

Optimization Algorithms

So, we started with Gradient Descent Algorithm for single Sigmoid Neuron & we saw how to extend to network of neurons with back propagation.

Goal

Finding a better way to traversing the Error Surface so that we can reach the minimum value quickly without resorting to brute force search.

↳ (Better way of navigating to Error Surface)

Gradient Descent Rule

So the better way we could derive it from Principled Solution from the Taylor Series.

So, we move in a direction Opposite to Gradient

Parameter Update Equations

$$w_{t+1} = w_t - n \nabla w_t$$

$$b_{t+1} = b_t - n \nabla b_t$$

Where,

$$\nabla w_t = \frac{\partial \alpha(w, b)}{\partial w}$$

$$\nabla b_t = \frac{\partial \alpha(w, b)}{\partial b}$$

$$x \rightarrow \theta \rightarrow y = f(x)$$

$$f(x) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

↳ So, we were just trying to adjust these weights & biases by hand

And we realized it's clearly not good and but we still try to do a guesswork, where we were driven by this loss function. (It tells whether the current guess is better than the previous guess or not).

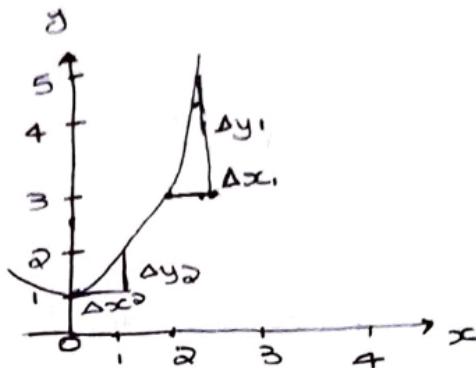
So, what we are actually doing is there is this Error Surface exists which can be plotted for all w, b values. So, we are finding a path over the surface so that we enter to the better region.

So, now we have a principled way of moving in the $w-b$ plane than our guess work algorithm.

So when we ran this Algorithm it was initially slow than Sgd.
Picked up & then it again became slow.

↓
This is because,

⇒ When the curve is steep the gradient ($\frac{\Delta y_1}{\Delta x_1}$) is large } ⇒ Larger in terms of Magnitude.



⇒ When the curve is gentle the gradient } ⇒ Smaller in terms of Magnitude.
($\frac{\Delta y_2}{\Delta x_2}$) is small

Algorithm

$t \leftarrow 0$

max_iterations $\leftarrow 1000$

while $t < \text{max_iteration}$ do

$w_{t+1} \leftarrow w_t - n \nabla w_t$

$b_{t+1} \leftarrow b_t - n \nabla b_t$

end

Momentum based Gradient Descent

Observation in Gradient Descent

- ↳ ① It takes a lot of time to navigate regions having a small gentle slope
- ↳ ② This is because gradient in these regions are very small.

So if we had set this to max_iteration Equal to 1000, so the integral happens to be such that you stuck in large flat region, then those 1000 iterations just keep moving around the flat region.

~~we will never enter into a valley. Valleys are what we're interested~~
that's because where we could find the minimal function.

Intuition

↳ If I am repeatedly being asked to move in the same direction then I should probably gain some confidence and start taking bigger steps in that direction.

Update rule

$$\text{update}_t = r \cdot \text{update}_{t-1} + n \nabla w_t$$

$$w_{t+1} = w_t - \text{update}_t$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = r \cdot \text{update}_0 + n \nabla w_1 = n \nabla w_1$$

$$\text{update}_2 = r \cdot \text{update}_1 + n \nabla w_2 = r \cdot n \nabla w_1 + n \nabla w_2$$

:

$$\text{update}_t = r \cdot \text{update}_{t-1} + n \nabla w_t = r^{t-1} \cdot n \nabla w_1 + r^{t-2} \cdot n \nabla w_1 + \dots + n \nabla w_t$$

* Is moving fast good? would there be a situation where momentum would cause us to run past our goal?

Ex:

So we have a very peculiar Error Surface. Two ~~saddles~~ at top and a very sharp valley. So the error is high on the either side of the minima Valley.

↗ Momentum based gradient descent oscillates in and out of the minima valley as the momentum carries it out of the valley.

- Takes a lot of u-turns before finally converging. Despite these u-turns it still converges faster than Vanilla gradient descent.

Nesterov Accelerated Gradient Descent (NAG)

So to reduce the oscillations we use ↗

Intuition

- Look before you leap
- Recall the update_t = r · update_{t-1} + n ∇w_t
- So, we know that we are going to move by at least by r · update_{t-1} & then a bit more by n ∇w_t
- Why not calculate the gradient ($\nabla w_{\text{look-ahead}}$) at this partially updated value of w ($w_{\text{look-ahead}} = w_t - r \cdot \text{update}_t$) instead of calculating it using the current value w_t .

