

## نوايا المذاكرة

- ١- رضا الله.
- ٢- رضا الوالدين وادخال السرور على قلوبهم.
- ٣- إجابة عن سؤاله وعن عمره فيما أفتنه ...
- ٤- نفع المسلمين بما تعلمت وسد ثغور الأمة.
- ٥- تحقيق الأهداف وما خلقت لأجلة.
- ٦- تعمير الأرض.
- ٧- رفع الجهل.
- ٨- المؤمن القوي خير وأحب إلى الله من المؤمن الضعيف.

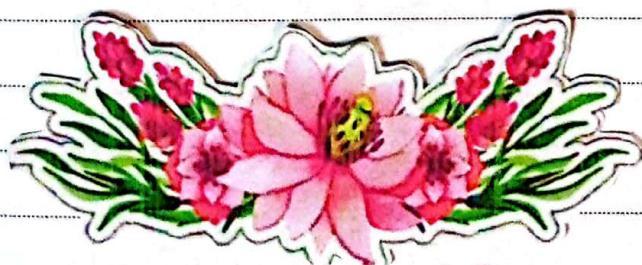
«مذاكرتك مسؤوليتك كمسلم فلا تخاذل!»

تعلم العلم فلا تدرى متى تحتاج للكتابة!

علمك = شركك = مقاومتك ❤

لا تتسرّى وأهلى وأهل غزة وال المسلمين المُلاجئ عن  
في الأرض من دعواكم ولكم بالمثل ❤

\* نرحل ويبقى الأثر  
\* أفعل الخير ولا تستصغر، فلا تدرى أي حسنة  
تُدخلك الجنة



## بصمة لد متطر

1- Quadratic function has Numerical & Analytical Solutions.

2- The derivative at the max & min points is zero.

3- Gradient descent idea is to take repeated steps

against the gradient direction to find the min. of the fun.

4- Linear Regression  $\Rightarrow Y = WX + b$  or  $Y_K = \sum_{i=1}^d W_{ki} X_i + b_K$

5- Objective function of Linear regression

Mean Square error "MSE" =  $\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

or  $\frac{1}{N} \sum_{n=1}^N E^n$  or  $\sum_{n=1}^N \sum_{K=1}^K (Y_K^n - t_K^n)^2$  or  $\sum_{n=1}^N \sum_{K=1}^K (\sum_{i=1}^d W_{ki} X_i + b_K - t_K^n)^2$

6- Activation function of Linear Regression is Linear Function.

7- The gradient of MSE loss function with respect to

Weight  $\Rightarrow (Y_r^n - t_r^n) X_s^n$  or  $\frac{1}{N} \sum_{n=1}^N (Y_r^n - t_r^n) X_s^n$

bias  $\Rightarrow (Y_r^n - t_r^n)$  or  $\frac{1}{N} \sum_{n=1}^N (Y_r^n - t_r^n)$

8- Update for weights & bias Vector

$W_{rs}^{t+1} = W_{rs}^t - \eta \frac{\partial E(w,b)}{\partial W_{rs}}$  or  $W_{rs}^{t+1} = W_{rs}^t - \eta \frac{1}{N} \sum_{n=1}^N (Y_r^n - t_r^n) X_s^n$

$b_r^{t+1} = b_r^t - \eta \frac{\partial E(w,b)}{\partial b_r}$  or  $b_r^{t+1} = b_r^t - \eta \frac{1}{N} \sum_{n=1}^N (Y_r^n - t_r^n)$

9- Objective function of Linear regression is Scalar. "MSE"

10- Predict house price  $\Rightarrow$  Application of Linear Regression.

11- Binary activation function isn't Continuous and not differentiable at zero.

12- Sigmoid is Continuous & differentiable.

