

Sequence Models

Dr. Venkataramana Veeramsetty
NVIDIA DLI Instructor & Faculty@CAIDL
Assistant Professor
Dept. of EEE
SR Engineering College

il : dr.vvr.research@gmail.com Venkataramana Veeramsetty)

December 31, 2020





Problems dealing with sequential data



√ Speech Recognition:

- Both input "X" and output "Y" are in sequence
- Able to convert speech into text
- X is an audio clip and so that plays out over time and Y, the output, is a sequence of words.
- Music Generation:
 - Input "X" is an empty or single integer and output "Y" are in sequence
- √ Sentiment Classification:
 - Input "X" is sequential data like "Movie is not good" and output "Y" may be rating in terms of an integer value





Problems dealing with sequential data



- √ DNA Sequence Analysis:
 - Input "X" is DNA sequence, and output "Y" is also sub sequence in DNA
- √ Machine Translation:
 - Input "X" is a sequence of words in one language and Output "Y" is a sequence of words in target language
- √ Video activity recognition:
 - Input "X" is sequence of frames in video, output "Y" tells the activity
- √ Name entity recognition:
 - Input "X" is a sentence which is sequence of words, output "Y" consists names in that sentence





Representation of Sequence: Name en identification

X(1)	Ram $X^{1,1}$	is <i>X</i> ^{1,2}	а <i>X</i> ,13	$X^{1,4}$	boy <i>X</i> ^{1,5}
Y(1)	1 Y ^{1,1}	0 Y ^{1,2}	-	0 Y ^{1,4}	0 Y ^{1,5}

 $InX^{i,t}$, "i" represents sample number and "t" represents time stamp.





One-Hot Vector



D En

esearch@gmail.com

Let Assume a dictionary as follows

									iail
a	an	are	boy	girl	good	is	Ram	Sitha	Zero
1	2	3	4	5	6	7	8	9	<u> 16</u>

Now X(1) can be represented as

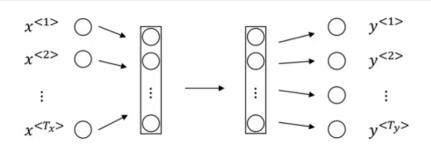
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$												
Ram	$X^{1,1}$	0	0	0	0	0	0	0	1	0	0	rams
is	$X^{1,2}$	0	0	0	0	0	0	1	0	0	0	setty
а	$X^{1,3}$	1	0	0	0	0	0	0	0	0	0	
good	$X^{1,4}$	0	0	0	0	0	1	0	0	0	0	
boy	$X^{1,5}$	0	0	0	1	0	0	0	0	0	0	

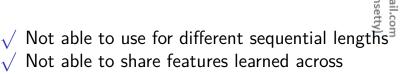




Can we ANN for Sequential Data **Problems**





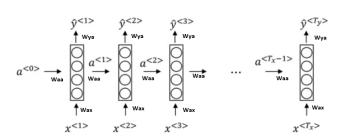


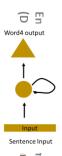
- different positions of text
- Network will be more complex with size of



Recurrent Neural Network

















$$a^{<0>} = 0$$
 $a^{} = g(Wax * X^{<1>} + Waa * a^{} + b_a)$ $y^{} = f(Wya * a^{} + b_v)$

Activation function to calculate $a^{<}t>$ is either Tanh or Relu. Where as for output it is depending on type of problem. For regression and binary classification can use sigmoid, for multiple classification can use softmax.









- Predict character "o" after feeding characters "h","e","I","I" sequentially
- one hot encoding dictionary

	h	е		0
h	1	0	0	0
е	0	1	0	0
-	0	0	1	0
0	0	0	0	1

Dr. Venkataramana Veeramsetty

Assume 3 hidden neurons in hidden layer, size of W_{ax} is 3X4 and W_{aa} is 3X3 and size of bias matrix $\frac{1}{5}$ 3X1





Dr. Venkataramana Veeramsetty)