Update Rule

$$w_{\text{look-ahead}} = w_t - r \cdot \text{update}_{t-1}$$

$$\text{update}_t = r \cdot \text{update}_{t-1} + n \nabla w_{\text{look-ahead}}$$

$$w_{t+1} = w_t - \text{update}_t$$

- Looking ahead helps NAG in correcting its course quicker than momentum based gradient descent.
- Hence, the oscillations are smaller and the chances of escaping the minima Valley also smaller.

pile

Stochastic & Mini Batch Gradient Descent

→ In Gradient Descent Algorithm it goes over the entire data once before updating the parameters.

So if we have million of point, we will go over all these million points & make this one update.

↓ Now imagine the consequence, when you are in plateau region. So imagine how much time it will take to converge. (One tiny update). → Very slow.

↓ To make it better ?? "Stochastic Gradient Descent"

↓
So this Algorithm updates the parameters for every single data point. So here we have million data points we update million times.

↓
This ^{is} not the True Gradient. True Gradient is the summation over all the points. Now this is no longer a true gradient this is just a point estimator, this is just a approximation of gradient.

↓ We see many oscillations why??

→ Because we are making greedy decisions.

↓
Because we are looking at one point, the point says to decrease the loss with respect to move in this direction & we blindly move in that direction. So all these points are actually trying to just make things better for themselves. They are not thinking about what is happening to all other points in my data right.

→ So some decision which I took with respect to where to move which was locally favourable for one of these points not good for another point.

↓

Can we reduce the oscillations by improving our stochastic estimates of the gradient?

⇒ Yes, let's look at mini Batch gradient Descent.

**

So looking at one data point is bad because it is very noisy, looking at the entire data is bad because it is very time consuming. So you need to do something in between which is mini batch gradient descent.

<u>Algorithm</u>	<u># of steps in 1 Epoch</u>
Vanilla (Batch) Gradient Descent	1
Stochastic Gradient Descent	N
Mini Batch Gradient Descent	N/B

Tips for Adjusting Learning Rate & Momentum

→ One could argue that we could have solved the problem of navigating gentle slopes by setting the learning rate high!!

→ So it is not that we choose high ' η ' and get away with it we actually want a adaptive ' η ' that somehow figures out that it is on gentle slope move ~~slowly~~ fast and now I am on fast slope move slow.

⇒ Tune learning rate [Try diff values on log scale, 0.001, 0.01, 0.1, 1, 10]

⇒ Run few Epoch with each of these figure out a learning rate which works best

⇒ Now do a finer search around this value around 0.05, 0.08, 0.1, 0.2 etc.

⇒ Tips for annealing learning rate

① Step decay: Halve the learning rate after Every 5 Epochs.

↳ after 5 epochs we expect closer to the solution

(Ex)

⇒ Halve the Learning rate after an Epoch if the validation error is more than what it was at the End of the epoch.

② Exponential decay: $n = n_0^{-kt}$ where n_0 & k are hyperparameters and t is the step number.

↳ where with each time step we just keep decreasing the learning rate.

③ 1/t Decay: $n = \frac{n_0}{1+kt}$ where n_0 & k are hyperparameters & t is the step number.

⇒ Tips for Momentum

$$\mu_t = \min(1 - \alpha^{-1 - \log_2(\lfloor t/250 \rfloor + 1)}, \mu_{\max})$$

Where μ_{\max} chosen from 0.999, 0.995, 0.99, 0.9, 0.7

↳ So as the Time step increasing gamma is increasing.

Gradient Descent with Adaptive Learning Rate

↳ Refer PPT & Video (5.9)

Explanation for why we need bias correction in adam

↳ Proof

*

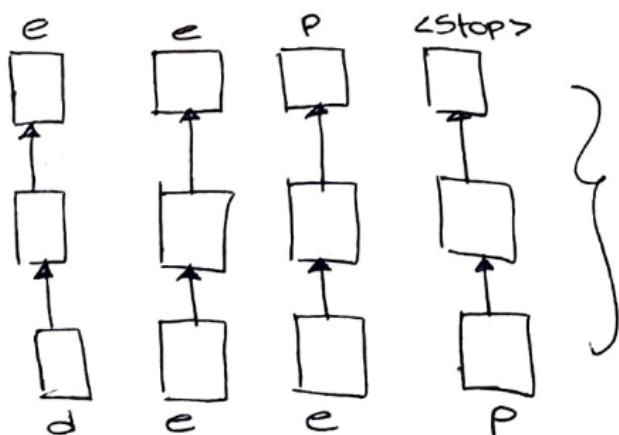
Sequence Learning Problems

↳ So far we have dealt with feedforward and convolutional neural networks and both these networks the input was always of a fixed size. Further, each input to the network was independent of the previous or future inputs.

→ But in many applications input is not of a fixed size. Further successive input may not be independent of each other.

Ex: Consider the task of auto completion

- ① Successive inputs are no longer independent
- ② the length of the input is not fixed
- ③ Each network performing the same task.



Each of this is a fully connected layer.

So these are known as Sequence learning problems. Where we have a sequence of inputs and then you need to produce some outputs and each input actually corresponds to one time step.