13- Sigmoid  $\Rightarrow F(z) = y = \frac{1}{1 + e^{-z}}$



- 14- Objective function of binary Classification is Cross entropy or Log maximum Likelihood
- 15- Binary Cross entropy (BCE) Loss function
- 16- Objective function of binary Classification  $\Rightarrow E^n = -[t^n \log y^n + (1-t^n) \log (1-y^n)]$
- 17- Range of Sigmoid Function is  $[0, 1]$ .
- 18- Output of Sigmoid function when Input is zero is 0.5
- 19- Sigmoid maximum Score is 1
- 20- Gradient of loss function with respect to weights =  $(y^n - t^n)x_i$   
with respect to bias =  $(y^n - t^n)$
- 21- Value of y is between 0 & 1, so to make decision we use Threshold (0.5).  $y = \begin{cases} 1 & \text{if } y \geq 0.5 \\ 0 & \text{if } y < 0.5 \end{cases}$
- 22- Threshold for decision  $\rightarrow y = \begin{cases} 1 & \text{if } y \geq 0.5 \\ 0 & \text{if } y < 0.5 \end{cases}$
- 23- Binary classification must have only one output scalar
- 24- Sigmoid function has a Probabilistic interpretation.
- 25- Minimum value of cross entropy objective function is zero.
- 26- Spam email detection is application of binary Classification.
- 27- Activation Function of Multiclass Classification is Softmax.
- 28- Softmax  $\Rightarrow y_k = \frac{e^z}{\sum_j e^z}, 0 < y_k < 1 \quad \sum_k y_k = 1$
- 29- Compute the Jacobian matrix g-  
 if  $k=1 \rightarrow y_i(1-y_k)$  or  $y_i - y_i y_k$   
 if  $k \neq 1 \rightarrow -y_i y_k$
- 30- Derivative of Softmax is Matrix (Jacobian matrix).
- 31- Objective function of Multiclass Classification is Categorical Cross entropy  $\Rightarrow E^n = -\sum_{k=1}^K t_k^n \log y_k^n$

- 32- The Summation of Softmax Vector elements equal to 1  
 33- Softmax Function is Vector, & its derivative is Matrix  
 34- Categorical Cross entropy objective function is Scalar.  
 35- In One-hot encoding Only a Single element Can have a Value of 1 . ex: [1,0,0] , [0,1,0,0]  
 36- Cross entropy is an objective function for binary & Multiclass classifications. "Scalar"  
 37- Recognition of handwritten digits (0-9) is an application of Multiclass Classification.  
 38- Objective Function of MultiLabel classification is K binary Cross entropy  $\rightarrow E = -t \log y - (1-t) \log(1-y)$   
 39- Multi-hot Vector more than one element can have a Value of One ex: [0,1,1,1] , [0,0,0,0] , [1,1,1,1]  
 40- Output of Linear regression is Scalar.  $K \geq 1$   
 41- Output of Binary Classification is SCalar  $K=1$   
 42- Output of Multiclass Classification is vector  $K \geq 3$   
 43- Output of MultiLabel Classification is Vector  $K \geq 2$   
 44- Range of tanh activation function is  $[-1, 1]$   
 45- Output of tanh function when input is zero = 0 , any +Ve no ( $\infty$ ) = 1 , any -ve no ( $-\infty$ ) = -1  
 46- Relu Function =  $F(X) = \max(0, X)$   
 47- Non linear activation Functions help the network to avoid Saturation or Vanishing gradients.  
 48- the role of hidden Layer(s) is to perform Computations & extract features From the data.  
 49- Information Flow Forward Only , from input to output in multilayer network  
 50- The main benefit of using multilayer network is the ability to Learn more Complex relationships in data.

الLinear regression ممكن تكون vector برضو ودا علي حسب التاسك.

في التاسك اللي اخذناه بتاع توقع سعر البيت كانت scalar

51- Convolution operation is  $\Rightarrow S(t) = \int X(a) W(t-a) da$   
 $= X(t) * W(t) = W(t) * X(t)$  "Commutative"

52- Discrete Convolution :  $S(n) = \sum_{a=-\infty}^{\infty} X(a) W(n-a)$

53- Two dimensional Convolution

$$S(i,j) = (X * W)(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} X_{k,l} W_{i-k, j-l}$$

54- Cross Correlation  $\Rightarrow$

$$S(i,j) = \sum_{k=0}^{m-1} \sum_{l=0}^{n-1} X_{i+k, l} W_{k,l} = (W \otimes X)(i,j)$$

55- Equation of the output size =  $(W - F + 2P)/S + 1$

$$\frac{W - F + 2P}{S} + 1 \Rightarrow \frac{\text{input size} - \text{Filter size} + 2 \text{ Padding}}{\text{Stride}} + 1$$

56- Dense Layers in CNN may be Linear regression , or  
Binary Classification or Multi-class Classification.

57- Main Purpose of Convolutional Layer in CNN is to extract  
Features from the input image.

58- The role of Pooling Layer in CNN is to downsample the  
Feature maps and reduce Computational Cost.

59- Padding helps to preserve spatial information.

60- Purpose of dropout layers in CNN is to prevent  
Overfitting.

61- Stride is the step size taken by the filter when moving  
across the input data.

62- The most Suitable activation function for hidden Layer  
is Relu.

63- The input of CNN is matrix , 2D matrix for grayscale  
image & 3D matrix for Color Image (RGB channels).

على حسب انا هسختم اي في  $\Rightarrow$  Dense Layer ١١

67- RNNs are nonlinear models.