Initialize Weights and Bias Parameters and activation states

$$W_{ax} = \begin{bmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ -0.1 & -0.2 & 0.3 & 0.1 \\ 0.1 & 0.2 & -0.3 & -0.4 \end{bmatrix}$$

$$W_{aa} = \begin{bmatrix} 0.1 \\ -0.1 \\ 0.1 \end{bmatrix}$$

$$\begin{bmatrix} b_a & b_y \\ 0.1 & 0.1 \\ 0.2 & 0.2 \\ 0.3 & 0.3 \\ 0.4 \end{bmatrix}$$



 $\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ -0.1 & -0.2 & 0.3 \end{bmatrix}$





Pass first input h=[1;0;0;0] and compute state of hidden layer using equation shown below.

$$a^{< t=1>} = g(Wax * X^{< t>} + Waa * a^{< t-1>} + b_a)$$

$$\tanh(\begin{bmatrix}0.1 & 0.2 & 0.3 & 0.4\\ -0.1 & -0.2 & 0.3 & 0.1\\ 0.1 & 0.2 & -0.3 & -0.4\end{bmatrix} * \begin{bmatrix}1\\0\\0\\0\end{bmatrix} + \begin{bmatrix}0.1 & 0.2 & 0.3\\ -0.1 & -0.2 & 0.3\\ 0.1 & 0.2 & -0.3\end{bmatrix} * \begin{bmatrix}0\\0\\0\\0\end{bmatrix} + \begin{bmatrix}0.1\\0.2\\0.3\end{bmatrix})$$

$$a(t=1) = \begin{bmatrix} 0.197 \\ 0.099 \\ 0.379 \end{bmatrix}$$

Email : dr.vyr.xesearch@gmail.com (Dr. VenkataEymana Veeramsetty)





$$y^{< t>} = f(Wya * a^{< t>} + b_y)$$

Compute output using
$$y^{< t>} = f(Wya*a^{< t>} + b_y)$$

$$\begin{cases} 0.1 & 0.2 & 0.3 \\ -0.1 & -0.2 & 0.3 \\ 0.1 & 0.2 & -0.3 \\ -0.1 & -0.2 & 0.3 \end{cases} * \begin{bmatrix} 0.197 \\ 0.099 \\ 0.379 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix})$$
 softmax(
$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ -0.1 & -0.2 & 0.3 \\ -0.1 & -0.2 & 0.3 \end{bmatrix}$$





Back Propagation Through Time



Email : dr.vvr./esea (Dr. Venkatarsman

$$L(Y^{P}, Y) = -Y log(Y^{P}) - (1 - Y) log(1 - Y^{P}) = \frac{T_{y}}{100} \left(-Y^{T} log(Y^{P,t}) - (1 - Y^{T}) log(1 - Y^{T}) \frac{1}{100} \right)$$

$$L(Y^{P}, Y) = \sum_{t=1}^{T_{y}} (-Y^{t} log(Y^{P,t}) - (1 - Y^{T}) log(1 - Y^{T}) \frac{1}{100} \left(-Y^{T} log(Y^{P,t}) - (1 - Y^{T}) log(1 - Y^{T}) \frac{1}{100} \right)$$

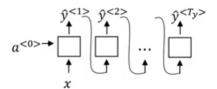




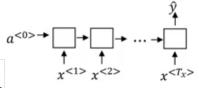
RNN Models



One To Many - Music Generation



Many To One - Sentiment Classification





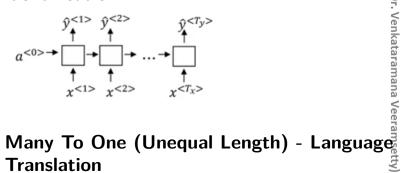




RNN Models

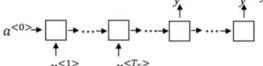


Many To Many (Equal Length) - Name entity identification Email : dr.vvr.research@gmail.com



