

Applied Natural Language Processing
Prof. Ramaseshan Ramachandran
Department of Computer Science and Engineering
Chennai Mathematical Institute, Madras

Lecture - 58
BPTT - Derivatives for W, V and U

(Refer Slide Time: 00:15)

BPTT - DERIVATIVE FOR V

$$h_t = Wx_t + Uh_{t-1}$$

$$s_t = \tanh(h_t)$$

$$z_t = Vs_t$$

$$\hat{y}_t = \text{softmax}(z_t)$$

$$E_t = -y_t \log(\hat{y}_t)$$

$$\frac{\partial E_t}{\partial V} = \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial V} \quad (6)$$

Let $\delta'_{out} = \frac{\partial E_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t}$

$$\frac{\partial E_t}{\partial V} = \delta'_{out} s_t \quad (7)$$

Here δ'_{out} is the loss for each of the units in the output layer

So, let us look at the derivatives now.

(Refer Slide Time: 00:20)

FORWARD PASS - NETWORK EQUATIONS

If the corpus contains T words, then $(x_1, x_2, x_3, \dots, x_T)$ are the corresponding word vectors

- ▶ $x_t \in R^{D_w}$ represents the input word at time t and D_w is the dimension of the word vector. If one-hot vector, it will be x^{D_w} 100 words
- ▶ $W \in R^{D_h \times D_w}$ is the weight matrix that conditions the input vector
- ▶ $U \in R^{D_h \times D_h}$ matrix that keeps the dependency of the word sequence
- ▶ $V \in R^{D_y \times D_h}$
- ▶ s_{t-1} is the output of the non-linear function (tanh) of the time step $t-1$
- ▶ $\hat{y}_t \in R^{|V|}$ is the probability distribution of the predicted word at time step t for the given context of $x_1, x_2, x_3, \dots, x_t$, where $|V|$ is the size of the vocabulary

Forward pass

$$h_t = Wx_t + Uh_{t-1} \quad (1)$$

$$s_t = \tanh(h_t) \quad (2)$$

$$z_t = Vs_t \quad (3)$$

$$\hat{y}_t = \text{softmax}(z_t) \quad (4)$$

$$E_t = -\sum_t y_t \log(\hat{y}_t) \quad (5)$$

So, what we want to find out is to minimize this error right this should be E_t . So, in order for us to find these smallest values of the derivative of E_t w.r.t. V , we have to do the derivation for all the parameters that we are looking at DV rather V W and U . So, let us start from the first parameter ok. So, we want to find out what is the change that we want to make when there is an error that is coming from the output layer ok. So, I have listed this for our convenience.

Let us look at $\frac{dE_t}{dV}$. So, we are taking these errors at state t right. So, we have to do $\frac{dE_t}{dV}$ which would be equal to $\frac{dE_t}{dy_t} \frac{dy_t}{dz_t} \frac{dz_t}{dV}$. So, we are going to be finding the error with respect to y_t right. So, we have $\frac{dE_t}{dy_t}$ and then we are going to be finding this $\frac{dy_t}{dz_t}$ with respect to E_t and then we are going to be finding the change with respect to V .

So, it is a chain rule that you can apply wherein you will have this $\frac{dE_t}{dy_t}$ there and then here you will write you y_t that by $\frac{dy_t}{dz_t}$ and then here you will have $\frac{dz_t}{dV}$.

So, we have to find these values individually and then multiply that you will get you $\frac{dE_t}{dV}$ ok. So, we know that when we differentiate with respect to this ok. So, this is nothing, but change in this value right. So, if we use a mean square value this is going to be $y_t - \hat{y}_t$ or it will be $1/y_t$ if you are going to be using a cross-entropy model this becomes 1 because we are going to be using a 1 vector if we use one vector, this will become 1 otherwise it will be y_t ok.

So, let us not really worry about what we are going to be taking on the right side. So, we will just have the notation for that and then use it in the computation of the derivative for V all right. So, what we have done here is. So, we are going to be taking these two and then calling it as $\delta_{out,t}$. So, this is going to be the loss for each of the units in the output layer here ok.

For each of the units in the output layer. So, we call it as $\delta_{out,t}$ and then when you differentiate this when you do the partial differentiation of $\frac{dE_t}{dy_t}$ by $\frac{dE_t}{dz_t}$ what you get? Oh, I am sorry. So, this is what we are calling it as $\delta_{out,t}$. So, we are going to be doing the partial derivative of $\frac{dE_t}{dz_t}$ with respect to V . So, what you get is s_t . So, the error with respect to the value here is $\delta_{out,t}$ into s_t all right. So, this is the derivative for V .

