

If we have to convert DNN to CNN -

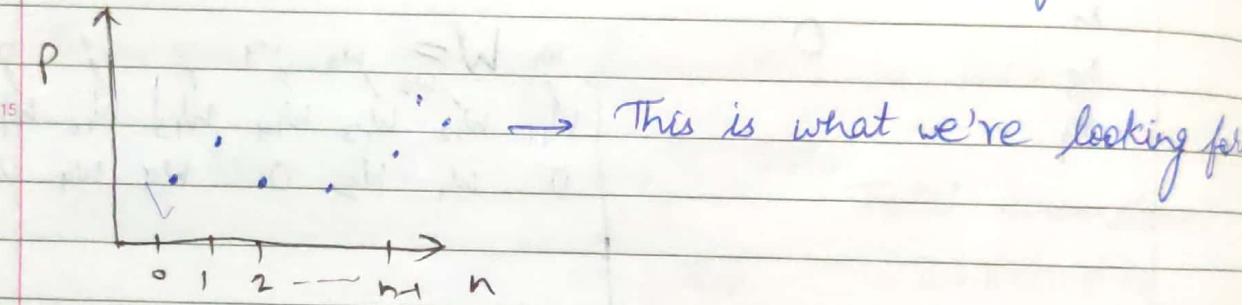
size of filter = I/p size

$$\begin{array}{c} \boxed{} \\ 3 \times 3 \end{array} * \begin{array}{c} \boxed{} \\ 3 \times 3 \end{array} = \boxed{}_{1 \times 1}$$

filter size Output

Softmax Activation fn

- For 2 classes, sigmoid is used. Softmax is used for multi-class classification.
- Softmax uses one-hot vector to give probability distribution.



$$\begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

1 is put where probab. is max.
Here, P_3 has the highest probab.

Probab. distri: is converted to One-Hot Vectors.

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g(Z^{[l]}).$$

→ ReLU is used most of the time
→ At last layer, sigmoid is used

$$Z^{[l]} \rightarrow (n^c \times 1) \quad n^c: H \text{ classes.}$$

Let $n^c = 3$,

$$z^{(l)} = \begin{bmatrix} z_1^{(l)} \\ z_2^{(l)} \\ z_3^{(l)} \end{bmatrix}, e^{z^{(l)}} = \begin{bmatrix} e^{z_1^{(l)}} \\ e^{z_2^{(l)}} \\ e^{z_3^{(l)}} \end{bmatrix}$$

$$a^{(l)} = \begin{bmatrix} a_1^{(l)} \\ a_2^{(l)} \\ a_3^{(l)} \end{bmatrix}, a_1^{(l)} = \frac{e^{z_1^{(l)}}}{\sum_{i=1}^3 e^{z_i^{(l)}}}, a_3^{(l)} = \frac{e^{z_3^{(l)}}}{\sum_{i=1}^3 e^{z_i^{(l)}}}$$

* Imagenet - 1 million images with 1000 classes.

* We're calculating prob. in order to back propagate. Else we could have just find the $z_a^{(l)}$ which is max & declare it as ans. inside one hot vector

15 G/P $\rightarrow 7 \times 7 \times 512 \rightarrow$ # repr of an image

Since, we're convolving so many times, we get diff. repr of base image

- With DNN, we get more # rich features for model to perform well

How powerful are DNNs as compared to CNNs?

has more no. of parameters (& don't take overfitting into account)

25 CNNs require order in images but DNNs don't

As DNNs are prone to overfitting, they perform poorly.

Date

15/10/19

Datasets

Camlin Page

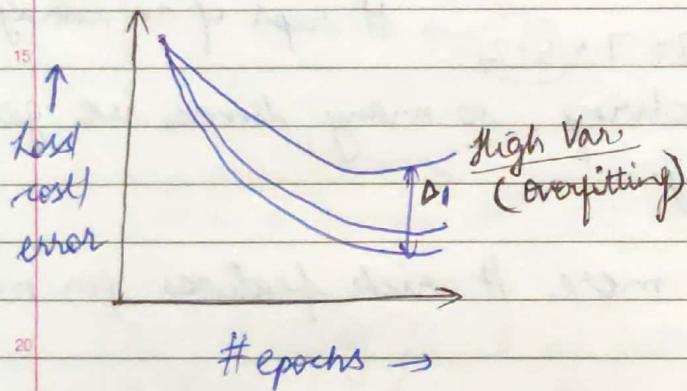
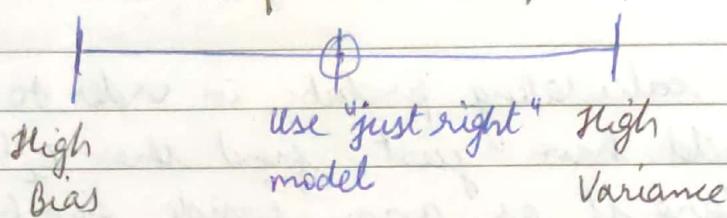
Date / /

- All 3 datasets - train / dev / test come from same distribution
- earlier we didn't have too many hyperparameters
no need as such to have validation set.