68-  $h_t = F_w(h_{t-1}, X_t)$

69- Hidden State activation  $\Rightarrow Z_t^h = W_h h_{t-1} + W_x X_t + b_h$   
 $h_t = \tanh(Z_t^h)$

70- Output activation  $\Rightarrow Z_t^y = W_y h_t + b_y$

$y_t = \text{Softmax}(Z_t^y)$

71- Norms of Jacobians ٨

•  $\left\| \frac{\partial h_{j+1}}{\partial h_j} \right\| \leq \|W_h\| \cdot \left\| \text{diag}(\phi'_h(h_j)) \right\|$  "الكتاب"

$\phi_h = \tanh$  activation function &  $\phi'_h$  is its derivative

•  $\left\| \prod_{i=K+1}^t \frac{\partial h_i}{\partial h_{i-1}} \right\| \leq \prod_{i=K+1}^t \|W_i\| \cdot \left\| \text{diag}(F'(h_{t-1})) \right\| \leq (\gamma_w \gamma_h)^{t-K}$  نوتس الدكتور

•  $\left\| \prod_{j=k}^t \frac{\partial h_{j+1}}{\partial h_j} \right\| \leq (\gamma_w \gamma_h)^{t-k}$  "الكتاب"  
الآن نفس الاتجاه بس عسان تكون عارف لما يجي في المقدمة

72- if the distance  $t-K$  is Large  $\Rightarrow \gamma > 1$  gradient

will Vanish, if  $\gamma > 1$  gradient will explode.

73- input gate  $\Rightarrow i_t = \sigma(W_i h_{t-1} + U_i X_t)$

- Forget gate  $\Rightarrow f_t = \sigma(W_f h_{t-1} + U_f X_t)$

- Output gate  $\Rightarrow o_t = \sigma(W_o h_{t-1} + U_o X_t)$

memory Cell  $\Rightarrow C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$

$\tilde{C}_t = \tanh(W_c h_{t-1} + U_c X_t)$

$C_t = f_t \odot C_{t-1} + i_t \odot \tanh(W_c h_{t-1} + U_c X_t)$

$h_t = o_t \odot \tanh(C_t) \Rightarrow$  hidden state



74- What function is typically used to activate the forget gate, input gate, and output gate in LSTMs ?  
Sigmoid for gates , tanh for cell state.

75- What's the main advantage of LSTMs Compared to Standard RNNs ? They Can Learn Long-term dependencies in data.

76- What Kind of neural network architecture is LSTM a Part of?  
Recurrent Neural Network (RNN)

77- Compared to a Standard Classification network with softmax activation, How does the output layer of multi-label Classification network differ ? It uses a Sigmoid activation function on each node.

78- what is Backpropagation Through Time (BPTT) ?

An extension of backpropagation that accounts for the recurrent nature of the network.

79- What is the Loss function used in RNN training ?

Categorical Cross-Entropy  $\Rightarrow E^n = - \sum_{k=1}^K t_k^n \log y_k^n$

80- What is the formula for the overall Loss E ?

$$\textcircled{1} \quad E = \sum_{t=1}^T E_t(L_t, y_t)$$

$$\textcircled{2} \quad E = - \sum_{t=1}^T L_t \log Y_t$$

$$\textcircled{3} \quad E = - \sum_{t=1}^T L_t \log [\text{softmax}(Y_t)]$$

81- Seq2Seq Translation, Input : Sequence of word vectors from the Source language (Matrix)

Output : Sequence of word vectors in the target language (Matrix)

سُبْحَانَ اللَّهِ، الْحَمْدُ لِلَّهِ  
خَلَقَ اللَّهُ، الْعَزِيزُ  
خَوْلٌ وَلَا قُوَّةٌ إِلَّا بِاللَّهِ  
الْفُعُولُ عَلَى سِنِّيَّةِ  
كُلِّ مُؤْمِنٍ وَسَعْيُهُ



82- Which Problem does the attention mechanism address in Seq2Seq model ? - Difficulty in processing Long Sequences.

83- How does the attention mechanism improve the performance of Seq2Seq models ?

By assigning different weights to different parts of the input sequence .

84- What's the role of the Context Vector in a Seq2Seq model with attention ?

Carry the weighted sum of the encoder hidden states .

85- what is the key difference between a traditional RNN & LSTM used in Seq2Seq models ?

LSTMs have gates to control the flow of information .

86- what is the benefit of using bidirectional RNNs in the encoder part of a Seq2Seq model ?

They provide context from both past & future tokens in the sequence .

87- What Problem can arise if Seq2Seq models without attention mechanisms are used for long sequences ?

Vanishing gradients .