} Refer it once.

- So in the task of Predicting the parts of speech tag of each word in a sentence. Once we see an adjective we are almost sure that next word is noun.
 - Thus the current output depends on the current input as well as the previous input.
 - Unlike the case of CNN where we feeded an apple it is no dependence on whether the previous input that I pass to the network was an apple (or) car (or) what nots.
- ⇒ Here we are interested in producing the output at each time step.

- ⇒ Sometimes we may not be interested in producing the output at each stage.
- Ex: task of predicting the polarity of a movie review.
(sentiment analysis)
↓
So after reading the entire review model should give the prediction otherwise it would be incomplete.

- ⇒ Even here the network is performing the same task but the only thing is it don't care about the intermediate outputs.

- Sequence could be composed of any things For Example, a video could be treated as a sequence of images.
- We may want to look at the entire sequence & then detect the activity being performed.

Ex: Surya Namaskar.

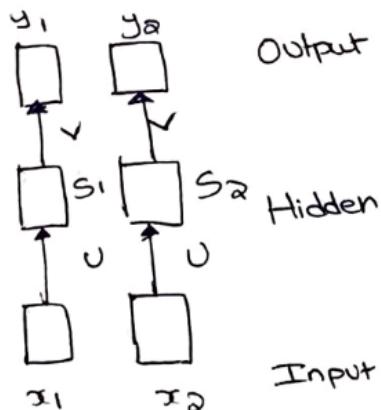
Recurrent Neural Networks (RNN)

→ So we have seen Sequence learning Problems. So the Ques is how to model this?

The Model will Come up with :

- ① Dependence b/w the inputs. (output actually depend on multiple inputs)
- ② Variable no of inputs.
- ③ It should also make sure the function Executed at each time step is the same.
(Because at Every time step we are doing the same activity).

⇒ What is the function Being Executed in each time step?



$$s_i = \sigma (Ux_i + b)$$

$$y_i = O (Vs_i + c)$$

i = timestep

→ Since we want same function to be Executed at each time step we would share the same network (Same parameters at each time step).

→ Since, we are simply going to compute the same function at each time step, the no of time steps doesn't matter

→ we just make a multiple copies of the network & Execute them at each time step.

⇒ How do we account for the case that the output actually depends on multiple inputs not on the current input?

Infeasible way of doing this

↳ At each time step we will feed all the previous inputs to the network.

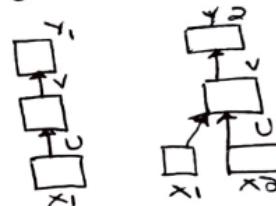


The problem with this is it violates the rules of the conditions

→ First, the function being executed at each time step now is different.

$$y_1 = f_1(x_1)$$

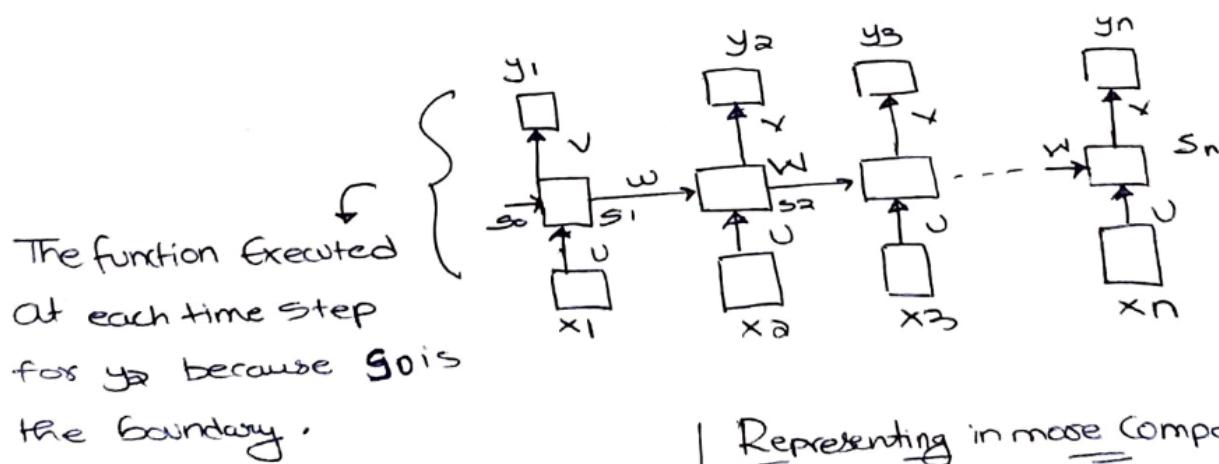
$$y_2 = f_2(x_1, x_2)$$



↳ So network is now sensitive to the length of the network.

↓
Solution

↳ Is to add a recurrent connection to the network.