High bias - Underfitting

High Variance - Overfitting

Infinite # classifiers possible.



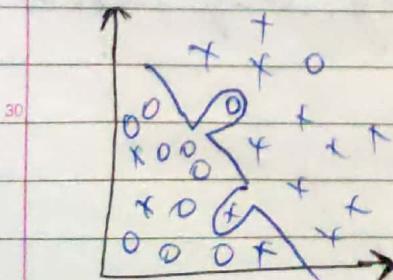
Underfitting - model performs poorly on both train & test set.
(High Bias)

e.g. 15y. on train set
16y. on test set

(High Bias & High Variance)

e.g. 15y. on train set [Performs bad on train set]
30y. on test set but even worse on test set

(High Variance) e.g. 1y. on train set
15y. on test set



Regularization

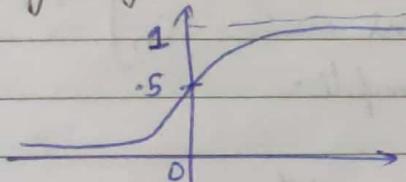
$$J(y, \hat{y}) = \left[-\frac{1}{m} \sum_{i=1}^m (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)) \right] + \frac{\lambda}{2m} \|w\|^2$$

\uparrow
parameters
 \uparrow
 $\frac{\lambda}{2m} \sum_{j=1}^p w_j^2$

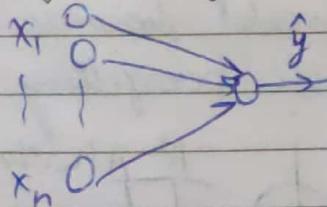
sq. of all parameters in our model

The model ultimately tries to min. cost fn. Now, addition of sq. of parameters will penalize if para. take very large values. (as goal is to min. J).

Apply regularization to avoid overfitting



* Logistic Regression



$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_{i=1}^n w_i^2 + \frac{\lambda}{2m} b^2$$

L₂ regularization: $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2$

$$= [W^T W]$$

Ignore as b is very small as compared to w

$$w_i = w_i - \alpha \frac{\partial J + J'}{\partial w} = w_i - \alpha \frac{\partial J}{\partial w} - \alpha \frac{\partial J'}{\partial w}$$

$$w_i = w_i - \alpha \frac{\partial J}{\partial w} - \alpha \frac{\lambda}{2m} (2w_i)$$

$$\alpha, \lambda > 0 \\ m > 0$$

$$w_i = \left(1 - \frac{\alpha \lambda}{m}\right) w_i - \alpha \frac{\partial J}{\partial w_i} \rightarrow$$

We're reducing w value even further

?? [Reg may contribute to overcome exploding grad. probm]

②: Regularization parameter

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(g_i, y_i) + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l+1)}} (w_{ij}^{(l)})^2$$

Very high bias - Regulari doesn't solve this.

[Weight-decay]

Reg turns complex classifier to simple classifier

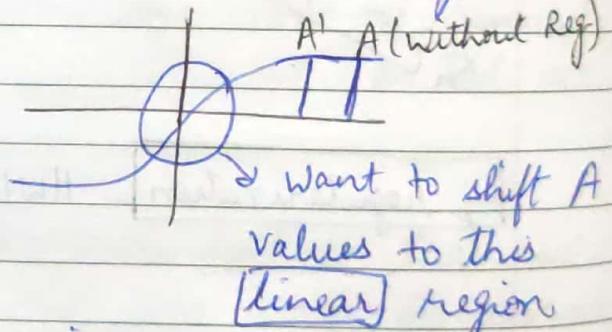
$C = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$ close to
Make non-linear to linear class
Regularization reduces $[w_0, w_1, \dots, w_n]$
to make classifier simpler

Regularization helps to decay weights.

(non-linear)

The more layers we add, our model becomes more complex & becomes more prone to overfitting.

- By reducing wts., we shift values of A by act. fn to linear region.
(less complex)



- Reg only penalizes high weights to reduce cost function

Q: How Reg deals with $[10, 0.01, 0.01, \dots, 0.01]$?
It tries to redistribute these weights so that all of them contribute almost equally.

λ value of $\left(\frac{\lambda}{2m} \sum w_i^2\right)$ is reduced

Effect of Reg: Tend to spread out weights. Doesn't make them zero.

Camlin	Page
Date	/ /

* L1 Regularization = $\frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^n \sum_{j=1}^m |w_{ij}^{(l)}|$

↓ may get sparse

5 [While L2 will give very low w] values but not sparse.

Dropout - Drop some of the neurons (keep prob = 0.5) makes n/w simpler randomly.

10 * While training, weights are distributed uniformly using drop out

15 * But at test time, we don't use dropout. As we should use all trained neurons in the n/w to contribute to the actual prediction.

* More data means less overfitting possible but still we can design very very huge NN that will overfit. Then we use dropout instead of using smaller NN.

