

1a. $-\sum_{m=1}^30 \log p(x_m)$

This is known as the log-loss of the model on training data. Note the negative sign. Minimizing log-loss is the same as maximizing log-likelihood.

You could also add a regularization term.

Remember that for an autoregressive language model, $p(x)$ is defined as a product of conditional probabilities, so $\log p(x) = \log p(x_1) + \log p(x_2|x_1) + \log p(x_3|x_1, x_2) + \dots$ where x_1, x_2, \dots here represent the tokens of sentence x .

1b. reduces BIAS, but at the cost of greater VARIANCE

1c. Option 1: Use a regularization term that penalizes the difference between the new parameters and the old ones. Thus, training will be reluctant to change parameters and will only do this if it substantially helps to predict the training examples of Trumpy English.

Option 2: Early stopping before the parameters have converged.

Either way: you have a hyperparameter to tune. In option 1, it is the regularization constant. In option 2, it is the number of epochs. Either way, you can select this hyperparameter using a dev corpus of Trumpy English.

Notice that choosing the hyperparameter for option 2 is very efficient -- you just keep training until the cross-entropy on the dev corpus starts getting worse. For example, you evaluate on the dev course after 15 epochs of training, 16 epochs of training, etc. It's cheap to find out whether 16 epochs gives a good model because you don't have to train a new model from scratch: just start at the 15-epoch model and train for one more epoch.

2a. The matrix V , the input embeddings and output embeddings for all word types in the vocabulary, and the initial hidden state vector h_0 .

2b. The number of rows is d . The number of columns is $(1 + d + \text{dimensionality of the word embeddings})$.

2c. We'll say that the hidden vectors h_j represent the upper layer, so (1) can be unchanged. Change equation (2) to replace w_j with g_j , where g_j is a hidden vector at the lower layer. Add equation (3):

$$g_j = \text{sigma}(U [1; g_{j-1}; w_j])$$

2d. This is a tricky question! The goal is to see how an RNN can track properties of the input over time, just like the FSA in the previous question. We saw in class how nodes in neural nets can implement AND/OR operations, and we'll do that here.

Note that the typesetting of this answer omits the vector arrow, so it does not distinguish properly between the word w_j and its embedding w_j (which should have a vector arrow).

In this answer, we will assume that a d -dimensional vector has indices $1\dots d$, as indicated in the last line of the question. This is common in mathematical notation, in contrast to Python's

indices $0 \dots (d-1)$. (We accepted either style in your answer, though.)

We can see that

$$h_j[3] = \sigma(v \cdot [1; h_{j-1}]; w_j)$$
if v denotes row 3 of matrix V .

Therefore, we need

$$v \cdot [1; h_{j-1}; w_j]$$
to be strongly negative if $w_j == \text{Trump}$ or if $h_{j-1}[3]$ is close to 0 (meaning that $w_i == \text{Trump}$ for some $i < j$),but it should be strongly positive otherwise.

We shouldn't pay attention to the other elements of h_{j-1} ,so we can set $v[2,3,5,6,\dots,d+1]$ to 0.

Then we have

$$v \cdot [1; h_{j-1}; w_j] = v[1] + v[4] * h_{j-1}[3] + v[d+2, \dots] \cdot w_j$$

We want $v[4]$ to be strongly positive so that if $h_{j-1}[3]$ is close to 1 (meaning that we haven't seen Trump yet),then the dot product will be strongly positive and thus $h_j[3]$ will also be close to 1.

However, this should be overridden if $w_j == t$, inwhich case we want $v[d+2, \dots] \cdot w_j$ to be negative enoughto drive the dot product negative. We can do this bysetting $v[d+2, \dots] = -c t$ for some large positive c .That ensures that $v[d+2, \dots] \cdot w_j$ is much more negativewhen $w_j = \text{Trump}$ than for any other w_j (becausewhen $w_j = \text{Trump}$, $t \cdot w_j$ is positive and larger thanfor any other w_j , according to the problem).

So, choose a large c , and then solve for $v[1]$ and $v[4]$ to ensure that the dot product is (for example) < -5 when it is supposed to be strongly negative andalso > 5 when it is supposed to be strongly positive.(This ensures that $h_j[3]$ will be < 0.01 and > 0.99 ,respectively.) This is a system of inequalitiesin two variables. If there is no solution, then make c larger so that there is a solution.

(Also, the initial hidden state vector h_0 should have $h_0[3] = 1$,to make the setup work.)

- 2e. To ensure that almost every sentence contains Trump, we need to ensure that $w_{j+1} == \text{EOS}$ is improbable if $h_0[3] = 1$. We can do this by setting $e_{j+1}[4]$ to be very negative.