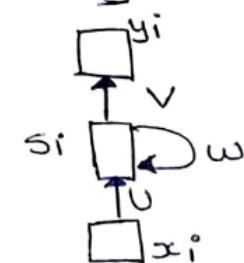
$$s_i = \sigma (Ux + Ws_{i-1} + b)$$

$$y_i = O (Vs_i + c)$$

(08)

$$y_i = f(x_i, s_i, W, U, V)$$

Representing in more compact way is



\Rightarrow s_i is the state of network at time step i

Back propagation through time

Dimensions of the Parameters :

$$\begin{aligned}x_i &\in \mathbb{R}^n && (\text{n-dimensional input}) \\s_i &\in \mathbb{R}^d && (\text{d-dimensional state}) \\y_i &\in \mathbb{R}^K && (\text{Say K classes}) \\U &\in \mathbb{R}^{n \times d} \\V &\in \mathbb{R}^{d \times K} \\W &\in \mathbb{R}^{d \times d}\end{aligned}$$

→ We train this network by using Back propagation.

Why not Standard back propagation Algorithm??

→ Ex: Task of Autocompletion

For simplicity we assume that there are only 4 characters in vocabulary ($\text{d, e, p, } \langle\text{stop}\rangle$)

→ At each time step we need to predict one of these 4 characters. So, suitable output function is Softmax and Suitable loss function is Cross Entropy.

→ Softmax is usually paired with Cross Entropy Loss Function. The combination simplifies gradients during back propagation & makes the training more efficient.

±

→ Suppose we initialize the U, V, W randomly. So the network predicts the probabilities and the true probabilities which we have.

So the two Questions ??

- ① Total loss made by the model ??
- ② How do we back propagate this loss & update the parameters ($\theta = \{U, V, W, b_i(y)\}$) of the network ??

① so the total loss is simply the sum of loss over all time steps.

$$\ell(\theta) = \sum_{t=1}^T \ell_t(\theta)$$

$$\ell_t(\theta) = -\log(y_{tc})$$

$y_{tc} \Rightarrow$ Predicted probability of true character at time step t

②

For Back propagation we need to compute the gradients w.r.t w, v, b

③ i) The loss function w.r.t to derivative $v \Rightarrow \frac{\partial \ell(\theta)}{\partial v}$

$$\frac{\partial \ell(\theta)}{\partial v} = \sum_{t=1}^T \frac{\partial \ell_t(\theta)}{\partial v}$$

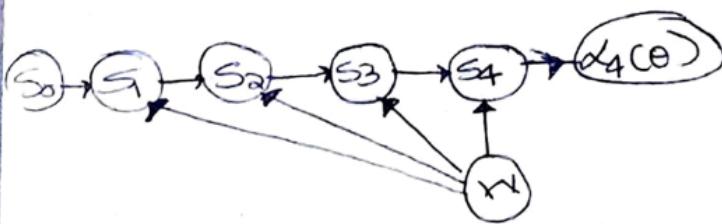
ii) The loss function w.r.t to derivative $w \Rightarrow \frac{\partial \ell(\theta)}{\partial w}$

$$\frac{\partial \ell(\theta)}{\partial w} = \sum_{t=1}^T \frac{\partial \ell_t(\theta)}{\partial w}$$

→ By the chain rule of derivatives we know that $\frac{\partial \ell_t(\theta)}{\partial w}$ is obtained by summing gradients along all paths from $\ell_t(\theta)$ to w

↓
What are the paths connecting $\ell_t(\theta)$ to w??

Consider $\ell_4(\theta)$



So we have an ordered network. In this each state is computed one at a time. (first S_0 , then S_1 & so on)

$$\frac{\partial L_4(\theta)}{\partial w} = \frac{\partial L_4(\theta)}{\partial s_4} \frac{\partial s_4}{\partial w}$$

\downarrow This is direct Let's compute this

$$s_4 = \sigma(w s_3 + b)$$

$$\frac{\partial s_4}{\partial w} = s_3$$

\hookrightarrow It is again depend on w .

~~* So, that's the problem with an ordered network. We cannot compute the gradient of s_4 w.r.t w assuming s_3 is a constant. s_3 is not a constant again it's a function of w & w is the parameter w.r.t computing derivative.~~

~~* In such network the total derivative $\frac{\partial s_4}{\partial w}$ has two parts.~~

Explicit : $\frac{\partial^+ s_4}{\partial w}$ treating all other inputs as constant

Implicit : Summing over all indirect paths from s_4 to w

Derivation

$$\frac{\partial s_4}{\partial w} = \underbrace{\frac{\partial^+ s_4}{\partial w}}_{\text{Explicit}} + \underbrace{\frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial w}}_{\text{Implicit}}$$

$$= \frac{\partial^+ s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \left[\frac{\partial^+ s_3}{\partial w} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w} \right]$$

$$= \cancel{\frac{\partial^+ s_4}{\partial w}} + \cancel{\frac{\partial s_4}{\partial s_3}} \cancel{\frac{\partial^+ s_3}{\partial w}} + \cancel{\frac{\partial s_4}{\partial s_3}} \cancel{\frac{\partial s_3}{\partial s_2}} \left[\cancel{\frac{\partial^+ s_2}{\partial w}} + \cancel{\frac{\partial s_2}{\partial s_1}} \cancel{\frac{\partial^+ s_1}{\partial w}} \right]$$

$$= \frac{\partial^+ s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial^+ s_3}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial s_2} \left[\frac{\partial^+ s_2}{\partial w} + \frac{\partial s_2}{\partial s_1} \frac{\partial^+ s_1}{\partial w} \right]$$

$$= \frac{\partial s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial s_1} \left[\frac{\partial s_1}{\partial w} \right]$$

