

Deep Learning
Prof. Mitesh M. Khapra
Department of Computer Science and Engineering
Indian Institute of Technology, Madras

Lecture – 111
How LSTMs avoid the problem of vanishing gradients (Contd.)

(Refer Slide Time: 00:12)

The diagram illustrates the internal operations of an LSTM cell at time t . It shows the flow of information and gradients through four gates: selective write, selective read, selective forget, and selective write. The input x_t and the previous hidden state h_{t-1} are combined with weights W and U and passed through a sigmoid function σ . The resulting gates are used to update the cell state s_t and produce the output h_t .

- If the state at time $t - 1$ did not contribute much to the state at time t (i.e., if $\|f_t\| \rightarrow 0$ and $\|o_{t-1}\| \rightarrow 0$) then during backpropagation the gradients flowing into s_{t-1} will vanish
- But this kind of a vanishing gradient is fine (since s_{t-1} did not contribute to s_t we don't want to hold it responsible for the crimes of s_t)
- The key difference from vanilla RNNs is that the flow of information and gradients is controlled by the gates which ensure that the gradients vanish only when they should (i.e., when s_{t-1} didn't contribute much to s_t)

29/33

Ok so, will start from where we left off. So, in the last class we started with this motivation that recurrent neural networks have this problem of Vanishing and Grade Exploding Gradients. And we wanted to arrive with some principle way of avoiding this. So, you have first started with this intuition that in many real life situations like for example, the human brain or the whiteboard. We tend to these to these three operations called selective read, selective write, and selective forget. And they essentially help us in dealing with these finite sized memories right or whether it is a whiteboard which is finite sized or your brain or whatever it is right.

So, can we is it possible to kind of improve RNN's which also suffer from this problem that they have this finite sized memory. And hence if you are trying to capture everything from time step one then by the time you reach time step t where say t is 30 or 40 or so on. It is quite natural that whatever you have learned earlier will get move off to an extent that it just is not recognisable anymore right.

So, you wanted to deal with this problem and with that we motivated selective read write and forget. And then we introduced some equations or converted this into a model and this is the diagram that you see is the model actually that is the LSTM cell data. And it has these three gates output gate, input gate, and forget gate and which perform these three functions of selective read write and forget. So, intuitively all these was fine, but we need to be more technical in terms of you trying to deal with a problem of vanishing and grade exploding gradients.

So, how does it solve that problem all that makes or the story seems fine, but how does this actually relate to the math. So, we saw some intuition for that and the intuition hinged on this observation that. During forward pass the gates control how much of information passes from one state to another. And in particular if you have the situation that from one time step to another say the forget gate tells you that keep forgetting point 5 of the previous state. Then by the time you reach say the 100 state you would have forgotten 0.5^{100} of the first state.

So; that means, even during forward pass the information from state 1 vanishes. So, if it vanishes during backward path that is also fine, because state 1 did not contribute to state 100. And that was the intuition that all this hinged on now we are not going to do much different from this intuition. We just going to see the corresponding equations for these intuitions and just make a more I would not call it rigorous but more mathematical proof on why LSTM solve the problem of vanishing gradients.

And we are also sure that they actually do not solve the problem of exploding gradients and then we will see a simple trick of dealing with exploding gradient. That is what we will do in the remainder of this particular lecture and then will move on to the next lecture in this lecture.

(Refer Slide Time: 03:02)

We will now see an illustrative proof of how the gates control the flow of gradients

So, we will now see an illustrative proof of how the gates control the flow of gradients right.

(Refer Slide Time: 03:11)

- Recall that RNNs had this multiplicative term which caused the gradients to vanish

$$\frac{\partial \mathcal{L}_t(\theta)}{\partial W} = \frac{\partial \mathcal{L}_t(\theta)}{\partial s_t} \sum_{k=1}^{t-1} \prod_{j=k}^{t-1} \left(\frac{\partial s_{j+1}}{\partial s_j} \right) \frac{\partial^+ s_k}{\partial W}$$

~~$(\lambda \gamma)^{t-k}$~~

So, we call that this is the control this is the flow diagram or the dependency diagram that you had for RNN's. And in particular because you are dealing with an ordered network we add this explicit and implicit derivatives and finally, you came up with this multiplicative form. And this term here is actually a matrix because it is a derivative of a vector with respect to a vector.