Translation

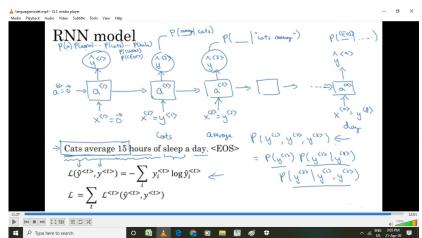






Language Modeling









Word Level Sequence Generation using Trained RNN



 $a^{<0>} \xrightarrow{g^{<1>}} y^{<2>} \xrightarrow{g^{<3>}} \cdots \xrightarrow{g^{<T_y>}} a^{<1>} \xrightarrow{g^{<T_y>}} y^{<1>} y^{<1>} y^{<2>} \cdots$

email: dr.vvr.research@gmail. Dr. Venkataramana Veeramse

 $\sqrt{\text{Vocabulary}} = [a, aaron, ..., zuzu, ... UNK, EOS]$

1/

 $y^{<1>}$ =np.random.choice(p(a),p(aaron),..,p(z $\frac{8}{2}$ z $\frac{1}{2}$),.

- $\sqrt{}$ Based on T_{ν} , sequence generation will stop
- \checkmark Can generate until "UNK" generated

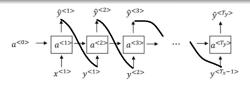




Character Level Sequence Generation using Trained RNN







$$\sqrt{\text{Vocabulary}=[a,b,..,z,-,*,...A,B,...Z]}$$

$$y^{<1>}$$
=np.random.choice(p(a),p(b),..,p(*),...,p

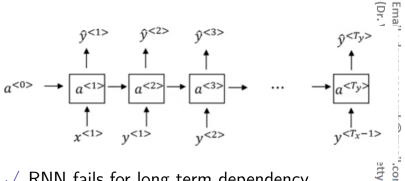
Limitations:

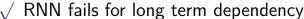
- √ Not recognize unknown word
- √ High complex
- √ More computation time
- √ High dependent previous time step.



Vanishing Gradient and Exploding Gradient







- The cat, which ate apple, banana,..., was full
- The cats, which ate apple, banana,..., were full





Vanishing Gradient and Exploding Gradient



main dr.vvr.research@gmail.con)r. Venkataraman Veeramsetty)

- √ Output at particular time stamp highly dependent
 on near time stamps and weakly effected by initial time stamps.

 Output at particular time stamp highly dependent
 on near time stamps.

 Output at particular time stamp highly dependent
 on near time stamps.

 Output at particular time stamp highly dependent
 on near time stamps.

 Output at particular time stamp highly dependent
 on near time stamps.

 Output at particular time stamp highly dependent
 on near time stamps.

 Output at particular time stamp highly dependent
 on near time stamps and weakly effected by initial time stamps.

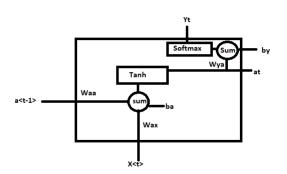
 Output at particular time stamp highly dependent
 output at particular time stamp hig
- √ Gradient of error will be have neglisible impact on initial weights (Gradient vanishing)
- √ Some times gradients may increase to high value (Exploding gradient), will lead high values of parameters.
- √ Exploding gradient can be overcome by using gradient clipping method





RNN- Graphical Representation





$$a_t = g(W_{ax}x_t + W_{aa}a_{t-1} + b_a)$$

 $Y_t = f(W_{va}a_t + b_v)$



g: tanh/Relu :Sigmoid/Softmax



Email: dr.vvr.research@gmail.com (Dr. Venkataramana Veeramsett)

RNN:Limitations



- RNN will estimate word/character/number in Vernkataramana Vernkata
- √ Due to deep structure, RNN has vanishing
- gradient problem
- Due to deep structure, the value at particular time stamp may forget initial values in RNN





Gated Recurrent Unit (GRU)



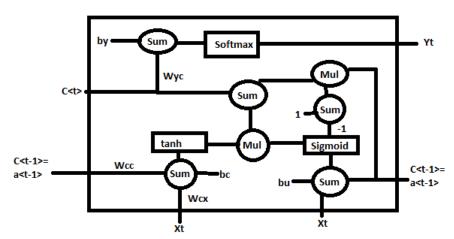
Email: dr.vvr.research@gmail.com Veeramsetty)

