

Feedback — Lecture 3 Quiz

[Help Center](#)

You submitted this quiz on **Tue 28 Jun 2016 2:24 AM CEST**. You got a score of **8.00** out of **8.00**. However, you will not get credit for it, since it was submitted past the deadline.

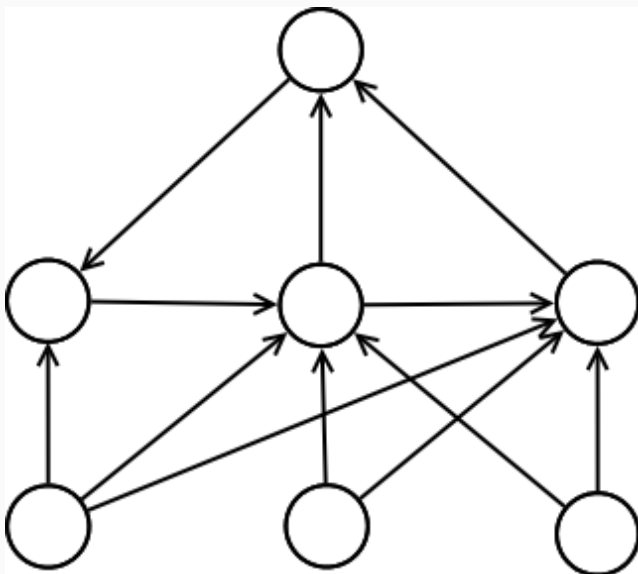
Question 1

Which of the following neural networks are examples of a feed-forward neural network?

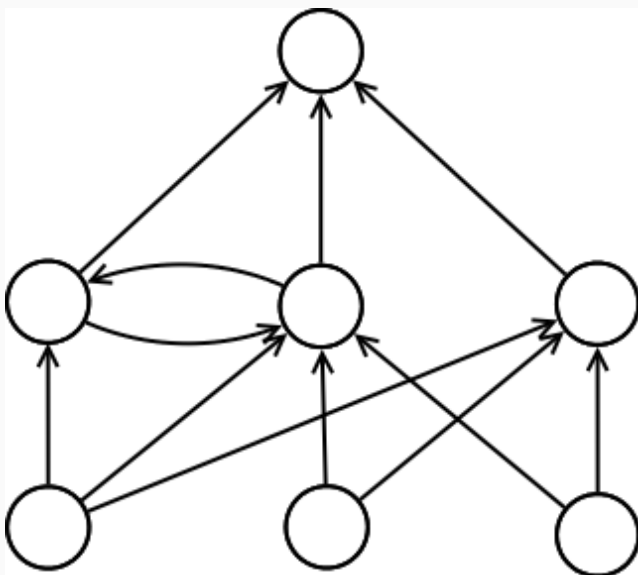
Your Answer

Score

Explanation



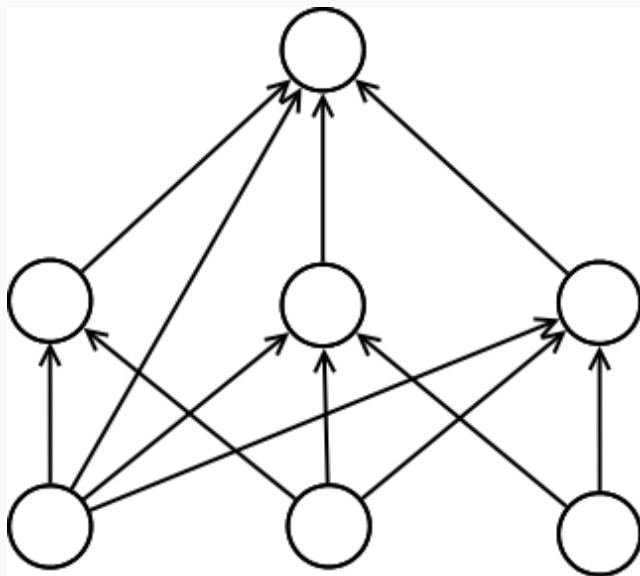
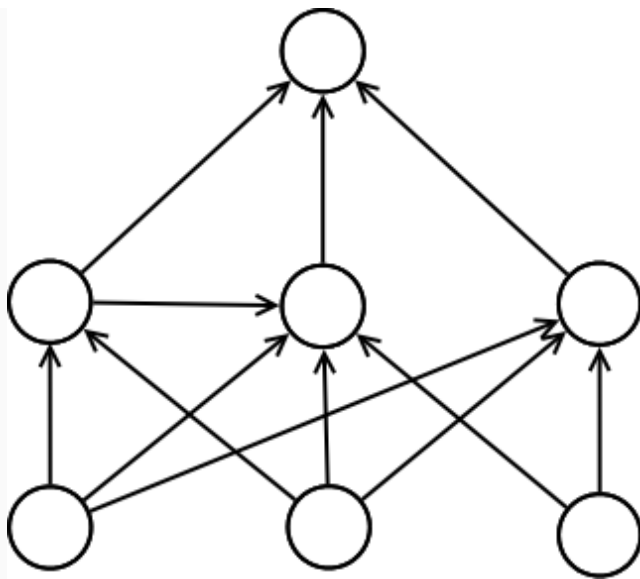
0.50