And then this same matrix was getting multiple times and then we did this proof it showed that this term is actually lambda gamma ok. It is actually proportional to this term right and as if lambda into gamma is greater than 1, then this will explode. If it is less than 1 then it will vanish given sufficient times that is.

(Refer Slide Time: 03:56)

- Recall that RNNs had this multiplicative term which caused the gradients to vanish

$$\frac{\partial \mathcal{L}_t(\theta)}{\partial W} = \frac{\partial \mathcal{L}_t(\theta)}{\partial s_t} \sum_{k=1}^t \prod_{j=k}^{t-1} \frac{\partial s_{j+1}}{\partial s_j} \frac{\partial^+ s_k}{\partial W}$$

- In particular, if the loss at $\mathcal{L}_t(\theta)$ was high because W was not good enough to compute s_1 correctly then this information will not be propagated back to W as the gradient $\frac{\partial \mathcal{L}_t(\theta)}{\partial W}$ along this long path will vanish

31/43

Mitesh M. Khapra CS7015 (Deep Learning)

Now, in particular what is happening here is the following that you have this loss at time step t you have the time step is 4. Now, what if this loss or this error occurred, because W was not good enough to compute a good value for S_1 right. So, W was at a certain configuration based on that you computed S_1 . And that S_1 was not good enough which eventually led to the error at time step 4 all of you if you can imagine this situation that you mean you not being not be able to do something well at S_1 . Now this needs to be told to W so that it can improve right. And that information has to come through S_1 that information is already going from here, but this information is about how badly it performed in computing S_4 .

This is not how badly it can perform in computing S_1 . So, that information has to travel to W all the way through S_1 and that was not happening because this path do not look at the bullets this path was actually vanishing. And that is what this multiplicative term says that as the number of times that increased that time that path would vanish ok. So, that is the actual problem that we are trying to deal.

(Refer Slide Time: 05:05)

• In general, the gradient of $\mathcal{L}_t(\theta)$ w.r.t. θ_i vanishes when the gradients flowing through **each and every path** from $\mathcal{L}_t(\theta)$ to θ_i vanish.

Mitesh M. Khapra CS7015 (Deep Learning)

So, now what is the general situation here right, the general principle is that the gradient of L theta at particular time step say here we are considering L_4 so I will just call it L_t with respect to any parameter theta i . The parameters at W u v b and c with respect to any parameter it would vanish if all the paths leading to that parameter if it vanishes. So, with respect to this particular path so that is the only path which leads to W through S_1 . If there were multiple paths if there was say one such direct path right if we had you some other kind of connection which gave us this direct path then it would still have been fine.

But there was only one path leading to W through S_1 at the gradient vanishes along that path then the gradient will vanish ok. If there were multiple paths then only if the gradient vanishes across all the paths then the gradient would vanish is it fine. What is the corresponding rule for exploding gradients? If there are multiple paths the gradient would explode if.

Student: (Refer Time: 06:06).

If it vanishes through any o[ne]- if it explodes though any one of the paths ok.

(Refer Slide Time: 06:11)

- In general, the gradient of $\mathcal{L}_l(\theta)$ w.r.t. θ_i vanishes when the gradients flowing through **each and every path** from $L_l(\theta)$ to θ_i vanish.
- On the other hand, the gradient of $\mathcal{L}_l(\theta)$ w.r.t. θ_i explodes when the gradient flowing through **at least one path** explodes.
- We will first argue that in the case of LSTMs there exists at least one path through which the gradients can flow effectively (and hence no vanishing gradients)

Mitesh M. Khapra CS7015 (Deep Learning) 32/43

So, these are the two things that we need to consider ok. So, to prove that in the case of LSTM this does not happen. For the first case will have to show that there are at least one path through which it does not vanish and for the second case because we are going to show that it explodes we just have to show that there is at least one path through which it can explode ok. So, these are the two things that we need to prove and the first thing that we are going to focus on; is the vanishing gradient problem.

