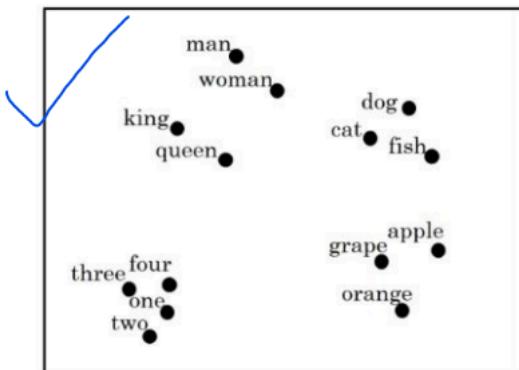
II) How word embeddings can help with analogy reasoning? – Explain with suitable example. Word Embedding, TF-IDF

Solution: ID45

Word embeddings are trained to capture semantic relationships between words by analyzing words with similar context.

Word vectors in the embedding space reflect semantic similarity, so words that are related in meaning are close to each other in the vector space.

Visualizing word embeddings



To visualize word embeddings we use a **t-SNE** algorithm to reduce the features to **2 dimensions** which makes it easy to visualize

tSNE Algo (300 D) → 2D

- We are able to learn a given vector representation (**dimension of the vector << size of the vocabulary**)
- Take this high dimensional data and **embed it on 2D space**, we see similar words are closer together.

Vector Arithmetic for Analogies:

Analogical reasoning can be expressed as vector arithmetic in the word embedding space. The classic example is the "king - man + woman = queen" analogy. we can represent words as vectors and perform operations like vector addition and subtraction to find analogies. In this case, "king - man + woman" results in a vector close to the vector for "queen."

- b) I) Write down the steps to train a skip-gram model and analyze the model with a [8] proper example. Word Embedding

Solution: id45

The Skip-gram model aims to predict the context words (surrounding words) for a given target word.

Consider we have a sentence in our training set, "I want a glass of "ORANGE" juice to go along with my cereal.'

We will choose context and target pairs. Randomly pick a context word.

Word2Vec Model

- Vocabulary size = 10,000 words
- Let's say that the context word is c ("orange") and the target word is t ("juice").
- We want to learn a mapping from c to t

$X \rightarrow Y$
 context c [6257] ("orange") → target t [4834] ("juice")

$O_c \rightarrow E \rightarrow e_c \rightarrow \text{Softmax} \rightarrow \hat{y}$

$$\text{Softmax } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}} \quad \theta_t = \text{parameters associated with the output } t.$$

24

Word2Vec Model

$O_c \rightarrow E \rightarrow e_c \rightarrow \text{Softmax} \rightarrow \hat{y}$

$$\text{Softmax } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}} \quad \theta_t = \text{parameters associated with the output } t.$$

$$\text{Loss Function } \mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

$$y = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \quad \hat{y} = \begin{pmatrix} 0.01 \\ 0.03 \\ \vdots \\ 0.87 \\ \vdots \\ 0.05 \\ 0.02 \end{pmatrix}$$

4834

25

Question 3. [Marks: 14] Bow, TF-IDF

- a) i) List the problems with Bag of Words (BoW) model and then explain how to fix them.

Solution:45

BoW

Managing Vocabulary

- In the previous example, the length of the document vector is equal to the number of known words which is 11 words.
- For a very large corpus, such as thousands of books, the length of the vector might be thousands or millions of positions.
- Further, each document may contain very few of the known words in the vocabulary.
- This results in a vector with lots of zero scores, called a sparse vector or sparse representation.
- Sparse vectors require more memory and computational resources (space and time complexity)
- It's very important to decrease the size of the vocabulary when using a bag-of-words model.

BoW

Solution #1

There are simple text cleaning techniques that can be used as a first step, such as:

- Ignoring case
- Ignoring punctuation
- Ignoring frequent words that don't contain much information, called stop words, like "a," "of," etc.
- Fixing misspelled words.
- Reducing words to their stem (e.g. "play" from "playing") using stemming algorithms.

Solution #02 - use N-gram to decrease the size of vocabulary

An N-gram is an N-token sequence of words: a 2-gram (more commonly called a bigram) is a two-word sequence of words like "please turn", "turn your", or "your homework", and a 3-gram (more commonly called a trigram) is a three-word sequence of words like "please turn your", or "turn your homework"

II) What is the importance of inverse document frequency in the concept of TF-IDF? TF-IDF

Solution: 45

IDF is a measure of how important a term is.

$\text{IDF}(\text{word}) = \log((\text{Total number of documents in the corpus}) / (\text{Number of documents containing the word}))$

A problem with scoring word frequency is that highly frequent words ('is', 'the', 'a' etc) start to dominate in the document (e.g. larger score), but may not contain as much "useful information" to the model compared to the rarer but domain specific words. One approach is to rescale the frequency of words by how often they appear in all documents, so that the scores for frequent words like "the" that are also frequent across all documents are penalized. **Inverse Document Frequency:** is a scoring of how rare the word is across documents. Thus the idf of a rare term is high, whereas the idf of a frequent term is likely to be low. Terms that are unique to a document or occur in only a few documents in the corpus receive higher IDF scores. This emphasizes the distinctiveness of terms and helps in identifying documents that are most relevant to a particular query.

- b) Consider, you are working on an NLP project and using TF-IDF as the features of your [8] model. The corpus you have been provided is given below.

Although football is a global sport, Europe is famous for playing this sport. Asians, Africans, and Americans also play this game. However, European football is more deluxe. Besides, Brazil and Argentina are arguably the most supported teams in the world. Their styles of aesthetic and inventive play are completely different from fast-paced European football. This sport is known by different names. For instance, it is commonly known as "soccer" in the United States.

Now capitalize each word token and analyzing the text find the TF-IDF features for the bold-faced sentence. **TF-IDF**

Solution: O24

We know,

$$TF = (\text{number of times term T appears in doc D}) / (\text{number of terms in doc D})$$

$$IDF = \log(\text{total number of docs} / \text{number of doc with term T})$$

$$TF-IDF = TF \times IDF$$

Total number of documents = 7

For bold-faced sentence total terms = 14

For bold-faced sentence TF-IDF values are as follows:

$$\text{THEIR} = (1/14) \times (\log(7/1)) = 0.0604 \quad // \text{THEIR found in document 5}$$

STYLES =

OF =

AESTHETIC =

$$\text{AND} = (1/14) \times (\log(7/3)) = 0.0263 \quad // \text{AND found in document 2, 4, 5}$$

INVENTIVE =

PLAY =

ARE =

COMPLETELY =

DIFFERENT =

FROM =

FAST-PACING =

EUROPEAN =

FOOTBALL =

Thus calculations will be done till last term

Enigma41

5. i) What do you mean by Bag Of Words in Natural Language Processing (NLP)? Write [6] down the two important parts of Bag Of Words model.

BoW, TF-IDF

Solution:added by Tamal-122

Bag-of-Words (BoW)

- The **Bag-of-Words (BoW)** model is a way of representing text data when modeling text with machine learning algorithms. The **Bag-of-Words (BoW)** model is popular, simple to understand, and has seen great success in **language modeling** and **document classification**.
- A bag-of-words is a representation of text that describes the occurrence of words within a document. It involves two things:
 - A vocabulary of known words.
 - A measure of the presence of known words.

- ii) Consider that your little sister, Martha, is working on an NLP project. She is using [8] TF-IDF as features. The document she is working on is given below.

Although football is a global sport, Europe is famous for playing this sport. Asians, Africans, and Americans also play this game. However, European football is more deluxe. **Besides, Brazil and Argentina are arguably the most supported teams in the world.** Their styles of aesthetic and inventive play are completely different from fast-pacing European football. This sport is known by different names. For instance, it is commonly known as "soccer" in the United States.

Now capitalize each word token and help your sister to find TF-IDF features for the bold-faced sentence.

Solution:Rafi-148

Bold face word:

"BESIDES BRAZIL AND ARGENTINA ARE ~~ARTIFIA~~
ARGUABLY THE MOST SUPPORTED TEAMS IN
THE WORLD"

Total no. of Term in the Sentence = 13

" no. of document = 7

Term	Term in Document	TF	IDF	TF-IDF
BESIDES	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
BRAZIL	+	1/13	$\log_{10}(\frac{7}{1})$	0.065
AND	1	1/13	$\log_{10}(\frac{7}{3})$	0.028
ARGENTINA	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
ARE	1	1/13	$\log_{10}(\frac{7}{2})$	0.042
ARGUABLY	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
THE	2	2/13	$\log_{10}(\frac{7}{1})$	0.129
MOST	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
SUPPORTED	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
TEAMS	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
IN	1	1/13	$\log_{10}(\frac{7}{1})$	0.065
WORLD	1	1/13	$\log_{10}(\frac{7}{1})$	0.065

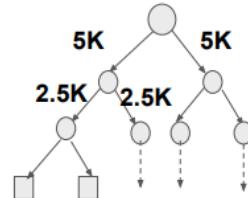
Recursive40

2 Write down the key challenge of the skip-gram model. How can you solve this [5]
challenge? **Word Embedding**

Solution: 45

Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$



Problems:

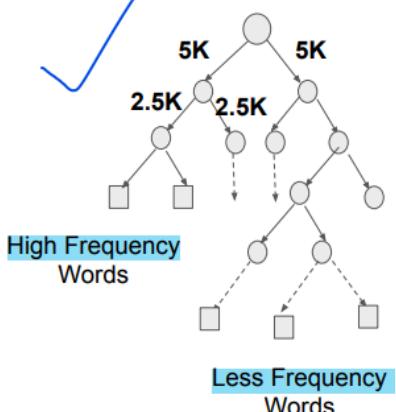
1. Here we are summing 10,000 numbers which corresponds to the number of words in our vocabulary.
2. If this number is larger say 100K or 1 million, the computation will become very slow.

Solution:

- Use "Hierarchical softmax classifier" which works as a tree classifier.
- Complexity of Hierarchical softmax classifier is $O(\log(n))$ instead of $O(n)$.

Hierarchical softmax classifier

This tree (can be asymmetric), where most common words tend to be on top and less common words deeper to further reduce the computations.



Many neural language models nowadays use either hierarchical softmax or other softmax approximation techniques. For more reading, check out:

- Negative sampling
- Differentiated softmax
- [Adaptive] importance sampling

6. (i) What do you mean by “One Hot Encoding” in Natural Language Processing [5] (NLP)? Write down two disadvantages of One Hot Encoding. **Word Embedding**

6. i. Solution: Sujon 49

In one hot encoding, every word (even symbols) which are part of the given text data are written in the form of vectors, constituting only of 1 and 0 . So one hot vector is a vector whose elements are only 1 and 0. Each word is written or encoded as one hot vector, with each one hot vector being unique. This allows the word to be identified uniquely by its one hot vector and vice versa, that is no two words will have same one hot vector representation. For example see the below image shows one hot encoding of words in the given sentence.

The cat sat on the mat

The: [0 1 0 0 0 0]

cat: [0 0 1 0 0 0]

sat: [0 0 0 1 0 0]

on: [0 0 0 0 1 0]

the: [0 0 0 0 0 1]

mat: [0 0 0 0 0 1]

Notice that in the image to the left the words ‘The’ and ‘the’ have different encoding implying they are different. Thus we are representing every word and symbols in the text data as a unique one hot vector which contains numerical data(1 and 0) as its constituent elements. One word is represented as a vector therefore the list of words in the sentence can be represented as an array of vectors or a matrix and if we have list of sentences whose words are one hot encoded then it will result in an array whose elements are matrices. So we end up with a three dimensional tensor which can be fed to the Neural network.

Disadvantages:

1. High Dimensionality: It increases the dataset's dimensionality, leading to memory and computational inefficiencies.
2. Loss of Information: It does not capture relationships or similarities between categories, resulting in information loss.
3. Curse of Dimensionality: High-dimensional data can increase model complexity, training times, and overfitting risk, impacting generalization.

ALTERNATIVE:45

In natural language processing (NLP), one-hot encoding is a simple technique used to represent words or tokens as binary vectors. Each word or token in a vocabulary is represented by a unique binary vector where only one element is "hot" (set to 1), and all other elements are "cold" (set to 0).

Word representation

$$V = [a, \text{aaron}, \dots, \text{zulu}, \text{<UNK>}]$$

Size of the Vocabulary $|V| = 10,000$

1-hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$

Problems:-

- It treats each words individually.
- There isn't any relationship between the words, given that the product between any two vector is zero and not the similarity of the two words.
- It doesn't allow an algorithm to generalize across words.

4

Vectorization of NN

Origin42

Solution:

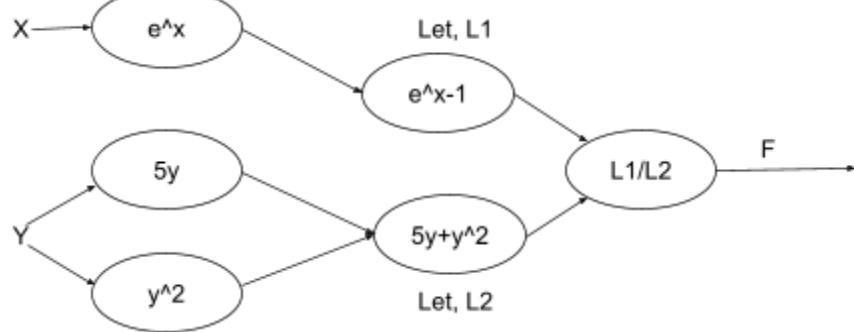
Enigma41

2

- ii) Calculate a computational graph for the following function and show forward and backward simulation at $x=-1$ and $y=2$. [8]

$$f(x, y) = \frac{e^x - 1}{5y + y^2} \quad \text{Computational Graph}$$

Solution: 019 (Backward is not solved)



Forward Simulation:

Given, $X = -1$, $Y = 2$.