For simplicity,

$$\frac{\partial S_A}{\partial w} = \frac{\partial S_A}{\partial S_4} \frac{\partial^+ S_4}{\partial w} + \frac{\partial S_A}{\partial S_3} \frac{\partial^+ S_3}{\partial w} + \frac{\partial S_A}{\partial S_2} \frac{\partial^+ S_2}{\partial w} + \frac{\partial S_A}{\partial S_1} \frac{\partial^+ S_1}{\partial w}$$

$$= \sum_{k=1}^4 \frac{\partial S_A}{\partial S_k} \frac{\partial^+ S_k}{\partial w} //$$

Finally we have,

$$\frac{\partial L_4(\theta)}{\partial w} = \frac{\partial L_4}{\partial s_4} \frac{\partial s_4}{\partial w}$$

$$\frac{\partial S_4}{\partial w} = \sum_{k=1}^4 \frac{\partial S_4}{\partial s_k} \frac{\partial^+ s_k}{\partial w}$$

In General,

$$\frac{\partial \Delta_k(\theta)}{\partial w} = \frac{\partial \Delta_k(\theta)}{\partial s_k} \sum_{k=1}^t \frac{\partial s_k}{\partial w}$$

→ This Algorithm is called Back propagation through Time (BPTT) as we back propagate over all these steps -- Previous time steps.

The Problem of Exploding & Vanishing

$$\frac{\partial S}{\partial t} = \frac{\partial S_t}{\partial S_{t-1}} \frac{\partial S_{t-1}}{\partial S_{t-2}} \dots \frac{\partial S_{k+1}}{\partial S_k}$$

$$= \prod_{j=k}^{t-1} \frac{\partial S_j + 1}{\partial S_j} \quad (\text{or}) \quad \frac{\partial S_j}{\partial S_j - 1}$$

↳ we are interested in this term.

$$a_j = Ws_{j-1} + b \quad \Rightarrow \text{Pre-activation}$$

$s_j = \sigma(a_j) \Rightarrow$ Hidden representation

$$\frac{\partial s_j}{\partial s_{j-1}} = \frac{\partial s_j}{\partial a_j} \frac{\partial a_j}{\partial s_{j-1}}$$

$$a_j = [a_{j1}, a_{j2}, \dots, a_{jd}]^T$$

$$s_j = [\sigma(a_{j1}), \sigma(a_{j2}), \dots, \sigma(a_{jd})]$$

$$\frac{\partial s_j}{\partial a_j} = \begin{bmatrix} \frac{\partial s_{j1}}{\partial a_{j1}} & \frac{\partial s_{j2}}{\partial a_{j1}} & \dots \\ \vdots & \ddots & \dots \end{bmatrix} = 0 \text{ because it doesn't depend on that}$$

$$= \begin{bmatrix} \sigma'(a_{j1}) & 0 & 0 & 0 \\ 0 & \sigma'(a_{j2}) & 0 & 0 \\ 0 & 0 & \ddots & \\ 0 & 0 & \dots & \sigma'(a_{jd}) \end{bmatrix}$$

$$= \text{diag}(\sigma'(a_j))$$

$$\frac{\partial s_j}{\partial s_{j-1}} = \frac{\partial s_j}{\partial a_j} \frac{\partial a_j}{\partial s_{j-1}}$$

$$\downarrow = \text{diag}(\sigma'(a_j))W$$

For some reason we are interested in the Magnitude.

$$\left\| \frac{\partial s_j}{\partial s_{j-1}} \right\| = \left\| \text{diag}(\sigma'(a_j))W \right\| \leq \left\| \text{diag}(\sigma'(a_j)) \right\| \|W\|$$

$\because \sigma'(a_j)$ is bounded function because we are using sigmoid tanh, $\sigma'(a_j)$ is bounded.