(Refer Slide Time: 06:38)

- We will start with the dependency graph involving different variables in LSTMs
- Starting with the states at timestep $k - 1$

$$o_k = \sigma(W_o h_{k-1} + U_o x_k + b_o)$$
- For simplicity we will omit the parameters for now and return back to them later

$$i_k = \sigma(W_i h_{k-1} + U_i x_k + b_i)$$

$$f_k = \sigma(W_f h_{k-1} + U_f x_k + b_f)$$

$$\tilde{s}_k = \sigma(W h_{k-1} + U x_k + b)$$

$$s_k = f_k \odot s_{k-1} + i_k \odot \tilde{s}_k$$

$$h_k = o_k \odot \sigma(s_k)$$

Mitesh M. Khapra CS7015 (Deep Learning)

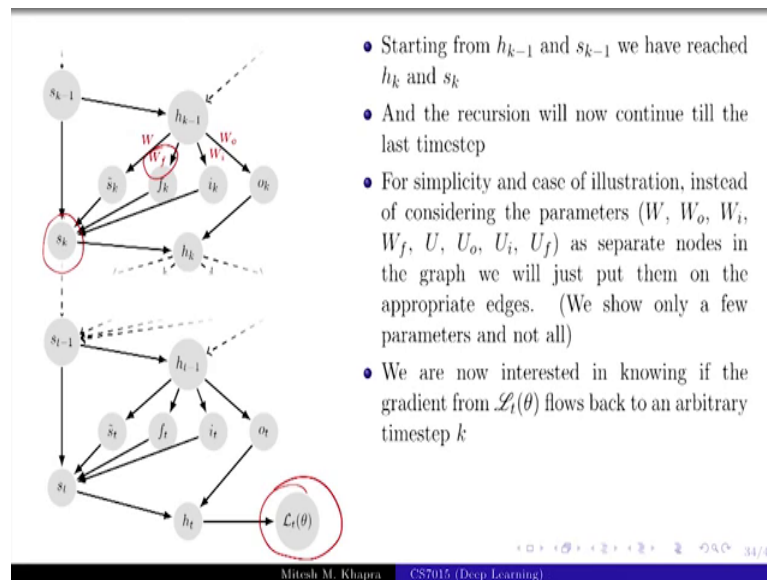
So, will start with the dependency graph for LSTM's; that means, I want to draw something similar for LSTM's involving all the different elements in LSTM. So, what are these different elements the two rhyming things one being gates states ok. So, gates and states that the two things that we care about. So, let us look at all these. So, starting with states at time step k minus so, at time step k minus 1 you have this two states s_{k-1} and h_{k-1} ok. Using h_{k-1} you are going to compute the output gate at time step k . And it is also depends on these parameters W_o , U_o and b_o right which is obvious from the equation.

Just to make sure that this diagram remains tractable, I am going to get rid of the parameters and I will come back to them later. So, right now will just focus on the states and the gates ok and then you have these other intermediate states and the other gates right. So, you had f_k you had i_k . So, add these three gates the temporary state and then what else what are the other two things at time step k . So, we saw this diagram about all the computations which happen at time step k right. How many computations happen? Three states and three gates right.

So, you seen the three gates and this one temporary state. So, which are the other two things? There is no selective forget with you guys is early everything forget. Hint look at the grey cells and change the time step. What will you get away are you all I mean we did LSTM's two days right I mean are you all with that or should I we need to revise something; mean I do not need to revise it, but we going to is it fine ok. So, s_k and the other thing h_k remember that s_k also depends on h_k just stare at this for 30 seconds and make sure that you are with it right. All the equations are there these are the see 6 equations that or the 6 computations which happen at time step k .

There are three gates and three states and the dependency graph is obvious from these equations; except for the fact that I have ignored the parameters. How many if you are comfortable with the equations and the graph corresponding graph please raise your hands high. So, I think it should be right we have these six equations and we have this dependency graph.

(Refer Slide Time: 08:50)



Now, starting so, what happened in the graph is, we started from s_{k-1} and h_{k-1} and we reached s_k and h_k which were the outputs at the next state. Now what will happen from here? We were looking at recurrent neural networks recursion is the answer. What will happen now?