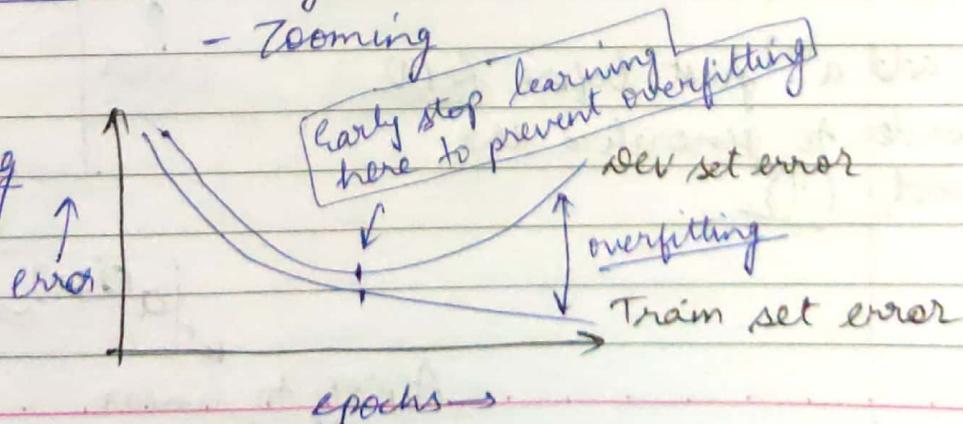
20 * Dropout is strongest form of Regularization that helps to curb over-fitting (Has same effect as L2-Reg)

Other regularization techniques

25 ① Data Augmentation - slight rotations

- zooming

30 ② Early stopping



Date
22/10/19

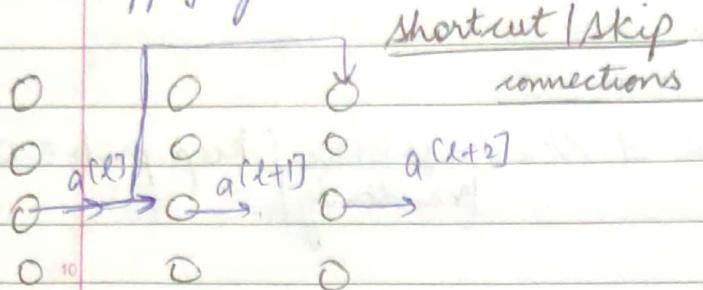
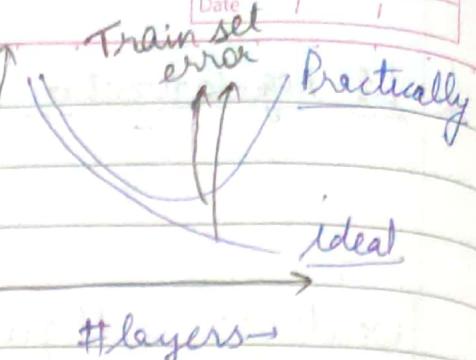
Camlin Page

Date / /

ResNets

152 layer architecture

O/P of l^{th} layer is directly fed to $(l+2)^{\text{th}}$ layer before applying ReLU.



- Very diff for optimizer to optimize weights (parameters) as our model becomes complex
- Learning becomes difficult

$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]}$$

$$a^{[l+1]} = g(z^{[l+1]})$$

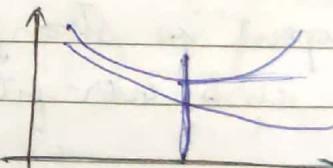
$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$$= g(W^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]}).$$

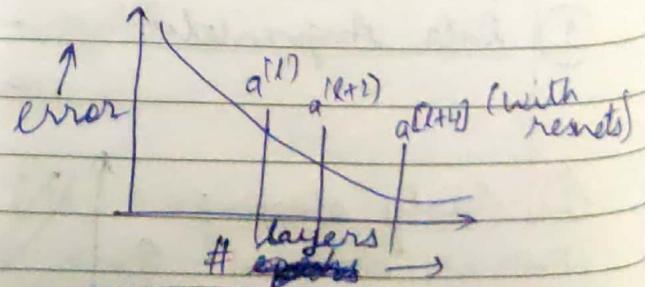
$$a^{[l+2]} \approx a^{[l]}$$

- ResNets help us overcome the above problem in plot
- It tries to make model learn at least identity fn so that at least our accuracy is maintained instead of decreasing.



Using regularization, we try to decay weights & biases so that output of $a^{l+2} \approx a^{l+1}$.

We add $a^{[l]}$ just before ReLU in order to generalize the model. (?)



$a^{[l+2]} \approx a^{[l]}$ so err reduces on increasing # layers

Approx. to linear identity fn.

Backprop in ResNets - I

- Helps us to overcome little bit of vanishing/exploding gradient probm as gradients are also back propagated from $n^{(l+2)}$ to $n^{(l)}$.

Why does a 1×1 convolution do?