- √ GRU can be helpful for long term dependency and also overcome vanishing gradient problem
 - The cat, which already ate apple,..., was full
 - The cats, which already ate apple,..., were full















$$C_t = a_t$$

$$C_t^{hat} = tanh(W_{cc}C_{t-1} + W_{cx}X_t + bc)$$

$$\Gamma_u = sigmoid(W_{uc}C_{t-1} + W_{ux}X_t + bu)$$

$$C_t = \Gamma_u * C_t^{hat} + (1 - \Gamma_u) * C_{t-1}^{hat}$$

nsail sm

 C_t represents memory cell C_t^{hat} represents candidate cell u will be either 0 or 1



Full GRU



$$C_t^{hat} = tanh(W_{cc} * \Gamma_r * C_{t-1} + W_{cx}X_t + bc)$$
 (Fig. 1) $\Gamma_u = sigmoid(W_{uc}C_{t-1} + W_{ux}X_t + bu)$ (Fig. 2)

$$\Gamma_r = sigmoid(W_{rc}C_{t-1} + W_{ux}X_t + br)$$
 (16)



$$C_t = \Gamma_u * C_t^{hat} + (1 - \Gamma_u) * C_{t-1}^{hat}$$

Long short Term Memory (LSTM)



$$C_t^{hat} = tanh(W_{cc}*a_{t-1} + W_{cx}X_t + bc)$$
 (18) $\Gamma_u = sigmoid(W_{ua}a_{t-1} + W_{ux}X_t + bu)$ (19) $\Gamma_f = sigmoid(W_{fc}a_{t-1} + W_{fx}X_t + bf)$ (20) $\Gamma_o = sigmoid(W_{oc}a_{t-1} + W_{ox}X_t + bo)$ (21)

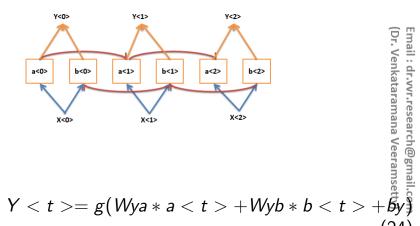


$$C_t = \Gamma_u * C_t^{hat} + \Gamma_f * C_{t-1}^{hat}$$



Bi-directional RNN





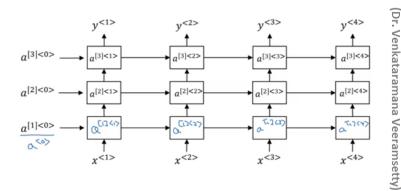
$$Y < t >= g(Wya * a < t > + Wyb * b < t > + by)$$
(24)



It can predict based on previous and upcoming

Deep-RNN





Email: dr.vvr.research @gmail.com

$$y < 1 > y < 2 > y < 3 > y < 4 > y < 5 > y < 4 > y < 6 > y < 7 > y < 7 > y < 7 > y < 8 > y < 8 > y < 8 > y < 8 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y < 9 > y$$





Word embedding



Let Assume a dictionary as follows

									ma Or.
6	3	an	Arjun	boy	girl	good	is	Ram	Zero ฐZoo
1	Ĺ	2	3	4	5	6	7	8	9 🔓 🗐 0
N	ow	X(r.resea raman					

												<u> </u>
Ram	$X^{1,1}$	0	0	0	0	0	0	0	1	0	0	ch@ Vee
is	$X^{1,2}$	0	0	0	0	0	0	1	0	0	0	gma
а	$X^{1,3}$	1	0	0	0	0	0	0	0	0	0	sett
Ram is a good	$X^{1,4}$	0	0	0	0	0	1	0	0	0	0	≤ ÿ
												-

• Ram is a good ? Laxman is a good





30 / 48

Indexing of words in embedding matrix

Х	LEARNING
_	

Word	Index
Α	O_1
An	O_2
Arjun	<i>O</i> ₃
Boy	O_4
Girl	O_5
Good	O_6
ls	O_7
Ram	<i>O</i> ₈
Zero	O_9
Zoo	O_{10}

Email: dr.vvr.research@gmail.com Dr. Venkataramana Veeramsetty)