The same graph will keep recursing right for the next time step and up to the last time step right, does that make sense ok? This is much more complicated than the dependency graph that we had for RNN's right by just because there are so many in RNN we just had this one state and no gates so here but we have these three states and three gates that is why this so many paths ok. Now for simplicity what I will do is, I will not draw separate nodes for the parameters all the in the case of the RNN dependency graph I had drawn them separately. What I am going to do is, I am just going to put the parameters on the corresponding edges right.

So, f_k actually depends on W_f , it also depends on U_f , and it also depends on that bias. But I am just going to take a small set of parameter I am only going to focus on the W 's not the U 's and the biases ok. There is only for illustration for no other reason right and whatever arguments or proof that we are going to see it holds for all the parameters, but we just need to prove it with respect to one parameter and the same story repeats for everything ok. So, this is the dependency graph and these are the parameters. Now what I

am interested in knowing is that, there was some loss at time step t and maybe that loss happened because W_f was not good enough to compute s_k .

Of course W_f computes f_k and then f_k helps in computing s_k , but maybe W_f was not I am just short it short circuiting it and saying that W_f was not good enough to compute s_k right. And that is why I want the gradient to reach to W_f through this s_k that is what I want do ok.

(Refer Slide Time: 10:42)

- For example, we are interested in knowing if the gradient flows to W_f through s_k
- In other words, if $\mathcal{L}_t(\theta)$ was high because W_f failed to compute an appropriate value for s_k then this information should flow back to W_f through the gradients
- We can ask a similar question about the other parameters (for example, W_i , W_o , W , etc.)
- How does LSTM ensure that this gradient does not vanish even at arbitrary time steps? Let us see

Mitesh M. Khapra CS7015 (Deep Learning)

And this exactly what I said I am interested in knowing that if this loss can reach W_f through s_k right. So, all the three highlighted things that what I am interested in, I am interested in the path to W_f through s_k . Of course, there are many other paths to W_f , but they do not account for the problem in s_k ; is that fine everyone is clear the setup ok?

Now and we can ask similar questions about all the other parameters the W 's the us the the input gate parameters the output gate parameters and so on right. There is nothing so special about W_f the same question holds for all these other parameters also ok. Now how does LSTM ensure that this does not vanish? So let us see that.

(Refer Slide Time: 11:20)

- It is sufficient to show that $\frac{\partial \mathcal{L}_t(\theta)}{\partial s_k}$ does not vanish (because if this does not vanish we can reach W_f through s_k)
- First, we observe that there are multiple paths from $\mathcal{L}_t(\theta)$ to s_k (you just need to reverse the direction of the arrows for backpropagation)
- For example, there is one path through s_{k+1} , another through h_k
- Further, there are multiple paths to reach h_k itself (as should be obvious from the number of outgoing arrows from h_k)
- So at this point just convince yourself that there are many paths from $\mathcal{L}_t(\theta)$ to s_k

Mitesh M. Khapra CS7015 (Deep Learning) 30/43

As I argued earlier it is sufficient to show that this gradient does not vanish ok. If I can show that this gradient does not vanish, then I am pretty sure there is only there is no recursive connection here because it just a single connection. So, there is no recursive connection here. So if I can show that the gradient reaches up to this point, then after that I can be sure that it is going to reach W_f everyone buys that set up right that is what I need to show?

So, to prove that the gradient reaches W_f I just need to show that it reaches s_k that is the only thing that I need to show. And the first thing I am going to observe is that there are multiple paths to reach to s_k which are these paths? One through s_{k+1} , because s_k contributes to s_{k+1} the other through h_k .

Student: h_k .

h_k which is visible, but now also notice that how many paths are there to reach h_k itself. Not four actually that is going to be combinate relate because there four outgoing edges from here, but then again there will be four next stage and four next stage and so on right. So, let us not count the number of paths, but let us just convince ourselves that there are many many paths to reach to s_k from $\mathcal{L}_t(\theta)$. Everyone is convinced about that we are not counting the exact number of paths that is not very hard to do. But all we are saying is that we know that there is one path through s_{k+1} , one path through h_k and h_k itself seems to have many incoming path during back propagation.