(1×1) filter helps reducing the depth of inputs. (# channels)

$$\begin{matrix} \boxed{} \\ | \\ \boxed{} \end{matrix} * \boxed{} = \text{Output}$$

$6 \times 6 \times 32$ $1 \times 1 \times 32$ $6 \times 6 \times \# \text{filters}$

We can reduce depth from 32 to (#filters).

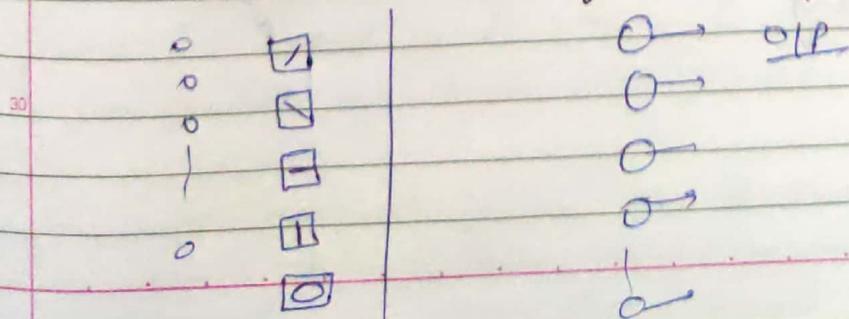
$$\begin{matrix} \boxed{} \\ | \\ \boxed{} \end{matrix} \xrightarrow[\substack{28 \times 28 \times 16 \\ 32, \text{ same}}]{\substack{\text{conv:} \\ (5 \times 5 \times 16)}} \begin{matrix} \text{O/P} \\ 28 \times 28 \times 32 \end{matrix}$$

$$\# \text{ multiplication operations} = \frac{(28 \times 28 \times 32) \times (5 \times 5 \times 16)}{\substack{\downarrow \\ \text{total O/P nos.}}} \times \frac{\substack{\uparrow \\ \text{computation done for each O/P no.}}}{\substack{\uparrow}}$$

≈ 10.0 M

In Inception nw, we've softmax Act. layers that have regularization effect (dropout) as we're capable of getting O/P in nw nw as well (dropout effect).

Transfer Learning



In case we've less data, we use transfer learning by freezing the weights of pretrained model.

We normalize inputs ($1/255$) so that large pixel values (that only repr color intensities) don't dominate.

Sequence Models

- Input size is not fixed in text/speech processing.
- for DL, we need mathematical input. It has to be a number (vector, matrix, etc). We can't input text here.

e.g. Harry Potter and Hermoine Granger invented a new spell

$$x \rightarrow x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(n)}$$

$$\hat{y} \rightarrow 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$$

$x^{(i), t}$ → 9/P, ith example's
 $y^{(i), t}$ → tth element

$T_{x^{(i)}} = 9$ (# inputs, words)

$T_{y^{(i)}}$

Output of tth element of
ith example.

Representing words :-

$\begin{bmatrix} & \\ & \\ & \end{bmatrix}$ $999 \rightarrow 0's$
1 → 1's

One-hot encoding

* Disadv

- Highly sparse.
- When our vocabulary changes, our representation changes.

Named Entity Problem

Q Stmt - "Hermione gives a new spell."

How to input it to a NN & find that whether a word is name or not?

- CNN shares features & parameters. So, they don't overfit as much that of DNN.

→ Use one-hot vector after constructing a corpus.

word 1 → 0
2 0
3 0
| |
n → 0

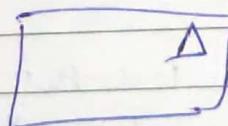
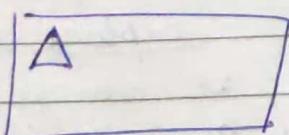
$0 \rightarrow 1$ → output layer should have as many neurons as in max-len sentence
 $0 \rightarrow 1$
 $0 \rightarrow 0$
|
 $0 \rightarrow 0$

Sentence with n words as IP.

* What if we've variable input length sentences?

We will O-pad for rest of the words.

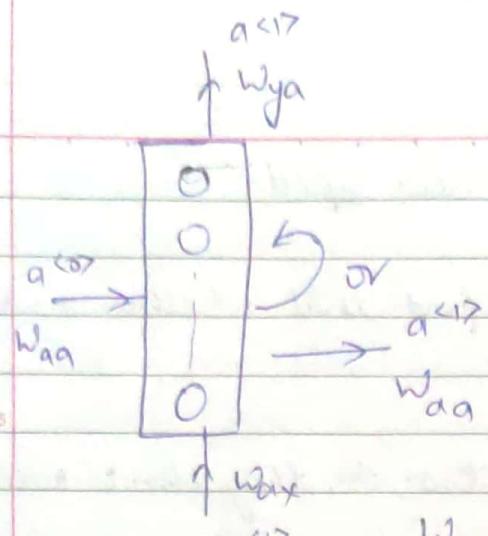
0 Take 1 sent. with max. words &
 0 take that many neurons.
 0 O-pad everyone else (not a very
 | good approach).



