

1. Which of the following are true? (Check all that apply.)

1 / 1 point

☒ $a_4^{[2]}$ is the activation output by the 4th neuron of the 2nd layer

✓ Correct

☐ $a_4^{[2]}$ is the activation output of the 2nd layer for the 4th training example

☒ $a^{[2]}$ denotes the activation vector of the 2nd layer.

✓ Correct

☒ X is a matrix in which each column is one training example.

✓ Correct

☐ X is a matrix in which each row is one training example.

☐ $a^{[2](12)}$ denotes activation vector of the 12th layer on the 2nd training example.

☒ $a^{[2](12)}$ denotes the activation vector of the 2nd layer for the 12th training example.

✓ Correct

↗ Expand

✓ Correct

Great, you got all the right answers.

2. The tanh activation is not always better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data, making learning complex for the next layer. True/False?

1 / 1 point

- ☒ False
- ☐ True

 Expand

 Correct

Yes. As seen in lecture the output of the tanh is between -1 and 1, it thus centers the data which makes the learning simpler for the next layer.

3. Which of these is a correct vectorized implementation of forward propagation for layer l , where $1 \leq l \leq L$?

1 / 1 point

- ☐ $Z^{[l]} = W^{[l]} A^{[l]} + b^{[l]}$
 $A^{[l+1]} = g^{[l+1]}(Z^{[l]})$
- ☐ $Z^{[l]} = W^{[l-1]} A^{[l]} + b^{[l-1]}$
 $A^{[l]} = g^{[l]}(Z^{[l]})$
- ☒ $Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$
 $A^{[l]} = g^{[l]}(Z^{[l]})$
- ☐ $Z^{[l]} = W^{[l]} A^{[l]} + b^{[l]}$
 $A^{[l+1]} = g^{[l]}(Z^{[l]})$

 Expand

 Correct

4. The use of the ReLU activation function is becoming more rare because the ReLU function has no derivative for $c = 0$. True/False?

1 / 1 point

4. The use of the ReLU activation function is becoming more rare because the ReLU function has no derivative for $c = 0$. True/False?

1 / 1 point

- ☐ True
- ☒ False

 Expand

 **Correct**

Yes. Although the ReLU function has no derivative at $c = 0$ this rarely causes any problems in practice. Moreover it has become the default activation function in many cases, as explained in the lectures.

5. Consider the following code:

1 / 1 point

```
#+begin_src python
```

```
x = np.random.rand(3, 2)
```

```
y = np.sum(x, axis=0, keepdims=True)
```

```
#+end_src
```

What will be `y.shape`?

- ☐ (2,)
- ☐ (3,)
- ☐ (3, 1)
- ☒ (1, 2)

 Expand

 **Correct**

Yes. By choosing the `axis=0` the sum is computed over each column of the array, thus the resulting array is a row vector with 2 entries. Since the option `keepdims=True` is used the first dimension is kept, thus (1, 2).

6. Suppose you have built a neural network with one hidden layer and tanh as activation function for the hidden layer. You decide to initialize the weights to small random numbers and the biases to zero. The first hidden layer's neurons will perform different computations from each other even in the first iteration. True/False?

1 / 1 point

- ☐ False No. Since the weights are most likely different, each neuron will do a different computation.
- ☒ True Yes. Since the weights are most likely different, each neuron will do a different computation.

 Expand

 Correct

7. Logistic regression's weights should be initialized randomly rather than to all zeros, because if you initialize to all zeros, then logistic regression will fail to learn a useful decision boundary because it will fail to "break symmetry", True/False?

1 / 1 point

- ☐ True
- ☒ False

 Expand

 Correct

Yes, Logistic Regression doesn't have a hidden layer. If you initialize the weights to zeros, the first example x fed into the logistic regression will output zero but the derivatives of the Logistic Regression depend on the input x (because there's no hidden layer) which is not zero. So at the second iteration, the weights' values follow x 's distribution and are different from each other if x is not a constant vector.

8. You have built a network using the tanh activation for all the hidden units. You initialize the weights to relatively large values, using `np.random.randn(...)*1000`. What will happen?

