

Recurrent Neural Network (RNN)

Sequence - series + time
length $\begin{cases} \text{Input} \\ \text{Output} \end{cases}$

Recurrent Neural Network

Architecture

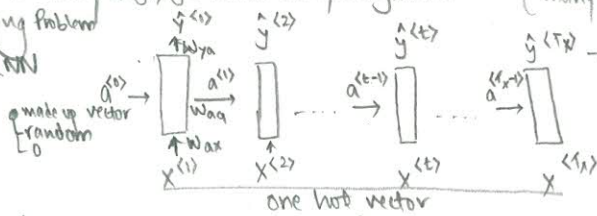
One to One - Normal NN

RNN $\begin{cases} \text{One to many} - \text{Music generation} \\ \text{many to one} - \text{Sentiment Analysis} \end{cases}$

attention based

$\begin{cases} \text{many to many} - \text{Named Entity Recognition} \\ \text{many to many} - \text{Machine Translation} \end{cases}$
→ binary classification of each word

Many to many → 0 → 0 - Named Entity Recognition
Supervised Learning Problem
Unidirectional RNN



Forward Propagation

weight w_{ax} input output

$$a^{(0)} = \vec{0}$$

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

$$\hat{y}^{(t)} = g_2(w_{ya} a^{(t)} + b_y) \quad \text{sigmoid}$$

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

sequence of word $\begin{cases} \text{vocabulary / Dictionary - Representative words} \\ \text{[padding n words]} \\ \text{[Unknown]} \end{cases}$

→ $w_a [a^{(t-1)}, x^{(t)}] + b_a$

$$\rightarrow w_y a^{(t)} + b_y$$

$$w_a = [w_{aa} | w_{ax}] \quad \text{stacked horizontally}$$

Standard NN

$\begin{cases} \text{Input \& output - different size} \\ \text{max-padding} \\ \text{feature/parameter sharing} \end{cases}$

Backward Propagation through time

$$L(\hat{y}, y) = \sum_t L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$L^{(t)} = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

Gradient Problem
gradient (layer)

Vanishing
Exploding

Backprop error → memory local influences

language model

gradient → numerical overflow (NaN)
→ gradient clipping

word - singular plural long range dependencies

Language Sequence

Model speech recognition

sentence → sequence of words

$$P(y^{(1)}, y^{(2)}, \dots, y^{(n)}) \quad \text{language model}$$

Training Set

tokenize - vocabulary - map → one hot vector

EOS - end of sentence

function

Uncommon words

Sample Novel Sequence

level - word

character

choose start

predict softmax → random sampling...

fill EOS

no of words

challenge - unknown words

Language model and sequence generation

Language model

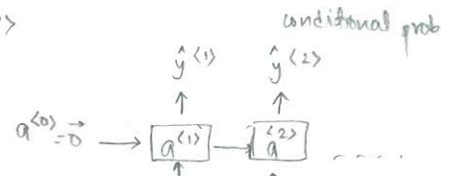
Speech recognition
Machine translation

probability (sentence)
next "

(sequence of words) $y^{(t)}$

Training set - large corpus of text

tokenize, \leftrightarrow OS \rightarrow vocabulary
unknown words



predicts one word at a time, going from left to right.

$$L(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

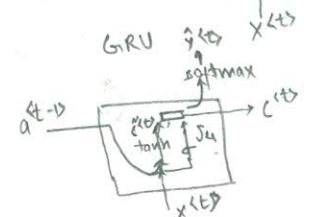
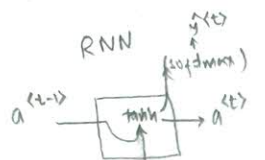
continue till EOS / # of words.

sampling novel sequences

Sample a sequence from trained RNN

$\hat{y}^{(t)}$ - probability (ith word) \rightarrow distⁿ
softmax distⁿ - reject unknown words
random sample

word level language model vocabulary - words
character " " " - letters...
longer sequences / complex computer extensive



Vanishing \uparrow gradients
Exploding \downarrow gradients \rightarrow NaN

capturing long term dependencies \rightarrow problem - linear
memorize

RNN - local influences backpropagate earlier - error - weakness

Exploding gradient

NaN / numerical overflow
gradient clipping - rescale gradient vector
max value

Gated Recurrent Unit (GRU) (2014)

long range connection - capture
vanishing gradient

memory cell $\rightarrow c^{(t)}$ $a^{(t)}$

$$\text{Candidate Value } \tilde{c}^{(t)} = \tanh(W_c [\Gamma_u^{(t-1)} \cdot c^{(t-1)} + x^{(t)}] + b_c)$$

vanishing gradient \leftarrow \uparrow -ve \leftarrow Update Gate $\rightarrow \Gamma_u$
 $\rightarrow c^{(t)}$ is maintained

$$\text{Update Gate } \Gamma_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u) \text{ mostly } 0 \text{ or } 1$$

$$\text{Gated fence } c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

$$\text{Relevance Gate } \Gamma_r = \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r) \quad a^{(t)} = c^{(t)}$$

STM (long short term) Memory unit - general (1997)

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

$$\text{Update Gate } \Gamma_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

$$\text{Forget Gate } \Gamma_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

$$\text{Output Gate } \Gamma_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

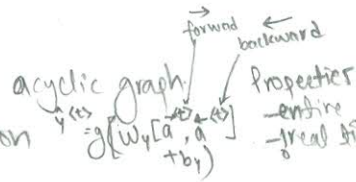
$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)}$$

$$a^{(t)} = \Gamma_o * \tanh c^{(t)}$$

peephole connection \rightarrow add $c^{(t-1)}$ term while calculating gate values.