If σ is logistic function then the bound is,

$$\sigma'(a_j) \leq \frac{1}{4} = r \quad [\text{If } \sigma \text{ is logistic}]$$

Max derivative we can get from the curve - 

For Tanh Function, $\leq 1 = r$ [If σ is tanh]

$$\leq \underbrace{\|\text{diag}(\sigma'(a_j))\|_1 \|w\|_1}_{\text{So this quantity is bounded and we call it as } 'r'}$$

weight matrix is also bounded because they are the real weights.

$$\leq r\lambda$$

$$\begin{aligned}\left\| \frac{\partial s_k}{\partial z_k} \right\| &= \left\| \begin{pmatrix} t \\ \prod_{j=k+1}^t r_j \end{pmatrix} \right\| \\ &\leq \prod_{j=k+1}^t r_j \\ &\leq (r\lambda)^{t-k}\end{aligned}$$

→ If $r\lambda < 1$ the gradient will vanish

→ If $r\lambda > 1$ the gradient could explode

$$\therefore \frac{\partial L_k(\theta)}{\partial w} = \frac{\partial L_k(\theta)}{\partial z_k} \sum_{k=1}^t \boxed{\frac{\partial s_k}{\partial z_k}} \frac{\partial^* z_k}{\partial w}$$

This vanishes if this vanishes the entire gradient could vanish.

→ And if gradients are vanished we have no updates we stuck where we are.

→ If the gradient explodes · If we have a very large gradient its gone way far away from where you are in the wB plane.

Because your update is $w = w - n \underline{\nabla w}$

→ this value is large so we are moving somewhere far.

→ to avoid this remember the $(t-k)$ term and the problem appears when $(t-k)$ is rather t is close to T so k is closed to 1.

Solution: Don't backpropagate through all the time steps. Use an approximation that if you are at time step n we just look at $(n-k)$ time steps and not all the way back.
↳ Exploding & Vanishing

→ Other thing to avoid Exploding

↓
In Gradient Descent we are always interested in the direction. So, we can just normalize it.

$$\frac{\nabla w}{\|\nabla w\|} \leq k \quad \begin{cases} \text{clip the Gradients} \end{cases}$$

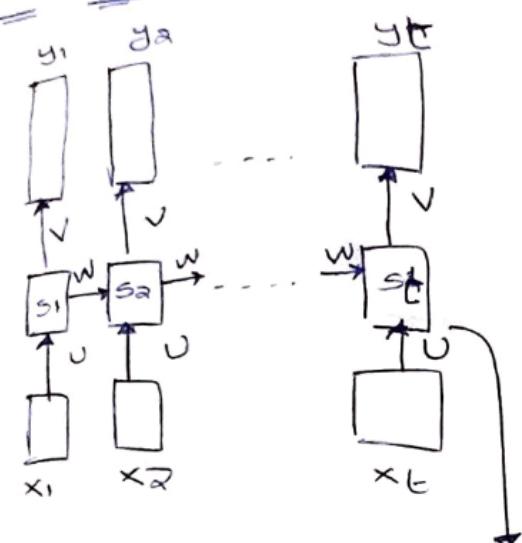
↳ Just a intuition for the Magnitude.

Some Groovy Details (13.5) → Refer once.

LSTM & GRU

Selective Read, Selective Write, Selective Forget - The WhiteBoard

Analogy



⇒ The state (s_t) of a RNN records information from all previous time steps

⇒ At each new time step the old information is accumulated by the current input.

But now the issue is that this state is going to be finite size dimension. Now as we keep writing information to that cell we are morphing the information that you had written in the earlier time steps.

↳ So, the problem is we want the information from all previous steps but the other hand we just have a finite amount of memory to deal with so, it is bound to leak over hidden state. So, it is completely impossible for the information to get morphed so much that completely impossible to say what was the original contribution at time step 1.

↳ So this is the problem with the **RE-RNN** (Recurrent Neural Network).

↳ Similar problem occurs when the information flows backwards (backpropagation).

↳ It is very hard to assign the responsibility of the error caused at the time step t to arbitrary time steps before it very far away time steps it is very hard so this is the vanishing gradient problem.

→ So both during the forward propagation information vanishes,
Even during the backward propagation information vanishes.

Analogy

Fixed Size White Board

→ we follow the following strategy at each time step

- ① Selectively write on Board
- ② Selectively read the already written Content
- ③ Selectively forget Some Content.

Selective write

→ There may be many steps in the derivation but we may just skip few which are not very imp.

Ex: $a=1 \quad c=5$

$$ac \Rightarrow \quad \Rightarrow \quad ac = 5 \quad \checkmark$$

we skip $a=1 \quad c=5$ $ac = 1 \times 5 \quad \checkmark$ These
 $= 5$ Steps

↓
Because we don't write everything for the reason
we run out of the memory.

Selective read

 = = =

→ While writing one step we typically read some of the previous steps we have already written & then what to write next.