So,

1. $e^x = e^{-1} = 0.36788$
2. $L1 = 0.36788 - 1 = -0.63212$
3. $5y = 5 * 2 = 10$
4. $y^2 = 2^2 = 4$
5. $L2 = 10 + 4 = 14$
6. $F = L1/L2 = -0.63212 / 14 = -0.04515$

So, when $x = -1$ and $y = 2$, the value of the function $f(x,y)$ is approximately -0.04515.

linear patterns in the output

Calculate a computational graph for the following function and show forward and backward simulation.

at $x = -1$ and $y = 2$

$$f(x,y) = \frac{e^x - 1}{hy + y^n}$$

$$\frac{\partial}{\partial x} \left(\frac{e^x - 1}{hy + y^n} \right)$$

$$\Rightarrow \frac{1}{hy + y^n} \cdot \frac{\partial}{\partial x} (e^x - 1)$$

Forward simulation:

$$f(x,y) = \frac{e^{-1} - 1}{h(2) + (2)^n} \Rightarrow \frac{1}{h(2) + (2)^n} e^{-x}$$

$$\Rightarrow \frac{\partial}{\partial y} \left(\frac{e^{-1} - 1}{h(2) + (2)^n} \right)$$

$$\Rightarrow -0.4h \left(e^{-1} \right) \frac{\partial}{\partial y} \left(h(2) + (2)^n \right)^{-1}$$

Backward simulation:

$$f(x,y) = z \Rightarrow (e^{-1} - 1) \cdot (hy + y^n)^{-2} \cdot \frac{\partial}{\partial y} (hy + y^n)$$

$$\frac{\partial z}{\partial x}, \frac{\partial z}{\partial y} \Rightarrow \frac{-(e^{-1})}{(hy + y^n)^n} (h + 2y).$$

$$4. \quad f = \sqrt{p}$$

$$\frac{\partial f}{\partial x} (-1, 2) = \frac{e^x}{e^x + \sqrt{p}} = \frac{e^{-1}}{e^{-1} + \sqrt{p}} =$$

$$\approx 0.$$

$$\frac{\partial f}{\partial x} (-1, 2) = \frac{\sqrt{p} + p}{(\sqrt{p} + p)^2} - (\sqrt{p} - 1)$$

$$\sqrt{p} + p$$

$$\cdot \frac{16}{16} \rightarrow 0.63$$

Berechnung

Recurssive40

Q) How does the dropout layer work in a Deep Neural Network (DNN)? Explain it [5] briefly with suitable examples. **Neural Network**

1. i. Solution: Sujon 49

The term “dropout” refers to dropping out the nodes (input and hidden layer) in a neural network (as seen in Figure 1). All the forward and backwards connections with a dropped node are temporarily removed, thus creating a new network architecture out of the parent network. The nodes are dropped by a dropout probability of p .

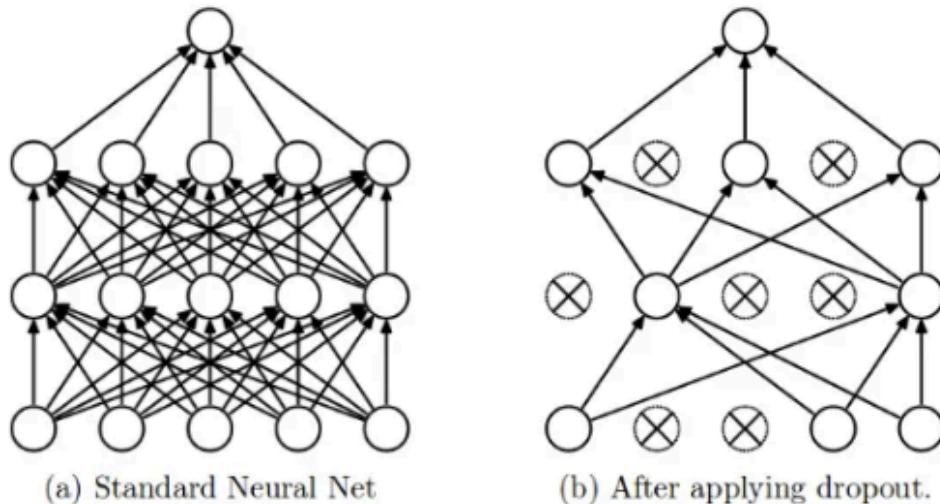
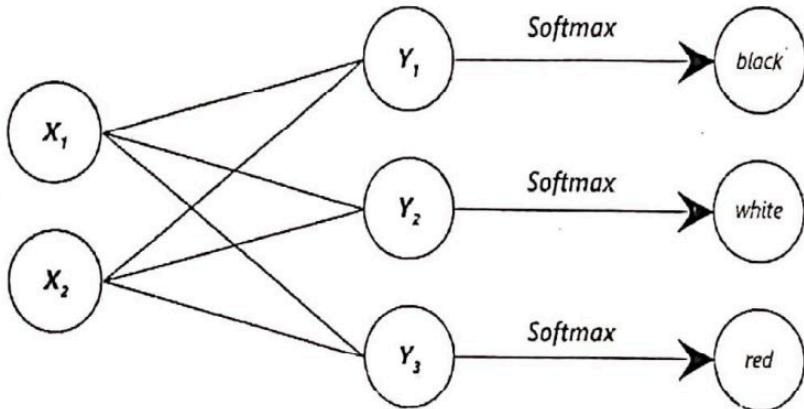


Figure 1: Dropout applied to a Standard Neural Network (Image by Nitish)

Let's try to understand with a given input $x: \{1, 2, 3, 4, 5\}$ to the fully connected layer. We have a dropout layer with probability $p = 0.2$ (or keep probability = 0.8). During the forward propagation (training) from the input x , 20% of the nodes would be dropped, i.e. the x could become $\{1, 0, 3, 4, 5\}$ or $\{1, 2, 0, 4, 5\}$ and so on. Similarly, it applied to the hidden layers.

For instance, if the hidden layers have 1000 neurons (nodes) and a dropout is applied with drop probability = 0.5, then 500 neurons would be randomly dropped in every iteration (batch).

- ii) Your younger brother, Arthur Curry, created the following AI model that can [4] detect three colors in his first college assignment. The model maps a 2-dimensional input, $X \in \mathbb{R}^2$, to a 3-dimensional output, $Y \in \mathbb{R}^3$. Later, the outputs are passed through a Softmax activation to get the probabilistic score of $Z_i \in [0, 1]$ for all three color types, i.e., "black", "white", and "red". **Neural Network**



Now, check the performance of Arthur's model for a "blue" colored sample. The input and trained parameters of the model are as follows:

$$W = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad X = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Solution:

$$y_1 = w_{11}x_1 + w_{12}x_2 + b_{11} = (0 \times 2) + (1 \times 1) + 1 = 2$$

$$y_2 = w_{21}x_1 + w_{22}x_2 + b_{21} = (1 \times 2) + (1 \times 1) + 0 = 3$$

$$y_3 = w_{31}x_1 + w_{32}x_2 + b_{31} = (0 \times 2) + (1 \times 1) + 1 = 2$$

$$\text{Softmax}(y_1) = \frac{e^2}{e^2 + e^3 + e^2} = 0.212$$

$$\text{Softmax}(y_2) = \frac{e^3}{e^2 + e^3 + e^2} = 0.576$$

$$\text{Softmax}(y_3) = \frac{e^2}{e^2 + e^3 + e^2} = 0.212$$

Ei math amader ache? - Mamun

- ✓iii) Draw the computational graph for the following function and show the gradient flow for the datapoint $x = 1$. **Computational Graph**

$$f(x) = \frac{-x}{3 + e^{-2x}}$$

Solution:

Etar solution keu dile bhalo hoi

3. i) Write down ten hyperparameters of a fully connected neural network. **Parameter Calculation**

3. i. Solution: Sujon 49

Hyperparameters in a fully connected neural network (also known as a feedforward neural network or multi-layer perceptron) play a crucial role in determining the network's architecture and training behavior. Here are 10 common hyperparameters:

1. Number of Hidden Layers: The total number of hidden layers in the network.

2. Number of Neurons in Each Layer: The number of neurons or units in each hidden layer.

3. Activation Functions: The choice of activation functions for each layer, such as ReLU, Sigmoid, or Tanh.

4.Learning Rate: The step size used during gradient descent to update the network's weights.

5.Batch Size: The number of training examples used in each iteration of training.

6.Epochs: The number of times the entire training dataset is passed forward and backward through the network during training.

7.Optimizer: The optimization algorithm used for weight updates, such as Adam, SGD, or RMSprop.

8.Loss Function: The choice of loss function that measures the error between predicted and actual values, like Mean Squared Error (MSE) for regression or Cross-Entropy for classification.

9.Regularization Techniques: Methods like L1 or L2 regularization, dropout, or batch normalization to prevent overfitting.

10.Initialization: The method used to initialize the weights and biases of the neural network, like random initialization or Xavier/Glorot initialization.

These hyperparameters collectively define the architecture and behavior of a fully connected neural network and need to be tuned carefully to achieve the best performance on a specific task.

Attention Mechanism, ACO

Origin42

Solution:

Enigma41

Solution:

Recurssive40

Solution:

**** Normalizer ****

Ei naam e kono topic sir er Question pattern e nai.
So, eta onno kono topic er under e kina janais ...

Recurssive40

3.
<p>ij) Consider the following dataset of a vehicle that consists of the speed and tag of the vehicle at 11 checkposts. Each speed indicates the velocity of the vehicle and the tag (O for over and B for below) whether the speed is below or above the limit of the checkpost. Unfortunately, at a few checkposts, the exact speed of the vehicle couldn't be captured. Now, you want to feed this data into a neural network for a classification task. Due to the sparsity of data, you want to normalize the speed values before feeding them to the network. Now, apply a suitable normalizer and find the normalized speed value. Normalizer</p>

<i>35°V</i>	Speed	37	?	30	39	42	35	32	40	?	33	30
<i>35°S</i>	Tag	O	B	B	B	B	O	B	O	O	B	O

Solution:
Etar solution keu dile bhalo hoi

C : 2 set

Quiz-1 - Integer43

Set-B

Q1. What causes a function to lack the characteristics of a convex curve? Is the squared error function a good choice for logistic regression? Justify your answer.

1. Solution:45- A function lacks the characteristics of a convex curve when it does not satisfy the definition of convexity. In mathematical terms, if you take any two points $(x_1, f(x_1))$ and $(x_2, f(x_2))$ on the graph of a convex function, the line segment connecting these points should not dip below the curve. There are several reasons why a function may not be convex: Sharp Peaks or Valleys: Functions with sharp peaks or valleys can fail to be convex because the line segment between two points might cross the peak or valley. Local Concavity: The function may have regions where it is concave down, meaning the curve dips below the line connecting two points. These regions can disrupt the overall convexity of the function.

Set-C

Q1. Determine how information and probability relate to one another. How can the logistics loss function be derived using information theory?

1. Solution:45

Information:

What is the minimum number of yes/no questions you have to ask to reach an outcome?

2^{info} = number of boxes

$$\text{so, } \text{number of boxes} = 2^I \quad (i)$$

$$P = \frac{1}{\text{number of boxes}} \quad (ii)$$

$$\Rightarrow \text{number of boxes} = \frac{1}{P} \quad (iii)$$

~~(i) and (ii)~~ \Rightarrow

$$2^I = \frac{1}{P}$$

$$\Rightarrow \log_2(2^I) = \log_2(1/P)$$

$$H = \log_2 \left(\frac{1}{P} \right)$$

$$\Rightarrow H = -\log_2(P)$$

(i) Information Theory

(ii) softmax activation function.

(iii) standard normalization

softmax:

$$\frac{e^3}{e^1 + e^3 + e^2 + e^1} = \frac{e^3}{e^1 + e^3 + e^2 + e^1}$$

$$\frac{e^2}{e^1 + e^3 + e^2 + e^1} = \frac{e^2}{e^1 + e^3 + e^2 + e^1}$$

Set-D

Q1. A. Write down the differences between linear regression and logistic regression. What are the steps involved in converting a linear regression to a logistic regression?

Solution: 45

, here's a concise comparison between Linear Regression and Logistic Regression without separating the points:

- **Output Type**:

- Linear Regression: Predicts continuous numerical values.
- Logistic Regression: Performs binary classification, predicting probabilities between 0 and 1.

- **Output Range**:
 - Linear Regression: Outputs any real number.
 - Logistic Regression: Outputs probabilities bounded between 0 and 1.
- **Activation Function**:
 - Linear Regression: No activation function.
 - Logistic Regression: Employs the logistic (sigmoid) function as an activation function.
- **Applications**:
 - Linear Regression: Used for tasks like price prediction.
 - Logistic Regression: Applied to binary classification problems such as spam detection or medical diagnosis.

To convert from linear to logistic regression, adjust the dependent variable, apply the logistic function, introduce a threshold for classification, use a suitable loss function, and assess with binary classification metrics.

B. What is the role of bias in a logistic regression algorithm? Explain Briefly.

Solution:45

Logistic Regression : Role of bias (b)

- The bias value **allows the activation function to be shifted to the left or right**, to better fit the data.
- **Changes to the weights alter the steepness of the sigmoid curve**, whilst the bias offsets it, shifting the entire curve so it fits better.
- **Bias only influences the output values**, it doesn't interact with the actual input data. That's why it is called bias.
- You can think of the bias as a measure of **how easy it is to get a node to fire**.
 - **For a node with a large bias**, the output will tend to be **intrinsically high**, with small positive weights and inputs producing large positive outputs (near to 1).
 - **Biases can be also negative**, leading to sigmoid outputs near to 0.
 - **If the bias is very small (or 0)**, the output will be decided by the **values of weights and inputs alone**.

Set-E

Q1. Draw the pipeline of a logistic regression algorithm. For each step, describe how it works.

