

1. (a) Give brief answers to the following. Each part carries 2 marks.

- (i) What is the total number of learnable parameters in a 2-layer network with 5 neurons in each layer and receiving 2-dimensional inputs.

Answer

Given

- layer = 2
- neuron = 5
- input = 2
- bias = 1

layer 1

- neuron * (input + bias) = $5 * (2 + 1) = 5 * 3 = 15$ parameters

layer 2

- neuron * (neuron + bias) = $5 * (5 + 1) = 5 * 6 = 30$ parameters

Total = $15 + 30 = 45$ parameters

- (ii) Write a tensor with the shape [2, 1, 3] and integer elements, as in Python.

Answer

[[[1., 2., 3.]], [[7., 8., 9.]]]

- (iii) One wants to initialize all the weights in a 3-layer deep neural networks to 0.4. Is this a good idea? And why?

Answer

- No, it is not a good idea.
- If all weights are the same, every neuron in a layer will compute the same gradient during backpropagation.

- (iv) You want to solve a classification problem with a training sample of 50 patterns. But the training loss was high. Then you decide to train with 1000 samples. Is this good approach? And why?

Answer

- Yes, it is a good approach.
- A larger training set gives the model more examples to learn from, improving its ability to generalize to unseen data.

1. (a) cont

- (v) You started training a neural network but the loss was neither decreasing or increasing. State two possible reasons.

Answer

- Learning rate is too low or too high.
- The gradients become too small or too large.

- (vi) State two ways to handle local minima problem in gradient descent learning.

Answer

- Use SGD or mini-batch GD
- Using adaptive learning rate methods (e.g., Adam)

- (vii) State one advantage and one disadvantage of using a small batch size for training.

Answer

- Helps the local minima problem
- Slower training, slower convergence

- (viii) State how you can fix underfitting in your neural network.

Answer

- Train longer (more epochs)
- Decrease learning rate or batch size

(16 marks)

1. (b) The output $\mathbf{y} \in \mathbf{R}^K$ of a linear neuron is given by $\mathbf{y} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$ where $\mathbf{x} \in \mathbf{R}^n$ is input, and \mathbf{W} and \mathbf{b} are the weight matrix and bias vector of the layer, respectively. Given a target vector $\mathbf{d} \in \mathbf{R}^K$, the square error loss J of the layer is given $J = (\mathbf{d} - \mathbf{y})^T (\mathbf{d} - \mathbf{y})$.

- (i) Write the expressions for $\nabla_{\mathbf{y}} J$, $\nabla_{\mathbf{W}} J$ and $\nabla_{\mathbf{b}} J$.

(5 marks)

Answer

$J = (\mathbf{d} - \mathbf{y})^T (\mathbf{d} - \mathbf{y})$ $= \ \mathbf{d} - \mathbf{y}\ ^2$ $\nabla_{\mathbf{y}} J = \frac{\partial J}{\partial \mathbf{y}}$ $= -2(\mathbf{d} - \mathbf{y})$	$\nabla_{\mathbf{W}} J = \frac{\partial J}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{W}}$ $\mathbf{y} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$ $\frac{\partial \mathbf{y}}{\partial \mathbf{W}} = \mathbf{x}$ $\nabla_{\mathbf{W}} J = \frac{\partial J}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{W}}$ $= -2(\mathbf{d} - \mathbf{y}) \mathbf{x}^T$	$\nabla_{\mathbf{b}} J = \frac{\partial J}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{b}}$ $\mathbf{y} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$ $\frac{\partial \mathbf{y}}{\partial \mathbf{b}} = \mathbf{1}$ $\nabla_{\mathbf{b}} J = \frac{\partial J}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{b}}$ $= -2(\mathbf{d} - \mathbf{y})$
---	---	--

Summary

$$\nabla_{\mathbf{y}} J = -2(\mathbf{d} - \mathbf{y})$$

$$\nabla_{\mathbf{W}} J = -2(\mathbf{d} - \mathbf{y}) \mathbf{x}^T$$

$$\nabla_{\mathbf{b}} J = -2(\mathbf{d} - \mathbf{y})$$

- (ii) If L_2 -norm weight regularization penalty $\beta \|\mathbf{W}\|^2$ is added to the loss J , write the expressions for $\nabla_{\mathbf{W}} J_1$ and $\nabla_{\mathbf{b}} J_1$ is the regularized loss. β is the penalty parameter.

