

✓

1 / 1 points

1. The squared error cost function with  $n$  linear units is equivalent to:

**Clarification:**

Let's say that a network with  $n$  linear output units has some weights  $w$ .  $w$  is a matrix with  $n$  columns, and  $w_i$  indexes a particular column in this matrix and represents the weights from the inputs to the  $i^{\text{th}}$  output unit.

Suppose the target for a particular example is  $j$  (so that it belongs to class  $j$  in other words).

The squared error cost function for  $n$  linear units is given by:

$$\frac{1}{2} \sum_{i=1}^n (t_i - w_i^T x)^2$$

where  $t$  is a vector of zeros except for 1 in index  $j$ .

The cross-entropy cost function for an  $n$ -way softmax unit is given by:

$$-\log \left( \frac{\exp(w_j^T x)}{\sum_{i=1}^n \exp(w_i^T x)} \right) = -w_j^T x + \log \left( \sum_{i=1}^n \exp(w_i^T x) \right)$$

Finally,  $n$  logistic units would compute an output of  $\sigma(w_i^T x) = \frac{1}{1 + \exp(-w_i^T x)}$  independently for each class  $i$ . Combined with the squared error the cost would be:

$$\frac{1}{2} \sum_{i=1}^n (t_i - \sigma(w_i^T x))^2$$

Where again,  $t$  is a vector of zeros with a 1 at index  $j$  (assuming the true class of the example is  $j$ ).

Using this same definition for  $t$ , the cross-entropy error for  $n$  logistic units would be the sum of the individual cross-entropy errors:

$$-\sum_{i=1}^n t_i \log(\sigma(w_i^T x)) + (1 - t_i) \log(1 - \sigma(w_i^T x))$$

For any set of weights  $w$ , the network with  $n$  linear output units will have some cost due to the squared error (cost function). The question is now asking whether we can define a new network with a set of weights  $w^*$  using some (possibly different) cost function such that:

a)  $w^* = f(w)$  for some function  $f$

b) For every input, the cost we get using  $w$  in the linear network with squared error is the same cost that we would get using  $w^*$  in the new network with the possibly different cost function.

☐ The cross-entropy cost function with an  $n$ -way softmax unit.

Un-selected is correct

☐ The cross-entropy cost function with  $n$  logistic units.

Un-selected is correct

☐ The squared error cost function with  $n$  logistic units.

Un-selected is correct

☒ None of the above.

Correct

✓

1 / 1 points

2. A logistic unit with the cross-entropy cost function is equivalent to:

**Clarification:**

In a network with a logistic output, we will have a single vector of weights  $w$ . For a particular example with target  $t$  (which is 0 or 1), the cross-entropy error is given by:

$$-t \log(\sigma(w^T x)) - (1 - t) \log(1 - \sigma(w^T x)) \text{ where } \sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}.$$

The squared error if we use a single linear unit would be:

$$\frac{1}{2} (t - w^T x)^2$$

Now notice that another way we might define  $t$  is by using a vector with 2 elements,  $[1, 0]$  to indicate the first class, and  $[0, 1]$  to indicate the second class. Using this definition, we can develop a new type of classification network using a softmax unit over these two classes instead. In this case, we would use a weight *matrix*  $w$  with two columns, where  $w_i$  is the column of the  $i^{\text{th}}$  class and connects the inputs to the  $i^{\text{th}}$  output unit.

Suppose an example belonged to class  $j$  (where  $j$  is 1 or 2 to indicate  $[1, 0]$  or  $[0, 1]$ ). Then the cross-entropy cost for this network would be:

$$-\log \left( \frac{\exp(w_j^T x)}{\exp(w_1^T x) + \exp(w_2^T x)} \right) = -w_j^T x + \log(\exp(w_1^T x) + \exp(w_2^T x))$$

For any set of weights  $w$ , the network with a logistic output unit will have some error due to the cross-entropy cost function. The question is now asking whether we can define a new network with a set of weights  $w^*$  using some (possibly different) cost function such that:

a)  $w^* = f(w)$  for some function  $f$

b) For every input, the cost we get using  $w$  in the network with a logistic output and cross-entropy error is the same cost that we would get using  $w^*$  in the new network with the possibly different cost function.

☒ A 2-way softmax unit (a softmax unit with 2 elements) with the cross entropy cost function.

Correct

With weights  $w_1$  and  $w_2$  we can write the softmax probability for an input  $x$  as:

$$P(x \text{ belongs to class 1}) = \frac{\exp(w_1^T x)}{\exp(w_1^T x) + \exp(w_2^T x)}.$$

Dividing the top and bottom by  $\exp(w_1^T x)$ :

$$P(x \text{ belongs to class 1}) = \frac{1}{1 + \exp(w_2^T x - w_1^T x)} = \frac{1}{1 + \exp(-(w_1 - w_2)^T x)}.$$

So this is equivalent to a logistic unit with weights  $w_1 - w_2$ . Since the cross-entropy is the negative log-probability, it will be the same for both representations.

☐ A linear unit with the squared error cost function.

Un-selected is correct

☐ A 2-way softmax unit (a softmax with 2 elements) with the squared error cost function.

Un-selected is correct

☐ None of the above.

Un-selected is correct

✓

1 / 1 points

3. The output of a neuro-probabilistic language model is a large softmax unit and this creates problems if the vocabulary size is large. Andy claims that the following method solves this problem:

At every iteration of training, train the network to predict the current learned feature vector of the target word (instead of using a softmax). Since the embedding dimensionality is typically much smaller than the vocabulary size, we don't have the problem of having many output weights any more. Which of the following are correct? Check all that apply.