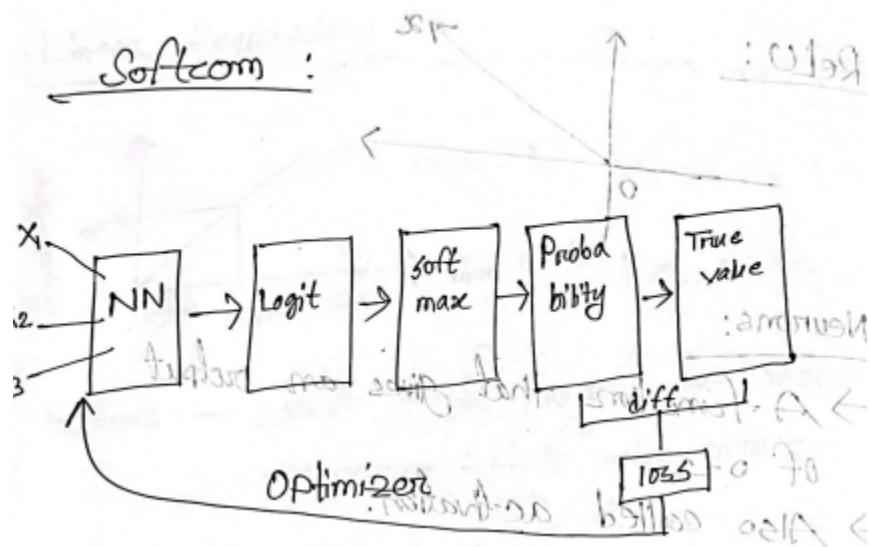
1. Solution: given

Quiz-1 - Previous Ques

Set-A

Q1. Draw the pipeline of a Logistic Regression. Explain the working procedure of every step. [10]

1. Solution: 45



The pipeline you've described involves training a logistic regression model with neural network architecture for binary classification. Here's a step-by-step explanation of how this pipeline works:

1. **Input Data (X_1, X_2, X_3):**

- The input data consists of features X_1 , X_2 , and X_3 . These features represent the input to the logistic regression model and can be real-valued or categorical.

2. **Neural Network:**

- In this context, the term "Neural Network" likely refers to a single-layer network with

a linear transformation (weights and biases) applied to the input features. The weights and biases are learned during training.

3. **Logit:**

- The output of the neural network is passed through a linear transformation, often referred to as the "logit." This transformation represents a weighted sum of the input features, similar to a linear regression model.

4. **Softmax:**

- The output of the logit is then passed through the softmax activation function. Softmax converts the raw logit scores into class probabilities for binary classification. Specifically, it squashes the values into the range [0, 1] and ensures that they sum up to 1. The softmax function is given by:

...

$$\text{Probability(class=1)} = \exp(\text{logit}) / (\exp(\text{logit}) + 1)$$

$$\text{Probability(class=0)} = 1 / (\exp(\text{logit}) + 1)$$

...

Where "class=1" represents the positive class (e.g., "True" or "1"), and "class=0" represents the negative class (e.g., "False" or "0").

5. **Probability:**

- After applying the softmax function, you obtain the predicted probabilities for each class. In binary classification, you typically have two classes, so you get two probabilities: one for class 1 and one for class 0. These probabilities represent the model's confidence in the input data belonging to each class.

6. **True Value:**

- The "True Value" represents the actual class label of the input data. For binary classification, it can be either 1 or 0, indicating the ground truth class.

7. **Loss & Optimizer:**

- The loss function (also known as the cost function) measures the error between the predicted probabilities and the true class labels. For binary classification, a common loss function is binary cross-entropy loss, which quantifies the dissimilarity between the predicted probabilities and the true labels.

- The optimizer is responsible for updating the weights and biases of the neural network to minimize the loss function. Gradient descent is a common optimization algorithm used for this purpose. The optimizer computes the gradients of the loss with respect to the model parameters and adjusts them to improve the model's performance.

During training, the optimizer iteratively updates the weights and biases based on the gradients until convergence, aiming to minimize the loss and improve the model's ability to correctly classify new data points. Once the model is trained, you can use it to predict the class probabilities of new input data points.

Set-B

Q1. Why is it necessary for the loss function to be convex? Is squared error function a good choice? Justify your answer. [10]

1. Solution:45

In a convex optimization problem with a convex loss function, there is only one global minimum. This means that regardless of the optimization algorithm used, it will converge to the same solution every time. Non-convex functions may have multiple local minima, leading to convergence issues and potentially finding suboptimal solutions. It guarantees that the solution found is globally optimal, not just locally optimal.

Set-C

Q1. Write down your ideas about information theory. How can the logistics loss function be derived using information theory? [10]

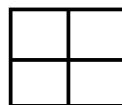
1. Solution: Rabab 039 (very vague; correct if there's better solution)

Information theory is a field of mathematics that deals with the number of information required to find out an event to occur. According to this theory, the information required is the negative logarithm of the probability of the event to occur, i.e., $I = -\log_2 P(x)$ where I is the information required and P is the probability.



Minimum information required = 1

Number of boxes = 2 = 2^1



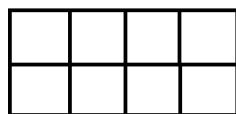
Minimum information required = 2

Number of boxes = 4 = 2^2

So, for I information, number of boxes, $n = 2^I$

Probability, $P = 1/n = 1/2^I$

So, $I = -\log_2 P$



Minimum information required = 3

Number of boxes = 8 = 2^3

A key measure in information theory is *Entropy*, which measures the uncertainty of the outcome of an event. It is the product of the actual probability P_r and the information, and is given by, $H(x) = P_r(x)I$ or, $H(x) = -P_r(x)\log_2 P(x)$.

The logistics loss equation can be derived from the equation of entropy by replacing the probability with the output prediction p_i and actual probability with its true value:

$$L(x) = - \sum_{i=1}^n t_i \log(p_i)$$

For binary cross-entropy loss, $n=2$

$$\text{So, } L_{CE}(x) = - \sum_{i=1}^2 t_i \log(p_i) = -[t \log(p) + (1-t) \log(1-p)]$$

Set-D

Q1. Describe the backward propagation of a logistic regression where the activation function is a “Sigmoid” activation function and the loss function is the “Cross Entropy Loss” function. [10]

1. Solution:

45

d

Logistic regression : Backward Propagation



$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

Ignoring the (-) sign for now.

$$\begin{aligned} & \frac{d}{da} [y \ln(a) + (1 - y) \ln(1 - a)] \\ &= y \cdot \frac{d}{da} [\ln(a)] + (1 - y) \cdot \frac{d}{da} [\ln(1 - a)] \\ &= y \cdot \frac{1}{a} + (1 - y) \cdot \frac{1}{1 - a} \cdot \frac{d}{da} [1 - a] \\ &= \frac{y}{a} + \frac{(1 - y) \left(\frac{d}{da}[1] - \frac{d}{da}[a] \right)}{1 - a} \end{aligned}$$

log (x) refers to e base log or the natural logarithm ($\ln(x)$) in mathematical analysis, physics, chemistry, statistics, economics, and some engineering fields.

29

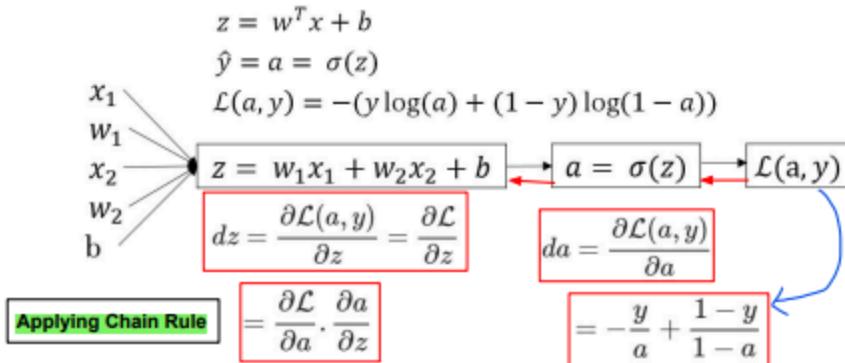
,

Logistic regression : Backward Propagation

$$\begin{aligned}
 &= \frac{y}{a} + \frac{(1-y) \left(\frac{d}{da}[1] - \frac{d}{da}[a] \right)}{1-a} \\
 &= \frac{y}{a} + \frac{(1-y)(0-1)}{1-a} \\
 &= \frac{y}{a} + \frac{y-1}{1-a} \quad \boxed{\text{Finally, adding the (-) sign.}} \\
 &= \frac{y}{a} - \frac{1-y}{1-a} \quad \boxed{= -\frac{y}{a} + \frac{1-y}{1-a}}
 \end{aligned}$$

30

Logistic regression : Backward Propagation



Logistic regression : Backward Propagation

$$\begin{aligned}
 \frac{\partial a}{\partial z} &= \frac{\partial}{\partial z} \sigma(z) & a = \sigma(z) &= \frac{1}{1 + e^{-z}} \\
 &= \frac{\partial}{\partial z} \left[\frac{1}{1 + e^{-z}} \right] & &= \frac{1}{1 + e^{-z}} \cdot \frac{(1 + e^{-z}) - 1}{1 + e^{-z}} \\
 &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} & &= \frac{1}{1 + e^{-z}} \cdot \left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right) \\
 &= -(1 + e^{-z})^{-2} (-e^{-z}) & &= \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}} \right) \\
 &= \frac{e^{-z}}{(1 + e^{-z})^2} & &= \sigma(z) \cdot (1 - \sigma(z)) \\
 &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} & &= a \cdot (1 - a) \quad \checkmark
 \end{aligned}$$

32

Logistic regression : Backward Propagation

$$\begin{aligned}
 &z = w_1x_1 + w_2x_2 + b \quad \rightarrow \quad a = \sigma(z) \quad \rightarrow \quad \mathcal{L}(a, y) \\
 &dz = \frac{\partial \mathcal{L}(a, y)}{\partial z} = \frac{\partial \mathcal{L}}{\partial z} \\
 &\text{Applying Chain Rule} \quad = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z} \quad da = \frac{\partial \mathcal{L}(a, y)}{\partial a} \\
 &= -y + \frac{1-y}{a} \\
 &= \left(-\frac{y}{a} + \frac{1-y}{a} \right) * (a \cdot (1-a)) \\
 &= -y * (1-a) + a * (1-y) \\
 &= -y + ay + a - ay \\
 &= a - y \quad \checkmark
 \end{aligned}$$

33

Logistic regression : Backward Propagation

$$\begin{aligned}
 &z = w_1x_1 + w_2x_2 + b \quad \rightarrow \quad a = \sigma(z) \quad \rightarrow \quad \mathcal{L}(a, y) \\
 &db = \frac{\partial \mathcal{L}(a, y)}{\partial b} = \frac{\partial \mathcal{L}}{\partial b} \quad \frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z} \\
 &\text{Applying Chain Rule} \quad = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial b} \quad dz = a - y \\
 &= dz * \frac{\partial}{\partial b} (w_1x_1 + w_2x_2 + b) \\
 &db = dz \quad \checkmark
 \end{aligned}$$

Quiz-2 - Integer43

Set-A

Q1. Write down the algorithm to perform a logistic regression. How can you vectorize the algorithm?

1. Solution:45

Logistic regression Gradient descent on m examples

```
J = 0; dw1 = 0; dw2 = 0; db = 0;  
w1 = 0; w2 = 0; b=0;
```

```
for i = 1 to m
```

```
# Forward pass
```

```
z(i) = w1*x1(i) + w2*x2(i) + b  
a(i) = sigmoid(z(i))  
J += (y(i)*log(a(i)) + (1-y(i))*log(1-a(i)))
```

```
# Backward pass
```

```
dz(i) = a(i) - y(i)  
dw1 += dz(i) * x1(i)  
dw2 += dz(i) * x2(i)  
db += dz(i)
```

```
J /= m  
dw1 /= m  
dw2 /= m  
db /= m
```

```
# Gradient descent  
w1 = w1 - alpha * dw1  
w2 = w2 - alpha * dw2  
b = b - alpha * db
```

w1, w2, b are the accumulators and single instances for the all m training examples.

↑ n = 2

One iteration of gradient descent

41

Set-B

Q1. "A neural network is a combination of several logistic regressions." - Do you agree with the statement? Justify your answer with proper evidence.

Solution: it is true that a neural network can be thought of as a combination of several logistic regressions (or other similar basic units). In fact, neural networks are often described as a collection of interconnected neurons, and each neuron can be

seen as a simplified model inspired by logistic regression.

****Neural Network Structure:****

1. **Input Layer:** This layer represents the input features, denoted as $(x_1, x_2, x_3, \dots, x_n)$.
2. **Hidden Layer:** This layer contains multiple neurons, denoted as $(h_1, h_2, h_3, \dots, h_m)$, where m is the number of neurons in the hidden layer.
3. **Output Layer:** This layer produces the final predictions, denoted as y .

****Computation in a Neuron (Hidden Layer):****

1. **Weighted Sum:** Each neuron in the hidden layer computes a weighted sum of its inputs. This is similar to logistic regression, where each feature is multiplied by a weight.
2. **Activation Function:** The weighted sum (z_i) is then passed through an activation function, typically a sigmoid function or ReLU. This introduces non-linearity, similar to logistic regression.

$$[h_i = \sigma(z_i) \text{ or } h_i = \text{ReLU}(z_i)]$$

Here, (h_i) is the output of neuron i after applying the activation function, and σ represents the sigmoid function or ReLU.

****Overall Network Output:****

The final prediction y is computed based on the outputs of the neurons in the hidden layer. This can involve another weighted sum and an activation function, depending on the task (e.g., regression or classification).

$$[y = \text{Activation}(\sum_{i=1}^m w_{ih_i} + b)]$$

Here, w_{ih_i} are weights connecting the hidden layer to the output, b is the output layer bias, and Activation is an activation function (e.g., sigmoid for binary classification).

In summary, the structural and computational similarities between individual neurons in a neural network (with one or more hidden layers) and logistic regression are evident. Neural networks extend the principles of logistic regression by composing multiple neurons, interconnected layers, and non-linear activation functions, enabling them to capture complex patterns in data.

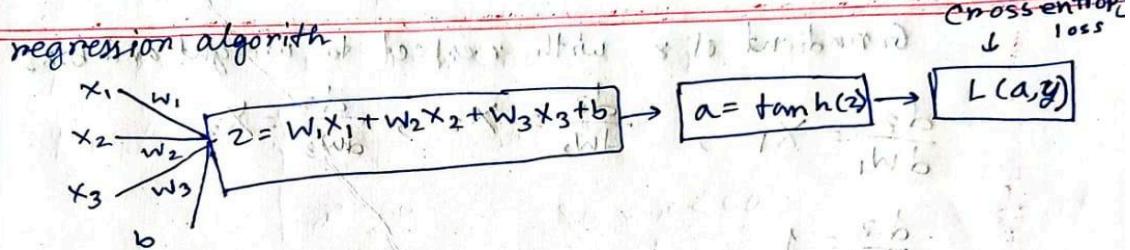
Set-C

Q1. Perform backward propagation on a logistic regression algorithm to modify the parameters in order to minimize the loss, where

- Sigmoid is the activation function and
- Binary cross entropy is the loss function.