(4 marks)

Answer

$$J_1 = (\mathbf{d} - \mathbf{y})^T (\mathbf{d} - \mathbf{y}) + \beta \|\mathbf{W}\|^2$$

Regularization $\beta \|\mathbf{W}\|^2$ does not depend on bias, gradients unchanged

$$\nabla_{\mathbf{W}} (\beta \|\mathbf{W}\|^2) = 2\beta \mathbf{W}$$

$$\nabla_{\mathbf{b}} J_1 = -2(\mathbf{d} - \mathbf{y})$$

$$\nabla_{\mathbf{W}} J_1 = -2(\mathbf{d} - \mathbf{y}) \mathbf{x}^T + 2\beta \mathbf{W}$$

Summary

$$\nabla_{\mathbf{W}} J_1 = -2(\mathbf{d} - \mathbf{y}) \mathbf{x}^T + 2\beta \mathbf{W}$$

$$\nabla_{\mathbf{b}} J_1 = -2(\mathbf{d} - \mathbf{y})$$

2. The 3-layer feedforward neural network shown in Figure Q2 receives two-dimensional inputs $(x_1, x_2) \in \mathbf{R}^2$ and produces one-dimensional output y . The first hidden layer consists of three neurons and the second hidden layer consists of two neurons. All hidden neurons have *Tanh* activations functions and the output neuron is a **logistic regression** neuron. The weights of the networks are initialized as indicated in the figure and all the biases are initialized to 0.5 (not shown).

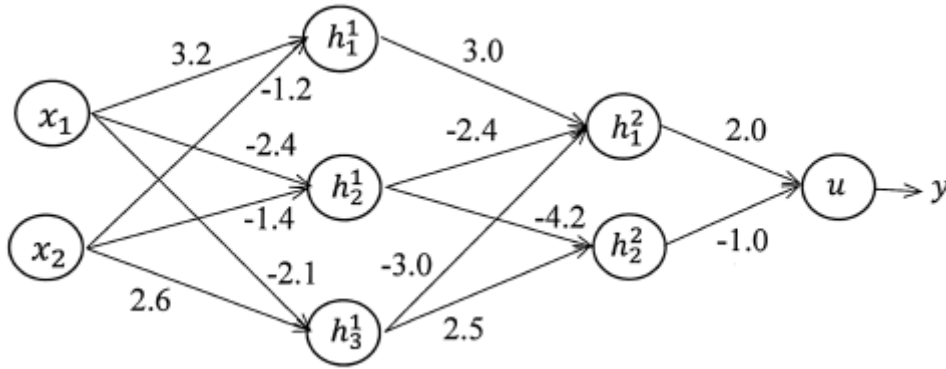


Figure Q2

The network is trained to produce a desired output $d = 0$ for an input $x = \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix}$. You are to perform one iteration of **stochastic gradient descent** learning with the example (x, d) . Give your answers rounded up to two decimal places.

2. (a) Write initial weight matrices \mathbf{W} and bias vectors \mathbf{b} , connected to the three layers. (3 marks)

Answer

$$W_1 = \begin{pmatrix} 3.2 & -2.4 & -2.1 \\ -1.2 & -1.4 & 2.6 \end{pmatrix} \quad b_1 = \begin{pmatrix} 0.5 \\ 0.5 \\ 0.5 \end{pmatrix}$$

$$W_2 = \begin{pmatrix} 3.0 & 0.0 \\ -2.4 & -4.2 \\ -3.0 & 2.5 \end{pmatrix} \quad b_2 = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

$$W_3 = \begin{pmatrix} 2.0 \\ -1.0 \end{pmatrix} \quad b_3 = 0.5$$

2. (b) Find the synaptic inputs \mathbf{u} and outputs \mathbf{h} of the two hidden layers.

(4 marks)

Answer

$$\mathbf{W}_1 = \begin{pmatrix} 3.2 & -2.4 & -2.1 \\ -1.2 & -1.4 & 2.6 \end{pmatrix} \quad \mathbf{b}_1 = \begin{pmatrix} 0.5 \\ 0.5 \\ 0.5 \end{pmatrix}$$

$$\mathbf{W}_2 = \begin{pmatrix} 3.0 & 0.0 \\ -2.4 & -4.2 \\ -3.0 & 2.5 \end{pmatrix} \quad \mathbf{b}_2 = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

Input $\mathbf{x} = \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix}$

Synaptic input to h_1 ,

$$\begin{aligned} u_1 = \mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1 &= \begin{pmatrix} 3.2 & -1.2 \\ -2.4 & -1.4 \\ -2.1 & 2.6 \end{pmatrix} \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \\ 0.5 \end{pmatrix} = \begin{pmatrix} (3.2)(1) + (-1.2)(2) \\ (-2.4)(1) + (-1.4)(2) \\ (-2.1)(1) + (2.6)(2) \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \\ 0.5 \end{pmatrix} \\ &= \begin{pmatrix} 0.8 \\ -5.2 \\ 3.1 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 1.3 \\ -4.7 \\ 3.6 \end{pmatrix} \end{aligned}$$

Output of h_1 ,

$$h_1 = \tanh(u_1) = \frac{e^u - e^{-u}}{e^u + e^{-u}} = \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix}$$

