

0.a

First let us denote by t the index in which y has the 1. For simplicity, we assume that the base of the log is e . We get that

$$CE(\hat{y}) = -\sum_i y_i \log(\hat{y}_i) = -\log(\hat{y}_t) = -\log\left(\frac{\exp(\theta_t)}{\sum_j \exp(\theta_j)}\right) = -\theta_t + \log\left(\sum_j \exp(\theta_j)\right)$$

Let us look at the right term of the last equation. We get that

$$\frac{\partial \log(\sum_j \exp(\theta_j))}{\partial \theta_i} = \frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)} = \hat{y}_i$$

For the term we that $\frac{\partial -\theta_t}{\partial \theta_i} = -1$ if $i = t$ and 0 otherwise. To summarize, we get

$$\frac{\partial CE}{\partial \theta_i} = \begin{cases} \hat{y}_i - 1, & i = t \\ \hat{y}_i, & i \neq t \end{cases}$$

Note that the partial derivative for θ_t is negative, which makes sense since when applying SGD we subtract the derivative from the current value of θ_t , hence we would enlarge θ_t and minimize the loss. For all other $i \neq t$ the derivative is positive.

0.d

The advantage of character level model stems from the fact that two words can be very similar, for example "walk" and "walked". With a word level model we don't leverage this similarity and ignore this information which may lead to poorer performance. Character level model can infer that the character sequence "walk" and the character sequence "walked" are similar. If we treat words with capital letters as different words, e.g. "Walk" and "walk" then it's even more confusing to word level model. The disadvantage of character level model is that it increases the length of the sequence. If we have a document of 100 words, and the average length of a word is for example 4, then with a word level model the sequence length is 100 and with a character level model the length is 400. Plus, with character level model there's an added difficulty - the model needs to learn what sequences of characters constitute as legitimate words. This is not a problem with word level model that get this information "for free".

0.e

By using basic power laws, we can rewrite the given expression as follows:

$$\begin{aligned} 2^{-\frac{1}{M} \sum \log_2 p(s_i | s_1, \dots, s_{i-1})} &= \\ &= (2^{\sum_i \log_2 p(s_i | s_1, \dots, s_{i-1})})^{-\frac{1}{M}} = \\ &= \left(\prod_i 2^{\log_2 p(s_i | s_1, \dots, s_{i-1})}\right)^{-\frac{1}{M}} = \\ &= \left(\prod_i p(s_i | s_1, \dots, s_{i-1})\right)^{-\frac{1}{M}} = \\ &= \left(\prod_i e^{\ln p(s_i | s_1, \dots, s_{i-1})}\right)^{-\frac{1}{M}} = \\ &= e^{-\frac{1}{M} \sum \ln p(s_i | s_1, \dots, s_{i-1})} \end{aligned}$$

References