These 2 imgs if we input to DNN will not be able to classify but will do with CNN.

What we want is to learn to recognise the data irrespective of where it's placed.

We need same thing in RNN. If Harry Potter appears somewhere in corpus (our sentence) same parameters should be able to learn to recognize the [name].



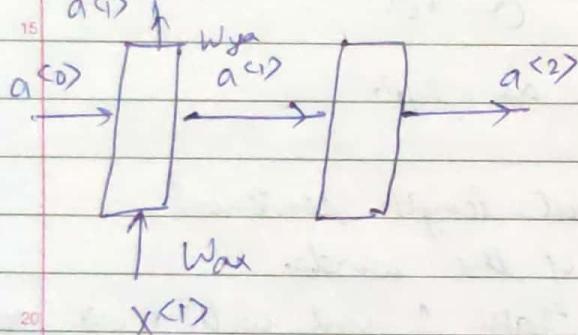
He said Teddy Roosevelt Name
(Teddy)

He said Teddy Bears.
Name X (Teddy)

$W_{xy} \rightarrow$ If p is of type x & output is of kind y.

Q. How to tell whether a sent. is positive or negative - like a movie rating → using Softmax fn

Now in this case, the architecture used for entity words will not work. RNNs have many variance among them.



This archi is Turing complete, but if you do not pass $a^{(1)}$ & pass $y^{(1)}$, it's not turing complete.

NLP & Word Embeddings

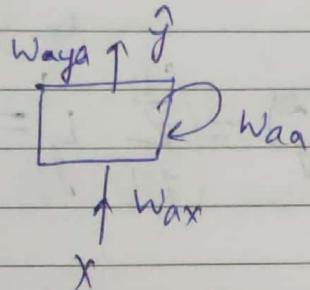
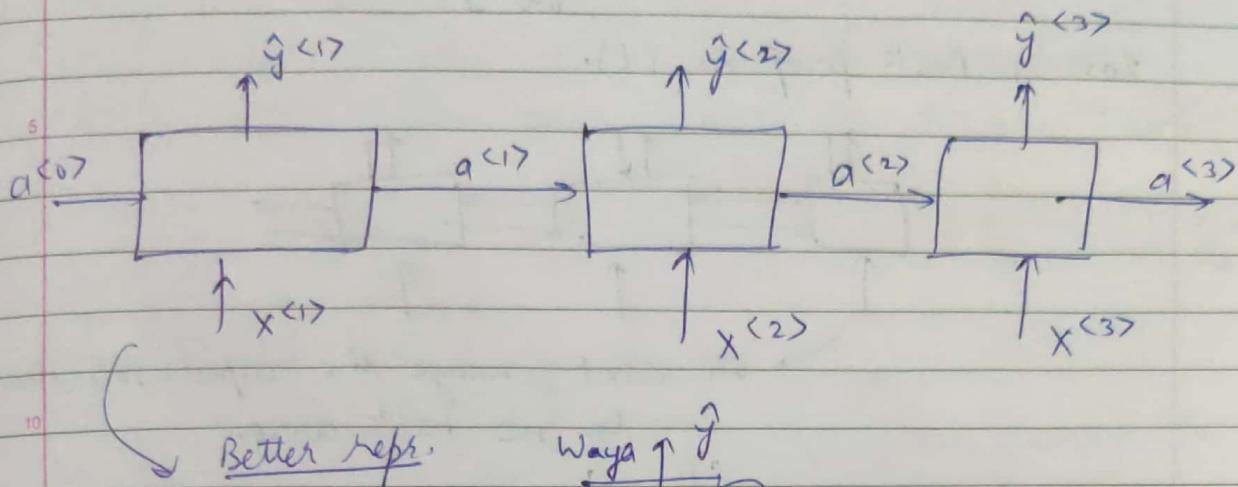
Bigger problem of One-hot Encoding than sparse - don't capture semantics.

Instead of using One-hot encoding, we can go with "Featureized representation": 'Word Embedding'.

Gender	Man, Woman	King, Queen
Age	↓ (denoting the correlation)	
Food		

Now each word is just repr by a vector of size [300]. Dot prod. will not be 0, it will be some value denoting correlation.

never
easily)
ly)
Problem of sparsity as well as semantics both are solved if we use 'Featurized representation'.



$$x \rightarrow (300 \times 1) \text{ dim.}$$

W_{xa} : Outputs a & inputs x .

equations: $z^{(t)} = b_a + W_{xa} x^{(t)} + W_{aa} a^{(t-1)}$

$$a^{(t)} = g_1(z^{(t)}) \leftarrow \text{tanh / ReLU}$$

$$z'^{(t)} = W_{ya} a^{(t)} + b_y$$

$$\hat{y}^{(t)} = g_2(z'^{(t)}) \leftarrow \text{sigmoid}$$

- RNN's have an equivalent power to that of a Turing Machine

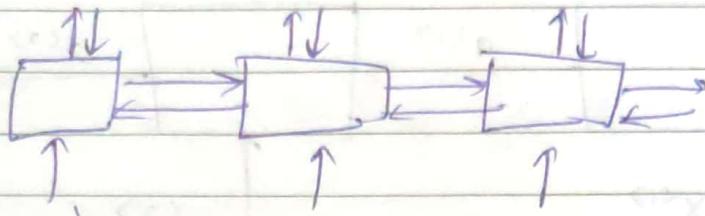
- Tanh can also be used unlike CNN that only uses ReLU.