Synaptic input to h_2 ,

$$\begin{aligned} u_2 = \mathbf{W}_2^T \mathbf{h}_1 + \mathbf{b}_2 &= \begin{pmatrix} 3.0 & -2.4 & -3.0 \\ 0.0 & -4.2 & 2.5 \end{pmatrix} \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \\ &= \begin{pmatrix} (3)(0.86) + (-2.4)(-1) + (-3.0)(1) \\ (0)(0.86) + (-4.2)(-1) + (2.5)(1) \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} \\ &= \begin{pmatrix} 1.98 \\ 6.7 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 2.48 \\ 7.20 \end{pmatrix} \end{aligned}$$

Output of h_2 ,

$$h_2 = \tanh(u_2) = \frac{e^u - e^{-u}}{e^u + e^{-u}} = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix}$$

Summary

synaptic input $u_1 = \begin{pmatrix} 1.3 \\ -4.7 \\ 3.6 \end{pmatrix}$

output $h_1 = \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix}$

synaptic input $u_2 = \begin{pmatrix} 2.48 \\ 7.20 \end{pmatrix}$

output $h_2 = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix}$

2. (c) Find the output y and the cross-entropy at the output layer.

(4 marks)

Answer

$$W_3 = \begin{pmatrix} 2.0 \\ -1.0 \end{pmatrix} \quad b_3 = 0.5$$

Desired output $d = 0$

$$h_2 = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix}$$

Synaptic input to y ,

$$\begin{aligned} u_3 &= W_3^T h_2 + b_3 = \begin{pmatrix} 2 & -1 \end{pmatrix} \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix} + (0.5) \\ &= ((2)(0.99) + (-1)(1)) + (0.5) \\ &= 0.98 + 0.5 = 1.48 \end{aligned}$$

Output of y ,

$$y = \text{sigmoid}(u_3) = \frac{1}{1 + e^{-u}} = \frac{1}{1 + e^{-1.48}} = 0.81$$

cross-entropy,

$$\begin{aligned} J &= -d \log(y) - (1 - d) \log(1 - y) \\ &= -0 \log(0.81) - (1 - 0) \log(1 - 0.81) \\ &= -\log(1 - 0.81) \\ &= -\log(0.19) = 1.66 \end{aligned}$$

Summary

output: $y = 0.81$

cross-entropy: $J = 1.66$

2. (d) Find the **derivatives** $f'(\mathbf{u})$ at the two hidden layers with respect to synaptic input \mathbf{u} , where f is the *Tanh* activation function.

(2 marks)

Answer

Derivatives of *Tanh*

$$y = f(u) = \tanh(u) = \frac{e^{+u} - e^{-u}}{e^{+u} + e^{-u}}$$

$$f'(u) = 1 - y^2$$

Hidden layer h_1 ,

$$h_1 = \tanh(u_1) = \frac{e^u - e^{-u}}{e^u + e^{-u}} = \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix}$$

$$f'(u_1) = 1 - h_1^2 = \begin{pmatrix} 1 - (0.86)^2 \\ 1 - (1.00)^2 \\ 1 - (-1.00)^2 \end{pmatrix} = \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix}$$

Hidden layer h_2 ,

$$h_2 = \tanh(u_2) = \frac{e^u - e^{-u}}{e^u + e^{-u}} = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix}$$

$$f'(u_2) = 1 - h_2^2 = \begin{pmatrix} 1 - (0.99)^2 \\ 1 - (1.00)^2 \end{pmatrix} = \begin{pmatrix} 0.02 \\ 0 \end{pmatrix}$$

Summary

derivatives $f'(u_1) = \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix}$

$$f'(u_2) = \begin{pmatrix} 0.02 \\ 0 \end{pmatrix}$$

2. (e) Find gradients $\nabla_{\mathbf{u}} J$ of the cost \mathbf{u} with respect to activations \mathbf{u} , at the three layers.
(6 marks)

Answer

Backpropagation for FFN:

Desired output $d = 0$

$$y = \text{sigmoid}(u_3) = 0.81$$

Outer layer u_3 ,

$$\nabla_{u_3} J = -(d - y) = -(0 - 0.81) = 0.81$$

Second hidden layer u_2 ,

$$W_3 = \begin{pmatrix} 2.0 \\ -1.0 \end{pmatrix} \quad f'(u_2) = \begin{pmatrix} 0.02 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \nabla_{u_2} J &= W_3 (\nabla_{u_3} J) \cdot f'(u_2) = \begin{pmatrix} 2.0 \\ -1.0 \end{pmatrix} (0.81) \cdot \begin{pmatrix} 0.02 \\ 0 \end{pmatrix} = \begin{pmatrix} (2)(0.81) \\ (-1)(0.81) \end{pmatrix} \cdot \begin{pmatrix} 0.02 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} 1.62 \\ -0.81 \end{pmatrix} \cdot \begin{pmatrix} 0.02 \\ 0 \end{pmatrix} = \begin{pmatrix} 0.03 \\ 0 \end{pmatrix} \end{aligned}$$