Ex: Step ② we computed bd

Step ③ $bd + a$

→ Information of Step ② is imp

Vanishes
in Vanished

Selective Read

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

=

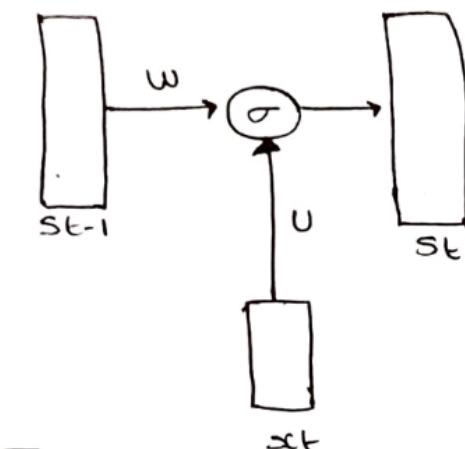
=

=

=

=

Selective Write



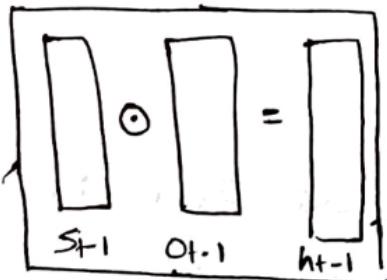
\Rightarrow So now instead of passing s_{t-1} as it is to s_t we want to pass only some portion of it to next state.

\Rightarrow So whenever we compute s_t we don't want to write the whole of s_{t-1} , just want to use selective positions.

\Rightarrow So what we do is we introduce something known as gate. (o_{t-1})

↓
which decides what fraction of each element of s_{t-1} should be passed to the next state.

\Rightarrow So we take original state s_{t-1} and do an element wise product with gate o_{t-1} (output gate) and then write the product to new vector which is h_{t-1} .



$\Rightarrow o_{t-1}$ is restricted by $0 \leq o_{t-1} \leq 1$

\Rightarrow But how do we compute o_{t-1} ? How does the RNN know what fraction of the state to pass on??



If we want to learn something we always introduce a parametric form of that quantity and then learn the parameters of that function.

→ RNN has to learn o_{t-1} along with other parameters (W, U, V)

↳ we compute o_{t-1} and h_{t-1} as

$$o_{t-1} = \sigma \left(\underbrace{W_o h_{t-2}}_{\text{Output at } t-2} + \underbrace{U_o x_{t-1}}_{\text{Input at timestep } t-1} + b_o \right)$$

↳ so we use Sigmoid, because we want the fraction to be bw 0 to 1

$$h_{t-1} = o_{t-1} \odot \sigma (s_{t-1})$$

o_t is called the Output gate as it decides how much to pass (write) to the next time step.

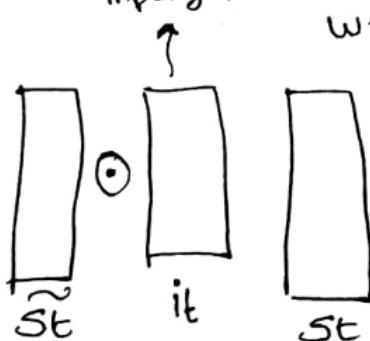
Selective Read

↳ we will now use h_{t-1} to compute the new state at the next time step

↳ we will also use x_t which is the new input at time step t .

$$\bar{s}_t = \sigma (W_h h_{t-1} + U_x x_t + b)$$

↓
This Eq form similar to RNN instead s_{t-1} we are using h_{t-1} because h_{t-1} have only selective written values.



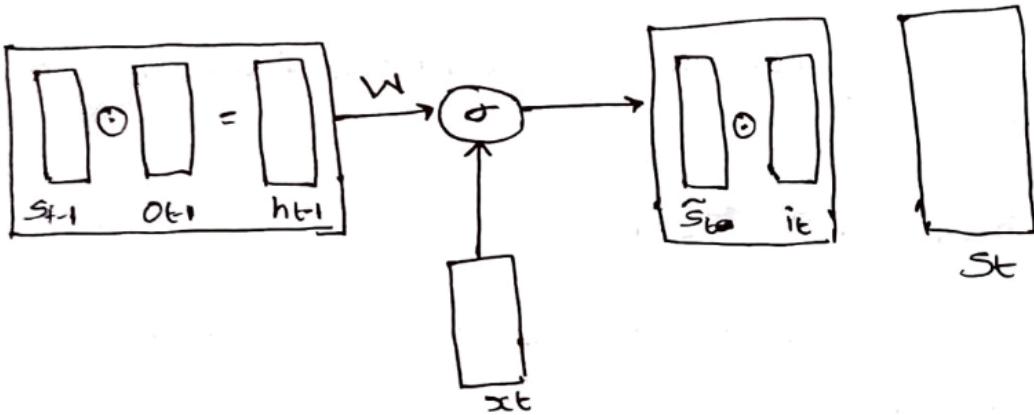
} we again need to selectively pass it or to. so again we introduce the gate if is called the Read gate.

↳ You don't want to pass it on as it is to s_t .

$$i_t = \sigma(W_{ht-1} + U_i x_t + b_i)$$

$i_t \odot \tilde{s}_t \Rightarrow$ Selective read state information.