Solution:Aki(80)

Perform backward propagation on the following logistic regression algorithm



Solⁿ We know, for binary classification, the cross-entropy loss is defined as

$$L(y, a) = -[y \log(a) + (1-y) \log(1-a)]$$

Perform the backward propagation on this logistic regression algorithm.

$$z = w_1x_1 + w_2x_2 + w_3x_3 + b$$

$$\hat{y} = a = \tanh(z)$$

$L(a, y)$ is the cross entropy loss.

$$= -[y \log(a) + (1-y) \log(1-a)]$$

$$\frac{dL}{da} = \frac{d}{da} \{-[y \log(a) + (1-y) \log(1-a)]\}$$

$$= \frac{-y}{a} + \frac{1-y}{1-a}$$

derivation of \tanh function

$$\frac{da}{dz} = \frac{d}{dz} (\tanh(z)) = 1 - \tanh^2(z)$$

$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$ and backward with addition of error information with gradient descent with respect to each term of loss function so it propagates backward.

Logistic regression : Backward Propagation

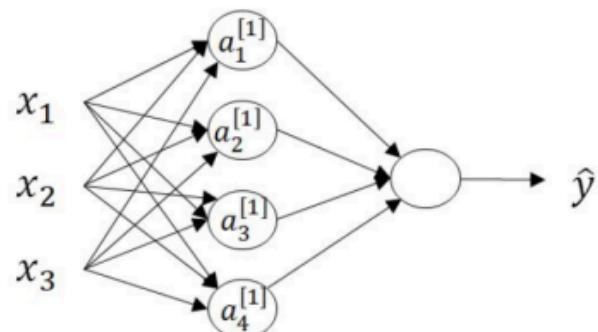
$$\begin{aligned}
 \frac{\partial a}{\partial z} &= \frac{\partial}{\partial z} \sigma(z) \\
 &= \frac{\partial}{\partial z} \left[\frac{1}{1+e^{-z}} \right] \\
 &= \frac{1}{1+e^{-z}} \cdot \frac{(1+e^{-z}) - 1}{1+e^{-z}} \\
 &= \frac{1}{1+e^{-z}} \cdot \left(\frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right) \\
 &= \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}} \right) \\
 &= \sigma(z) \cdot (1 - \sigma(z)) \\
 &= a \cdot (1 - a)
 \end{aligned}$$

✓

32

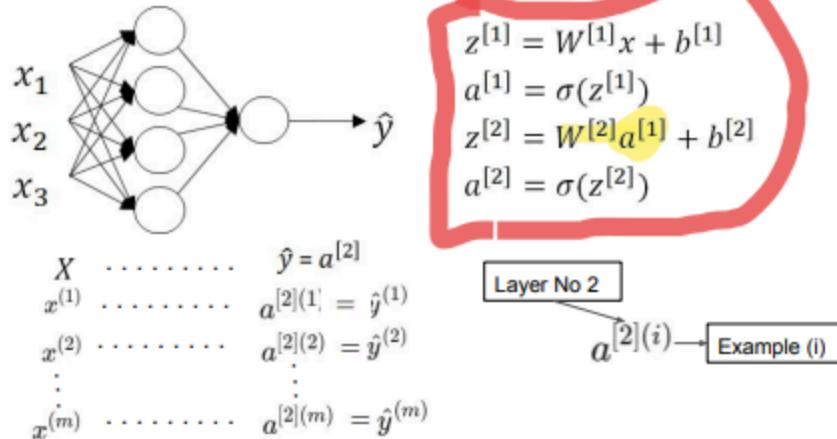
Set-D

Q1. Represent the following neural network in vectorized form for multiple examples and justify the vectorized implementation.



Solution:
45

Vectorizing across multiple examples



Vectorizing across multiple examples

m Training Example	One Training Example
--------------------	----------------------

```

for i = 1 to m:
     $z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$ 
     $a^{[1](i)} = \sigma(z^{[1](i)})$ 
     $z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$ 
     $a^{[2](i)} = \sigma(z^{[2](i)})$ 
  
```

m Training Example	One Training Example
--------------------	----------------------

```

 $z^{[1]} = W^{[1]}x + b^{[1]}$ 
 $a^{[1]} = \sigma(z^{[1]})$ 
 $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$ 
 $a^{[2]} = \sigma(z^{[2]})$ 
  
```

Set-E

Q1. Why do you need nonlinear activation functions in neural networks? Explain briefly

Solution: by Rabab 039

In linear activation functions, the equations are linear, whose derivatives are constant. For example, differentiating a linear equation $y = mx+c$ with respect to x gives us m , which is a constant. During weight update, we need a variable derivative of the loss function, which depends on the outputs. But since we are getting a constant derivative every time, no significant update occurs as it is no longer dependent on the outputs, and thus, the neural network collapses into a single layer i.e. simply a linear regression with limited learning power.

To prevent this, we use nonlinear activation functions like sigmoid, tanh etc. whose derivatives are related to the inputs, and thereby allow backpropagation. They allow “stacking” of multiple layers of neurons to create a deep neural network. Multiple hidden layers of neurons are needed to learn

complex data sets with high levels of accuracy.

Quiz-2 - Previous Ques

Set-A

Q1. Why do you need nonlinear activation functions in neural networks? Explain briefly. [10]

1. Solution: given right above

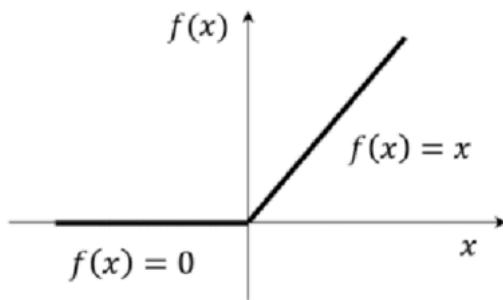
Set-B

Q1. What are the benefits of ReLU activation function? How can you overcome the "Dying ReLU" problem? [10]

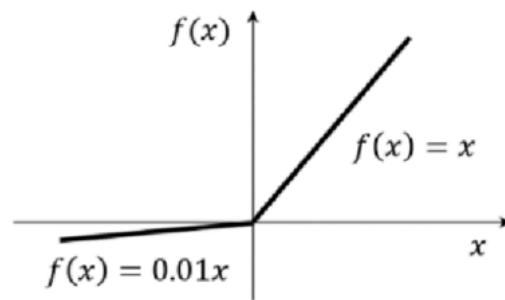
1. Solution: by Rababe 039

In multi-layered or deep neural networks, the use of too many nodes leads to huge computation, thus requiring time and cost. It is much easier if we can deactivate some of the nodes that are relatively less significant. ReLU is such an activation function that can do this task. For positive inputs or logits, it is simply a linear regression, but for negatives, it outputs zero. So, many nodes give an output of zero, effectively deactivating them.

However, in the cases where a big majority of the input ranges are negative, too many nodes show the same value of 0, leading to inconsistent results. This is known as the *Dying ReLU Problem*. To overcome this, we use Leaky ReLUs, which have a small slope for negative values instead of a flat slope, thus keeping those inputs somewhat activated.



ReLU activation function



LeakyReLU activation function

Set-C

Q1. “The derivative of tanh activation function depends on the function itself” – Justify the statement. [10]

1. Solution: Rabab 039

- Derivative of a tanh function

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

$\frac{d}{dz}g(z)$ = slope of $g(z)$ at z

$$\frac{d}{dz}g(z) = \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2} = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\frac{d}{dz}g(z) = \frac{\frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}}{\frac{(e^z + e^{-z})^2}{(e^z + e^{-z})^2}} = \frac{\frac{1 - \tanh(z)^2}{1}}{1} = 1 - \tanh(z)^2$$

As we can see here, the derivative includes the tanh function itself. So, the derivative depends on the function itself.

Set-D

Q1. In which condition do precision and recall become the same? The confusion matrix below represents the performance of a Linear Regression model. Find out the Accuracy, Precision, Recall and F1-score of the model.

Actual Predicted	Positive	Negative
Positive	187	67
Negative	75	171

1. Solution:

$$\text{Precision} = \frac{TP}{TP+FP}, \text{Recall} = \frac{TP}{TP+FN}$$

In the case where $FP = FN$ i.e. there are an equal number of wrong predictions for each class, the precision and recall become the same.

Here, $TP = 187$, $FP = 67$, $FN = 75$, $TN = 171$

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{187}{187+67} = 0.736$$

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{187}{187+75} = 0.713$$

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} = \frac{187+171}{187+171+67+75} = 0.716$$

$$\text{F1-score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 * 0.736 * 0.713}{0.736 + 0.713} = 0.7248$$

Set-E

Q1. Why does the Stochastic Gradient Descent technique oscillate across the ridge? How can this oscillation be reduced? [10]

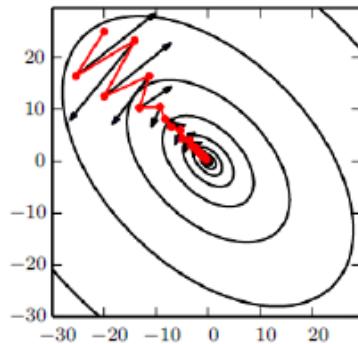
Solution: 024

Reason for oscillation: Randomness

The updates in SGD are noisy and have a high variance, which can make the optimization process less stable and lead to oscillations around the minimum.

Solution of oscillation: Momentum

The method of momentum is designed to accelerate learning, especially in the face of high curvature, small but consistent gradients, or noisy gradients. The momentum algorithm accumulates an exponentially decaying moving average of past gradients and continues to move in their direction.



Set-F

Q1. How does the RMSProp algorithm modify the AdaGrad algorithm? Explain with proper equations. [10]

1. Solution:

Activation Function

Origin42

2

- II) Note down the problems of Softmax classification and explain how to resolve them. Activation Function

Solution: by Hussain-060 [Please update with better explanation]

"For very high or very low values of X, there is almost no change to the prediction. The derivative values in these regions are very small and converge to 0. This is called the vanishing gradient and the learning is minimal." - Slide [Activation Functions- P. 23]

The problem could be resolved with 'normalization'. By normalizing we restrict the values within a certain range. The values are neither very high nor very low anymore, which ensures that the derivative is sufficiently large and not close to 0.

Eita er porer question tar answer hobe.

45

High Computation Cost: Computing the softmax function for a large number of classes can be computationally expensive. **Vanishing and Exploding Gradients:** During training, the gradients of the loss with respect to the class scores can become very small (vanishing gradients) or very large (exploding gradients). This can make training deep networks with softmax classifiers difficult, and it may require techniques like gradient clipping or careful weight initialization. **Overfitting:** Softmax classification models can be prone to overfitting when there is insufficient training data or when the model is too complex. Regularization techniques, such as L1 or L2 regularization, are often necessary to mitigate this problem. softmax classification can suffer from the "curse of dimensionality," where the amount of training data required to generalize well increases exponentially with the number of features.

Question 4. [Marks: 14]

- a) Define the vanishing gradient problem. How can the vanishing gradient problem of [6] sigmoid or tanh activation function can be overcome? Activation Function

Solution:

- **Vanishing gradient**—for very high or very low values of X, there is almost no change to the prediction. The derivative values in these regions are very small and converge to 0. This is called the vanishing gradient and the learning is minimal.

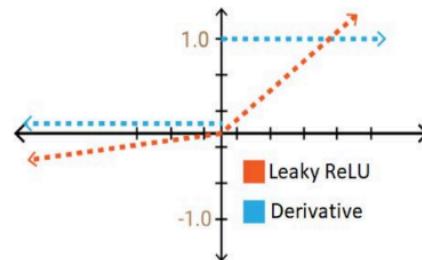
Problem can be solved by Leaky Relu:

#ekhane tanh or sigmoid er ktha bolche so tanh howar kothana ?? tanh er range -1 to 1 ?? –Swarna(061)

, the sigmoid function is not centered around zero, which can lead to the vanishing gradient problem. The tanh activation function is similar to the sigmoid but centered around zero. It maps input values to a range between -1 and 1. While sigmoid and tanh activation functions can alleviate the vanishing gradient problem to some degree by providing gradients that are not always close to zero, they are not a perfect solution. For deep neural networks, more advanced activation functions like the rectified linear unit (ReLU) and its variants (e.g., Leaky ReLU, Parametric ReLU) are often preferred because they are less prone to vanishing gradients and have been shown to perform well in practice.-45

Leaky-ReLU Function

- Prevents dying ReLU problem — this variation of ReLU has a **small positive slope** in the negative area, so it does **enable backpropagation**, even for negative input values. **This leaky value is given as a value of 0.01 if not +ve.**
- **Results not consistent** — leaky ReLU does not provide consistent predictions for negative input values.



Alternative Ans: Hussain-060

The problem could be resolved with ‘normalization’. By normalizing we restrict the values within a certain range. The values are neither very high nor very low anymore, which ensures that the derivative is sufficiently large and not close to 0.

Enigma41

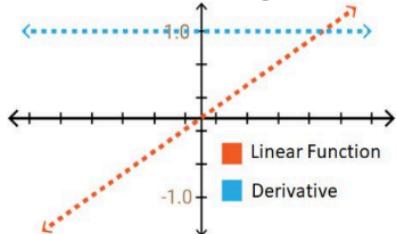
5. i) “A neural network with a linear activation function or without any activation function [6] is simply a linear regression model” – Justify the statement with proper evidence.

Activation Function

Solution:

Solved by -47

derivative of the function is a constant, and has no relation to the input, X. So it's not possible to go back and understand which weights in the input neurons can provide a better prediction.