Sidirectional RNN
BRNN

- Unidirectional / forward activation
+ backward (from future) - backward connection



Projector
entire sequence of data.
real time speech recognition

deep RNNs

complex function - stack multiple layers
computationally expensive to train

Natural Language Processing & Word Embeddings

NLP - word representation

One hot encoding - Vocabulary Index - words are independent

Word Embeddings - words are embedded in high dimension space

- Get \rightarrow train
 - Pretrained - Transfer learning
 - Named Entity Recognition
 - Text summarization
 - Co-reference resolution
 - Parsing
 - Language model
 - Machine translation

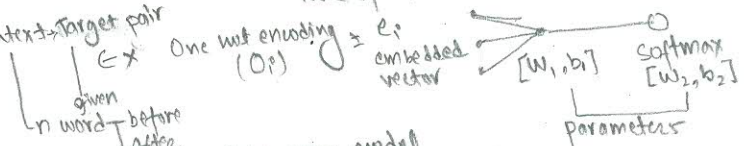
distance - Euclidean
 inner product
 features - parallel in n dimension
 $\arg \max_w (\text{similarity}(e_w, e_{f_1} - e_{f_2} + e_{w_1}))$
 cosine $\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$

- Analogy reasoning \rightarrow similarity
 - Visualize 0D \rightarrow 2D
 non-linear transformation
 t-SNE algorithm

- Embedding Matrix - learn from gradient descent
 to convert one hot encoding \rightarrow embedded representation
 multiplication
 look-up

Learn Word Embeddings - E

1) 2003 Supervised learning Problem - Context \rightarrow Target pair



2) 2013 Word2Vec Algorithm

① Skip gram model
 Context random
 Target Pair random + n word window
 Heuristic balance
 Common words
 Uncommon words

$O_c \rightarrow E \rightarrow e_c \rightarrow \text{softmax } \hat{y}$
 Embedding matrix
 Loss function $L(y, \hat{y}) = \sum y_i \log \hat{y}_i$
 $P(\hat{y}|e_c) = \frac{e^{e_c^T \theta_{\hat{y}}}}{\sum e^{e_c^T \theta_{\hat{y}}}}$
 $\theta_{\hat{y}} \rightarrow$ parameter - output + included bias
 [slow \rightarrow Hierarchical Softmax - tree computational complexity increases by log not linearly]

② CBOW model Continuous backwards context - surrounding context

3) 2013 Negative Sampling Context Target Pair \rightarrow Target (0,1)

select \rightarrow 1
 select \rightarrow 0 - negative sampling
 random
 positive to negative ratio
 observed frequency - stop words
 $\frac{f(w_i)}{\sum f(w_j)^{3/4}}$
 Empirical

Basic logistic regression model
 $P(y=1|e, \theta) = \sigma(\theta^T e_c)$
 softmax func sigmoid

4) 2014 GloVe Algorithm Global vector

$X_{ij} = \# \text{ of } i \text{ appears in context } j$
 symmetric - nearby

minimize $\sum_i \sum_j f(X_{ij}) (\theta_i^T e_j + b_i + b_j - \log X_{ij})^2$
 Symmetric
 weighting term \rightarrow frequent words
 $X_{ij} = \log$
 starts with uniformly random
 Optimize \rightarrow GD
 $e_w^{\text{final}} = \frac{e_w + \theta_w}{2}$

Application

Sentiment classification

- Use word embedding
 - ignore word order
 - Use RNN many to one architecture
 $O \rightarrow E \rightarrow e \xrightarrow{\text{sum/average}} \hat{y} \xrightarrow{\text{softmax}}$

Debiasing Word Embedding - 2016

bias
 gender
 ethnicity
 sexual orientation
 college admission
 job loan
 criminal justice

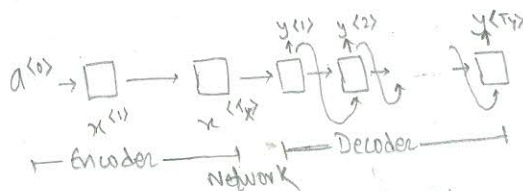
- cannot guarantee that individual component of feature representation \rightarrow interpretable
 $\theta_i^T e_j = (A \theta_j)^T (A^{-1} e_j)$

- 1) Identify bias direction difference [average 2D singular value decomposition (SVD)]
- 2) Neutralization - word not definitional - intrinsic bias
- 3) Equalize pair - difference - pairwise