Q. What problem is with $W_{ya}^{(2)}$ \rightarrow No use of parameter sharing
then, for some words, we want features to be shared instead of another parameter.

* Back Propagation

$$w = w - \alpha \frac{\delta L}{\delta w}$$

Loss in Back prop is L .

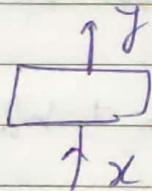


we can't change the input, so there can be no back arrow

$$L \leftarrow (\hat{y}^{(t)}, \hat{y}) = \left(\sum_{j=1}^t L_j \right) = (\hat{y}^{(t)} \log \hat{y}^{(t)} + (1 - \hat{y}^{(t)}) \log (1 - \hat{y}^{(t)}))$$

- Many-to-one practical application \rightarrow [sentiment Analysis]

One-to-one



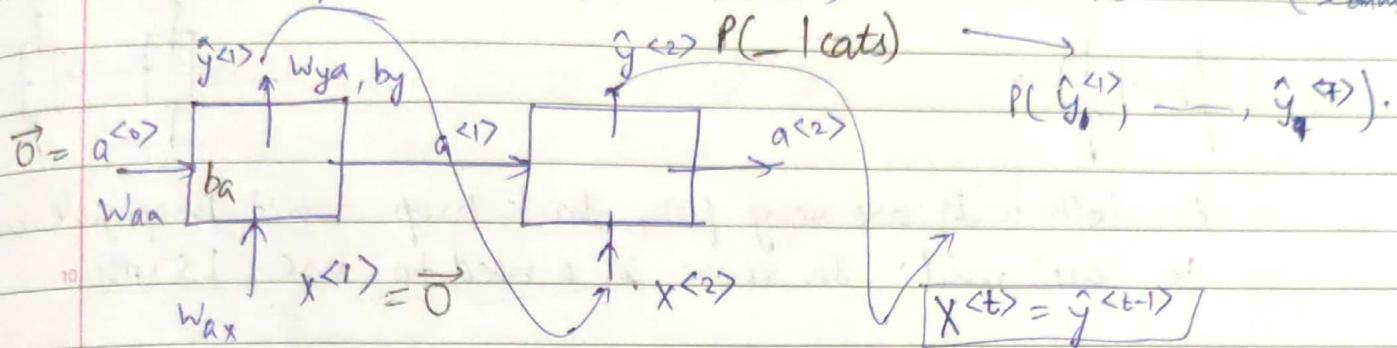
DNN's

(No timestamp means whatever we did in DNN's)

* Language Modelling - What is the probab. associated with the sentence or word?

$$P(\text{The apple & the pair salad}) \approx 3.3 \times 10^{-13}$$

$$P(\text{The apple & the pear sadad}) \approx 5.7 \times 10^{-10} \quad \checkmark \text{ (more common)}$$



$a^{<0>}$ is initialised to $\vec{0}$ (No input)

$x^{<1>} = \vec{0}$ [most frequently occurring words is not a good idea as all an generated sent. by RNNs will start with 'The'.

$$a^{<t>} = g_1(w_{ax} x^{<t>} + w_{aa} a^{<t-1>} + b_a) \quad \xrightarrow{\text{tanh}}$$

$$\hat{y}^{<t>} = g_2(w_{ya} a^{<t>} + b_y) \quad \xrightarrow{\text{softmax}}$$

Idea: We select the 1st word, based on the probab of words in the corpus. eg- $P(\text{the}) = 0.07$, there are 70% chances of choosing 'the' out of all possible words/choice

$$P(\text{cats}) = 0.0004$$

$$P(\hat{y}^{<2>} | \text{cats})$$

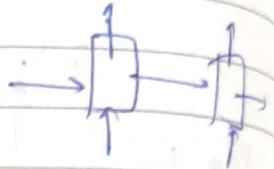
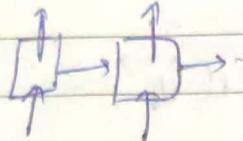
$$P(\hat{y}^{<3>} | \text{cats average})$$

- <EOS>: end of sent. marker
- To deal with <UNK>, we can even have a model ~~is that~~ at character level as no chr(ASCII) can be unknown

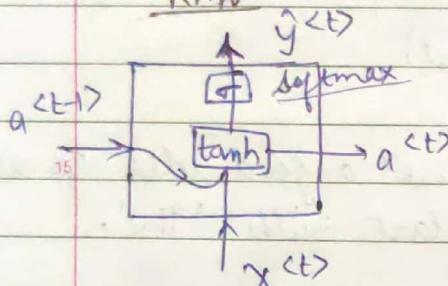
LSTMs

RNNs suffer from Vanishing Gr. prob. much more than DNNs or CNNs.