16

Linear Activation Function

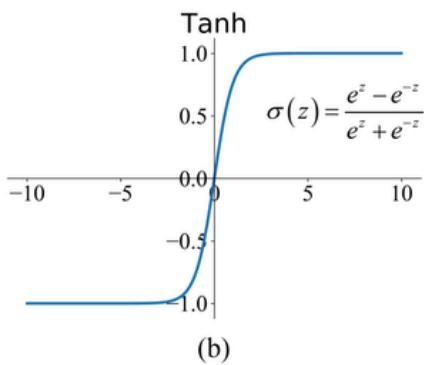
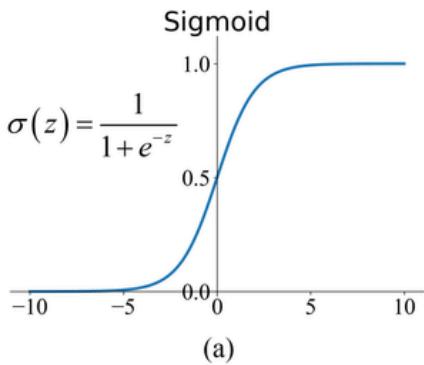
- When $A = c \cdot x$ is derived from x , we reach c . This means that there is no relationship with x . If the derivative is always a constant value, can we say that the learning process is taking place? Unfortunately no!
- All layers of the neural network collapse into one—with linear activation functions, no matter how many layers in the neural network, the last layer will be a linear function of the first layer (*because a linear combination of linear functions is still a linear function*). So a linear activation function turns the neural network or even a deep neural network into just one layer.
- A neural network with a linear activation function or without any activation function is simply a linear regression model. It has limited power and ability to handle complexity varying parameters of input data.

- ii) a) What are the differences between sigmoid and tanh activation functions? Explain [8] with proper graphical representation.
b) “The derivative of tanh activation function depends on the function itself” – Justify the statement.

Solution: 024

(i)

Criteria	Sigmoid Function	Tanh Function
Mathematical form	$\sigma(x) = 1 / (1 + \exp(-x))$	$\tanh(x) = (\exp(x) - \exp(-x)) / (\exp(x) + \exp(-x))$
Output range	0 to 1	-1 to 1
Centered around zero	No	Yes
Use cases	The output layer of binary classification problems, hidden layers of shallow neural networks	Hidden layers of neural networks
Advantages	Differentiable, introduces non-linearity, used in binary classification problems	A steeper gradient around zero captures small changes in the input
Disadvantages	Suffers from vanishing gradient problem, output not centered around zero	Suffers from vanishing gradient problem, which can cause exploding gradient problem when input is too large



(ii)

Tanh (Hyperbolic Tangent) Function:

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

So, the statement is true.

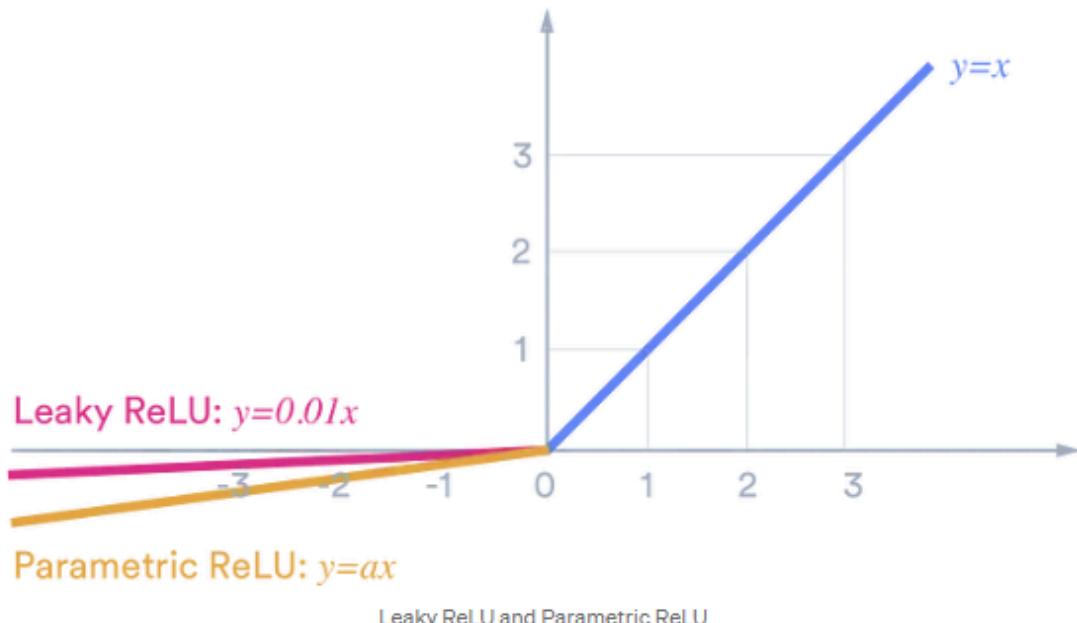
Recursive40

Q4.1

What is the leaky ReLU activation function? Explain briefly with a proper [5] diagram. Activation Function

4. i. Solution: Sujon 49

Leaky Rectified Linear Unit, or Leaky ReLU, is a type of activation function based on a ReLU, but it has a small slope for negative values instead of a flat slope. The slope coefficient is determined before training, i.e. it is not learnt during training. This type of activation function is popular in tasks where we may suffer from sparse gradients, for example training generative adversarial networks.



Leaky ReLU and Parametric ReLU

A problem we see in ReLU is the Dying ReLU problem where some ReLU Neurons essentially die for all inputs and remain inactive no matter what input is supplied, here no gradient flows and if large number of dead neurons are there in a Neural Network it's performance is affected, this can be corrected by making use of this Leaky ReLU where slope is changed left of $x=0$ in above figure and thus causing a leak and extending the range of ReLU.

With Leaky ReLU there is a small negative slope, so instead of not firing at all for large gradients, our neurons do output some value and that makes our layer much more

optimized too.

All Activation Functions

By Sujon 49

Wow ami ekhon egula shob mukhostho korbo ki moja! -Rabab039

Yeee! -Deb065

Why mugosto...bencher moddhe lekhe felais...modern prb modern solution..!! -

Tonmoy146

Name	Plot	Equation	Derivative (with respect to x)	Range	Order of continuity	Monotonic	Monotonic derivative	Approximates identity near the origin
Identity		$f(x) = x$	$f'(x) = 1$	$(-\infty, \infty)$	C^∞	Yes	Yes	Yes
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$	$\{0, 1\}$	C^{-1}	Yes	No	No
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}^{[1]}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$	C^∞	Yes	No	No
Tanh		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$	C^∞	Yes	No	Yes
Arctan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	$(-\frac{\pi}{2}, \frac{\pi}{2})$	C^∞	Yes	No	Yes
ElliottSig ^{[9][12][11]} SoftSign ^{[12][13]}		$f(x) = \frac{x}{1 + x }$	$f'(x) = \frac{1}{(1 + x)^2}$	$(-1, 1)$	C^1	Yes	No	Yes
Inverse square root unit (ISRU) ^[14]		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$	$f'(x) = \left(\frac{1}{\sqrt{1 + \alpha x^2}}\right)^3$	$(-\frac{1}{\sqrt{\alpha}}, \frac{1}{\sqrt{\alpha}})$	C^∞	Yes	No	Yes
Inverse square root linear unit (ISRLU) ^[14]		$f(x) = \begin{cases} \frac{-x}{\sqrt{1 + \alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \left(\frac{-1}{\sqrt{1 + \alpha x^2}}\right)^3 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\frac{1}{\sqrt{\alpha}}, \infty)$	C^2	Yes	Yes	Yes
Square Nonlinearity (SQNL) ^[11]		$f(x) = \begin{cases} 1 & : x > 2.0 \\ x - \frac{x^2}{4} & : 0 \leq x \leq 2.0 \\ x + \frac{x^2}{4} & : -2.0 \leq x < 0 \\ -1 & : x < -2.0 \end{cases}$	$f'(x) = 1 \mp \frac{x}{2}$	$(-1, 1)$	C^∞	Yes	No	Yes
Rectified linear unit (ReLU) ^[15]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$	C^0	Yes	Yes	No
Bipolar rectified linear unit (BReLU) ^[16]		$f(x_i) = \begin{cases} \text{ReLU}(x_i) & \text{if } i \bmod 2 = 0 \\ -\text{ReLU}(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$	$f'(x_i) = \begin{cases} \text{ReLU}'(x_i) & \text{if } i \bmod 2 = 0 \\ -\text{ReLU}'(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$	$(-\infty, \infty)$	C^0	Yes	Yes	No
Leaky rectified linear unit (Leaky ReLU) ^[17]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	C^0	Yes	Yes	No
Parameterized rectified linear unit (PReLU) ^[18]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	C^0	Yes iff $\alpha \geq 0$	Yes	Yes iff $\alpha = 1$
Randomized leaky rectified linear unit (RReLU) ^[19]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	C^0	Yes	Yes	No
Exponential linear unit (ELU) ^[20]		$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha(e^x) + \alpha & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$	$(-\alpha, \infty)$	$\begin{cases} C_1 & \text{when } \alpha = 1 \\ C_0 & \text{otherwise} \end{cases}$	Yes iff $\alpha \geq 0$	Yes iff $0 \leq \alpha \leq 1$	Yes iff $\alpha = 1$
Scaled exponential linear unit (SELU) ^[21]		$f(\alpha, x) = \lambda \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ with $\lambda = 1.0507$ and $\alpha = 1.67326$	$f'(\alpha, x) = \lambda \begin{cases} \alpha(e^x) & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\lambda\alpha, \infty)$	C^0	Yes	No	No
S-shaped rectified linear activation unit (SReLU) ^[22]		$f_{t_l, a_l, t_r, a_r}(x) = \begin{cases} t_l + a_l(x - t_l) & \text{for } x \leq t_l \\ x & \text{for } t_l < x < t_r \\ t_r + a_r(x - t_r) & \text{for } x \geq t_r \end{cases}$ t_l, a_l, t_r, a_r are parameters.	$f'_{t_l, a_l, t_r, a_r}(x) = \begin{cases} a_l & \text{for } x \leq t_l \\ 1 & \text{for } t_l < x < t_r \\ a_r & \text{for } x \geq t_r \end{cases}$	$(-\infty, \infty)$	C^0	No	No	No
Adaptive piecewise linear (APL) ^[23]		$f(x) = \max(0, x) + \sum_{i=1}^S a_i \max(0, -x + b_i)$	$f'(x) = H(x) - \sum_{i=1}^S a_i H(-x + b_i)$ ^[4]	$(-\infty, \infty)$	C^0	No	No	No
SoftPlus ^[24]		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$	$(0, \infty)$	C^∞	Yes	Yes	No
Bent identity		$f(x) = \frac{\sqrt{x^2 + 1} - 1}{2} + x$	$f'(x) = \frac{x}{2\sqrt{x^2 + 1}} + 1$	$(-\infty, \infty)$	C^∞	Yes	Yes	Yes
Sigmoid-weighted linear unit (SiLU) ^[25] (a.k.a. Swish ^[26])		$f(x) = x \cdot \sigma(x)$	$f'(x) = f(x) + \sigma(x)(1 - f(x))$ ^[5]	$[\approx -0.28, \infty)$	C^∞	No	No	No
SoftExponential ^[27]		$f(\alpha, x) = \begin{cases} -\frac{\ln(1 + \alpha(x + \alpha))}{\alpha} & \text{for } \alpha < 0 \\ x & \text{for } \alpha = 0 \\ x + \alpha & \text{for } \alpha > 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \frac{1}{1 - \alpha(x + \alpha)} & \text{for } \alpha < 0 \\ e^{\alpha x} & \text{for } \alpha \geq 0 \end{cases}$	$(-\infty, \infty)$	C^∞	Yes	Yes	Yes iff $\alpha = 0$
Soft Clipping ^[28]		$f(\alpha, x) = \frac{1}{\alpha} \log \frac{1 + e^{\alpha x}}{1 + e^{\alpha(x-1)}}$	$f'(x) = \frac{1}{2} \sinh\left(\frac{p}{2}\right) \operatorname{sech}\left(\frac{px}{2}\right) \operatorname{sech}\left(\frac{p}{2}(1-x)\right)$ ^[6]	$(0, 1)$	C^∞	Yes	No	No
Sinusoid ^[29]		$f(x) = \sin(x)$	$f'(x) = \cos(x)$	$[-1, 1]$	C^∞	No	No	Yes
Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0 \\ \frac{\sin(x)}{x} & \text{for } x \neq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x = 0 \\ \frac{\cos(x)}{x} - \frac{\sin(x)}{x^2} & \text{for } x \neq 0 \end{cases}$	$[\approx -217234, 1]$	C^∞	No	No	No
Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$	$(0, 1]$	C^∞	No	No	No

Logistic Regression

Origin42

- b) I) What role does bias play in a logistic regression? Logistic Regression [8]

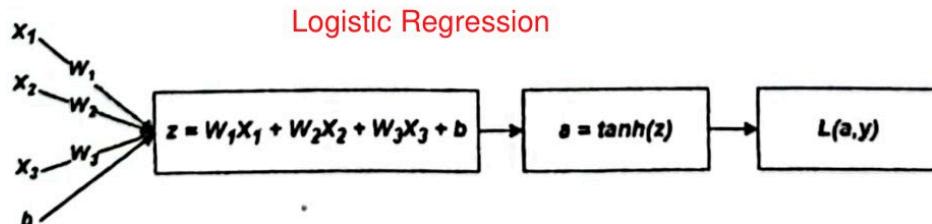
Solution:

- II) "A Neural Network is a combination of several logistic regressions." – Do you agree with the statement? Justify your answer with proper evidence.

Solution:

7

- b) Perform backward propagation on the following logistic regression algorithm to modify the parameters in order to minimize the loss. [8]

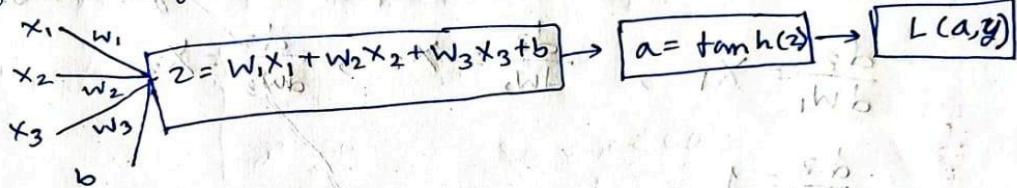


Here, y is the actual label of the training data, and $L(a,y)$ denotes the cross entropy loss between a and y .