So, there are many many paths which are reaching from $L_t(\theta)$ to s_k . Everyone is convinced about that anyone who has a problem with that? Now to show that the gradient does not vanish what do I need to show of all the paths the set there exist at least one path through which the gradient can flow that is what I need to show ok. Even if I vanishes across all the other paths I am still fine with it ok.

(Refer Slide Time: 13:05)

- Consider one such path (highlighted) which will contribute to the gradient
- Let us denote the gradient along this path as t_0

$$t_0 = \frac{\partial \mathcal{L}_t(\theta)}{\partial h_t} \frac{\partial h_t}{\partial s_t} \frac{\partial s_t}{\partial s_{t-1}} \dots \frac{\partial s_{k+1}}{\partial s_k}$$

- The first term $\frac{\partial \mathcal{L}_t(\theta)}{\partial h_t}$ is fine and it doesn't vanish (h_t is directly connected to $\mathcal{L}_t(\theta)$ and there are no intermediate nodes which can cause the gradient to vanish)
- We will now look at the other terms $\frac{\partial h_t}{\partial s_t} \frac{\partial s_t}{\partial s_{t-1}} (\forall t)$

Mitabh M. Khapra CS7015 (Deep Learning) 37/43

So, now consider one such path which is this highlighted path that is a valid path to reach to s_k . Now let us denote the gradient along this path to be t_0 and the total gradient is going to be a sum of many such paths right. So, I am calling this path as t_0 and this is what the gradient look like ok. So, this is simple just this red path the next red path and then the series of problematic multiplications right you have this recursive multiplications again. So, everyone agrees that red is good, the red path there is no recursion the gradient will flow right we just need to focus on the blue path everyone is convinced about that right ok.

So, that is good the first term is fine as I said because it directly connected to $L_t(\theta)$ there is no recursive or no other intermediate nodes. So, the gradient will just flow through that there is not a problem there and now we look at the other terms which is first is $\frac{\partial h_t}{\partial s_t}$ and the other is this ok.

(Refer Slide Time: 14:03)

- Let us first look at $\frac{\partial h_t}{\partial s_t}$
- Recall that

$$h_t = o_t \odot \sigma(s_t)$$

$$h_t = [h_{t,1} \quad h_{t,2} \quad \dots \quad h_{t,d}]$$

$$o_t = [o_{t,1} \quad o_{t,2} \quad \dots \quad o_{t,d}]$$

$$s_t = [s_{t,1} \quad s_{t,2} \quad \dots \quad s_{t,d}]$$

i, j

38/43

Mitesh M. Khapra CS7015 (Deep Learning)

So, let us look at $\frac{\partial h_t}{\partial s_t}$ by s_t what is this going to be? Tensor, vector, matrix scalar at this point in the course I want a unanimous answer.

Student: Matrix

Matrix right and recall that in particular the equation was of this form ok. So, what is the derivative going to look like even without computing can you tell me something profound about it; it will be a dash matrix. Big matrix, how many if you say diagonal matrix? How many if you do not think it is a diagonal matrix please raise your hands total sum is never one. So, remember that h_t is equal to $h_{t,1}, h_{t,2}$ up to $h_{t,d}$. And you have o_t equal to $o_{t,1}, o_{t,2}, o_{t,d}$ and s_t equal to $s_{t,1}, s_{t,2}, s_{t,d}$. So, $h_{t,2}$ depends only on $o_{t,2}$ and $s_{t,2}$ right it does not depend on in particular does not depend on any of the other s_t 's.

So, we have already seen this before in such cases what is the ij th entry of this matrix of the gradient matrix derivative of $h_{t,i}$ with respect to $s_{t,j}$ which of these terms are going to be zero wherever?

Student: I not equal to 0.

I is not equal to 0; that means, it results in a.

Student: Diagonal matrix

Diagonal matrix.