First hidden layer u_1 ,

$$W_2 = \begin{pmatrix} 3.0 & 0.0 \\ -2.4 & -4.2 \\ -3.0 & 2.5 \end{pmatrix} \quad f'(u_1) = \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \nabla_{u_1} J &= W_2 (\nabla_{u_2} J) \cdot f'(u_1) = \begin{pmatrix} 3.0 & 0.0 \\ -2.4 & -4.2 \\ -3.0 & 2.5 \end{pmatrix} \begin{pmatrix} 0.03 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} (3.0)(0.03) + (0.0)(0) \\ (-2.4)(0.03) + (-4.2)(0) \\ (-3.0)(0.03) + (2.5)(0) \end{pmatrix} \cdot \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} 0.09 \\ -0.07 \\ -0.09 \end{pmatrix} \cdot \begin{pmatrix} 0.26 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0.02 \\ 0 \\ 0 \end{pmatrix} \end{aligned}$$

Summary

$$\nabla_{u_3} J = 0.81$$

$$\nabla_{u_2} J = \begin{pmatrix} 0.03 \\ 0 \end{pmatrix}$$

$$\nabla_{u_1} J = \begin{pmatrix} 0.02 \\ 0 \\ 0 \end{pmatrix}$$

2. (f) Find gradients $\nabla_W J$, and $\nabla_b J$ of the cost J with respect to weights W , and biases b , respectively, at three layers.

(6 marks)

Answer

Outer layer u_3 ,

$$h_2 = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix} \quad \nabla_{u_3} J = 0.81$$

$$\nabla_{W_3} J = h_2 (\nabla_{u_3} J)^T = \begin{pmatrix} 0.99 \\ 1.00 \end{pmatrix} 0.81 = \begin{pmatrix} (0.99)(0.81) \\ (1.00)(0.81) \end{pmatrix} = \begin{pmatrix} 0.80 \\ 0.81 \end{pmatrix}$$

$$\nabla_{b_3} J = \nabla_{u_3} J = 0.81$$

Second hidden layer u_2 ,

$$h_1 = \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix} \quad \nabla_{u_2} J = \begin{pmatrix} 0.03 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \nabla_{W_2} J &= h_1 (\nabla_{u_2} J)^T = \begin{pmatrix} 0.86 \\ -1.00 \\ 1.00 \end{pmatrix} \begin{pmatrix} 0.03 & 0 \end{pmatrix} = \begin{pmatrix} (0.86)(0.03) & (0.86)(0) \\ (-1.00)(0.03) & (-1.00)(0) \\ (1.00)(0.03) & (1.00)(0) \end{pmatrix} \\ &= \begin{pmatrix} 0.03 & 0 \\ -0.03 & 0 \\ 0.03 & 0 \end{pmatrix} \end{aligned}$$

$$\nabla_{b_2} J = \nabla_{u_2} J = \begin{pmatrix} 0.03 \\ 0 \end{pmatrix}$$

First hidden layer u_1 ,

$$x = \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix} \quad \nabla_{u_1} J = \begin{pmatrix} 0.02 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \nabla_{W_1} J &= x (\nabla_{u_1} J)^T = \begin{pmatrix} 1.0 \\ 2.0 \end{pmatrix} \begin{pmatrix} 0.02 & 0 & 0 \end{pmatrix} = \begin{pmatrix} (1)(0.02) & (1)(0) & (1)(0) \\ (2)(0.02) & (2)(0) & (2)(0) \end{pmatrix} \\ &= \begin{pmatrix} 0.02 & 0 & 0 \\ 0.04 & 0 & 0 \end{pmatrix} \end{aligned}$$

$$\nabla_{b_1} J = \nabla_{u_1} J = \begin{pmatrix} 0.02 \\ 0 \\ 0 \end{pmatrix}$$

Summary

$$\nabla_{W_3} J = \begin{pmatrix} 0.80 \\ 0.81 \end{pmatrix} \quad \nabla_{b_3} J = 0.81$$

$$\nabla_{W_2} J = \begin{pmatrix} 0.03 & 0 \\ -0.03 & 0 \\ 0.03 & 0 \end{pmatrix} \quad \nabla_{b_2} J = \begin{pmatrix} 0.03 \\ 0 \end{pmatrix}$$

$$\nabla_{W_1} J = \begin{pmatrix} 0.02 & 0 & 0 \\ 0.04 & 0 & 0 \end{pmatrix} \quad \nabla_{b_1} J = \begin{pmatrix} 0.02 \\ 0 \\ 0 \end{pmatrix}$$

3. Figure Q3 depicts a network that consists of two convolutional layers and a fully connected layer. The size of input or output volume is represented as $D \times H \times W$, where D is the number of channels, and $H \times W$ is the spatial size.