So far we have,



Previous State : s_{t-1}

Output gate : $o_{t-1} = \sigma(W_{oo} h_{t-1} + U_{oo} x_t + b_o)$

Selective Write : $h_{t-1} = o_{t-1} \odot \sigma(s_{t-1})$

Current / Temporary State : $\tilde{s}_t = \sigma(W_{ht-1} + U_{ht} x_t + b_h)$

Input gate : $i_t = \sigma(W_{hi-1} + U_i x_t + b_i)$

Selective Read : $i_t \odot \tilde{s}_t$

Selective Forget
= - =

↳ How do we combine s_{t-1} & \tilde{s}_t to get new state.

Simple Way

$$\hookrightarrow s_t = \underbrace{s_{t-1} + i_t \odot \tilde{s}_t}_{\downarrow}$$

We should
forget something do not
use directly

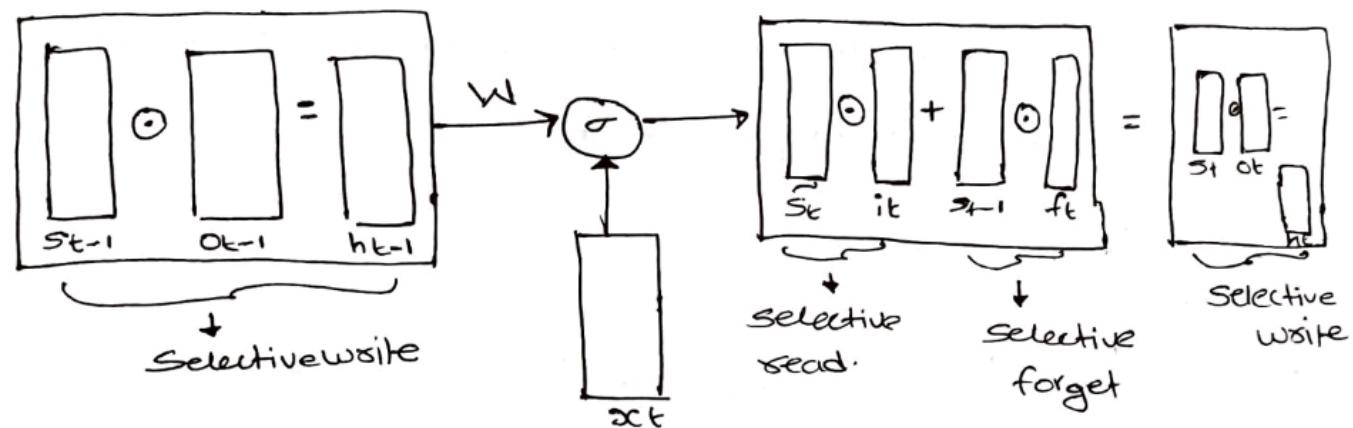
$$f_t = \sigma(w_f h_{t-1} + u_f x_t + b_f)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

Gates

$$\{ \begin{aligned} o_t &= \sigma(w_o h_{t-1} + u_o x_t + b_o) \\ i_t &= \sigma(w_i h_{t-1} + u_i x_t + b_i) \\ f_t &= \sigma(w_f h_{t-1} + u_f x_t + b_f) \end{aligned}$$

$$\} \quad \left. \begin{aligned} \tilde{s}_t &= \sigma(w_h h_{t-1} + u_x x_t + b) \\ s_t &= f_t \odot s_{t-1} + i_t \odot \tilde{s}_t \\ h_t &= o_t \odot \sigma(s_t) \end{aligned} \right\}$$



$s_t \rightarrow C_t$ ↴ current state

→ LSTM has many variants which include different no of gates and also different arrangement of gates.

→ This is the most popular variant of LSTM.

Another popular variant in LSTM is Gated Recurrent Units (GRU).

Gates :

$$o_t = \sigma(w_o s_{t-1} + u_o x_t + b_o)$$

$$i_t = \sigma(w_i s_{t-1} + u_i x_t + b_i)$$

States : $\tilde{s}_t = \sigma(w(o_t \odot s_{t-1}) + u_x x_t + b)$

$$s_t = (1 - i_t) \odot s_{t-1} + i_t \odot \tilde{s}_t$$

- No Explicit forget gate (Input gates are tied)
- The gates depend directly on s_{t-1} & not intermediate h_{t-1}

How LSTM avoid the problem of Vanishing Gradients??

- ↪ In Rnn due to the Recurrent network Structure there is the Vanishing Gradient Problem.
- ↪ Even in LSTM we have Recurrent Structure. But Why it doesn't vanishing Gradient??



Intuition

- ↪ During forward propagation the gates control the flow of information.
- ↪ They prevent any irrelevant information from being written to the state.
- ↪ Similarly during back propagation they control the flow of gradients

$$s_t = f_t \odot s_{t-1} + i_t \odot h_t$$

↑
If this vanishes in the worst case. We still have the Gradient f_t

↓
This is good Why is this OK Because f_t decides how much flew in the forward direction.

⇒ Refer the Video (14.3)

↓
Part-1

⇒ The key difference from Vanilla RNN is that flow of information and gradients is controlled by the gates which ensures that the gradient vanish only when they should.

Illustrative Proof ⇒ Refer the video (14.3) ⇒ Part 2