

Internal Assessment Test 1 Answer scheme & Solutions – October 2024

Sub:	DEEP LEARNING				Sub Code:	21CS743	Branch:	AInDS		
Date:	16/10/2024	Duration:	90 minutes	Max Marks:	50	Sem	VII		OBE	
<u>Answer any FIVE Questions</u>								MARKS	CO	RBT
1	<p>Explain the different parts of an AI system relate to each other within different AI disciplines.</p> <p>A computer can reason about statements in these formal languages automatically using logical inference rules. This is known as the knowledge base. (2 Marks)</p> <p>AI systems need the ability to acquire their own knowledge, by extracting patterns from raw data. This capability is known as machine learning.(2 Marks)</p> <p>A representation learning algorithm can discover a good set of features for a simple task in minutes.(2 Marks)</p> <p>Deep learning allows the computer to build complex concepts out of simpler concepts.(2 Marks)</p> <p>Diagram (2 Marks)</p>						[10]	1	L3	
2	<p>Explain in Detail about Probabilistic Supervised Learning and non-probabilistic supervised learning.</p> <p>Probabilistic Supervised Learning (5 Marks)</p> <p>To find the best parameter vector θ for a parametric family of distributions $p(y x; \theta)$. We have already seen that linear regression corresponds to the family</p> $p(y x; \theta) = N (y; \theta = xT, I)$ <p>A distribution over a binary variable is slightly more complicated, because its mean must always be between 0 and 1. One way to solve this problem is to use the logistic sigmoid function to squash the output of the linear function into the interval (0, 1) and interpret that value as a probability:</p> $p(y = 1 x; \theta) = \sigma(\theta Tx)$ <p>This approach is known as logistic regression</p> <p>non-probabilistic supervised learning. (5 Marks)</p> <p>k-nearest neighbors - there is not even really a training stage or learning process. Instead, at test time, when we want to produce an output y for a new test input x, we find the k-nearest neighbors to x in the training data X. We then return the average of the corresponding y values in the training set</p> <p>decision tree - breaks the input space into regions and has separate parameters for each region is the decision tree</p>						[10]	2	L3	
3	<p>List The Historical Trends in Deep Learning.</p> <p>1.Deep learning has had a long and rich history, but has gone by many names reflecting different philosophical viewpoints, and has waxed and waned in popularity.</p> <p>a 2.Deep learning has become more useful as the amount of available training data has increased.</p> <p>3.Deep learning models have grown in size over time as computer infrastructure(both hardware and software) for deep learning has improved.</p> <p>4.Deep learning has solved increasingly complicated applications with increasing accuracy over time. (4 *1 = 4)</p>						[4]	1	L2	
	<p>Explain the Key Trends i) Increasing Dataset Sizes ii) Increasing Model Sizes over time in Deep Learning</p> <p>i) Deep learning algorithm will generally achieve acceptable performance with around 5,000 labeled examples per category (2 marks)</p> <p>b Diagram (1 Mark)</p> <p>ii) Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years.(2 marks)</p> <p>Diagram (1 Mark)</p>						[6]			

What is Unsupervised Learning Algorithms ? Explain PCA.

- iii) Unsupervised learning algorithms perform other roles, like clustering, which consists of dividing the dataset into clusters of similar examples. (2 marks)
- iv) PCA learns a representation that has lower dimensionality than the original input. It also learns a representation whose elements have no linear correlation with each other. (2 Marks)

PCA finds a representation (through linear transformation) $\mathbf{z} = \mathbf{a}^T \mathbf{x}$ such that $\text{Var}[\mathbf{z}]$ is diagonal.