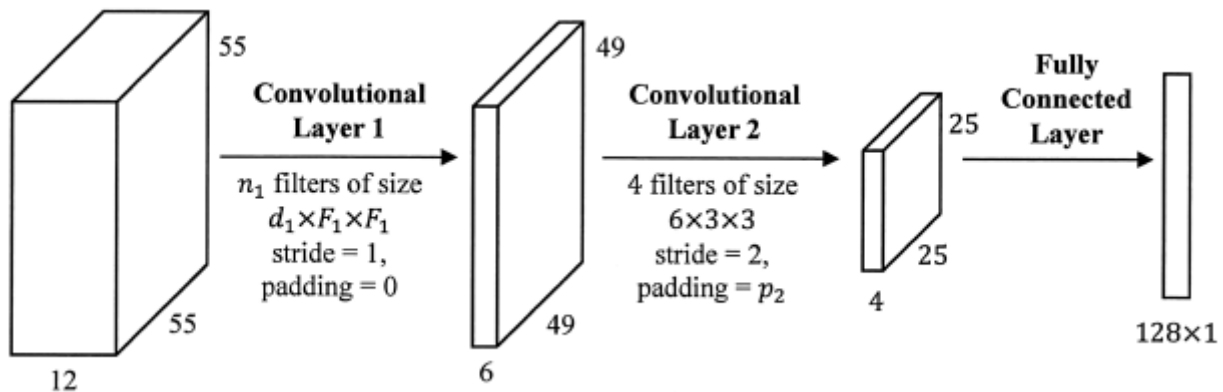


Figure Q3

3. (a) (i) Give the values of n_1, d_1, F_1 and p_2 .

(4 marks)

Answer

Convolution Layer 1

- Since output size: $6 \times 49 \times 49$, $n_1 = 6$
- Since input size: $12 \times 55 \times 55$, $d_1 = 12$

$$\text{Output} = \frac{N - F + 2P}{S} + 1$$

$$49 = \frac{55 - F_1 + 2(0)}{1} + 1$$

$$48 = 55 - F_1$$

$$F_1 = 7$$

- Output = 49
- $N = 55$
- $S = 1$
- $P = 0$

Convolution Layer 2

$$\text{Output} = \frac{N - F + 2P}{S} + 1$$

$$25 = \frac{49 - 3 + 2(p_2)}{2} + 1$$

$$48 = 46 + 2(p_2)$$

$$p_2 = 1$$

- Output = 25
- $N = 49$
- $S = 2$
- $F = 3$

Summary

- $n_1 = 6$
- $d_1 = 12$
- $F_1 = 7$
- $p_2 = 1$

3. (a) (ii) Calculate the total number of parameters in each layer, namely the Convolutional Layer 1, Convolutional Layer 2 and Fully Connected Layer. Be reminded to account for the bias terms.

(4 marks)

Answer

Convolution Layer 1

- filter size = $d_1 \times F_1 \times F_1$
- $d_1 = 12$
- $F_1 = 7$

- convolution filters = 6
- $n_1 = 6$

$$\begin{aligned}\text{Parameters per filter} &= (d_1 * F_1 * F_1) + 1 = (12 * 7 * 7) + 1 = 589 \\ \text{Total parameters} &= \text{param} * n_1 = 589 * 6 = 3534\end{aligned}$$

Convolution Layer 2

- $n_2 = 4$
- $d_2 = 6$
- $F_2 = 3$

$$\begin{aligned}\text{Parameters per filter} &= (d_2 * F_2 * F_2) + 1 = (6 * 3 * 3) + 1 = 55 \\ \text{Total parameters} &= \text{param} * n_2 = 55 * 4 = 220\end{aligned}$$

Fully Connected Layer

- input flattened = $4 * 25 * 25 = 2500$
- output = 128
- bias = 128

$$\text{Total parameters} = (\text{input} * \text{output}) + \text{bias} = (2500 * 128) + 128 = 320128$$

Total Parameters Summary

$$\begin{aligned}\text{Convolution Layer 1} &= 3534 \\ \text{Convolution Layer 2} &= 220 \\ \text{Fully Connected Layer} &= 320128\end{aligned}$$

3. (b) (i) Assume the same input volume $12 \times 55 \times 55$ and output volume $6 \times 49 \times 49$, replace the Convolution Layer 1 with a depthwise separable convolution, i.e., "depthwise + pointwise" convolutions. Assume stride = 1 and padding = 0. State the number of filters and the size of filter for each of the depthwise and pointwise convolution layers.

(4 marks)

Answer

Depthwise convolution layer

- Since input size: $12 \times 55 \times 55$, number of filters = 12
- Since filter size: $6 \times 7 \times 7$, size of filter = $1 \times 7 \times 7$

Pointwise convolution layer

- Since output size: $6 \times 49 \times 49$, number of filters = 6
- Since input size: $12 \times 55 \times 55$, size of filters = $12 \times 1 \times 1$

Summary

Depthwise convolution layer

- number of filters = 12
- size of filter = $1 \times 7 \times 7$

Pointwise convolution layer

- number of filters = 6
- size of filters = $12 \times 1 \times 1$

3. (b) (ii) Calculate the FLOPs of the original Convolution Layer 1, and the FLOPs of the new depthwise separable convolution designed by you in Q3(b)(i). Finally, compute the ratio of the FLOPs between the new and original layers.