(Refer Slide Time: 15:31)

- Let us first look at $\frac{\partial h_t}{\partial s_t}$
- Recall that

$$h_t = o_t \odot \sigma(s_t)$$
- Note that h_{ti} only depends on o_{ti} and s_{ti} and not on any other elements of o_t and s_t
- $\frac{\partial h_t}{\partial s_t}$ will thus be a square diagonal matrix $\in \mathbb{R}^{d \times d}$ whose diagonal will be $o_t \odot \sigma'(s_t) \in \mathbb{R}^d$ (see slide 35 of Lecture 14)
- We will represent this diagonal matrix by $\mathcal{D}(o_t \odot \sigma'(s_t))$

[]

38/43

Mitesh M. Khapra CS2015 (Deep Learning)

So, that is exactly what is written here and the diagonal elements are going to be this is that fine everyone with this ok. So, now, this diagonal matrix which contains this on the diagonal I am going to represented by the following notation is that fine. So, this is a diagonal matrix where every element is, I mean this is actually a vector right everyone agrees this is a vector. So, this diagonal is this vector is along the diagonal of this matrix how many if you get this notation? If you do not get this you will not understand anything else.

(Refer Slide Time: 16:07)

- Now let us consider $\frac{\partial s_t}{\partial s_{t-1}}$
- Recall that

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$
- Notice that \tilde{s}_t also depends on s_{t-1} so we cannot treat it as a constant
- So once again we are dealing with an ordered network and thus $\frac{\partial s_t}{\partial s_{t-1}}$ will be a sum of an explicit term and an implicit term (see slide 37 from Lecture 14)
- For simplicity, let us assume that the gradient from the implicit term vanishes (we are assuming a worst case scenario)
- And the gradient from the explicit term (treating \tilde{s}_t as a constant) is given by $D(f_t)$.

39/43

Now let us consider $\frac{\partial s_t}{\partial s_{t-1}}$ this is what s_t is equal to. So, what is the derivative of $\frac{\partial s_t}{\partial s_{t-1}}$? f_t right f_t right what else, why no, why are you rebelling, what the I mean s_t only right. If it is can you treat this as a constant no why? Because this is a dashed network.

Student: (Refer Time: 16:34).

So in an ordered network the derivative will have.

Two terms which are those?

Student: Explicit.

Explicit and implicit In the explicit term what you assume? The other terms to be a constant right fine. So, s_t I mean \tilde{s}_t also depends on s_{t-1} . So, we cannot treated as a constant so once again this derivative is going to contain an explicit term and an implicit term. Now I am going to make a worst case assumption. I making this assumption I making this assumption that actually the implicit term vanishes. Notice that this not favourable to me I am trying to prove that the gradient does not vanish the gradient is a sum of two terms I am saying it let the worst case be that one of these terms vanishes ok.

So, this is not a favourable assumption this is a unfavourable assumption which I am making. So, let us fine. So, I making the assumption that the implicit term vanishes. So, what is the explicit term actually?

Student: F t.

F t and what kind of a matrix is that?

Student: Diagonal matrix.

If you agree that it is a matrix first of all. It is a diagonal matrix again and what is the diagonal?

Student: F t.

F t right. So, I am going to represent it as D of f t with that fine?

(Refer Slide Time: 17:46)

The diagram illustrates an unrolled LSTM cell. It shows a sequence of hidden states $h_{k-1}, h_k, h_{t-1}, h_t$ and states $s_{k-1}, s_k, s_{t-1}, s_t$. The states s_k, s_{t-1}, s_t and hidden states h_k, h_{t-1}, h_t are highlighted with blue circles. The loss function $\mathcal{L}_t(\theta)$ is shown at the bottom right. The diagram is annotated with red and blue lines to indicate the flow of gradients through the forget gates.

- We now return back to our full expression for t_0 :

$$t_0 = \frac{\partial \mathcal{L}_t(\theta)}{\partial h_t} \frac{\partial h_t}{\partial s_t} \frac{\partial s_t}{\partial s_{t-1}} \dots \frac{\partial s_{k+1}}{\partial s_k}$$

$$= \mathcal{L}'_t(h_t) \cdot \mathcal{D}(o_t \odot \sigma'(s_t)) \mathcal{D}(f_t) \dots \mathcal{D}(f_{k+1})$$

$$= \mathcal{L}'_t(h_t) \cdot \mathcal{D}(o_t \odot \sigma'(s_t)) \mathcal{D}(f_t \odot \dots \odot f_{k+1})$$

$$= \mathcal{L}'_t(h_t) \cdot \mathcal{D}(o_t \odot \sigma'(s_t)) \mathcal{D}(\odot_{i=k+1}^t f_i)$$
- The red terms don't vanish and the blue terms contain a multiplication of the forget gates
- The forget gates thus regulate the gradient flow depending on the explicit contribution of a state (s_t) to the next state s_{t+1}