In section 2.12, we saw that the principal components of a design matrix \mathbf{X} are given by the eigenvectors of $\mathbf{X}^T \mathbf{X}$. From this view,

$$\mathbf{X}^T \mathbf{X} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T.$$

In this section, we exploit an alternative derivation of the principal components. The principal components may also be obtained via the singular value decomposition. Specifically, they are the right singular vectors of \mathbf{X} . To see this, let \mathbf{W} be the right singular vectors in the decomposition $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T$. We then substitute \mathbf{W} into the original eigenvector equation with \mathbf{W} as the eigenvector basis:

$$\mathbf{X}^T \mathbf{X} = (\mathbf{U} \mathbf{\Sigma} \mathbf{W}^T)^T \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T = \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^T.$$

The SVD is helpful to show that PCA results in a diagonal $\text{Var}[\mathbf{z}]$. Using the SVD of \mathbf{X} , we can express the variance of \mathbf{X} as:

$$\begin{aligned} \text{Var}[\mathbf{x}] &= \frac{1}{m-1} \mathbf{X}^T \mathbf{X} \\ &= \frac{1}{m-1} (\mathbf{U} \mathbf{\Sigma} \mathbf{W}^T)^T \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T \\ &= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T \\ &= \frac{1}{m-1} \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^T, \end{aligned}$$

where we use the fact that $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ because the \mathbf{U} matrix of the singular value decomposition is defined to be orthogonal. This shows that if we take $\mathbf{z} = \mathbf{W}^T \mathbf{x}$, we can ensure that the covariance of \mathbf{z} is diagonal as required:

$$\begin{aligned} \text{Var}[\mathbf{z}] &= \frac{1}{m-1} \mathbf{Z}^T \mathbf{Z} \\ &= \frac{1}{m-1} \mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} \\ &= \frac{1}{m-1} \mathbf{W}^T \mathbf{W} \mathbf{\Sigma}^2 \mathbf{W}^T \mathbf{W} \\ &= \frac{1}{m-1} \mathbf{\Sigma}^2, \end{aligned}$$

(6 Marks)

What is Task ? List most common machine learning tasks.

- a) Learning is our means of attaining the ability to perform the task (1 Mark)
- Examples : Classification ,Regression,Transcription,Machine Translation etc (4 Marks)

[10]

2

L2

4

5

a

[5]

1

L2

	<p>Explain Deep Feedforward Networks in detail. Deep feedforward networks, also often called feedforward neural networks, or multilayer perceptrons (MLPs), are the quintessential deep learning models. (2 Marks)</p> <p>b Linear models also have the obvious defect that the model capacity is limited to linear functions, so the model cannot understand the interaction between any two input variables. To extend linear models to represent nonlinear functions of x, we can apply the linear model not to x itself but to a transformed input $\phi(x)$, then how to choose the mapping ϕ. (3 Marks)</p>	[5]	2	L2
	<p>What is Gradient-Based Learning ? Explain about Cost Function. For feedforward neural networks, it is important to initialize all weights to small random values. The biases may be initialized to zero or to small positive values. The iterative gradient-based optimization algorithms used to train feedforward networks and almost all other deep models. train models such as linear regression and support vector machines with gradient descent too, and in fact this is common when the training set is extremely large. (3 Marks)</p> <p>a In most cases, our parametric model defines a distribution $p(y x; \theta)$ and we simply use the principle of maximum likelihood. This means we use the cross-entropy between the training data and the model's predictions as the cost function, where rather than predicting a complete probability distribution over y, we merely predict some statistic of y conditioned on x. Specialized loss functions allow us to train a predictor of these estimates. The total cost function used to train a neural network will often combine one of the primary cost functions described here with a regularization term. (2 Marks)</p>	[5]	2	L3
6	<p>Explain Softmax Units for Multinoulli Output Distributions. linear layer predicts unnormalized log probabilities: $z = WT h + b$ where $z_i = \log \tilde{P}(y = i x)$. The softmax function can then exponentiate and normalize z to obtain the desired \hat{y}. Formally, the softmax function is given by $\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$ to maximize $\log P(y = i; z) = \log \text{softmax}(z)_i$. Defining the softmax in terms of \exp is natural because the log in the log-likelihood can undo the exp of the softmax: $\log \text{softmax}(z)_i = z_i - \log(\sum_j \exp(z_j))$ To see that the softmax function responds to the difference between its inputs, observe that the softmax output is invariant to adding the same scalar to all of its inputs: $\text{softmax}(z) = \text{softmax}(z + c)$ Using this property, we can derive a numerically stable variant of the softmax: $\text{softmax}(z) = \text{softmax}(z - \text{imax}z_i) \quad (5 \text{ Marks})$</p>	[5]	2	L3