(4 marks)

Answer

From Q3(b)(i),

Depthwise convolution layer

- number of filters $D_1 = 12$
- size of filter $F = 7$

Pointwise convolution layer

- Since output size: $6 \times 49 \times 49$, $D_2 * H_2 * W_2 = 6 * 49 * 49$

Original Convolution Layer 1

$$\begin{aligned} \text{FLOPs}_{\text{standard}} &= (2 * D_1 * F^2) * D_2 * H_2 * W_2 \\ &= 2 * 12 * 7^2 * 6 * 49 * 49 = 16,941,456 \end{aligned}$$

Depthwise separable convolution

$$\text{FLOPs}_{\text{depthwise}} = (2 * F^2) * D_1 * H_2 * W_2 = (2 * 7^2) * 12 * 49 * 49 = 2,823,576$$

$$\text{FLOPs}_{\text{pointwise}} = (2 * D_1) * D_2 * H_2 * W_2 = (2 * 12) * 6 * 49 * 49 = 345,744$$

$$\begin{aligned} \text{FLOPs}_{\text{separable}} &= \text{FLOPs}_{\text{depthwise}} + \text{FLOPs}_{\text{pointwise}} \\ &= 2,823,576 + 345,744 = 3,169,320 \end{aligned}$$

Ratio of the FLOPs

$$\text{Ratio} = \frac{1}{D_2} + \frac{1}{F^2} = \frac{1}{6} + \frac{1}{7^2} = \frac{55}{294} = 0.187$$

Summary

FLOPs of the original Convolution Layer 1 = 16,941,456

FLOPs of the new depthwise separable convolution = 3,169,320

Ratio of the FLOPs (new/original) layers = 0.187

3. (b) (iii) You want to apply batch normalization in your network. Explain why you should not choose a very small mini-batch size during your training.

(3 marks)

Answer

- Unstable estimates of mean and variance
- Inconsistent behaviour across batches
- Poor model generalization

3. (c) The statements below are all related to autoencoders. Answer "TRUE" or "FALSE" to the following statements. Each part carries 1 mark.

Answer

- F** (i) The primary goal of an autoencoder is to achieve high classification accuracy on labeled data.
- T** (ii) Autoencoders can be used as a pre-training step for deep networks in situations where labeled data is limited.
- F** (iii) Stacked autoencoders are created by training multiple autoencoders in parallel, with each one independently processing the input data.
- T** (iv) Overcomplete autoencoders have a latent space dimension that is larger than the input dimension.
- F** (v) Sparse autoencoders utilize a loss function that encourages the network to have many activations close to one in the bottleneck layer.
- F** (vi) Autoencoders with a symmetrical architecture have different layers and neuron counts in the encoder and decoder parts, up to the bottleneck.

(6 marks)

4. (a) Consider a Jordan-type recurrent neural network (RNN) that receives 2-dimensional input patterns $\mathbf{x} \in \mathbf{R}^2$ and has one hidden layer. The RNN has two neurons in the hidden layer (which are initialized to zeros) and one neuron in the output layer. The hidden layer neurons have *Tanh* activation functions and the output layer neurons use *Sigmoid* activation functions.

The weight matrices \mathbf{U} connecting the input to the hidden layer, the top-down recurrence weight matrix \mathbf{W} , and \mathbf{V} connecting the hidden output to the output layer are given by

$$\mathbf{U} = \begin{pmatrix} 0.2 & 0.3 \\ 0.8 & 0.9 \end{pmatrix}, \mathbf{W} = \begin{pmatrix} 2.0 & 1.0 \end{pmatrix} \text{ and } \mathbf{V} = \begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix}$$

All **bias** connections to neurons are set to **0.2**. The **output** layer is initialized to an output of **1.0**.

Find the output of the network for an input sequence $(\mathbf{x}(1), \mathbf{x}(2))$. Provide your answers rounded to four decimal places.

$$\mathbf{x}(1) = \begin{pmatrix} 1.5 \\ 1.0 \end{pmatrix} \text{ and } \mathbf{x}(2) = \begin{pmatrix} -3.0 \\ 2.0 \end{pmatrix}$$

(6 marks)