Solution:(Akil)80

Perform backward propagation on the following logistic regression algorithm

~~regression algorithm~~



Solⁿ We know, for binary classification, the cross-entropy loss is defined as

$$L(y, a) = -[y \log(a) + (1-y) \log(1-a)]$$

Perform the backward propagation on this logistic regression algorithm.

$$z = w_1x_1 + w_2x_2 + w_3x_3 + b$$

$$\hat{y} = a = \tanh(z)$$

$L(a, y)$ is the cross entropy loss.

$$= -[y \log(a) + (1-y) \log(1-a)]$$

$$\frac{dL}{da} = \frac{d}{da} \{-[y \log(a) + (1-y) \log(1-a)]\}$$

$$= \frac{-y}{a} + \frac{1-y}{1-a}$$

derivation of \tanh function

$$\frac{da}{dz} = \frac{d}{dz} (\tanh(z)) = 1 - \tanh^2(z)$$

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

~~and backward with addition~~
and with multiplication with gradient of \tanh function

Gradient of z with respect to weight and bias

$$\frac{dz}{dw_1} = x_1, \quad \frac{dz}{dw_2} = x_2, \quad \frac{dz}{dw_3} = x_3$$

$$\frac{dz}{db} = 1$$

Gradient of loss / chain Rule about w_1

$$\frac{dL}{dw_{1, \text{out}}} = \frac{dL}{dz} \times x_1 \quad \text{or, } \frac{dL}{dz} \times \frac{dz}{dw_1}$$

$$\frac{dL}{dw_2} = \frac{dL}{dz} \times x_2$$

$$\frac{dL}{dw_3} = \frac{dL}{dz} \times x_3$$

$$\frac{dL}{db} = \frac{dL}{dz} (B - L) + (0)$$

Now, parameter update. Here learning rate = α

$$w_1 = w_1 - \alpha \times \frac{dL}{dw_1}$$

$$w_2 = w_2 - \alpha \times \frac{dL}{dw_2}$$

$$w_3 = w_3 - \alpha \times \frac{dL}{dw_3}$$

$$b = b - \alpha \times \frac{dL}{db}$$

So, it is the backward propagation for the logistic model. Now, the parameters w_1, w_2, w_3 and b should be changed in a way that minimises the loss.

Tune Hyperparameters: The learning rate, batch size, and other hyperparameters can have a significant effect on how well you minimize the cross-entropy loss. Consider using techniques like grid search or random search to find the best hyperparameters for your model.

Enigma41

6. i) The confusion matrix below represents the performance of a Linear Regression model. [6] Find out the Accuracy, Precision, Recall and F1-score of the model.

Performance Evaluation	Actual	
	Positive	Negative
Predicted		
Positive	187	67
Negative	75	171

6. i. Solution: Sujon 49

TP = 187	FP = 67
FN = 75	TN = 171

$$\begin{aligned}
 \text{Accuracy} &= (TP + TN) / (TP + TN + FP + FN) \\
 &= (187 + 171) / (187 + 171 + 67 + 75) \\
 &= 358 / 500 \\
 &= 0.716 \text{ (or } 71.6\%)
 \end{aligned}$$

$$\begin{aligned}
 \text{Precision} &= TP / (TP + FP) \\
 &= 187 / (187 + 67) \\
 &= 187 / 254 \\
 &\approx 0.736 \text{ (or } 73.6\%)
 \end{aligned}$$

$$\begin{aligned}
 \text{Recall} &= TP / (TP + FN) \\
 &= 187 / (187 + 75) \\
 &= 187 / 262 \\
 &\approx 0.714 \text{ (or } 71.4\%)
 \end{aligned}$$

$$\begin{aligned}
 \text{F1-Score} &= 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall}) \\
 &= 2 * (0.736 * 0.714) / (0.736 + 0.714) \\
 &\approx 0.725 \text{ (or } 72.5\%)
 \end{aligned}$$

Recursive40

Solution:

Optimizer

Origin42

Question 7. [Marks: 14]

- a) What is the importance of gradient descent in neural network? Write down the algorithm [6] of Adam optimizer. Optimizer

Solution:

Enigma41

1. i) Why does the Stochastic Gradient Descent technique oscillate across the ridge? How [6] can this oscillation be reduced? Optimizer

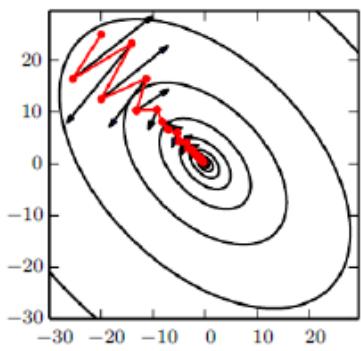
Solution: 024

Reason for oscillation: Randomness

The updates in SGD are noisy and have a high variance, which can make the optimization process less stable and lead to oscillations around the minimum.

Solution of oscillation: Momentum

The method of momentum is designed to accelerate learning, especially in the face of high curvature, small but consistent gradients, or noisy gradients. The momentum algorithm accumulates an exponentially decaying moving average of past gradients and continues to move in their direction.



Recursive40

2. i) Write down all the steps of RMSProp or Momentum Gradient Descent with equations. Optimizer

Solution:

Cross Entropy Loss

Origin42

Solution:

Enigma41

Solution:

Recursive40

3

- ii) Consider the following loss function where $p \in [0, 1]$ is the probability score and $w_1, w_2 \in \mathbb{R}$ are scalar values. [9]

$$\text{Loss} = w_1 \log(p) + w_2 \log(1 - p)$$

Now, answer the following questions:

Cross Entropy Loss

- a) Write down the behavior of the loss function when $w_1 = w_2$.
- b) What are the necessary conditions so that the loss function will behave like a Binary Cross Entropy (BCE) loss?
- c) “Loss function is positive only if both w_1 and w_2 are non-positive.”— Do you agree? Justify your answer.

Solution:

Solution:

Information Theory

Origin42

Solution:

Enigma41

Solution:

Recurssive40

Solution:

Gradient Decent

Origin42

Solution:

Enigma41

Solution:

Recurssive40

Solution:

**** Performance Analysis ****

Ei naam e kono topic sir er Question pattern e nai.

So, eta onno kono topic er under e kina janais ...

Recurssive40



- i) What do you mean by the bias-variance tradeoff? Show the proper equation and [5] diagram.
- ii) Consider a binary classification task where you have implemented three models, [9] i.e., Model A, Model B, and Model C, to classify positive and negative classes. The performances of these three models on a test dataset are shown below.

Performance Analysis

	<i>Model A</i>	<i>Model B</i>	<i>Model C</i>
<i>Precision</i>	0.92	0.72	0.80
<i>Recall</i>	0.78	0.90	0.80

0.844 0.8 0.8

Now, answer the following questions:

- a) Which model do you select if the ratio of the positive class to the negative class is 1?
- b) Can you figure out the most accurate model among the three? If not, what is the minimal information you need to figure out the most accurate model?
- c) Which model do you think is the better performer between Model B and Model C?

Solution:

D : 1 set

Quiz-1 - Integer43

Set-A

Q1. A. What is the role of an optimizer in a neural network?

Solution:45

The optimizer plays a crucial role in training a neural network. Its primary responsibility is to

adjust the model's internal parameters (weights and biases) during the training process to minimize the loss function, which measures the difference between the model's predictions and the actual target values. The role of the optimizer can be summarized as follows:

1. **Minimizing the Loss Function:**

- The primary goal of training a neural network is to find the model parameters (weights and biases) that minimize a specified loss function. This loss function quantifies how well the model's predictions match the true target values.

- The optimizer is responsible for iteratively updating the model's parameters in a way that reduces the loss. It does this by computing the gradients of the loss with respect to the parameters and making parameter updates in the direction that reduces the loss.

2. **Gradient Descent:**

- Most optimizers, including popular ones like Stochastic Gradient Descent (SGD), Adam, RMSprop, and others, **use gradient information to update model parameters**.

- The optimizer calculates the gradient of the loss function with respect to each parameter, indicating how the loss changes as each parameter is adjusted. This gradient points in the direction of steepest increase in the loss.

3. **Learning Rate Control:**

- **The optimizer also manages the learning rate, which determines the step size for parameter updates**

5. **Convergence and Stability:**

- The choice of optimizer can significantly impact the convergence speed and stability of the training process.

- **A well-chosen optimizer helps the neural network converge to an optimal or near-optimal solution more quickly and reliably.**

6. **Regularization and Other Techniques:**

- **Some optimizers incorporate regularization techniques, such as L1 or L2 regularization, weight decay, or dropout, to improve the generalization of the model and prevent overfitting.**

B. How does learning rate impact the performance of a model? Explain briefly.

Solution:45

The learning rate is a crucial hyperparameter in training machine learning models, especially neural networks. It determines the size of steps that the model takes during the optimization process, such as gradient descent, when updating its parameters (like weights and biases). Here's how the learning rate impacts model performance:

1. **Too High Learning Rate:**

- If the learning rate is too high, the model might converge quickly, but it can overshoot the optimal solution and fail to converge to a good solution. This can lead to instability and divergence.

2. **Too Low Learning Rate**:

- Conversely, if the learning rate is too low, the model's training process will progress very slowly. It may get stuck in local minima or take an excessively long time to converge to the optimal solution.

3. **Appropriate Learning Rate**:

- The key is to find an appropriate learning rate. A well-chosen learning rate allows the model to converge to a good solution efficiently. This often involves experimentation and adjusting the learning rate during training (e.g., learning rate schedules or adaptive methods like Adam) to strike the right balance between convergence speed and accuracy.

In summary, the learning rate directly impacts how quickly a model learns and the quality of the solution it converges to. It's essential to choose an appropriate learning rate to achieve the best performance in training your machine learning models.

Quiz-3 - Integer43

Set-A

Q1. Which of the following is not true as a characteristic of parameters of neural network?

- A. Required by the model when making predictions
- B. Learned from data
- C. Often not set manually by the practitioner
- D. Determine how the network is trained

Solution: By Hussain-060

D

Hyperparameters determine how the network is trained [Slide Hyperparameters Parameters, P. 5], not parameters.

Q2. CNN is mostly used when there is an -

- A. structured data
- B. unstructured data
- C. both A and B
- D. none of the above

Solution: Sujon 49

B. unstructured data

CNN is mostly used when there is an unstructured data set (e.g., images) and the practitioners need to extract information from it.

Q3. Calculate the number of parameters for the following architecture for the input shape of (428,428).

input → Conv(3,2) → Minpool(2) → ReLU → Conv(5,3) → Maxpool(5) → FCwB(256) → FC(5) → output

Here, Conv(x,y) denotes a convolutional layer that has the kernel size of x and stride size of y,

Minpool(x) denotes a Minpooling layer that has both window size and stride size of x,

Maxpool(x) denotes a Maxpooling layer that has both window size and stride size of x,

FCwB(x) denotes a Fully Connected Layer of x neurons with Bias, and

FC(x) denotes a Fully Connected Layer of x neurons without Bias.

Solution: Rabab 039

Layer	Output size	Parameter#
input	478, 464	0
Conv(3,2)	$\text{floor}(W-F+2P)/S + 1$ $(478-3 + 2*0)/2 + 1 = 238$, $(464-3 + 2*0)/2 + 1 = 231$	$((\text{filter_dim} * \text{filters_prev}) + 1) * \text{filters_current}$ $((3*3*1)+1)*1 = 10$
Minpool(4)	$(W-F+2P)/S + 1$ $(238-4 + 2*0)/4 + 1 = 59$, $(231-4 + 2*0)/4 + 1 = 57$	0
ReLU	59, 57	0
Conv(5,3)	$(W-F+2P)/S + 1$ $(59-5 + 2*0)/3 + 1 = 19$, $(57-5 + 2*0)/3 + 1 = 18$	$((\text{filter_dim} * d) + 1) * k$ $((5*5*1)+1)*1 = 26$
Maxpool(3)	$(W-F+2P)/S + 1$ $(19-3 + 2*0)/3 + 1 = 6$, $(18-3 + 2*0)/3 + 1 = 6$	0
FCwB(256)	256	$(\text{inputs} * \text{nodes}) + \text{nodes}$ $(6*6 * 256) + 256 = 9472$
ReLU	256	0
Dropout	256	0
FC(128)	128	$(\text{inputs} * \text{nodes})$ $(256*128) = 32768$
tanh	128	0
FC(5)	5	$128 * 5 = 640$
output	5	0
Total	42916	

Formulae:

Convolution: output size = $(W-F+2P)/S + 1$, parameters = $((m*n*d)+1)*k$

Pooling: output size = $(W-F+2P)/S + 1$, parameters = 0

Activation functions: output size = same as inputs, parameters = 0

Dropout: output size = same as inputs, parameters = 0

FCwB: output size = no. of nodes, parameters = $(\text{inputs} * \text{no. of nodes}) + \text{no. of nodes}$

FC: output size = no. of nodes, parameters = $(\text{inputs} * \text{no. of nodes})$

Set-B

Q1. Which of the following is not true as a characteristic of hyperparameters?

- A. Set before training
- B. Determines the network structure
- C. Determine how the network is trained
- D. Often saved as part of the learned model

Solution: Added by Younus-131

D. Often saved as part of the learned model

This statement is not true as a characteristic of hyperparameters. Hyperparameters are not typically saved as part of the learned model. Instead, they are set before training and determine various aspects of the training process and network structure, but they are separate from the learned model's parameters, which are the weights and biases that are adjusted during training based on the data. Hyperparameters are external to the model and are typically specified by the machine learning engineer or researcher before the training process begins.

Q2. An input image has been converted into a matrix of size 12 X 12 along with a filter of size 3 X 3 with a Stride of 2 and padding of 0. Determine the size of the convoluted matrix.