Mitlesh M. Khapra CS7015 (Deep Learning) 40/43

So, remember that the original equation ah had three terms all of these the last blue once for all identical. So, this is not problematic because this is a directly the last layer this we have already derived a form this is sum diagonal. And now for each of these we have a form; do you get that these are the three paths that we have done so far. So, let me just substitute them, this is what it looks like ok. Now this is a product of diagonal matrices what will the product look like?

Student: Diagonal matrix.

A diagonal matrix and each element would be each element on the diagonal would be a.

Student: Product

A product of all those things right. So, is it fair if I write it as this right which I can write it as this ok. Now just stare at this equation and tie it back to the intuition that we developed something about the gates regulating the flow of information you have a multiplicative term here right. Whenever there is a multiplicative term we have a problem, because remember these gates are between 0 to 1.

So, there is a chance of vanishing agencies that? You are multiplying t terms all of which are between 0 to 1. So, there is a chance of vanishing. But I am going to end this proof by saying that the gradient does not vanish. So, by what am I going to do now ok. I make the statement the gradient could vanish, but this kind of vanishing is fair what do you mean by that now when will the gradient vanish?

Student: Product

At this product of the forget gates vanishes, but if the product of the forget gate vanishes; that means, what would have happened during the forward pass that information was not carried all the way back, all the way front two times step t right do you see that ok. So, that is the main reason here right. So, the red term does not vanish, the red term time zone vanish the blue term can vanish, but it will vanish only if during the forward pass also this multiplicative term at cause the information to vanish by the time you are reach the time step t . how many if you get this and this exactly what I meant earlier by saying that.

(Refer Slide Time: 19:45)

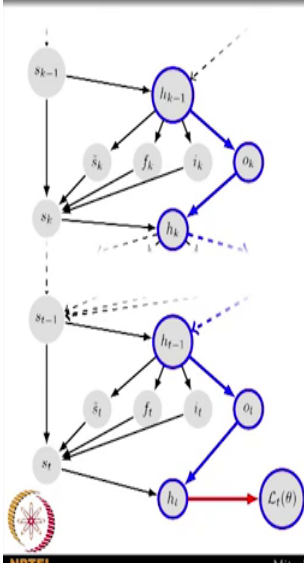
- If during forward pass s_t did not contribute much to s_{t+1} (because $f_t \rightarrow 0$) then during backpropagation also the gradient will not reach s_t
- This is fine because if s_t did not contribute much to s_{t+1} then there is no reason to hold it responsible during backpropagation (f_t does the same regulation during forward pass and backward pass which is fair)
- Thus there exists this one path along which the gradient doesn't vanish when it shouldn't
- And as argued as long as the gradient flows back to W_f through one of the paths (t_0) through s_k we are fine !
- Of course the gradient flows back only when required as regulated by f_i 's (but let me just say it one last time that *this is fair*)

Mitesh M. Khapra CS7015 (Deep Learning) 41/43

If during the forward pass s_t did not contribute much to s_{t+1} . Because the forget gate was tending to 0. Then during backward pass there is no need to pass this information back to s_t right because during forward pass you did not contribute. So, during backward pass why should I hold you responsible right? And this is absolutely fine to do this and this is exactly what the equation tells us that they gauge the gradient will vanish only if things vanished in the forward pass ok. And the gates are doing the same regulation in the forward pass as they will do in the backward pass so everything is fair is that ok?

And does there exist one path along which the gradients will not vanish when they do not need to vanish. So, if during forward pass all the gates were on; that means, the information from state one was actually carried all the way up to state t then during backward pass what will happen? The information will go all the way back right is that fine. So, the gradients flow back only when required as regulated by the forget gates and this is fair because if you are regulating the same thing in the forward as well as the backward direction then you are not doing anything wrong ok.