1 / 1 point

- ☐ This will cause the inputs of the tanh to also be very large, thus causing gradients to also become large. You therefore have to set α to a very small value to prevent divergence; this will slow down learning.
- ☐ So long as you initialize the weights randomly gradient descent is not affected by whether the weights are large or small.
- ☒ This will cause the inputs of the tanh to also be very large, thus causing gradients to be close to zero. The optimization algorithm will thus become slow.
- ☐ This will cause the inputs of the tanh to also be very large, causing the units to be “highly activated” and thus speed up learning compared to if the weights had to start from small values.

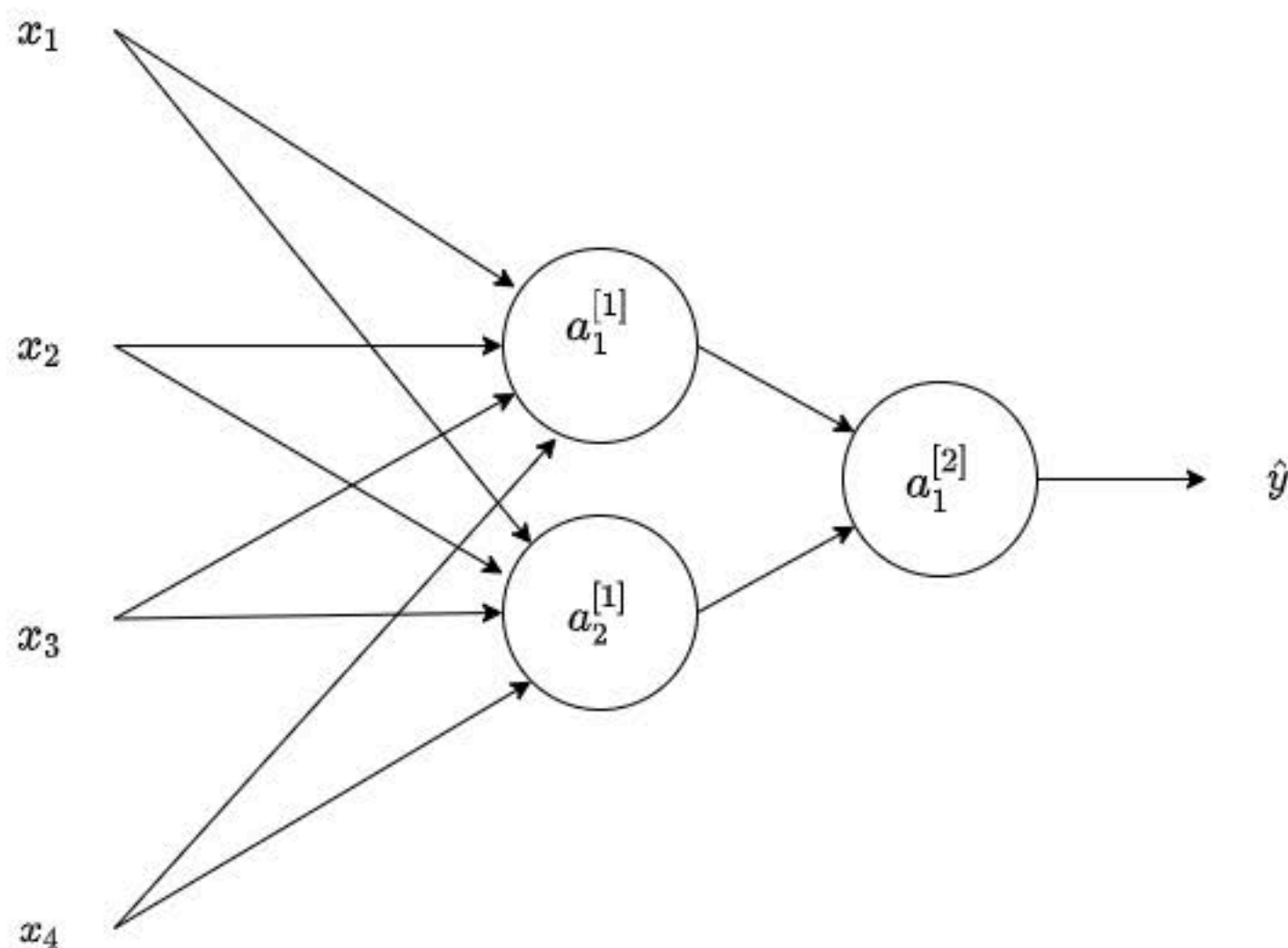
[Expand](#)

✓ **Correct**

Yes. tanh becomes flat for large values; this leads its gradient to be close to zero. This slows down the optimization algorithm.