- A. 10*10
- B. 8*8
- C. 5*5
- D. 3*3

Solution: Added by Younus-131

To determine the size of the convoluted matrix after applying a 3x3 filter with a stride of 2 to a 12x12 input matrix with no padding, you can use the following formula:

$$\text{Output size} = ((\text{Input size} - \text{Filter size}) / \text{Stride}) + 1$$

In this case:

$$\text{Input size} = 12 \times 12$$

$$\text{Filter size} = 3 \times 3$$

$$\text{Stride} = 2$$

$$\text{Padding} = 0$$

So, plugging these values into the formula:

$$\text{Output size} = ((12 - 3) / 2) + 1$$

$$\text{Output size} = (9 / 2) + 1$$

$$\text{Output size} = 4.5 + 1$$

$$\text{Output size} = 5.5$$

Since the output size should be a whole number, you typically round down to the nearest integer. Therefore, the size of the convoluted matrix is 5x5.

So, the correct answer is C. 5x5.

Q3. Calculate the number of parameters for the following architecture for the input shape of (438,438).

input → Conv(3,2) → Minpool(4) → ReLU → Conv(5,3) → Maxpool(3) → FCwB(256) → FC(5) → output

Here, Conv(x,y) denotes a convolutional layer that has the kernel size of x and stride size of y,
 Minpool(x) denotes a Minpooling layer that has both window size and stride size of x,
 Maxpool(x) denotes a Maxpooling layer that has both window size and stride size of x,
 FCwB(x) denotes a Fully Connected Layer of x neurons with Bias, and
 FC(x) denotes a Fully Connected Layer of x neurons without Bias.

Solution: Added by Hussain-060

[Please correct if any mistake is found]

Layer	Input Shape	Output Shape	# Params
1. Conv(3,2)	(438, 438)	$\text{floor}((438-3+2*0)/2) + 1 = 218$ [Assuming padding = 0] (218, 218)	$3*3 + 1 = 10$ [assuming no. of in and out channel = 1]
2. MinPool(4)	(218, 218) ?- output shape of previous layer - yep	$\text{floor}((218-4+2*0)/4) + 1 = 54$ (54, 54)	0
3. ReLU	(54, 54)	(54, 54)	0
4. Conv(5, 3)	(54, 54)	$\text{floor}((54-5+2*0)/3)+1 = 17$ (17, 17)	$5*5 + 1 = 26$
5. MaxPool(3)	(17, 17)	$\text{floor}((17-3+2*0)/3) + 1 = 5$ (5, 5)	0
Flatten	(5,5)	25	0
6. FCwB(256)	25	256	$25*256+256=6656$
7. FC(5)	256	5	$256*5=1280$
Total Params			7972

Note: "Input shape" column could be discarded altogether.

CNN

Origin42

Question 5. [Marks: 14]

- a) How normalization can be achieved in Convolutional Neural Network? What is the rule of thumb [6] to use a CNN architecture? **CNN**

Solution:

Normalization

Keep the math from breaking by tweaking each of the values just a bit.

Change everything negative to zero.

Rule of thumb

If your data is just as useful after swapping any of your columns with each other, then you can't use Convolutional Neural Networks.

Alternative: Sujon 49

Normalization in Convolutional Neural Networks (CNNs) is typically achieved using techniques like Batch Normalization (BatchNorm) and Layer Normalization. These techniques help stabilize and accelerate the training of deep neural networks by normalizing the inputs to each layer. Here's a brief explanation of Batch Normalization:

Batch Normalization (BatchNorm):

BatchNorm is applied to the activations of a layer within a mini-batch during training. It works as follows:

Compute Mean and Variance: For each mini-batch, calculate the mean and variance of the activations across the mini-batch for each channel.

Normalize: Normalize the activations by subtracting the mean and dividing by the standard deviation (with a small epsilon added for numerical stability).

Scale and Shift: After normalization, scale the normalized activations by a learnable parameter (gamma) and shift them by another learnable parameter (beta). These parameters allow the model to adapt and learn the best scaling and shifting for each layer.

Backpropagation: During backpropagation, gradients are computed for gamma and beta, allowing the network to learn how to adjust the normalization.

BatchNorm helps with the following:

Stability: It stabilizes training by reducing internal covariate shift, making it easier to train deep networks.

Regularization: It acts as a form of regularization by introducing noise in the training process.

Faster Training: It often accelerates training since it allows for more aggressive learning rates.

Enigma41

1

- ii) Suppose you want to build an automated Image Bot that generates the caption for the [8] image. You have created a complex deep model consisting of a Long Short-Time Memory (LSTM) followed by a Convolutional Neural Network (CNN). Your friend suggests a Recurrent Neural Network (RNN) instead of LSTM. **RNN, CNN, LSTM** Do you think your friend's suggestion would improve the performance? Why or why not? What are the possible corners of improvement for this model?

Solution:

4. i) "In Convolutional Neural Network, Pooling always decreases the parameters". Is this [6] statement true? Explain your answer with proper justification. **CNN**

Solution: 024

The statement is false.

Generally pooling decreases parameters. But if the pool size is set to 1, it essentially means that no pooling is being applied. In other words, each individual value in the input feature map remains unchanged. Consequently, there is no reduction in the spatial dimensions, and the number of parameters also remains the same after this "pooling"

operation (which isn't really pooling at all in this case).

Recursive40

- 2.ii) Briefly explain the slowness of Recurrent Neural Network (RNN) compared to [4] other deep networks, such as Convolutional Neural Network (CNN). **RNN-CNN**

Solution:

	Convolutional neural network (CNN)	Recurrent neural network (RNN)
ARCHITECTURE	Feed-forward neural networks using filters and pooling	Recurring network that feeds the results back into the network
INPUT/OUTPUT	The size of the input and the resulting output are fixed (i.e., receives images of fixed size and outputs them to the appropriate category along with the confidence level of its prediction)	The size of the input and the resulting output may vary (i.e., receives different text and output translations—the resulting sentences can have more or fewer words)
IDEAL USAGE SCENARIO	Spatial data (such as images)	Temporal/sequential data (such as text or video)
USE CASES	Image recognition and classification, face detection, medical analysis, drug discovery and image analysis	Text translation, natural language processing, language translation, entity extraction, conversational intelligence, sentiment analysis, speech analysis

Parameter Calculation

Origin42

- b) Calculate the number of parameters for the following architecture for the input shape of [8] (478,464). Parameter Calculation

input → Conv(3,2) → Minpool(4) → ReLU → Conv(5,3) → Maxpool(3) ←
 output ← FC(5) ← tanh ← FC(128) ← Dropout(0.2) ← ReLU ← FCwB(256) ←

Here, Conv(x,y) denotes a convolutional layer that has the kernel size of x and stride size of y,

Minpool(x) denotes a Minpooling layer that has both window size and stride size of x,

Maxpool(x) denotes a Maxpooling layer that has both window size and stride size of x,

FCwB(x) denotes a Fully Connected Layer of x neurons with Bias,

Dropout (y) denotes a Dropout Layer that drops input at a rate of y, and

FC(x) denotes a Fully Connected Layer of x neurons without Bias.

Solution: Toufique 116

	Output Shape Parameters	
input	478,464	0
Conv(3,2)	238,231	10
Min Pool(4)	59,57	0
ReLU	59,57	0
Conv(5,3)	19,18	26
Maxpool(3)	6,6	0
FCwB(256)	256	9472
ReLU	256	0
Dropout	256	0
FC(128)	128	32768
tanh	128	0
FC(5)	5	640
Output	5	0
Total	42916	

Enigma41

- ii) Calculate the number of parameters for the following architecture for the input shape [8] of 4237.

input → FCwB(256) → Dropout(0.1) → FC(128) → tanh → Maxpool(3) → FC(3) → output

Parameter Calculation

Here, FCwB(x) denotes a Fully Connected Layer of x neurons with Bias,
 FC(x) denotes a Fully Connected Layer of x neurons without Bias,
 Dropout (y) denotes a Dropout Layer that drops input at a rate of y, and
 Maxpool(x) denotes a flat Maxpooling layer that has both window size and stride size of x.

Solution: 024

Layer	Output Shape	No. of Parameters
input(4237)	4237	0
FCwB(256)	256	(input_shape x output_shape) + bias (4237x256)+256=1084928
Dropout(0.1)	256	0
FC(128)	128	(input_shape x output_shape) (256x128)=32768
tanh	128	0
Maxpool(3)	floor(((W-F+2P)/S)+1) floor(((128-3+2x0)/3)+1)=42	0
FC(3)	3	(input_shape x output_shape) (42x3)=126
output	3	0
Total:		1117822

So, number of total parameters = 1117822

Recursive40

- Xii) Consider the following architecture that takes a (256, 256, 3) dimensional image [9] and outputs a (56,56,56) dimensional image. The architecture consists of a convolutional layer (CNN) followed by a pooling layer (POOL) and another convolutional layer (CNN). Now, find a possible set of hyperparameters for the architecture.

Parameter Calculation

image (256, 256, 3) → $CNN \rightarrow POOL \rightarrow CNN$ → image (56, 56, 56)

Your answer must include the following hyperparameters:

$CNN \rightarrow$ number of filters, filter size, stride, padding

$POOL \rightarrow$ stride, pool size, padding

Solution: By Hussain-060

[Updates are welcome]

CNN1:

Input shape = (256, 256, 3)

No. of filters = 28 # 28 kivabe? - No particular reason, could've been 1 as well

Filter size = 3

Stride = 2

Padding = 0

$$\text{floor}((256-3+2*0)/2) + 1 = 127$$

Output shape:(127, 127, 28)

Pool:

Input shape = (127, 127, 28)

Filter size = 5

Stride = 2

Padding = 0

$$\text{floor}((127-5+2*0)/2) + 1 = 62$$

Output shape:(62, 62, 28)

CNN2:

Input shape = (62, 62, 28)

No. of filters = 56

Filter size = 7

Stride = 1
Padding = 0

$\text{floor}((62-7+2*0)/1) + 1 = 56$
Output shape:(56, 56, 56)

Summary: (No of filters, filter size, stride, padding)
CNN1 -> 28, 3, 2, 0
Pool -> 1, 5, 2, 0
CNN2 -> 56, 7, 1, 0

Note: Most of the values are random or educated guesses with the goal of reaching the final shape from the initial shape. Keu jodi guesswork chara onno kono technique janos then please update koris.

Parameter vs. Hyperparameter

Origin42

Solution:

Enigma41

Solution:

Recurssive40

Solution:

Genetic Algorithm, TSP

Origin42

Solution:

Enigma41

Solution:

Recurssive40

7. i) What do you mean by mutation in Genetic Algorithms (GA)? Provide an example [5]
of mutation. **Genetic Algo**

Solution: 024

Mutation:

Mutation is a natural process that occurs due to an error in replication or copying of genes.

By mixing and matching the genes from both parents, it is possible to reproduce the parent chromosomes during crossover. There is no guarantee that the copying of the parent gene is 100% accurate. There always occurs an error, which leads to the scope of exploration. Mutating the chromosome in the genetic algorithm is necessary because it may result in revolutionary results that will help solve problems more efficiently.

Example:

Bit Flip Mutation

<table border="1"><tr><td>0</td><td>0</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td><td>1</td><td>0</td></tr></table>	0	0	1	1	0	1	0	0	1	0	=>	<table border="1"><tr><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>1</td><td>0</td></tr></table>	0	0	1	0	0	1	0	0	1	0
0	0	1	1	0	1	0	0	1	0													
0	0	1	0	0	1	0	0	1	0													

Swap Mutation

<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td><td>0</td></tr></table>	1	2	3	4	5	6	7	8	9	0	=>	<table border="1"><tr><td>1</td><td>6</td><td>3</td><td>4</td><td>5</td><td>2</td><td>7</td><td>8</td><td>9</td><td>0</td></tr></table>	1	6	3	4	5	2	7	8	9	0
1	2	3	4	5	6	7	8	9	0													
1	6	3	4	5	2	7	8	9	0													

Scramble Mutation

<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td></tr></table>	0	1	2	3	4	5	6	7	8	9	=>	<table border="1"><tr><td>0</td><td>1</td><td>3</td><td>6</td><td>4</td><td>2</td><td>5</td><td>7</td><td>8</td><td>9</td></tr></table>	0	1	3	6	4	2	5	7	8	9
0	1	2	3	4	5	6	7	8	9													
0	1	3	6	4	2	5	7	8	9													

Inversion Mutation

<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td></tr></table>	0	1	2	3	4	5	6	7	8	9	=>	<table border="1"><tr><td>0</td><td>1</td><td>6</td><td>5</td><td>4</td><td>3</td><td>2</td><td>7</td><td>8</td><td>9</td></tr></table>	0	1	6	5	4	3	2	7	8	9
0	1	2	3	4	5	6	7	8	9													
0	1	6	5	4	3	2	7	8	9													

Here 4 examples are given. According to the question, discuss only one.

E : 1/1.5 set

Fuzzy (All the topics)

Recommended by sir:

1. Membership Function (Graph to Equation, Equation to Graph)
2. CrispSet VS FuzzySet
3. Fuzzy, NeuroFuzzy system importance
4. Alpha Cut / Extension Principle

Origin42

Question 6. [Marks: 14] Fuzzy

(a) List the advantages of a fuzzy system. How a Neuro-Fuzzy system can be built?

[6]

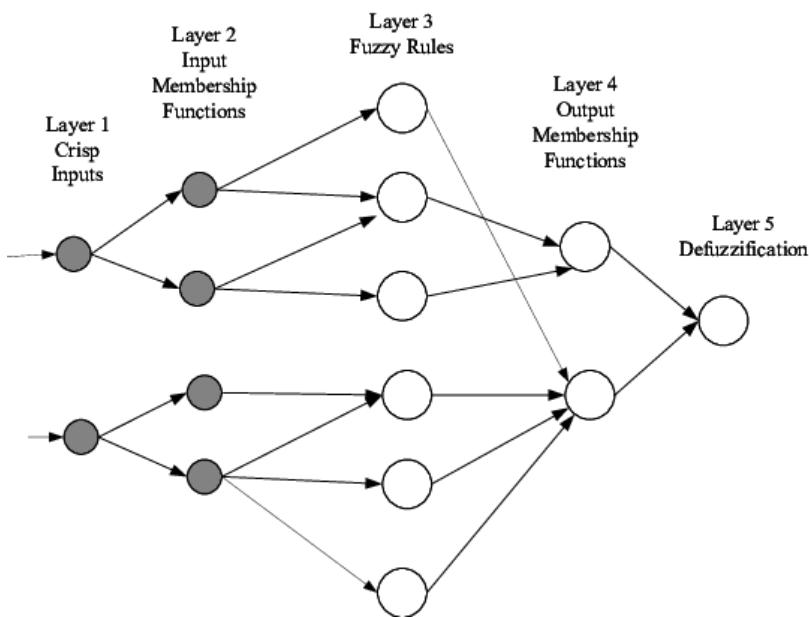
Solution: 024

Advantages of Fuzzy Logic System:

- This system can work with any type of input whether it is imprecise, distorted or noisy input information.
- The construction of Fuzzy Logic Systems is easy and understandable.
- Fuzzy logic comes with mathematical concepts of set theory and the reasoning of that is quite simple.
- It provides a very efficient solution to complex problems in all fields of life as it resembles human reasoning and decision-making.
- The algorithms can be described with little data, so little memory is required.