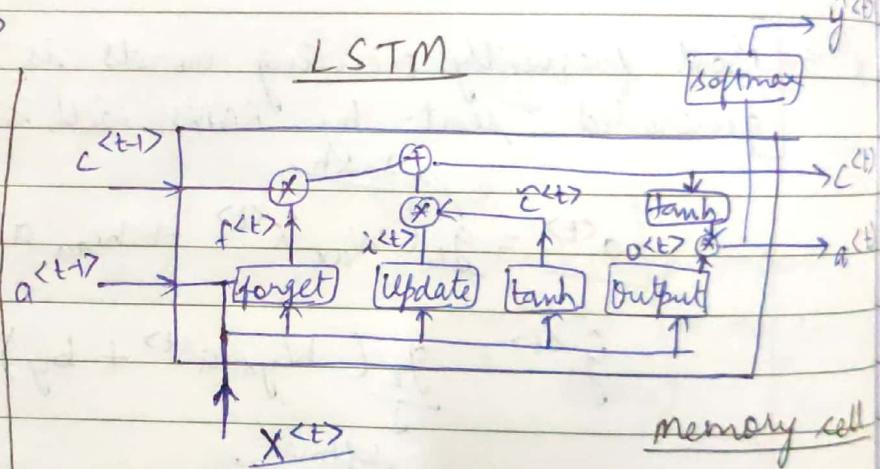
The cat, which already ate fruits like apple, pear, etc., was full.

RNNs

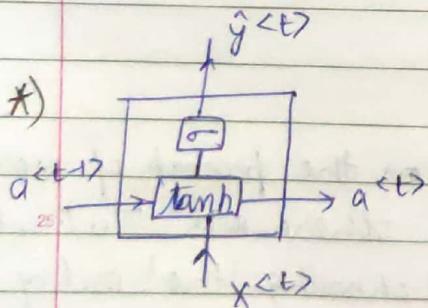
As both units are very far, back prop. won't propagate to 'cat' unit. So, there is a need to use LSTMs.

RNN

Doesn't capture long term dependencies.

LSTM

memory cell



(forget gate)

$$f_f = \sigma(W_f [a^{t-1}, x^t] + b_f).$$

$$i_u = \sigma(W_u [a^{t-1}, x^t] + b_u).$$

$$a^t = g_a (W_{ax} x^t + W_{aa} a^{t-1} + b_a)$$

$$o_o = \sigma(W_o [a^{t-1}, x^t] + b_o).$$

$$\hat{y}^t = g_y (W_{ya} a^t + b_y).$$

$$c^t = i_u * z^t + f_f * c^{t-1}.$$

update a^{t-1}
& change it to
 a^t .

$$z^t = \tanh (W_c [a^{t-1}, x^t] + b_c).$$

$$a^t = o_o * z^t.$$

In our terminology, $w_f = [w_{fa}, w_{fx}]$. Composed of 2 wts.

(Andrew Ng combined both in 1 slide).

$c^{<t>}$: cell state (for memorizing)

$a^{<t>}$: Activation of previous unit

Before LSTMs, people tried skip connections to maintain long term dep. But there were lots of extra conn. resulting in very high # parameters. Hence, LSTMs were invented.

GRUs (much simpler than LSTMs).

$$\tilde{c}^{<t>} = \tanh(W_c [r_y * c^{<t-1>}, x^{<t>}] + b_c).$$

$$r_u = \sigma(W_u [c^{<t-1>}, x^{<t>}] + b_u).$$

(Relevance gate)

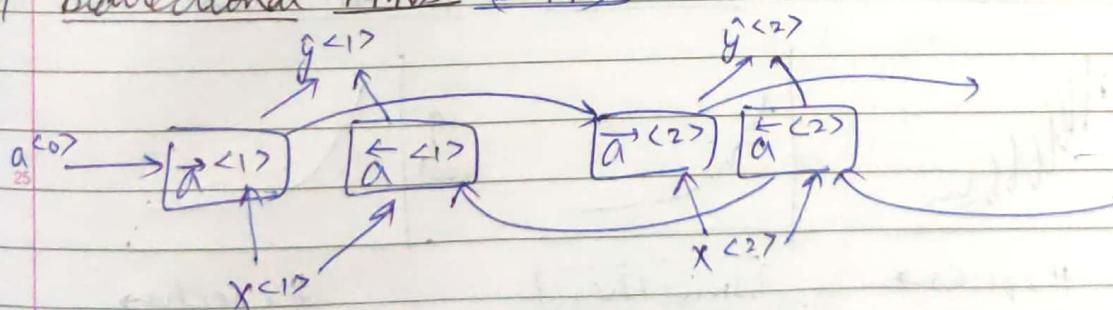
$$r_y = \sigma(W_y [c^{<t-1>}, x^{<t>}] + b_y).$$

$$c^{<t>} = r_u * \tilde{c}^{<t>} + (1 - r_u) * c^{<t-1>}.$$

$$a^{<t>} = c^{<t>}.$$

{ has only 2 gates - update and Relevance .