9. Consider the following 1 hidden layer neural network:

1 / 1 point



Which of the following statements are True? (Check all that apply).

☐ $W^{[1]}$ will have shape (4, 2).

Which of the following statements are True? (Check all that apply).

☐ $W^{[1]}$ will have shape (4, 2).

☒ $W^{[2]}$ will have shape (1, 2)

✓ **Correct**

Yes. The number of rows in $W^{[k]}$ is the number of neurons in the k-th layer and the number of columns is the number of inputs of the layer.

☐ $W^{[2]}$ will have shape (2, 1)

☒ $W^{[1]}$ will have shape (2, 4).

✓ **Correct**

Yes. The number of rows in $W^{[k]}$ is the number of neurons in the k-th layer and the number of columns is the number of inputs of the layer.

☒ $b^{[1]}$ will have shape (2, 1).

✓ **Correct**

Yes. $b^{[k]}$ is a column vector and has the same number of rows as neurons in the k-th layer.

☐ $b^{[1]}$ will have shape (4, 2)

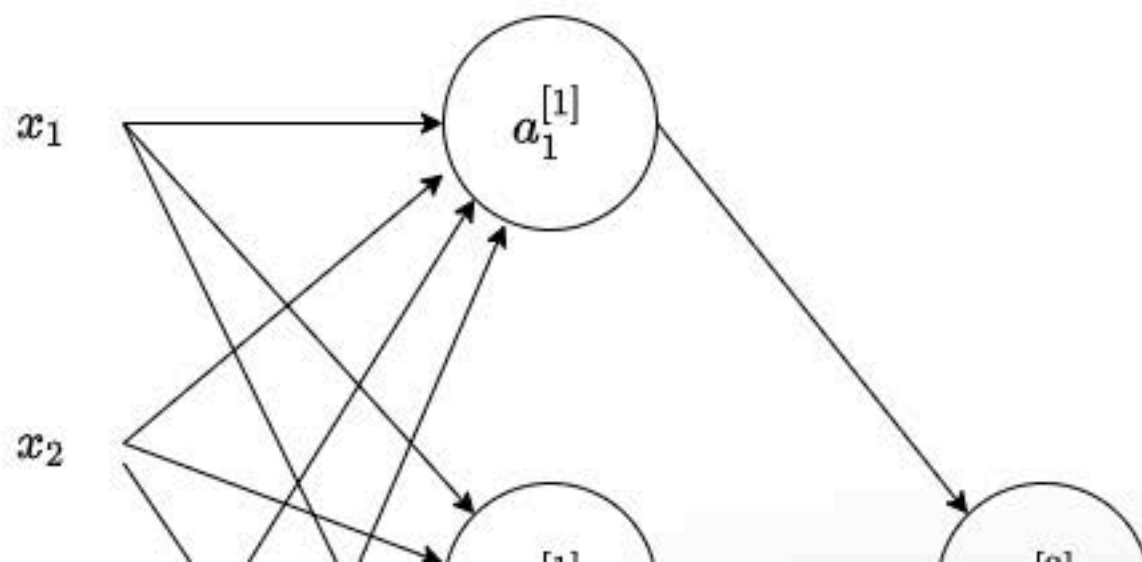
↗ **Expand**

✓ **Correct**

Great, you got all the right answers.