0.50



0.50



✓ 0.50

Total

2.00 / 2.00

Question Explanation

A feed-forward network does not have cycles.

Question 2

Consider a neural network with only one training case with input $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$ and correct output t . There is only one output neuron, which is linear, i.e. $y = \mathbf{w}^\top \mathbf{x}$ (notice that there are no biases). The loss function is squared error. The network has no hidden units, so the inputs are directly connected to the output neuron with weights $\mathbf{w} = (w_1, w_2, \dots, w_n)^\top$. We're in the process of training the neural network with the backpropagation algorithm. What will the algorithm add to w_i for the next iteration if we use a step size (also known as a learning rate) of ϵ ?

Your Answer

Score

Explanation

<input type="radio"/>	$\epsilon(\mathbf{w}^\top \mathbf{x} - t)x_i$		
<input checked="" type="radio"/>	$\epsilon(t - \mathbf{w}^\top \mathbf{x})x_i$	✓	1.00
<input type="radio"/>	x_i		
<input type="radio"/>	x_i if $\mathbf{w}^\top \mathbf{x} > t$ $-x_i$ if $\mathbf{w}^\top \mathbf{x} \leq t$		
Total		1.00 / 1.00	

Question Explanation

There are multiple components to this, all multiplied together: the learning rate, the derivative of the loss function w.r.t. the state of the output unit, and the derivative of the input to the output unit w.r.t. w_i .

Question 3

Suppose we have a set of examples and Brian comes in and duplicates every example, then randomly reorders the examples. We now have twice as many examples, but no more information about the problem than we had before. If we do not remove the duplicate entries, which one of the following methods will *not* be affected by this change, in terms of the computer time (time in seconds, for example) it takes to come close to convergence?

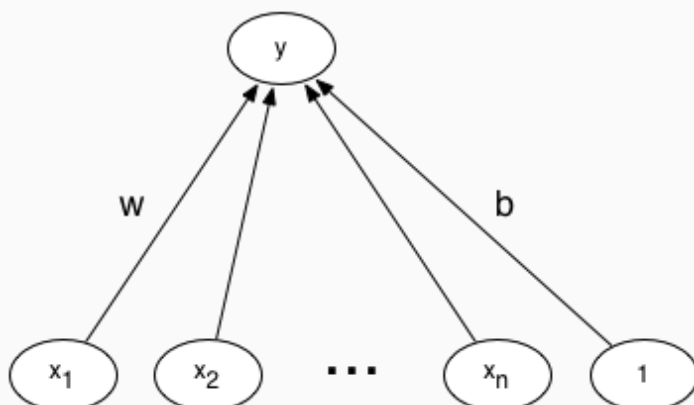
Your Answer	Score	Explanation
<input checked="" type="radio"/> Online learning, where for every iteration we randomly pick a training case.	✓ 1.00	
<input type="radio"/> Full-batch learning.		
<input type="radio"/> Mini-batch learning, where for every iteration we randomly pick 100 training cases.		
Total	1.00 / 1.00	

Question Explanation

Full-batch learning needs to look at every example before taking a step, therefore each step will be twice as expensive. Online learning only looks at one example at a time so each step has the same computational cost as before. On expectation, online learning would make the same progress after looking at half of the dataset as it would have if Brian has not intervened. Although this example is a bit contrived, it serves to illustrate how online learning can be advantageous when there is a lot of redundancy in the data.

Question 4

Consider a linear output unit versus a logistic output unit for a feed-forward network with *no hidden layer* shown below. The network has a set of inputs x and an output neuron y connected to the input by weights w and bias b .



We're using the squared error cost function even though the task that we care about, in the end, is binary classification. At training time, the target output values are 1 (for one class) and 0 (for the other class). At test time we will use the classifier to make decisions in the standard way: the class of an input x according to our model **after training** is as follows:

$$\text{class of } x = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Note that we will be training the network using y , but that the decision rule shown above will be the same at *test* time, regardless of the type of output neuron we use for training. Which of the following statements is true?