Answer

$$\mathbf{U} = \begin{pmatrix} 0.2 & 0.3 \\ 0.8 & 0.9 \end{pmatrix}, \mathbf{W} = \begin{pmatrix} 2.0 & 1.0 \end{pmatrix} \text{ and } \mathbf{V} = \begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}, \text{ and } c = b = 0.2$$

$$y = 1.0$$

$$\mathbf{x}(1) = \begin{pmatrix} 1.5 \\ 1.0 \end{pmatrix} \text{ and } \mathbf{x}(2) = \begin{pmatrix} -3.0 \\ 2.0 \end{pmatrix}$$

$$\phi(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

$$\sigma(u) = \text{sigmoid}(u) = \frac{1}{1 + e^{-u}}$$

4. (a) cont

Answer

$$\text{At } t = 1, \mathbf{x}(1) = \begin{pmatrix} 1.5 \\ 1.0 \end{pmatrix},$$

$$\mathbf{h}(t) = \phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{y}(t-1) + \mathbf{b})$$

$$\mathbf{h}(1) = \tanh(\mathbf{U}^T \mathbf{x}(1) + \mathbf{W}^T \mathbf{y}(0) + \mathbf{b})$$

$$= \tanh\left(\begin{pmatrix} 0.2 & 0.8 \\ 0.3 & 0.9 \end{pmatrix} \begin{pmatrix} 1.5 \\ 1.0 \end{pmatrix} + \begin{pmatrix} 2.0 \\ 1.0 \end{pmatrix} 1.0 + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right)$$

$$= \tanh\left(\begin{pmatrix} (0.2)(1.5) + (0.8)(1) \\ (0.3)(1.5) + (0.9)(1) \end{pmatrix} + \begin{pmatrix} (2)(1) \\ (1)(1) \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right)$$

$$= \tanh\left(\begin{pmatrix} 1.10 \\ 1.35 \end{pmatrix} + \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \tanh\left(\begin{pmatrix} 3.30 \\ 2.55 \end{pmatrix}\right) = \begin{pmatrix} 0.9973 \\ 0.9879 \end{pmatrix}$$

$$\mathbf{y}(t) = \sigma(\mathbf{V}^T \mathbf{h}(t) + c)$$

$$\mathbf{y}(1) = \text{sigmoid}(\mathbf{V}^T \mathbf{h}(1) + c)$$

$$= \text{sigmoid}\left(\begin{pmatrix} 0.2 & -0.2 \end{pmatrix} \begin{pmatrix} 0.9973 \\ 0.9879 \end{pmatrix} + 0.2\right)$$

$$= \text{sigmoid}((0.2)(0.9973) + (-0.2)(0.9879) + 0.2) = \text{sigmoid}(0.2019) = 0.5503$$

$$\text{At } t = 2, \mathbf{x}(2) = \begin{pmatrix} -3.0 \\ 2.0 \end{pmatrix},$$

$$\mathbf{h}(t) = \phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{y}(t-1) + \mathbf{b})$$

$$\mathbf{h}(2) = \tanh(\mathbf{U}^T \mathbf{x}(2) + \mathbf{W}^T \mathbf{y}(1) + \mathbf{b})$$

$$= \tanh\left(\begin{pmatrix} 0.2 & 0.8 \\ 0.3 & 0.9 \end{pmatrix} \begin{pmatrix} -3.0 \\ 2.0 \end{pmatrix} + \begin{pmatrix} 2.0 \\ 1.0 \end{pmatrix} 0.5503 + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right)$$

$$= \tanh\left(\begin{pmatrix} (0.2)(-3) + (0.8)(2) \\ (0.3)(-3) + (0.9)(2) \end{pmatrix} + \begin{pmatrix} (2)(0.5503) \\ (1)(0.5503) \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right)$$

$$= \tanh\left(\begin{pmatrix} 1.0 \\ 0.9 \end{pmatrix} + \begin{pmatrix} 1.1006 \\ 0.5503 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \tanh\left(\begin{pmatrix} 2.3006 \\ 1.6503 \end{pmatrix}\right) = \begin{pmatrix} 0.9801 \\ 0.9289 \end{pmatrix}$$

$$\mathbf{y}(t) = \sigma(\mathbf{V}^T \mathbf{h}(t) + c)$$

$$\mathbf{y}(2) = \text{sigmoid}(\mathbf{V}^T \mathbf{h}(2) + c)$$

$$= \text{sigmoid}\left(\begin{pmatrix} 0.2 & -0.2 \end{pmatrix} \begin{pmatrix} 0.9801 \\ 0.9289 \end{pmatrix} + 0.2\right)$$

$$= \text{sigmoid}((0.2)(0.9801) + (-0.2)(0.9289) + 0.2) = \text{sigmoid}(0.2102) = 0.5524$$

Summary

- Output of the network:
 - $y(1) = 0.5503$
 - $y(2) = 0.5524$

4. (b) Select the correct option (A, B, C or D) for each question.

- A. Both statements are TRUE.
- B. Statement I is TRUE, but statement II is FALSE.
- C. Statement I is FALSE, but statement II is TRUE.
- D. Both statements are FALSE.