(Refer Slide Time: 04:19)

BPTT - DERIVATIVE FOR W

$$\frac{\partial E_t}{\partial W} = \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial z_t} \frac{\partial z_t}{\partial s_t} \frac{\partial s_t}{\partial h_t} \frac{\partial h_t}{\partial W} \quad (8)$$

$$= \delta_{out}^t V \sigma'(h_t) x_t \quad (9)$$

Since the hidden layer activation depends on the previous time state, we have another similar term $\hat{\delta}_{-1}^t$ that get added to ()

9 / 27
NPTEL

And then now let us look at the derivative for W, it is in the similar fashion we are going to be using the chain rule to find out what is the change that we require. So, that we can update W with that particular change and this is the change that you want to find out with using which we will update the W.

Again use the chain rule when I am going to leave this to you as an exercise to find out how these values are obtained. So, in this case I am not completely changing these values I am just using a sigma dash this is a derivative I am just leaving it as is. So, if you can actually expand.

When you want to write your computer application using this ok. So, $\frac{\partial E_t}{\partial W}$ is obtained using the chain rule and then the value would be $\delta_{out}^t V \sigma'(h_t) x_t$ and next is the x_t comes from here and this is what you get from $\frac{\partial s_t}{\partial h_t}$ what you get from here. So, you may want to expand this and then see what the actual value is ok.

(Refer Slide Time: 06:03)

BPTT - DERIVATIVE FOR U

$$\frac{\partial E_t}{\partial U} = \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial z_t} \frac{\partial z_t}{\partial s_t} \frac{\partial s_t}{\partial h_t} \frac{\partial h_t}{\partial U} \quad (10)$$

$$= \delta_{out}^t V \sigma'(h_t) h_{t-1} \quad (11)$$

Since we are back propagating the error from the current state to the previous state, $\delta_{next} = \sigma(h_t) U \delta_{out}^t V \sigma'(h_t)$ needs to be added

9 / 27

Let us move on to the next one. So, now, we have to find the derivative with respect to U so, that we can update the matrix U here. Again apply the chain rule. So, we are going to be getting values similar to this when you are doing the backpropagation you will notice that this particular state depends on the previous state right. So, there is some contribution that is coming from the previous state. So, we need to make sure that the value of the previous state also is included in this particular computation.

So, that is what I write here as delta next this value will turn out to be and this needs to be added to $\frac{\partial E_t}{\partial U}$ all right.

(Refer Slide Time: 06:59)

The error for the entire duration of T for all the vocabulary is the sum of all the error across the layers

$$E(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{|V|} y_{t,j} \log(\hat{y}_{t,j}) \quad (12)$$

This term (\cdot) is known as the perplexity. Lower the value of $2^{E(\theta)}$ better is the confidence of the network in predicting the next word

19 / 27
NPTEL

So, when I unroll this we have done only for one small piece in the RNN. So, now, we unroll it. So, now, this is how the entire network looks ok. So, we have from the timestamp of 1 to time slice t ok. So, we have all the states now available.

So, how do you do the backpropagation in this case? We have done it for one now we have to see how can be extended for the entire talk. So, since they are unrolled now when they are enrolled you do not see any of those, but when you when they are not enrolled you see like this, when they are unrolled you see what you are seeing on the screen right like this and then every unit or every state you should do the operation of what we are done earlier the error for the entire duration T is obtained using this relationship ok.

So, that means, the error is summed across. So, every one that you find here they are summed and then the error with respect to those parameters should be minimized ok. So, that is the actual aim of this backpropagation through time all right. So, this has to be continued until we reach a state where there cannot be any more updates possible we are able to do the backpropagation through time in this fashion.

Once in neural networks settles we will stop or when the error becomes very small we stop or when the error does not change beyond a point we will stop again if you look at this parameter when this become smaller and smaller the confusion in terms of what is the pattern that I want to identify you know within the network, I am talking about

network looking at the patterns right. The confusion for it is going away when the value of E becomes smaller and smaller ok.

So, that is what we call it as perplexity. So, when this becomes smaller we know that network now has a lot of confidence in terms of identifying or predicting the next word if you are using it for natural language processing application ok.

(Refer Slide Time: 09:41)

PERPLEXITY

Perplexity is a measurement of how well a model predicts a sample. Perplexity is defined as

For bigram model, $PP(W_N) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i|w_{i-1})}}$ ✓ (13)

For trigram model $PP(W_N) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i|w_{i-1}w_{i-2})}}$ ✓ (14)

A good model gives maximum probability to a sentence or minimum perplexity to a sentence

$P(w_{10} | w_1, w_2, w_3, \dots, w_9)$

11 / 27
NPTEL

So, let us text look at what this is in terms of probability. See if I am sure you remember this right probability of predicting this word I given I minus 1 or this is for the bigram right and this is for the trigram and then if you use a character-based model that we are looking at right now, let us say I am going to be looking at the 10th character and I want to predict the 10th character. So, in this model what we are actually computing is this correct this is what we are to computing.

So, if in order for you to compute this and then identify that the machine has a lot of confidence in terms of protecting the right value, we need some kind of a mechanism to measure that right. So, and that what we use here. So, it is very similar to what we are here finding is nothing but right this is what we are trying to find out. So, once the prediction is right we move on to the next job correct. So, the perplexity is computed using this formula the lower the perplexity better it is all right.