Feature matrix



								(E
Fe	a	an	Arjun	boy	girl	good	is	Ram₹ Z ero
	O_1	O_2	O_3	O_4	O_5	O_6	O_7	$O_8 \stackrel{4}{=} O_9$
Age	0	0	0.93	0.95	0.92	0	0	0.96 🛊 🔞 0
Gen	0	0	0.93	0.95	0.92	0	0	0.96 2 50
Place	0	0	0	0	0	0.1	0	o o o
Verb	0.9	0.9	0	0	0	0.1	0.9	0 arc0
Qual	0	0	0	0	0	0.9	0	0 👸 👰 0

As both Ram and Arjun have similar feature matrix. Will get same one hot vector for both below words



- Ram is a good boy
- Laxman is a good boy



Analogies



	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

Email : dr.vvr.research@gmail.com eramsetty)

Man-Woman :: King-? $e_m an = [-1,0.01,0.01,0.09]$ $e_w oman = [1,0.02,0.02,0.01]$ $e_m an - e_w oman = [-2,0,0,0]$



Find e_w which suits in question mark in previous slide



$$arg_{max}^{w} = sim(e_{w}, e_{king} - e_{man} + e_{woman})$$

Cosine Similarity

$$sim(u,v) = \frac{u^T v}{|u||v|}$$

/enkataran /eeram:

Euclidean distance similarity

$$sim(u, v) = ||u - v||^2$$

Applications

- Man:Woman as Boy:Girl
- Ottawa:Canada as Delhi:India
- Rig-Rigger as Tall-Taller



Embedding matrix (E)



Email (Dr. Ve

							<u>n</u>
а	an	Arjun	boy	girl	good	is	Ram 🛣 Zero
O_1	O_2	O_3	O_4	O_5	O_6	O_7	$O_8 \stackrel{\triangleleft}{=} \stackrel{\triangleright}{=} O_9$
0	0	0.93	0.95	0.92	0	0	0.96 0
0	0	0.93	0.95	0.92	0	0	0.96 🖁 🖁 0
0	0	0	0	0	0.1	0	o √ee
0.9	0.9	0	0	0	0.1	0.9	0 🖁 🖁 0
0	0	0	0	0	0.9	0	o nset
	O ₁ 0 0 0 0 0.9	$\begin{array}{c cc} O_1 & O_2 \\ \hline 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.9 & 0.9 \\ \end{array}$	$\begin{array}{c ccc} O_1 & O_2 & O_3 \\ \hline 0 & 0 & 0.93 \\ 0 & 0 & 0.93 \\ 0 & 0 & 0 \\ 0.9 & 0.9 & 0 \\ \end{array}$	$\begin{array}{c cccc} O_1 & O_2 & O_3 & O_4 \\ \hline 0 & 0 & 0.93 & 0.95 \\ 0 & 0 & 0.93 & 0.95 \\ 0 & 0 & 0 & 0 \\ 0.9 & 0.9 & 0 & 0 \\ \end{array}$	$\begin{array}{c ccccc} O_1 & O_2 & O_3 & O_4 & O_5 \\ \hline 0 & 0 & 0.93 & 0.95 & 0.92 \\ 0 & 0 & 0.93 & 0.95 & 0.92 \\ 0 & 0 & 0 & 0 & 0 \\ 0.9 & 0.9 & 0 & 0 & 0 \\ \end{array}$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$

$$e_w = E.O_w$$

(29)





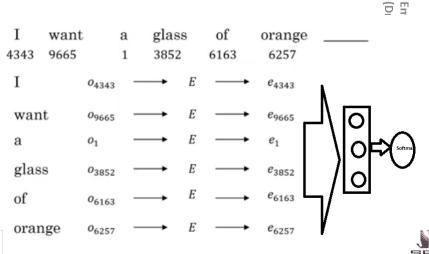
One-Hot Vector (O)



•	а	an	Arjun	boy	girl	good	is	Ram	Zero	⊋Zoo
	1	2	3	4	5	6	7	8	9	√en ₹0
	1	0	0	0	0	0	0	0	0	kata
	0	1	0	0	0	0	0	0	0	r.wr.resear
	0	0	1	0	0	0	0	0	0	
	0	0	0	1	0	0	0	0	0	ch@gMail.con Veeramsetty)
	0	0	0	0	1	0	0	0	0	gmai rams
	0	0	0	0	0	1	0	0	0	iil.com setty)
	0	0	0	0	0	0	1	0	0	0 3
	0	0	0	0	0	0	0	1	0	0
	0	0	0	0	0	0	0	0	1	0
	0	0	0	0	0	0	0	0	0	ENGINEERING COLLEGE