Your Answer	Score	Explanation
-------------	-------	-------------

☐ For a logistic unit, the derivatives of the error function with respect to the weights can have unbounded magnitude, while for a linear unit they will have bounded magnitude.

☐ At the solution that

minimizes the error, the learned weights are always the same for both types of units; they only differ in how they get to this solution.

☒ Unlike a logistic unit, using a linear unit will penalize us for getting the answer right too confidently.

✓ 1.00

If the target is 1 and the prediction is 100, the logistic unit will squash this down to a number very close to 1 and so we will not incur a very high cost. With a linear unit, the difference between the prediction and target will be very large and we will incur a high cost as a result, despite the fact that we get the classification decision correct.

☐ The error function (the error as a function of the weights) for both types of units will form a quadratic bowl.

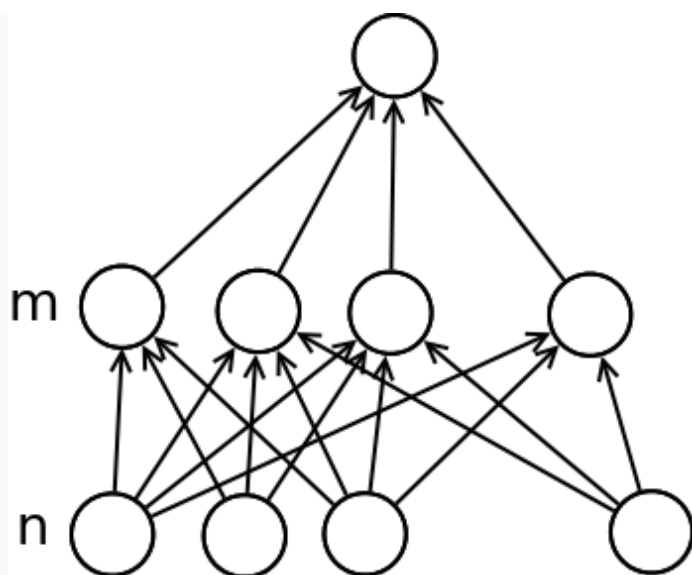
Total 1.00 / 1.00

Question Explanation

The squared error cost function is $\frac{1}{2} (t - y)^2$. For a linear unit y is a real number whereas for a logistic unit, y will be between 0 and 1. Try to picture the case where $w^T x + b$ is very large. What will the derivative be? Will the errors be the same?

Question 5

Consider a neural network with one layer of **logistic** hidden units (intended to be fully connected to the input units) and a linear output unit. Suppose there are n input units and m hidden units. Which of the following statements are true? Check all that apply.



Your Answer

Score

Explanation

☒ If $m > n$, this network can learn more functions than if m is less than n (with n being the same).

✓ 0.50

☒ As long as $m \geq 1$, this network can learn to compute any function that can be learned by a network without any hidden layers (with the same inputs).

✓ 0.50

If the weights into the hidden layer are very small, and the weights out of it are large (to compensate), then the hidden units behave like linear units, which makes lots of things possible.

☒ A network with $m > n$ has more learnable parameters than a network with $m \leq n$ (for a fixed value of n).

✓ 0.50

The bulk of the learnable parameters is in the connections from the input units to the hidden units. There are $m \cdot n$ learnable parameters there.

☐ Any function that can be learned by such a network can also be learned by a network without any hidden layers (with the same inputs).

✓ 0.50

Total

2.00 /
2.00

Question Explanation

This is quite a flexible model. It can learn many functions that cannot be learned without the use of a hidden layer. The nonlinearity in the hidden layer is essential.

Question 6

Brian wants to make his feed-forward network (with no hidden units) using a **logistic** output neuron more powerful. He decides to combine the predictions of two networks by averaging them. The first network has weights w_1 and the second network has weights w_2 . The predictions of this network for an example x are therefore:

$$y = \frac{1}{2} \frac{1}{1+e^{-z_1}} + \frac{1}{2} \frac{1}{1+e^{-z_2}} \text{ with } z_1 = w_1^T x \text{ and } z_2 = w_2^T x.$$

Can we get the exact same predictions as this combination of networks by using a single feed-forward network (again with no hidden units) using a **logistic** output neuron and weights $w_3 = \frac{1}{2} (w_1 + w_2)$?

Your Answer	Score	Explanation
<input type="radio"/> Yes		
<input checked="" type="radio"/> No	✓ 1.00	
Total	1.00 / 1.00	