☐ Andy is correct: this is equivalent to the serialized version of the model discussed in the lecture.

Un-selected is correct

☐ In theory there's nothing wrong with Andy's idea. However, the number of learnable parameters will be so far reduced that the network no longer has sufficient learning capacity to do the task well.

Un-selected is correct

☒ If we add in extra derivatives that change the feature vector for the target word to be more like what is predicted, it may find a trivial solution in which all words have the same feature vector.

Correct

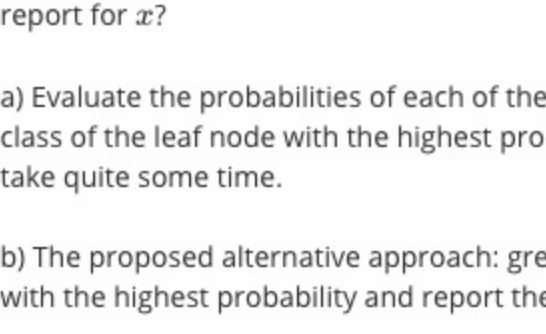
☒ The serialized version of the model discussed in the slides is using the current word embedding for the output word, but it's optimizing something different than what Andy is suggesting.

Correct

✓

1 / 1 points

4. We are given the following tree that we will use to classify a particular example  $x$ :



In this tree, each  $p$  value indicates the probability that  $x$  will be classified as belonging to a class in the right subtree of the node at which that  $p$  was computed. For example, the probability that  $x$  belongs to Class 2 is  $(1 - p_1) \times p_2$ . Recall that at training time this is a very efficient representation because we only have to consider a single branch of the tree. However, at test-time we need to look over all branches in order to determine the probabilities of each outcome.

Suppose we are not interested in obtaining the exact probability of every outcome, but instead we just want to find the class with the maximum probability. A simple heuristic is to search the tree greedily by starting at the root and choosing the branch with the highest probability at each node on our way from the root to the leaves. That is, at the root of this tree we would choose to go right if  $p_1 \geq 0.5$  and left otherwise.

If  $p_1 = 0.45$ ,  $p_2 = 0.6$ , and  $p_3 = 0.95$ , then which class will the following methods report for  $x$ ?

a) Evaluate the probabilities of each of the four classes (the leaf nodes) and report the class of the leaf node with the highest probability. This is the standard approach but may take quite some time.

b) The proposed alternative approach: greedily traverse the tree by choosing the branch with the highest probability and report the class of the leaf node that this finds.

☐ Method a) will report class 4

Method b) will report class 1

☐ Method a) will report class 2

Method b) will report class 3

☐ Method a) will report class 1

Method b) will report class 4

☐ Method a) will report class 4

Method b) will report class 4

☐ Method a) will report class 3

Method b) will report class 2

☒ Method a) will report class 4

Method b) will report class 2

Correct

✓

1 / 1 points

5. Brian is trying to use a neural network to predict the next word given several previous words. He has the following idea to reduce the amount of computation needed to make this prediction.

Rather than having the output layer of the network be a 100,000-way softmax, Brian says that we can just encode each word as an integer: 1 will correspond to the first word, 2 to the second word and so on up to 100,000. For the output layer, we can then simply use a single linear neuron with a squared error cost function. Is this a good idea?

☐ Yes. With this method, there are fewer parameters, while the network can still learn equally well.

☒ No. Brian is implicitly imposing the belief that words with similar integers are more related to each other than words with very different integers, and this is usually not true.

Correct

☐ Partly. Brian's method should only be used at the input layer. The output layer must always report probabilities, and squared error loss is not appropriate for that.

✓

1 / 1 points

6. In the Collobert and Weston model, the problem of learning a feature vector from a sequence of words is turned into a problem of:

☐ Learning to predict the next word in an arbitrary length sequence.

☐ Learning to reconstruct the input vector.

☒ Learning a binary classifier.

Correct

The model is made to answer this binary question: "Is the word in the middle correct, or is it just some random word?"

☐ Learning to predict the middle word in the sequence given the words that came before and the words that came after.

✓

1 / 1 points

7. Suppose that we have a vocabulary of 3 words, "a", "b", and "c", and we want to predict the next word in a sentence given the previous two words. Also suppose that we don't want to use feature vectors for words: we simply use the local encoding, i.e. a 3-component vector with one entry being 1 and all other two entries being 0.

In the language models that we have seen so far, each of the context words has its own dedicated section of the network, so we would encode this problem with two 3-dimensional inputs. That makes for a total of 6 dimensions; clearly, the more context words we want to include, the more input units our network must have. Here's a method that uses fewer input units:

We could instead encode the **counts** of each word in the context. So a context of "a a" would be encoded as input vector  $[2 \ 0 \ 0]$  instead of  $[1 \ 0 \ 0 \ 1 \ 0 \ 0]$ , and "b c" would be encoded as input vector  $[0 \ 1 \ 1]$  instead of  $[0 \ 1 \ 0 \ 0 \ 1]$ . Now we only need an input vector of the size of our vocabulary (3 in our case), as opposed to the size of our vocabulary times the length of the context (which makes for a total of 6 in our case). Are there any significant problems with this idea?

☐ Yes: the neural networks shown in the course so far cannot deal with integer inputs (as opposed to binary inputs).

☒ Yes: the network loses the knowledge of the location at which a context word occurs, and that is valuable knowledge.

Correct

☐ Yes: although we could encode the context in this way, we would then need a smaller bottleneck layer than we did before, thereby lowering the learning capacity of the model.

☐ Yes: even though the input has a smaller dimensionality, each entry of the input now requires more bits to encode, because it's no longer just 1 or 0. Therefore, there would be no significant advantage.