Note that this is just a demonstration that the architecture has the ability to achieve the desired behavior, with appropriate parameters. In practice, we will rely on training to find good parameters.

- 3a. $\sigma(W_{\{S \rightarrow NP VP\}} [1; h_{NP}; h_{VP}])$
Here $[1; h_{NP}; h_{VP}]$ is a column vector.

An equivalent expression is
 $\sigma([1, h_{NP}, h_{VP}] W_{\{S \rightarrow NP VP\}})$
where $[1; h_{NP}; h_{VP}]$ is a row vector.

σ denotes the logistic function, $\lambda x / (1 + \exp(-x))$.

- 3b. For the version of the formula shown in the previous answer, it must have d rows and $2d+1$ columns, so that it can multiply

by the $(2d+1)$ -dimensional column vector $[1; h_{NP}; h_{VP}]$ to yield a d -dimensional embedding h_S .

- 3c. I'll write out the full encoding as a sum of terms, and then point to the one that was requested as an answer.

Remember that we use $,$ to concatenate matrices horizontally and use $;$ to concatenate them vertically.

If we drop the σ , then for the parse tree shown, we have

$$\begin{aligned}h_S &= W_{\{S \rightarrow NP VP\}} [1; h_{NPs}; h_{VP}] \\h_{VP} &= W_{\{VP \rightarrow V NP\}} [1; h_V; h_{NPo}] \\h_{NPo} &= W_{\{NP \rightarrow Det N\}} [1; h_{Det}; h_N]\end{aligned}$$

where h_{NPs} is the encoding of the subject and h_{NPo} is the encoding of the object.

To make this easier to read, let's use numeric suffixes instead of rule subscripts:

$$\begin{aligned}h_S &= W1 [1; h_{NPs}; h_{VP}] \\h_{VP} &= W2 [1; h_V; h_{NPo}] \\h_{NPo} &= W3 [1; h_{Det}; h_N]\end{aligned}$$

Express each binary rule's matrix W as $[b, L, R]$ where b is the first column, L is the next d columns, and R is the final d columns. ("b" = bias, "L" = left child, "R" = right child)

This allows us to rewrite the above as

$$\begin{aligned}h_S &= b1 + L1 h_{NPs} + R1 h_{VP} \\h_{VP} &= b2 + L2 h_V + R2 h_{NPo} \\h_{NPo} &= b3 + L3 h_{Det} + R3 h_N\end{aligned}$$

Substituting within these equations, we get

$$\begin{aligned}h_S &= b1 + L1 h_{NPs} \\&\quad + R1 (b2 + L2 h_V \\&\quad\quad + R2 (b3 + L3 h_{Det} \\&\quad\quad\quad + R3 h_N)) \\&= b1 + L1 h_{NPs} \\&\quad + R1 b2 + R1 L2 h_V \\&\quad\quad + R1 R2 b3 + R1 R2 L3 h_{Det} \\&\quad\quad\quad + R1 R2 R3 h_N\end{aligned}$$

which is a sum of vectors as claimed.

The vectors associated with the 4 preterminal nodes are

L1 h_{NPs} for (NP Papa)
R1 L2 h_V for (V ate)
R1 R2 L3 h_{Det} for (Det the)
R1 R2 R3 h_N for (N caviar)

There are also vectors associated with the binary nodes, which the question forgot to mention. These would go away if we dropped the "1" from the formula in part (a).

b1 for the root S node
R1 b2 for the VP node
R1 R2 b3 for the object NP node

Notice the pattern. The embedding at the root, h_S , is a sum over terms corresponding to the nodes of the tree. The path from S to each node determines the sequence of matrices that are

multiplied together and by the node's embedding.

In particular, to answer the question: the node (Det the) is reached from the root by taking

the right child of S \rightarrow NP VP,
then the right child of VP \rightarrow V NP,
then the left child of NP \rightarrow Det N

so it contributes the vector

R1 R2 L3 h_Det

or to write that out in full,

$R_{\{S \rightarrow NP VP\}} R_{\{VP \rightarrow V NP\}} L_{\{NP \rightarrow Det N\}} h_{Det}$

Remark: Is this a *good* encoder? It basically just adds up all of the words (and internal nodes) of the tree, representing them by their role in the tree. The fact that the roles are represented is good, but there is no interaction among the words. For example, the representation of "caviar" is not altered by the fact that the caviar is being eaten, only by the fact that it is an object. The σ function does allow interactions -- you could have a dimension of the encoding that detects when caviar is being eaten.