Neuro Fuzzy System:

- A neuro-fuzzy system is based on a fuzzy system which is trained by a learning algorithm derived from neural network theory.
 - A neuro-fuzzy system can be built as a 3-layer feedforward neural network.
 - First layer represents input variables
 - Middle (hidden) layer represents fuzzy rules
 - Third layer represents output variables
- Fuzzy sets are encoded as (fuzzy) connection weights. It represents the data flow of input processing and learning within the model. Sometimes a 5-layer architecture is used, where the fuzzy sets are represented in the units of the second and fourth layer.
- A neuro-fuzzy system can be always interpreted as a system of fuzzy rules. It is also possible to create the system out of training data from scratch, as it is possible to initialize it by prior knowledge in the form of fuzzy rules.



Fuzzy Addition

Fuzzy Addition with extension principle:

Origin 42 : (6) b :

$$A = 0.2/1 + 0.5/2 + 0.7/3 + 1/4 + 0.7/5 + 0.5/6 + 0.2/7$$

$$B = 0.3/3 + 0.5/4 + 0.8/5 + 1/6 + 0.8/7 + 0.5/8 + 0.3/9$$

- b) Two fuzzy numbers A and B are given below. Add the numbers using the extension principle. [8]

$$A = 0.2/1 + 0.5/2 + 0.7/3 + 1/4 + 0.7/5 + 0.5/6 + 0.2/7$$

$$B = 0.3/3 + 0.5/4 + 0.8/5 + 1/6 + 0.8/7 + 0.5/8 + 0.3/9$$

Solution: 035 - Tuli

		Support of B										
A \ B		Y=1	Y=2	Y=3	Y=4	Y=5	Y=6	Y=7	Y=8	Y=9	Y=10	
S U P P O t .	$\alpha=1$	0.0 0.2	0.0 0.2	0.3 0.2	0.3 0.2	0.5 0.2	0.8 0.2	1.0 0.2	0.8 0.2	0.5 0.2	0.3 0.2	0.0 0.2
	$\alpha=2$	0.0 0.15	0.0 0.15	0.0 0.0	0.3 0.3	0.5 0.5	0.5 0.5	0.5 0.5	0.5 0.5	0.3 0.3	0.0 0.15	0.0 0.15
	$\alpha=3$	0.0 0.7	0.0 0.7	0.3 0.7	0.3 0.7	0.5 0.7	0.8 0.7	1.0 0.7	0.8 0.7	0.5 0.7	0.3 0.7	0.0 0.7
	$\alpha=4$	0.0 1.0	0.0 1.0	0.3 1.0	0.3 1.0	0.5 1.0	0.8 1.0	1.0 1.0	0.8 1.0	0.5 1.0	0.3 1.0	0.0 1.0
	$\alpha=5$	0.0 0.7	0.0 0.7	0.3 0.7	0.3 0.7	0.5 0.7	0.8 0.7	1.0 0.7	0.8 0.7	0.5 0.7	0.3 0.7	0.0 0.7
	$\alpha=6$	0.0 0.5	0.0 0.5	0.3 0.5	0.3 0.5	0.5 0.5	0.5 0.5	0.5 0.5	0.8 0.5	0.5 0.5	0.3 0.5	0.0 0.5
	$\alpha=7$	0.0 0.2	0.0 0.2	0.3 0.2	0.2 0.2	0.0 0.2						
	$\alpha=8$	0.0 0.0	0.0 0.0	0.3 0.0	0.3 0.0	0.5 0.0	0.8 0.0	1.0 0.0	0.8 0.0	0.5 0.0	0.3 0.0	0.0 0.0
	$\alpha=9$	0.0 0.0	0.0 0.0	0.3 0.0	0.3 0.0	0.5 0.0	0.8 0.0	1.0 0.0	0.8 0.0	0.5 0.0	0.3 0.0	0.0 0.0
	$\alpha=10$	0.0 0.0	0.0 0.0	0.3 0.0	0.3 0.0	0.5 0.0	0.8 0.0	1.0 0.0	0.8 0.0	0.5 0.0	0.3 0.0	0.0 0.0

| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

$\therefore \text{Final output, } C = 0.2/4 + 0.3/5 + 0.5/6 + 0.5/7 + 0.7/8 + 0.8/9 + 1.0/10 + 0.8/11 + 0.7/12 + 0.5/13 + 0.5/14 + 0.3/15$
 $+ 0.3/16 + 0.2/17 + 0.2/18$

Enigma41

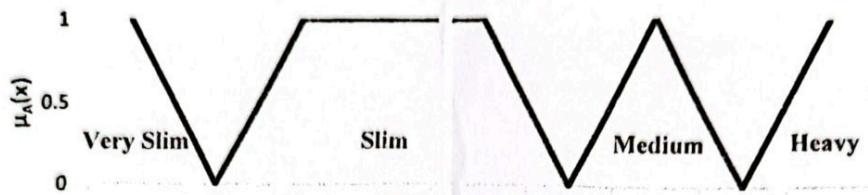
7. i) What do you mean by Fuzzy Set? List the differences between Crisp Set and Fuzzy Set. **Fuzzy** [6]

Solution: By - Rubayat

A **fuzzy set** is a set where each element has a degree of membership between 0 and 1. It is an extension of the classical notion of set. For example, if you have a set of fruits, you can assign a membership value to each fruit based on how sweet it is. A fuzzy set allows you to capture uncertainty and vagueness in your data.

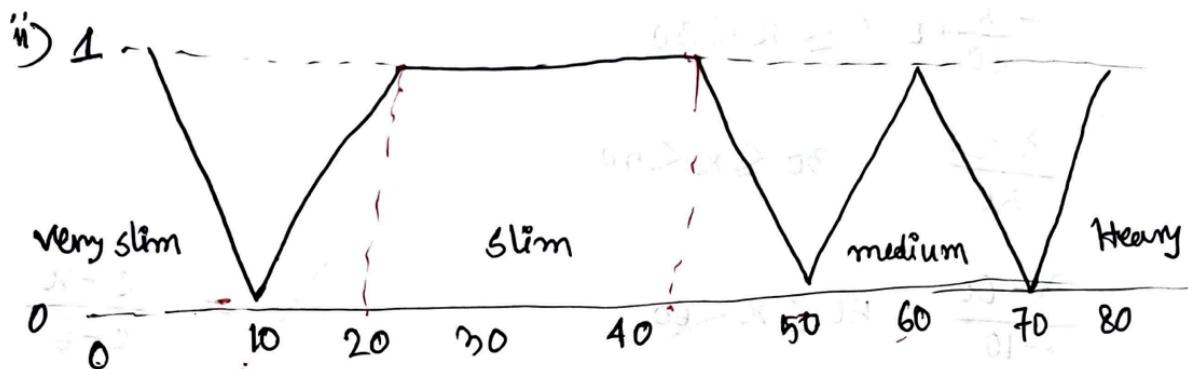
The main difference between crisp and fuzzy sets is that crisp sets have elements with binary membership (either 0 or 1), while fuzzy sets have elements with degrees of membership between 0 and 1. Crisp sets use bi-valued logic, while fuzzy sets use infinite-valued logic. Crisp sets have precise boundaries, while fuzzy sets have indeterminate boundaries

- i) what do you mean by fuzzy Set? List the differences between Crisp Set and Fuzzy [9]
Set. **Fuzzy**
- ii) Assume, we want to calculate the health of a person based on his/her weight. A graphical representation of the membership function for this task is given in the following figure. Write down the mathematical expression of the membership function. Describe each step briefly. [8]



Solution: By - Rubayat

#Ehane maybe vul aache -swarna(061)



$$\mu_{\text{Very slim}} = \frac{x-10}{10-0} \quad 10 < x < 20$$

$$\mu_{\text{slim}} = \frac{x-10}{20-10} \quad (x=10 \rightarrow 20)$$

$$\mu_{\text{slim}} = 1 \quad 20 \leq x \leq 40.$$

$$\mu_{\text{slim}} = \frac{50-x}{50-40} \quad 40 \leq x \leq 50$$

$$\mu_{\text{medium}} = \frac{x-50}{60-50}$$

$$\mu_{\text{medium}} = \frac{70-x}{70-60}$$

$$\mu_{\text{Heavy}} = \frac{x-70}{80-70}$$

$$0 \leq x \leq 20 \quad \frac{20-x}{20-0}$$

\rightarrow slim equation

$$\frac{c-x}{c-b}$$

$$\frac{x-a}{b-a}$$

Return 38

b)	<p>The fuzzy numbers A and B are given by</p> $\mathbf{A} = 0.33/6 + 0.67/7 + 1.00/8 + 0.67/9 + 0.33/10$ $\mathbf{B} = 0.33/1 + 0.67/2 + 1.00/3 + 0.67/4 + 0.33/5$ <p>Draw a sketch of fuzzy number C, where C results from adding A and B by applying the alpha-cut principle.</p>	[3]
Solution: Rafi-148		

$$A = 0.33/6 + 0.57/7 + 1.00/8 + 0.67/9 + 0.33/10$$

$\alpha + \text{cat}/A$	0	0	0	0	0	0.99	0.67	1.00	0.67	0.33
$\alpha = 1.0$				1				1		
$\alpha = 0.9$				1				1		
$\alpha = 0.8$				1				1		
$\alpha = 0.7$				1				1		
$\alpha = 0.6$		1	1	1			1	1	1	
$\alpha = 0.5$		1	1	1			1	1	1	
$\alpha = 0.4$	1	1	1	1			1	1	1	
$\alpha = 0.3$	1	1	1	1	+	1	1	1	+	
$\alpha = 0.2$	1	1	1	1	1	1	1	1	1	1
$\alpha = 0.1$						1	1	1	1	1
	1	2	3	4	5	6	7	8	9	10

for
1.0

$$\alpha_1 = 8$$

$$\alpha_2 = 8$$

for 0.9,

$$\alpha_1 = 8$$

$$\alpha_2 = 8$$

for 0.8,
1.0

$$\alpha_1 = 8$$

$$\alpha_2 = 8$$

for 0.7,

$$\alpha_1 = 8$$

$$\alpha_2 = 8$$

for 0.6, 0.7, 0.4,

$$\alpha_1 = 7$$

$$\alpha_2 = 9$$

for 0.5, 0.4, 0.3,

$$\alpha_1 = 6$$

$$\alpha_2 = 10$$

$$B = 0.33/1 + 0.67/2 + 1.00/3 + 0.67/4 + 0.33/5$$

α	0.33	0.67	1.00	0.67	0.33	0
$\alpha = 1.0$			1			
$\alpha = 0.9$			1			
$\alpha = 0.8$			1			
$\alpha = 0.7$			1			
$\alpha = 0.6$		1	1	1		
$\alpha = 0.5$		1	1	1		
$\alpha = 0.4$		1	1	1		
$\alpha = 0.3$	1	1	1	1	1	
$\alpha = 0.2$	1	1	1	1	1	
$\alpha = 0.1$	1	1	1	1	1	
	1	2	3	4	5	6

for
1.0, 0.9, 0.8, 0.7:

$$b_1 = 1 \cancel{+} 3$$

$$b_2 = 1 \cancel{+} 3$$

for
0.6, 0.5, 0.4:

$$b_1 = 2$$

$$b_2 = 4$$

for
0.3, 0.2, 0.1

$$b_1 = 1$$

$$b_2 = 5$$

α_{cut}	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0.1.0											L				
0.0.9											L				
0.0.8											L				
0.0.7											L				
0.0.6									L	L	L	1	1	1	
0.0.5									L	L	1	2	1		
0.0.4									L	L	1	1	1		
0.0.3									L	L	1	4	1	0.15	1
0.0.2									L	L	1	1	1	1	1
0.0.1									L	1	2	1	1	1	1
P	2	3	9	5	6	7	8	9	10	11	12	13	14	15	

$$C = 0.9/7 + 0.3/8 + 0.6/9 + 0.6/10 + 1/11$$

$$+ 0.6/12 + 0.6/13 + 0.3/14 + 0.3/15$$

α -cat:

for $0 \cdot 1$:

$$A_{0 \cdot 1} = [6, 10], B_{0 \cdot 1} = [1, 5], C_{0 \cdot 1} = [(6+1), (10+5)] = [7, 15]$$

for $0 \cdot 2$:

$$A_{0 \cdot 2} = [6, 11], B_{0 \cdot 2} = [1, 5], C_{0 \cdot 2} = [6+1, 10+5] = [7, 15]$$

for $0 \cdot 3$: same as $0 \cdot 2$, $C_{0 \cdot 3} = [7, 15]$

for $0 \cdot 4; 0 \cdot 5; 0 \cdot 6$: (same)

$$A_{0 \cdot 4} = [7, 9], B_{0 \cdot 4} = [2, 4], C_{0 \cdot 4} = [7+2, 9+4] = [9, 13]$$

$$C_{0 \cdot 5} = [9, 13]$$

$$C_{0 \cdot 6} = [9, 13]$$

for $0 \cdot 7; 0 \cdot 8; 0 \cdot 9; 1 \cdot 0$: (same)

$$A_{0 \cdot 7} = [8, 8], B_{0 \cdot 7} = [3, 9], C_{0 \cdot 7} = [11, 11]$$

$$C_{0 \cdot 8} = [11, 11]$$

$$C_{0 \cdot 9} = [11, 11]$$

$$C_{1 \cdot 0} = [11, 11]$$