Sequence Model

Sequence to sequence model

2014



2014/2015 Image Captioning CNN+RNN

Pick up most likely sentence

Language Model + condition → Machine Translation

output $P(y^{(1)}, \dots, y^{(n)} | x)$

random sample
most likely translation

$$\arg \max_{y^{(1)}, \dots, y^{(n)}} P(y^{(1)}, \dots, y^{(n)} | x)$$

Beam Search Algorithm

beam width - n words

- keep track of (n x no. of words) in vocabulary at each evaluation
- instantiate n copy of network
- till EOS - end of sentence

- 1 word greedy search

Length Normalization

Beam Search → ① $\arg \max_y \prod_{t=1}^T P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$

$$\text{product} \rightarrow \text{numerical underflow}$$

$$\text{sum} \rightarrow \text{numerical stable}$$

shorter sentence is preferred

② Divide by $\prod_{t=1}^T \alpha^{y^{(t)}}$

softer approach
 $\alpha = 0.7 \rightarrow$ heuristic

Error Analysis

Error | Beam Search
RNN

$$P(y^* | x) > P(\hat{y} | x)$$

$$P(y^* | x) \leq P(\hat{y} | x)$$

y^* - Human
 \hat{y} - Algorithm

Bleu Score

2002

Bilingual evaluation understudy (take role of senior actor)
Evaluate machine translation → multiple equally good translation (human)

Precision = $\frac{\text{present in reference}}{\# \text{ words in machine translated output}}$
repeated words → clip numerator → max # of words present in reference

words in 2-gram - Unigram

Bleu Score < - Bigram
- n-gram

$$p_n = \frac{\sum_{n\text{-gram} \in \hat{y}} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram}} \text{count}(n\text{-gram})}$$

Combined Bleu Score

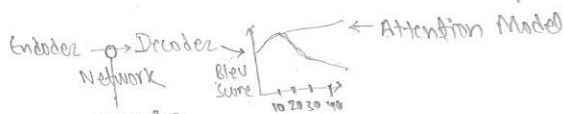
$$= B.P. + \exp\left(\frac{1}{N} \sum_{n=1}^N p_n\right)$$

↳ brevity penalty - penalize short sentence

$$B.P. = \begin{cases} 1 & \text{if (machine translation) > reference output length} \\ \exp(1 - \text{ref len} / \text{hyp len}) & \text{otherwise} \end{cases}$$

Attention Model

2014

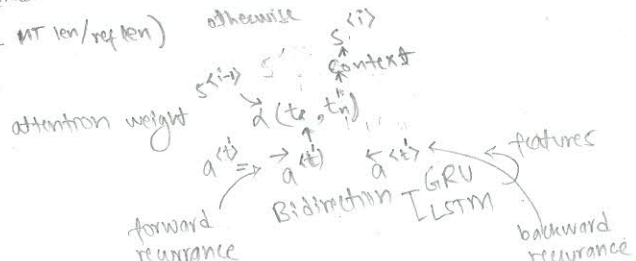


Runtime - Quadratic Cost

↳ input & output sentence is not long.

Application - Image captioning

Visualization of attention weights



$$\sum_t \alpha^{(t,t')} = 1 \quad \text{sum to 1}$$

$$c^{(i)} = \sum_t \alpha^{(t,i)} a^{(t)}$$

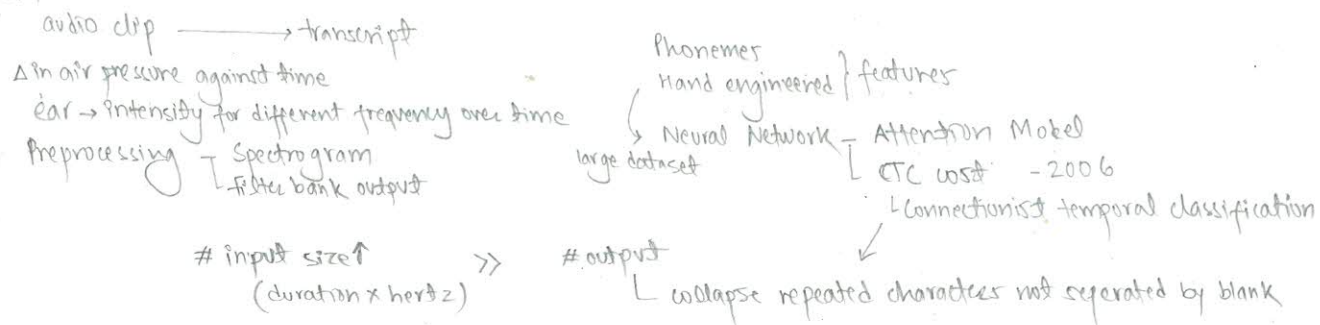
$$a^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^T \exp(e^{(t,t')})}$$

softmax

weight amount of attention $y^{(t)}$ should pay to $a^{(t')}$

small neural network
 $s^{(t-1)}, a^{(t)}$ → $e^{(t,t')}$
function learn

Speech Recognition



Trigger Word Detection