(Refer Slide Time: 20:51)



- Now we will see why LSTMs do not solve the problem of exploding gradients
- We will show a path through which the gradient can explode
- Let us compute one term (say t_1) of $\frac{\partial \mathcal{L}_t(\theta)}{\partial h_{k-1}}$ corresponding to the highlighted path

$$t_1 = \frac{\partial \mathcal{L}_t(\theta)}{\partial h_t} \left(\frac{\partial h_t}{\partial o_t} \frac{\partial o_t}{\partial h_{t-1}} \right) \dots \left(\frac{\partial h_k}{\partial o_k} \frac{\partial o_k}{\partial h_{k-1}} \right)$$

$$= \mathcal{L}'_t(h_t) (\mathcal{D}(\sigma(s_t) \odot o'_t) \cdot W_o) \dots$$

$$(\mathcal{D}(\sigma(s_k) \odot o'_k) \cdot W_o)$$

$$\|t_1\| \leq \|\mathcal{L}'_t(h_t)\| (\|K\| \|W_o\|)^{t-k+1}$$
- Depending on the norm of matrix W_o , the gradient $\frac{\partial \mathcal{L}_t(\theta)}{\partial h_{k-1}}$ may explode
- Similarly, W_i , W_f and W can also cause the gradients to explode

NPTEL Mitesh M. Khapra CS7015 (Deep Learning) 42/43

Now that is a proof for LSTM solve the vanishing gradient problems or in other words the gradients vanish only when required and not unnecessarily or arbitrary as is to happen in the case RNN's. Now we will show there exist one path along which the gradients can explode right. So, let us show that path. So, consider this path now this path is again also active so I if we consider the path to h_k there is going to be active for all the gates and all the states right. So, in whatever gates or states you are considering this paths would be there.

And this is what this path looks like you have the derivative with respect to the last layer and then you have these guys ok. These pairs h_t by o_t to h_{t-1} and so on again fine with this so, far what is the derivative of h_t with respect to o_t ? We do not remember the equations. So, I will just tell you directly. So, based on whatever we have done so, far just trust me that this is what each of the terms in the bracket looks like. We can go back and check this is just comes directly from the equations. Now what is happening here does this look very similar to the situation that we had with RNN's.

We had a diagonal matrix and a weight matrix and a repeated multiplication of these right and again the diagonal matrix is bounded the weight matrix is bounded. So, now, the repeated multiplication could explode is that fine. So, it does not solve the problem of exploding gradients. But it solves the problem of vanishing gradients, but now still this is

bad for us right whether the gradients explode or vanish our training is going to get messed up. So, how do we deal with this for exploding gradients what will you do?

Student: (Refer Time: 22:27) clipping.

We will do?

Student: Clipping.

Clipping right.

(Refer Slide Time: 22:32)

- So how do we deal with the problem of exploding gradients ?
- One popular trick is to use gradient clipping
- While backpropagating if the norm of the gradient exceeds a certain value, it is scaled to keep its norm within an acceptable threshold*

*Pascanu, Razvan, Tomas Mikolov, and Yoshua Bengio. "On the difficulty of training recurrent neural networks." ICML(3)28(2013):1310-1318

NPTEL Mitesh M. Khapra CS7015 (Deep Learning)

So in fact, is the way of dealing with this is gradient clipping if the norm of the gradient exceeds the certain value, when we are going to just clip it to a certain threshold. And this is fine because we care about the gradients only for the direction and not for the magnitude. Anyways when we introduce a learning rate, we are doing some kind of scale down for the gradient magnitude. So, this is just being more explicit and being careful that if the gradients are non manageable in terms of their magnitude.

Then we just going to keep them to some manageable value while being faithful to the direction. And the direction is what? Matters so that is why exploding gradients is easy, but in the case of vanishing gradients you do not have direction also because the entire gradient becomes 0, so there is no direction there right. So, that is why vanishing

gradients is more serious than exploding gradients. And as long as LSTM solve that they are fine with it is that fine ok. So, that is the end of this lecture.