88- What is the main advantage of using an LSTM or GRU over a traditional RNN in Seq2Seq model ?

They mitigate the vanishing gradient problem .

89- What is a significant benefit of multi-head attention in transformer models ?

It allows the model to focus on different parts of the input simultaneously .

قال رسول الله صلى الله عليه وسلم: إِنَّ اللَّهَ يُحِبُّ إِذَا عَمِلَ أَخْذَهُ عَمَلًا" أَنْ يُتَقْتَلَهُ .



90- What does the 'alignment model' in an attention based Seq2Seq System do?

It calculates the attention weights.

91- What type of Function is Commonly Used to calculate the alignment scores in attention mechanisms?

Softmax Function.

92- What is the main function of the "query", "key" & "Value" Vectors in the attention mechanism?

They are used to compute the alignment scores for attention.

93- How does the transformer model handle the Vanishing gradient problem that is common in RNNs?

By employing residual connections.

\* Transformer model uses residual connections to mitigate the vanishing gradient problem, allowing gradients to flow more easily through the network during backpropagation.

94- Which mechanism allows transformers to better capture long dependencies compared to RNNs?

Self - attention

95- What is the purpose of layer normalization in Transf. Stabilize & accelerates training.

96- What is the purpose of the Softmax function applied to the alignment scores?

To convert the scores into probabilities.

97- In transformer model, what type of input is provided to the encoder & decoder? Matrices.

98- What is purpose of key matrix K?

To compute the dot product with query matrix.



99- In multi-head attention mechanism, what type of output is produced by each individual attention head?

Sequence of Vectors (Matrix)

100- Which equation describes the scaling factor applied in the scaled dot-product attention?

$$\text{Scale} = \sqrt{d_k}$$

101- Equation represents the attention scores before scaling  $\Rightarrow \text{Attention}(Q, K, V) = \text{Softmax}(QK^T)V$

102- Equation represents the scaling step in the scaled dot-product attention mechanism?

$$\text{Scaled Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

103- Equation represents the calculation of the attention output in the multi-head attention mechanism after concatenation?

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_0$$

104- Equation represents the transformation applied to each attention head in the multi-head mechanism?

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

105- Equation represents the output of the feed-forward neural network (FFN) for a given input  $x$ ?

$$\text{FFN}(x) = W_2(\text{ReLU}(W_1x + b_1)) + b_2$$

علاقة استقرارية السعى بتخمين اليمى  
... بذلك يمكن حتى حقن لوما ومحاسن في الآخر.

- 106- RNNs are Sequential models, operate in a recurrent manner.
- 107- Transformers are non-Sequential models that process data in Parallel.
- 108- Transformers rely on Self attention mechanisms. They don't have recurrent Connections or hidden states.
- 109- Transformers Can handle both short & long Sequences due to their parallel processing nature.
- 110- Model architecture is best suited for image Classification task? CNN

«فِي زَاهِدِ الْأَيَادِ لَا تَنْسِي الغَارِيَةَ الَّتِي خَلَقْتَ لِأَجْلِهَا،  
إِنْتَ هُنَا وَدِيْعَةٌ سُرُّدٌ إِلَى رِبْهَا يَوْمًا فَأَحْسَنَ»



لَوْ أَنْتَ عَلَى بِإِخْرَاصِ فِي أَيِّ مَكَانٍ، حَتَّى لَوْ عَدَّشَ سَمْعَ  
بَيْكَ أَوْ شَاقَّكَ، هَلْ تَلْقَى الرِّزْقَ بِيَدِكَ عَلَيْكَ، رَبِّنَا  
بِيَحْبَبِ عِبَادِهِ الْمُخْلِصِينَ، وَالَّذِي رَبَّنَا صَرِحَ بِجَبَهِهِ مِنْ تَخَافِشِ  
عَلَيْهِ ❤

لَيْسَ كُلُّ فُجُورٍ هُدُّى بِالْخُضُورَةِ، وَلَكِنْ فُجُورٌ مَأْجُورٌ  
فَالْحَمْدُ لِلَّهِ الَّذِي يَجْازِي نَاساً عَلَى السُّعْيِ وَالْمَحَاوِلَاتِ،  
لَا عَلَى النَّتْيُوجَةِ وَالْوُصُولِ وَالْكَمَالِ.

أَذْكُرْنَفْسِي وَطَيَاكِمْ بِقُولِ الْحَبِيبِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ  
«هُنَّ سَلَّكُ طَرِيقًا يُلْتَهِسُ فِيهِ عَنْهَا، سَعْلَ اللَّهِ عَلَيْهِ  
طَرِيقًا إِلَى الْجَنَّةِ ❤»