Bidirectional RNNs (BRNN)

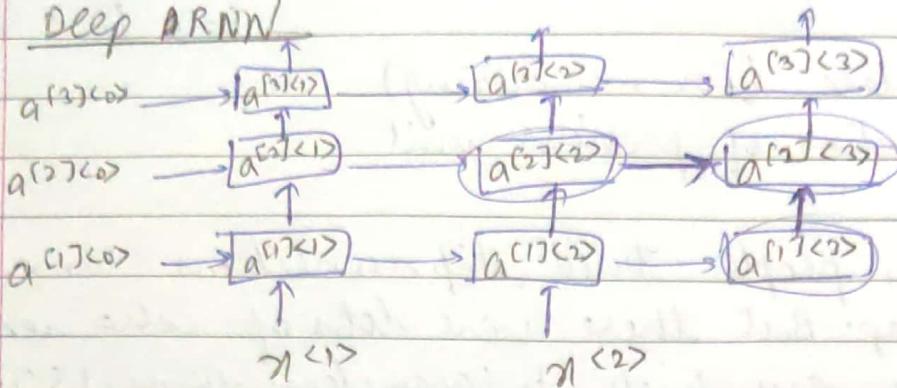


$$\hat{y}^{<t>} = g(w_y a^{<t>} + b_y) \quad - \text{one dir}$$

$$= g(\underbrace{w_y \vec{a}^{<t>}}_{\text{fwd}} + \underbrace{w_y \vec{a}^{<t>}}_{\text{Bwd}} + b_y)$$

$a^{(1)<2>} \rightarrow$ means 1st layer, 2nd unit | timestamp.

* Deep ARNN



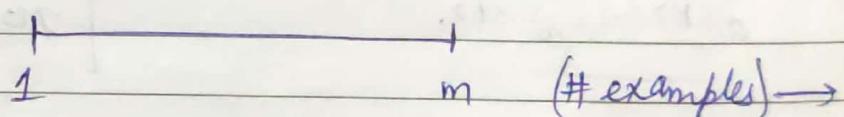
$$a^{(2)<3>} = g(W_a^{(2)} [a^{(2)<2>} \ a^{(1)<3>}] + b_a^{(2)}).$$

Turing complete archi

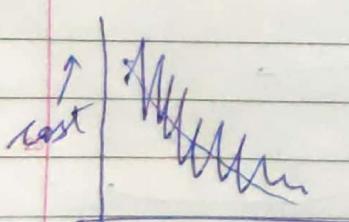
Date 15
21/11/19

SGD (Mini-batch Gr. Descent)

We can consider following # examples in 1 Gr. desc update -

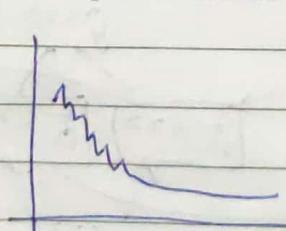


for 1 ex.



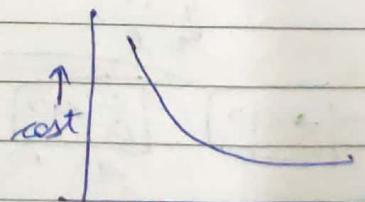
epochs →
(Stochastic GD)

$K < m$



epochs →
(Mini-Batch GD)

for m exs.



epochs →
(Batch GD)

$x^{(1)} \rightarrow$ Mini-batch 1.
 $x^{(1)} \rightarrow$ 1st tr. ex.

$x^{(2)} \rightarrow$ Mini-batch 2.
 $x^{(2)} \rightarrow$ 2nd tr. ex.

Stochastic Gr. Descent → When batch size = 1. We train on just 1 ex. before updating the parameters (ext. chosen randomly) epoch-wise

With Batch-GrD, learning is very slow

with SGD, learning is very fast

while Mini-Batch GrD takes good of both the things.

We require:

- 1) Learning should be moderately fast
- 2) Should reach to min as fast as possible.

Moving Avg. : Taking avg of obs. within a window.

Assign higher wts. to nearer obs. & so on

$$\begin{cases} V_0 = 0 \\ V_1 = 0.9 V_0 + 0.1 \theta_1 \\ V_2 = 0.9 V_1 + 0.1 \theta_2 \\ \vdots \\ V_t = 0.9 V_{t-1} + 0.1 \theta_t \end{cases}$$

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

exponentially weighted average

$\beta=1$, just considering fast avg

$\beta=0$, just considering current value



	$\beta = 0.98$
	$\beta = 0.9$
	$\beta = 0.5$

#³⁰ momentum