Answer

Answer

- B** (i) Statement I: The self-attention mechanism in Transformers allows each token in the input sequence to focus on different parts of the sequence when producing the output. T
Statement II: In Transformers, the number of self-attention heads is always fixed at one for all models and applications. F (2 marks)
- C** (ii) Statement I: The Transformer model uses convolutional layers to capture local patterns within the sequence data. F
Statement II: Layer normalization is a critical component in the Transformer's architecture, helping stabilize the activations throughout the network. T (2 marks)
- A** (iii) Statement I: The "query", "key", and "value" matrices in the self-attention mechanism of Transformers are all derived from the same initial input embeddings. T
Statement II: In the multi-head self-attention mechanism, different heads can potentially learn to attend to different parts or aspects of the input sequence. T (2 marks)

4. (c) The Transformer architecture employs a mechanism known as "positional encoding" to account for the order of words or tokens in a sequence.
- (i) Explain how sinusoidal positional encoding is computed and justify why it might be beneficial over other potential positional encoding methods.

(5 marks)

Answer

- *Sinusoidal positional encoding*
 - *interweaves the two signals (sine for even indices and cosine for odd indices)*

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

where pos is the position and i is the dimension,

- Beneficial
 - *Each position is uniquely encoded and the encoding can deal with sequences longer than any sequence seen in the training time.*

4. (ii) Based on the sinusoidal positional encoding, calculate the positional encoding values for dimension [0, 10] when the position of a word in the sequence is 5 (assuming the initial position is 0) and the dimension of the embeddings is 512. Provide your answers rounded to four decimal places. (4 marks)

Answer

Given $pos = 5$ and $d_{model} = 512$

At dimension 0, $2i = 0$ thus $i = 0$, therefore $PE_{(pos,2i)} = PE_{(5,0)} = \sin\left(\frac{5}{10000^{\frac{0}{512}}}\right) = -0.9598$

At dimension 1, $2i + 1 = 1$ thus $i = 0$, therefore $PE_{(pos,2i+1)} = PE_{(5,1)} = \cos\left(\frac{5}{10000^{\frac{0}{512}}}\right) = 0.2837$

At dimension 2, $2i = 2$ thus $i = 1$, therefore $PE_{(pos,2i)} = PE_{(5,2)} = \sin\left(\frac{5}{10000^{\frac{2}{512}}}\right) = -0.9939$

At dimension 3, $2i + 1 = 3$ thus $i = 1$, therefore $PE_{(pos,2i+1)} = PE_{(1,3)} = \cos\left(\frac{5}{10000^{\frac{2}{512}}}\right) = 0.1107$

At dimension 4, $2i = 4$ thus $i = 2$, therefore $PE_{(pos,2i)} = PE_{(1,4)} = \sin\left(\frac{5}{10000^{\frac{4}{512}}}\right) = -0.9982$

At dimension 5, $2i + 1 = 5$ thus $i = 2$, therefore $PE_{(pos,2i+1)} = PE_{(1,5)} = \cos\left(\frac{5}{10000^{\frac{4}{512}}}\right) = -0.0595$

At dimension 6, $2i = 6$ thus $i = 3$, therefore $PE_{(pos,2i)} = PE_{(1,6)} = \sin\left(\frac{5}{10000^{\frac{6}{512}}}\right) = -0.9750$

At dimension 7, $2i + 1 = 7$ thus $i = 3$, therefore $PE_{(pos,2i+1)} = PE_{(1,7)} = \cos\left(\frac{5}{10000^{\frac{6}{512}}}\right) = -0.2221$

At dimension 8, $2i = 8$ thus $i = 4$, therefore $PE_{(pos,2i)} = PE_{(1,8)} = \sin\left(\frac{5}{10000^{\frac{8}{512}}}\right) = -0.9277$

At dimension 9, $2i + 1 = 9$ thus $i = 4$, therefore $PE_{(pos,2i+1)} = PE_{(1,9)} = \cos\left(\frac{5}{10000^{\frac{8}{512}}}\right) = -0.3733$

At dimension 10, $2i = 10$ thus $i = 5$, therefore $PE_{(pos,2i)} = PE_{(1,10)} = \sin\left(\frac{5}{10000^{\frac{10}{512}}}\right) = -0.8600$

4. (d) Explain the concept of "mode collapse" in the context of GANs and discuss its implications on the quality of the generated data.

(4 marks)

Answer

- Mode collapse is a common training problem in GANs where the generator learns to produce limited varieties of outputs, even when the input (noise vector) varies.
- Instead of capturing the full diversity of the training data, the generator repeatedly outputs identical or nearly identical samples, effectively "collapsing" to a few modes.
- Implications
 - Generated samples lack variation, even when input noise is different.
 - The generator becomes overly specialized to a few patterns, reducing its ability to generalize to unseen data or edge cases.
 - Mode collapse often correlates with oscillations or divergence in GAN training, where the generator and discriminator fail to reach equilibrium.