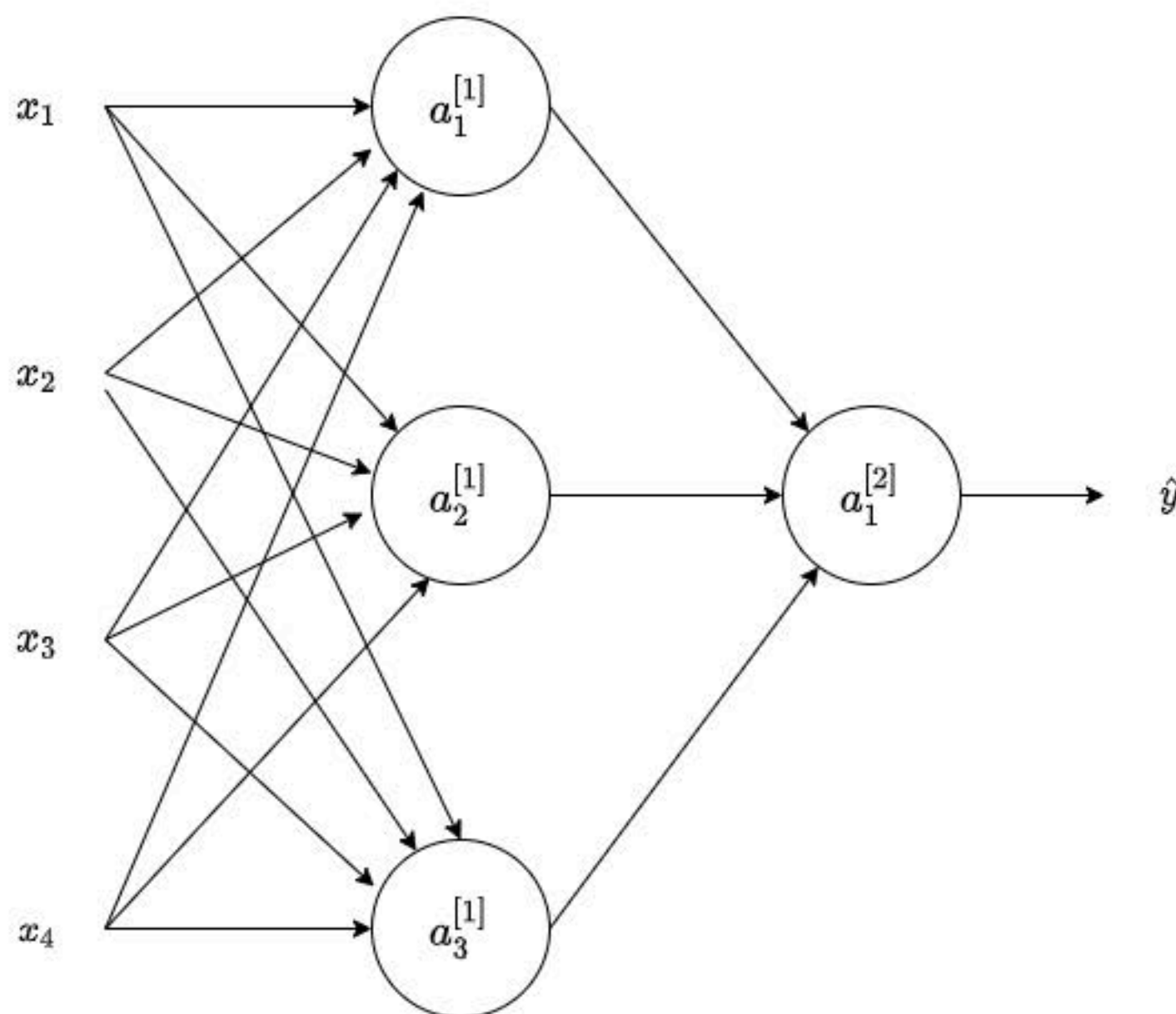
10. Consider the following 1 hidden layer neural network:

1 / 1 point



10. Consider the following 1 hidden layer neural network:

1 / 1 point



What are the dimensions of $Z^{[1]}$ and $A^{[1]}$?

- ☐ $Z^{[1]}$ and $A^{[1]}$ are (4, m)
- ☐ $Z^{[1]}$ and $A^{[1]}$ are (4, 1)
- ☐ $Z^{[1]}$ and $A^{[1]}$ are (3, 1)
- ☒ $Z^{[1]}$ and $A^{[1]}$ are (3, m)

[Expand](#)

✓ **Correct**

Yes. The $Z^{[1]}$ and $A^{[1]}$ are calculated over a batch of training examples. The number of columns in $Z^{[1]}$ and $A^{[1]}$ is equal to the number of examples in the batch, m. And the number of rows in $Z^{[1]}$ and $A^{[1]}$ is equal to the number of neurons in the first layer.