Natural Language Model







Skip-grams model



Ram is a good boy unlike rahul

• Choose context word and corresponding target word, And apply supervised learning, loss function categorical cross entropy

Context	Target
Ram	ls
Ram	good
Ram	boy

mail : dr.vvr.research@gmail.cc Dr. Venkataramana Veeramsett

Ram O8 E E e8 (size 5) Layer



Negative Sampling



• Similar to Skip-gram model but more efficient

Ram is a good boy unlike rahul

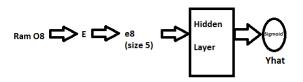
 Choose context word and corresponding target word, and set output as 1, similarly choose negative target words for given sampling word and set output as zero

X		Υ
Context	Target word	Output
Ram	ls	1
Ram	good	0
Ram	boy	0



• No. of negative sample generally between 5-10







Note:
Number of neurons in output layer depends on size of vocabulary, however will update K+1 neurons. K= is number of negative words. Negative examples will choose based on below equation



$$P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum Voc size f(w_i)^{\frac{3}{4}}}$$

GloVe - Global Vectors for word representation



Minimize

$$\sum_{i=1}^{Voc_Size}\sum_{j=1}^{Voc_Size}f(X_{ij})(\Theta_i^Te_j+b_i+b_j-log(X_{ij}))^2$$

 X_{ij} represent how many number of time j^{th} target comes with i^{th} context (X_{ij}) if X_{ij} =0

$$f(X_{ij}) = \text{if } X_{ij} = 0$$



Sentiment Classification



x

The dessert is excellent.

Service was quite slow.

Good for a quick meal, but nothing special.

Completely lacking in good taste, good service, and good ambience.

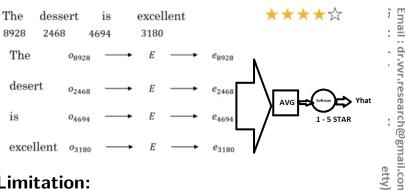






Sentiment Classification Model





Limitation:

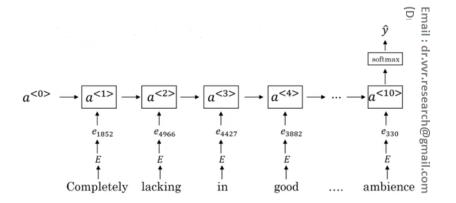
Ram is not a good learner, good listener and good human





Sentiment Classification RNN model





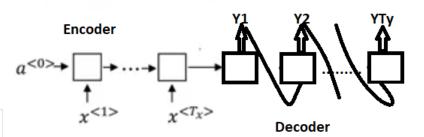




Sequence To Sequence Models: Machi

$$x^{<1>}$$
 $x^{<2>}$ $x^{<3>}$ $x^{<4>}$ $x^{<5>}$
Jane visite l'Afrique en septembre

Jane is visiting Africa in September. $v^{<1}>v^{<2}>v^{<3}>v^{<4}>v^{<5}>v^{<6}>$

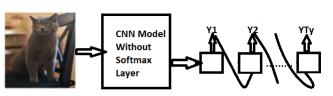




Sequence To Sequence Models: Image Captioning



Email:dr.vvr.research@gmail.com Dr. Venkataramana Veeramsetty)



Y1 Y2 . . . YTy: : A Cat Sitting On The Chair





Machine Translations using Conditiona Load Modeling



Email:dr.vvr.research@gmail.com Dr. Venkataramana Veeramsetty

Jane visite l'Afrique en septembre.

$$P(y^{<1>}, ..., y^{< T_y>} | x)$$

- → Jane is visiting Africa in September.
- → Jane is going to be visiting Africa in September.
- → In September, Jane will visit Africa.
- → Her African friend welcomed Jane in September.







Language model:

û<1> Machine translation:



