



RECURRENT NEURAL NETWORKS/LONG-SHORT TERM MEMORY FOR NATURAL LANGUAGE PROCESSING

AND APPLICATIONS IN CLINICAL NOTE ANALYTICS

Yanshan Wang, PhD
Mayo Clinic
wang.yanshan@mayo.edu

AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

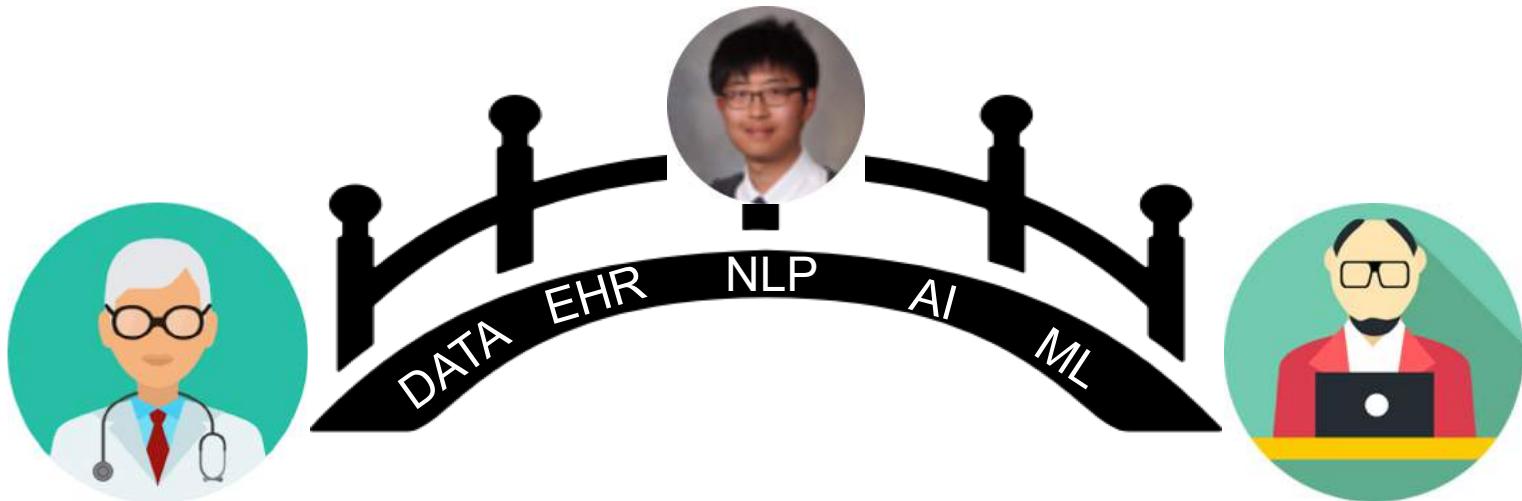
AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

INTRODUCTION

■ Who am I?

Computer scientist &
Informatician



■ How can you contact me?

- Email: wang.yanshan@mayo.edu
- LinkedIn, Twitter (yanshan_wang)

WHY TAKE THIS TUTORIAL

- Artificial Intelligence (AI) is one of the most interesting fields of research today.
- The growth of interest in AI is due to the recent advances in deep learning.
- Language is the most compelling manifestation of intelligence.



From left, Yann LeCun, Geoffrey Hinton and Yoshua Bengio. The researchers worked on key developments for neural networks, which are reshaping how computer systems are built.

GOAL OF THIS TUTORIAL

- To get an understanding of some core concepts in deep learning (RNN/LSTM) and natural language processing.
- To connect theoretical with practical knowledge.
- To get an introduction into programming of RNN/LSTM in Python.

PREREQUISITES

■ Maths

- Linear Algebra
- Probability

■ Machine learning

- Evaluating ML models (train/validation/test, cross validation, etc.)
- Overfitting, generalization
- Optimization (objective function, stochastic gradient descent)

■ Programming

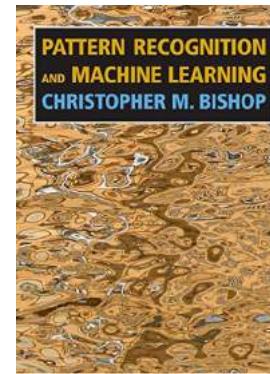
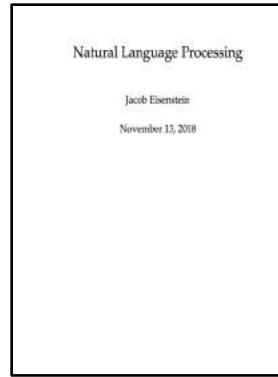
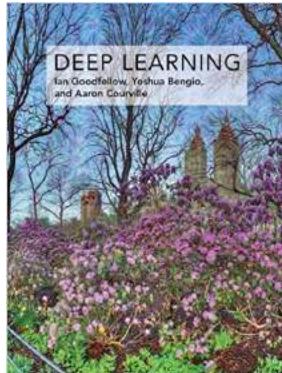
- Python (e.g., Keras, NumPy, TensorFlow, etc.)

SUGGESTED READING

■ Deep learning for NLP Courses

- Oxford Deep Learning for NLP class: <https://github.com/oxford-cs-deepnlp-2017/lectures>
- Stanford CS 224D: Deep Learning for NLP class: <http://cs224d.stanford.edu/syllabus.html>

■ Books



■ Resources

- <https://github.com/andrewt3000/DL4NLP>
- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

SUGGESTED READING

■ Papers (general NLP)

- [A neural probabilistic language model.](#) Bengio et al. 2003
- [Efficient estimation of word representations in vector space.](#) Mikolov et al. 2013.
- [Long short-term memory.](#) Hochreiter and Schmidhuber. 1997.
- [Recurrent neural network based language model.](#) Mikolov et al. 2010.
- [Generating sequences with recurrent neural networks.](#) Graves. 2013.
- [On the state of the art of evaluation in neural language models.](#) Melis et al. 2018

■ Paper (clinical NLP)

- [Clinical information extraction applications: a literature review.](#) Wang et al. 2018
- [Using clinical natural language processing for health outcomes research: Overview and actionable suggestions for future advances.](#) Velupillai et al. 2018.

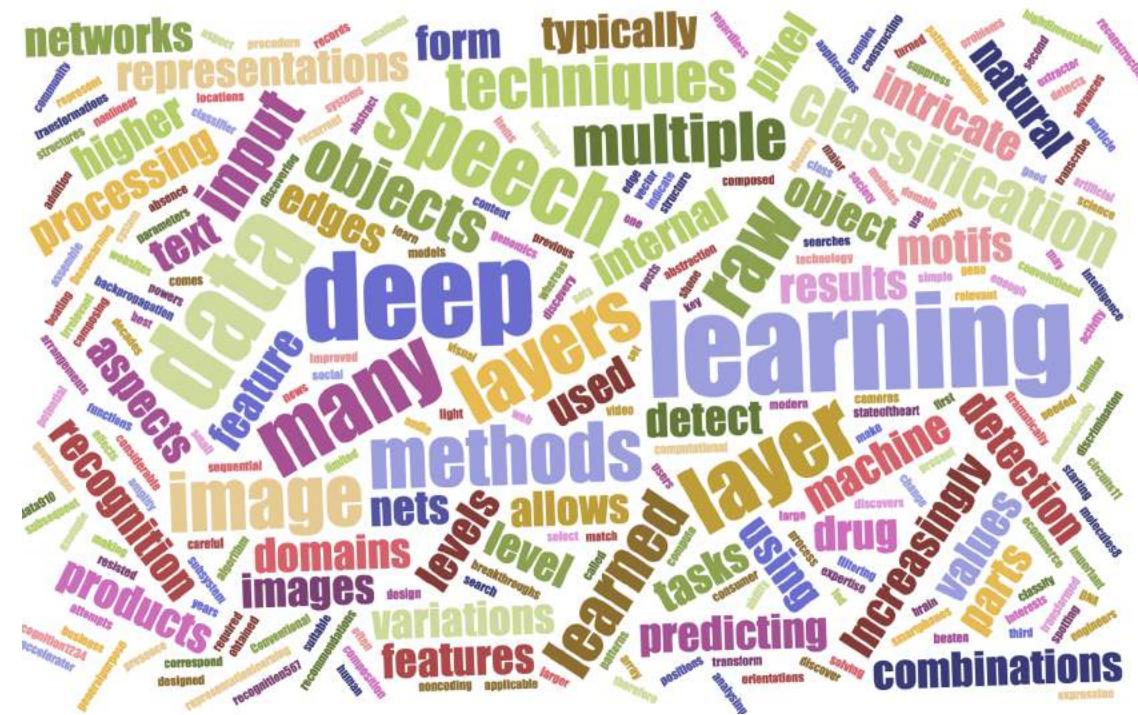
AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

NATURAL LANGUAGE PROCESSING

■ What is NLP?

- "Natural language processing (NLP) is a subfield of computer science, information engineering, and artificial intelligence concerned with the interactions between computers and human (natural) languages, in particular how to program computers to process and analyze large amounts of natural language data." (Wikipedia)



NATURAL LANGUAGE PROCESSING

■ What is NLP?

- "Natural language processing (NLP) is a subfield of computer science, information engineering, and artificial intelligence concerned with the interactions between computers and human (natural) languages, in particular how to program computers to process and analyze large amounts of natural language data." (Wikipedia)



QUESTION ANSWERING



IBM Watson

Voice Assistant



Source: <https://www.youtube.com/watch?v=BkpAro4zlwU>

INFORMATION EXTRACTION

A presenilin 1 mutation (*Ser169Pro*) associated with *early-onset AD* and *myoclonic seizures*.

INFORMATION EXTRACTION

A presenilin 1 mutation (*Ser169Pro*) associated with *early-onset AD* and *myoclonic seizures*.

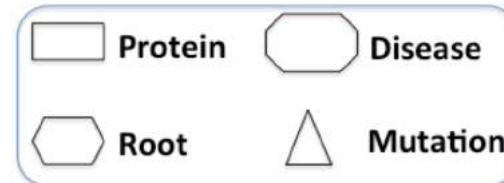
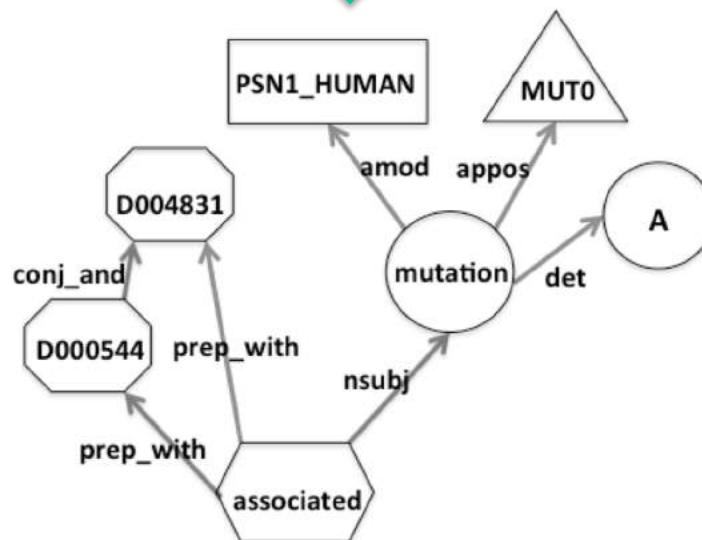


Entity normalization

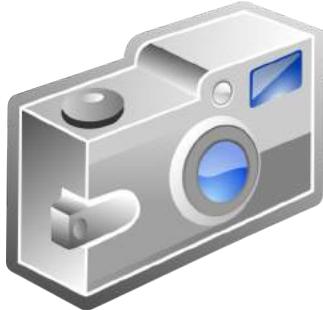
A PSN1_HUMAN mutation (*MUTO*) associated with *D000544* and *D004831*.



Dependency parser



SENTIMENT ANALYSIS



Attributes:

zoom

affordability

size and weight

flash

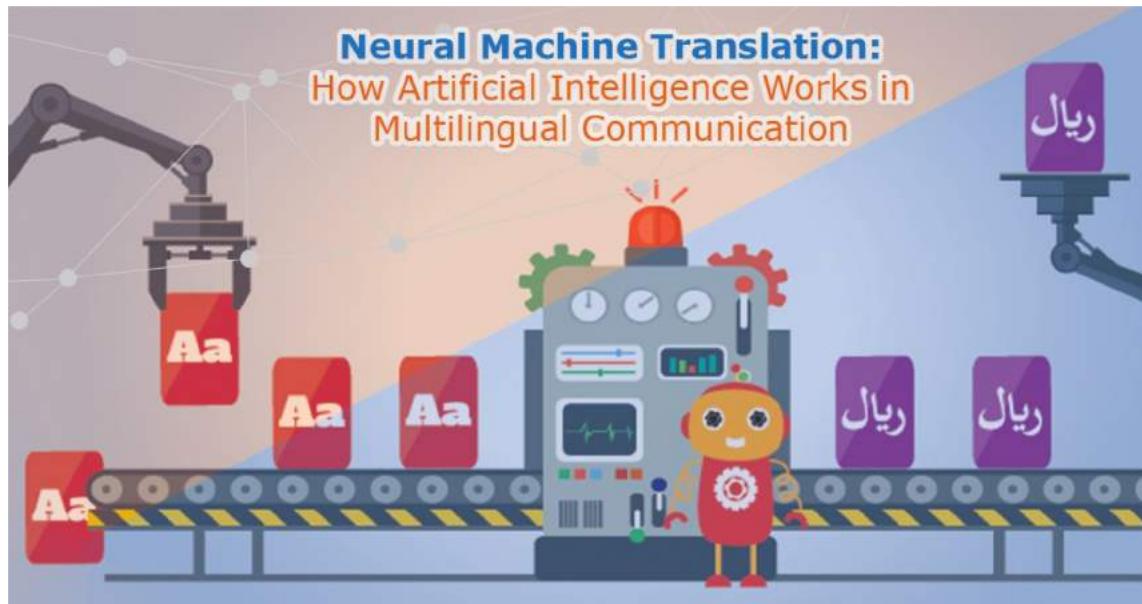
ease of use



Reviews:

- ✓ ■ nice and compact to carry!
- ✓ ■ since the camera is small and light, I won't need to carry around those heavy, bulky professional cameras either!
- ✗ ■ the camera feels flimsy, is plastic and very light in weight
you have to be very delicate in the handling of this camera

MACHINE TRANSLATION



Source: <https://chatbotslife.com/machine-translation-by-artificial-intelligence-9d7a732a4ee1>



Source: <http://www.tagxp.com/m/html/2372.html>

AGENDA

- Introduction
- Background
 - Natural Language Processing
 - **Language Models**
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

How TO REPRESENT NATURAL LANGUAGE

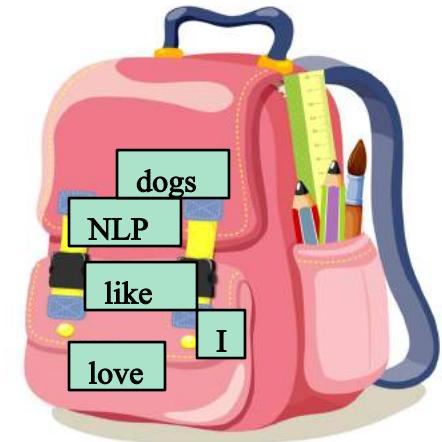
- Natural language text = sequences of discrete symbols (e.g. words).

I love NLP and I like dogs

- Vector representations of words: Vector Space Model
 - Bag-of-words

	I	love	NLP	and	like	dogs
I	1	0	0	0	0	0
love	0	1	0	0	0	0
NLP	0	0	1	0	0	0
like	0	0	0	0	1	0

Vocabulary list



HOW TO REPRESENT NATURAL LANGUAGE

Sparse representation

I	=	[1	0	0	0	0	0]
love	=	[0	1	0	0	0	0]
NLP	=	[0	0	1	0	0	0]
like	=	[0	0	0	0	1	0]

■ Drawbacks

- $\text{love}=[0,1,0,0,0,0]$ AND $\text{like}=[0,0,0,0,1,0] = 0!$
- Using such an encoding, there's no meaningful (semantic) comparison we can make between word vectors.

How TO REPRESENT NATURAL LANGUAGE

■ Semantic representations

- Count-based methods
 - Define a basis vocabulary C of context words.
 - Define a word window size w .
 - Count the basis vocabulary words occurring w words to the left or right of each instance of a target word in the corpus.
 - Form a vector representation of the target word based on these counts.

... and the cute **kitten** purred and then ...
... the cute fury **cat** purred and miaowed ...
... that the small **kitten** miaowed and she ...
... the loud fury **dog** ran and bit ...

Example **basis vocabulary**: {bit, cute, furry, loud, miaowed, purred, ran, small}.

kitten context words: {cute, purred, small, miaowed}.

cat context words: {cute, fury, miaowed}.

dog context words: {loud, fury, ran, bit}.



$$\text{kitten} = [0, 1, 0, 0, 1, 1, 0, 1]^T$$

$$\text{cat} = [0, 1, 1, 0, 1, 0, 0, 0]^T$$

$$\text{dog} = [1, 0, 1, 1, 0, 0, 1, 0]^T$$

HOW TO REPRESENT NATURAL LANGUAGE

Use inner product or cosine as **similarity kernel**. E.g.:

$$sim(\text{kitten}, \text{cat}) = cosine(\text{kitten}, \text{cat}) \approx 0.58$$

$$sim(\text{kitten}, \text{dog}) = cosine(\text{kitten}, \text{dog}) = 0.00$$

$$sim(\text{cat}, \text{dog}) = cosine(\text{cat}, \text{dog}) \approx 0.29$$

Reminder: $cosine(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \times \|\mathbf{v}\|}$

Cosine has the advantage that it's a *norm-invariant* metric.

LANGUAGE MODELS

- The simple objective of modelling the next word given the observed history of natural language

$$p(\cdot | \text{There she built } a)$$

- With more context we are able to use our knowledge of both language and the world to heavily constrain the distribution over the next word:

$$p(\cdot | \text{Alice went to the beach. There she built } a)$$

- Language modelling is a time series prediction problem in which we must be careful to train on the past and test on the future.

LANGUAGE MODELS

- Using language models, we learn “meaning” of a word using dense representation

Dense representation

I	= [0.99	0.05	0.1	0.87	0.1	0.1]
love	= [0.1	0.85	0.99	0.1	0.83	0.09]
NLP	= [0.67	0.23	0.01	0.02	0.01	0.81]
like	= [0.1	0.73	0.99	0.05	1.79	0.09]

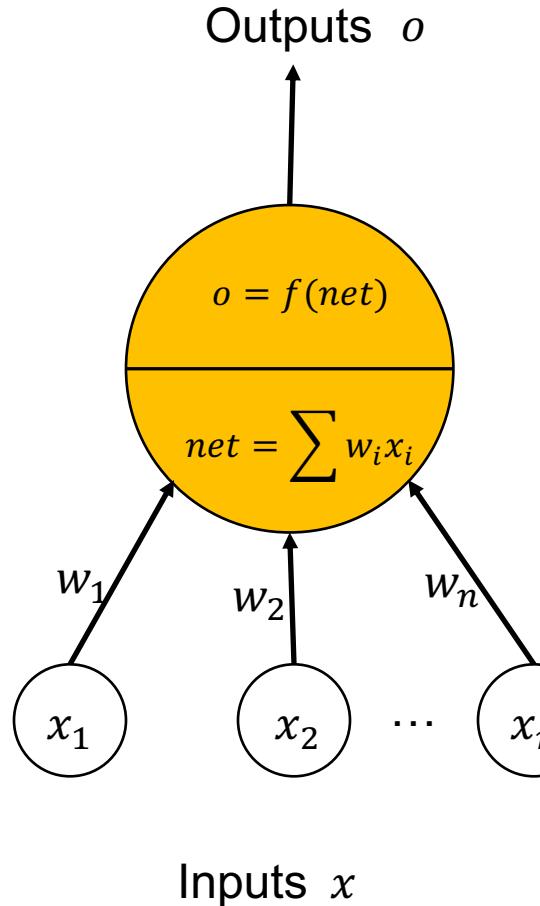
- Simply examining a large corpus it's possible to learn word vectors that are able to capture the relationships between words in a surprisingly expressive way.
- Popular language models
 - Topic models: LSA, pLSA, LDA, etc.
 - Neural language models

AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

DEEP LEARNING BASICS

- Artificial neuron is a nonlinear processing unit

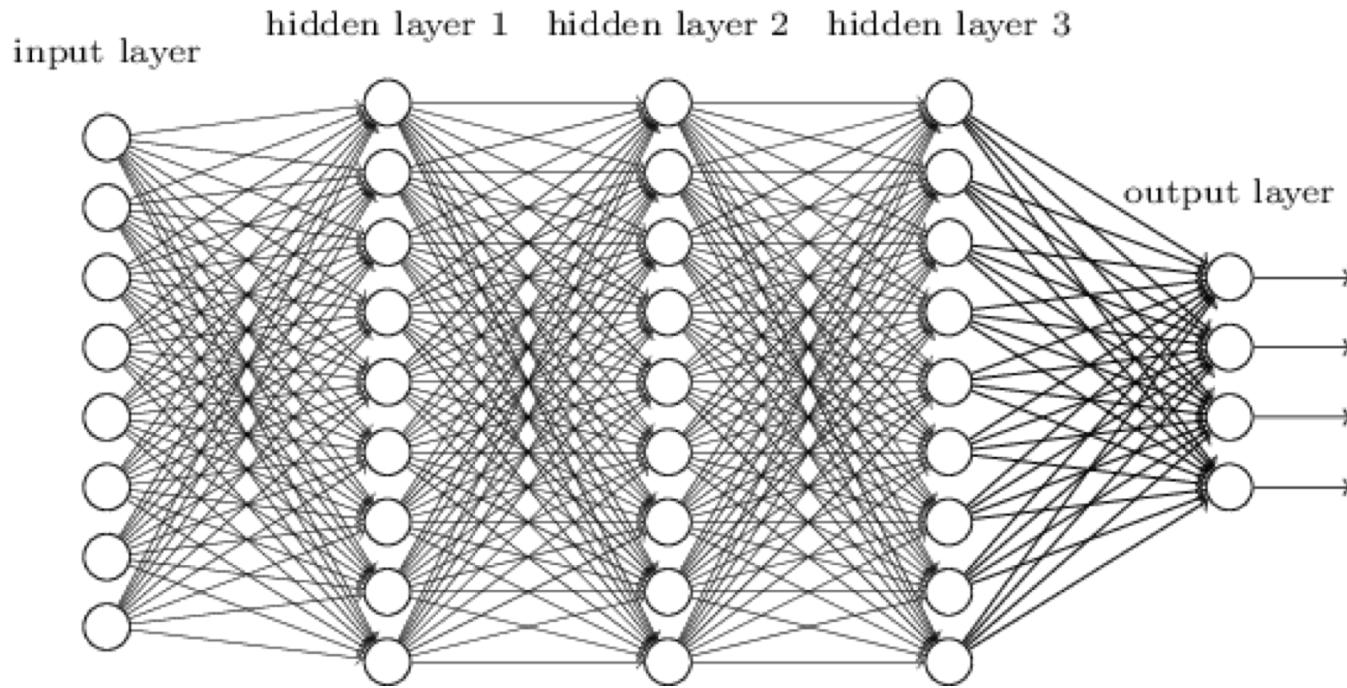


Activation functions f

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Arctan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU) ^[2]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) ^[2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) ^[3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

DEEP LEARNING

- Deep learning is a class of machine learning algorithms that use a cascade of multiple layers of nonlinear processing units for feature extraction and transformation.



AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

AGENDA

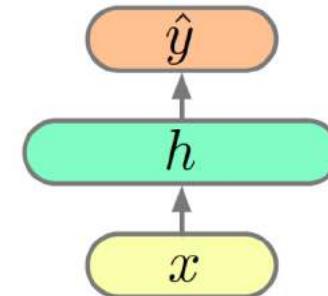
- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

NEURAL LANGUAGE MODELS

■ Feed forward network

$$h = g(Vx + c)$$

$$\hat{y} = Wh + b$$



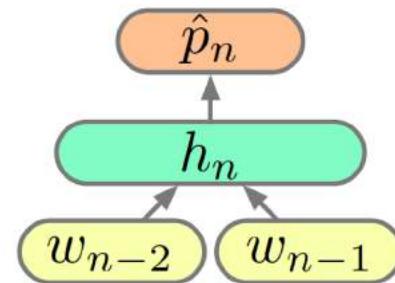
■ Trigram neural language model

$$h_n = g(V[w_{n-1}; w_{n-2}] + c)$$

$$\hat{p}_n = \text{softmax}(Wh_n + b)$$

- w_i are one hot vectors, \hat{p}_n is distribution
- $|w_i|$ is the size of vocabulary which is usually very large.

Computational expensive.

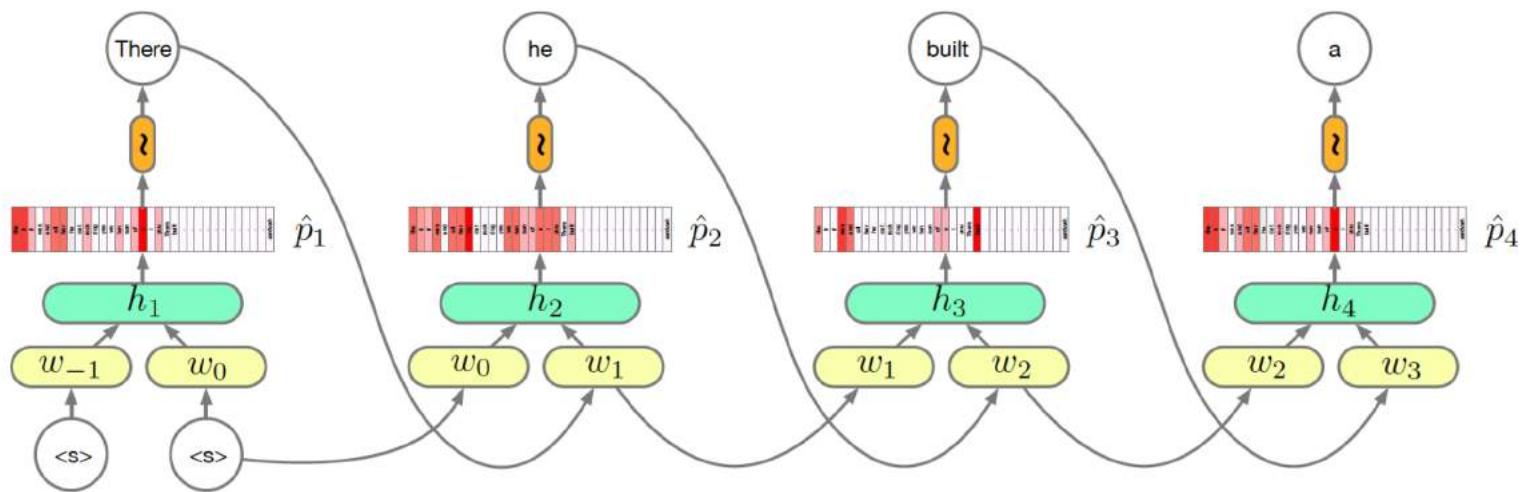


NEURAL LANGUAGE MODELS

■ Sampling.

- When learning, the output is a probability distribution instead of one word. When generating text we choose only one of the words ourselves given the probabilities and feed that back into the network. This is called sampling

$$w_n | w_{n-1}, w_{n-2} \sim \hat{p}_n$$



NEURAL LANGUAGE MODELS

■ Objective function

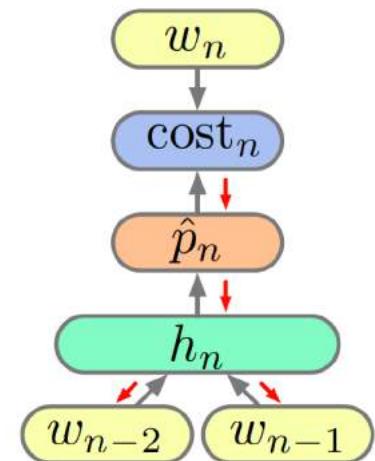
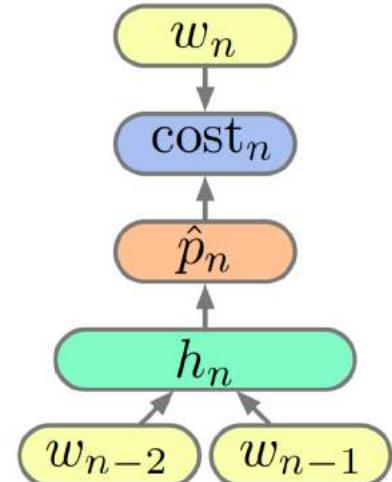
- The usual training objective is the cross entropy of the data given the model (MLE):

$$L = -\frac{1}{N} \text{cost}_n(w_n, \hat{p}_n)$$

- Calculating the gradients is straightforward with back propagation:

$$\frac{\partial L}{\partial W} = -\frac{1}{N} \sum_n \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial W}$$

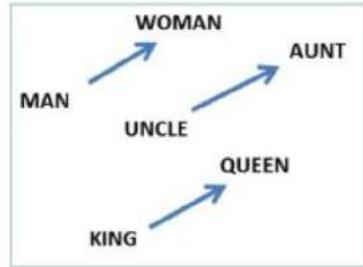
$$\frac{\partial L}{\partial V} = -\frac{1}{N} \sum_n \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial h_n} \frac{\partial h_n}{\partial V}$$



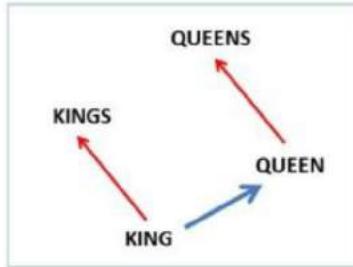
NEURAL LANGUAGE MODELS

- The insight behind neural language models is to treat word prediction as a discriminative learning task.
- The goal is to compute the probability $p(\text{word}|\text{context})$. Rather than directly estimating the word probabilities from (smoothed) relative frequencies, we can treat language modeling as a machine learning problem, and estimate parameters that maximize the log conditional probability of a corpus.
- The vector representation of words is also called word embeddings.

WORD2VEC



Source: <http://nlp.yvespeirsman.be/blog/visualizing-word-embeddings-with-tsne/>



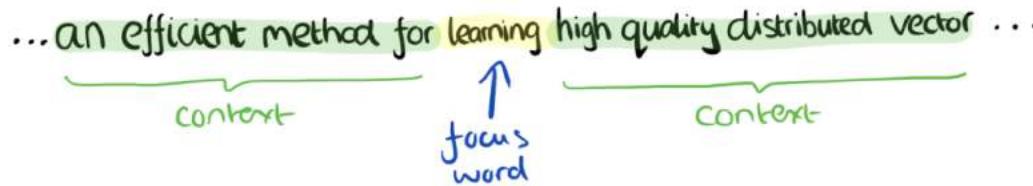
Source: <http://colah.github.io/posts/2014-07-NLP-RNNs-Representations/>

1. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems (pp. 3111-3119).
2. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781.

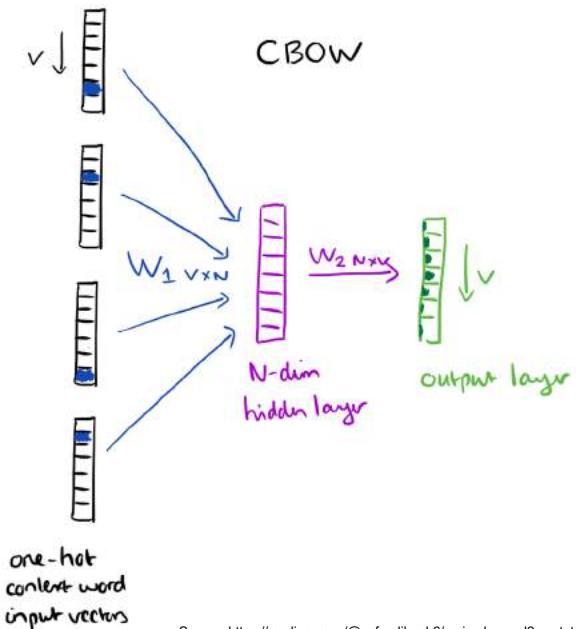
- Mikolov et al. weren't the first to use continuous vector representations of words, but they did show how to reduce the computational complexity of learning such representations – making it practical to learn high dimensional word vectors on a large amount of data.
- Two new architectures are proposed in their study: a Continuous Bag-of-Words (CBOW) model, and a Skip-gram model.

CBOW

- Imagine a sliding window over the text, that includes the central word currently in focus, together with the four words (a.k.a. word window =4) and precede it, and the four words that follow it:



- The context words form the input layer. Each word is encoded in one-hot form, so if the vocabulary size is V these will be V -dimensional vectors with just one of the elements set to one, and the rest all zeros. There is a single hidden layer and an output layer.



CBOW

- The training objective is to maximize the conditional probability of observing the actual output word (the focus word) given the input context words, with regard to the weights.

$\max \text{Prob}(\text{"learning"} | \text{"an"}, \text{"efficient"}, \text{"method"}, \text{"for"}, \text{"high"}, \text{"quality"}, \text{"distributed"}, \text{"vector"})$

- Since our input vectors are one-hot, multiplying an input vector by the weight matrix W_1 amounts to simply selecting a row from W_1 . Given C input word vectors, the activation function for the hidden layer h amounts to simply summing the corresponding ‘hot’ rows in W_1 , and dividing by C to take their average.

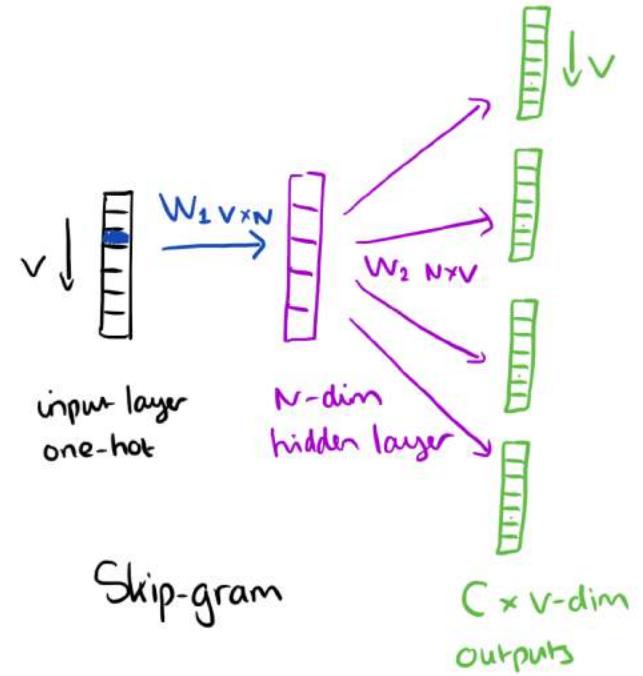
$$\begin{array}{c} \text{input} \\ \text{1} \times V \end{array} \quad \begin{array}{c} W_1 \\ V \times N \end{array} \quad \begin{array}{c} \text{hidden} \\ 1 \times N \end{array}$$
$$\begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & j & k & l \end{bmatrix}_{W_1} = \begin{bmatrix} e & f & g & h \end{bmatrix}$$

- From the hidden layer to the output layer, the second weight matrix W_2 can be used to compute a score for each word in the vocabulary, and softmax can be used to obtain the posterior distribution of words.

SKIP-GRAM

- The skip-gram model is the opposite of the CBOW model. It is constructed with the focus word as the single input vector, and the target context words are now at the output layer:
- At the output layer, we now output C multinomial distributions instead of just one. The training objective is to minimize the summed prediction error across all context words in the output layer.

- Optimization
- Having to update every output word vector for every word in a training instance is very expensive.
- Mikolov et al. used Negative Sampling. It is simply the idea that we only update a sample of output words per iteration. The target output word should be kept in the sample and gets updated, and we add to this a few (non-target) words as negative samples.



MOTIVATION

- Neural word embeddings have been widely used in biomedical machine learning applications.
- Different word embeddings have been utilized
 - Local corpus
 - External data resources, such as Wikipedia
 - Public pre-trained word embeddings.
- Questions remain unanswered
 - Do we need to train word embeddings for a specific biomedical machine learning task since there are a number of public pre-trained word embeddings? Or which word embeddings should we use for a biomedical machine learning task?

MATERIALS

- Clinical notes

Clinical notes for a cohort of 113k patients receiving their primary care at Mayo Clinic, spanning a period of 15 years from 1998 to 2013. The vocabulary size of this corpus is 103k.

- Biomedical publications

1.25 million PubMed Central articles. The vocabulary size of this corpus is 2 million.

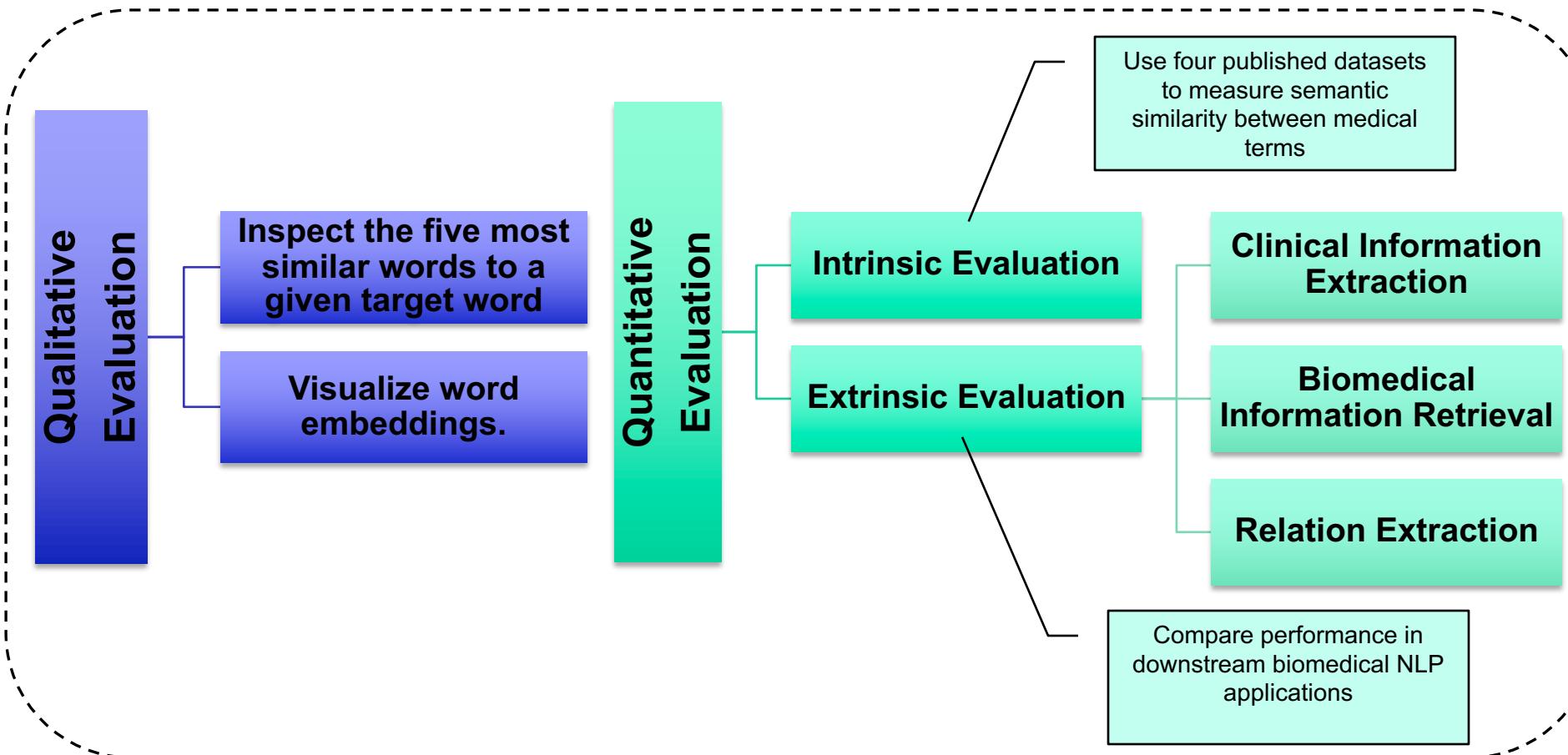
Pre-trained GloVe embeddings¹.

- Wikipedia

Pre-trained Google News embeddings².

- Public news

EVALUATION METHODS



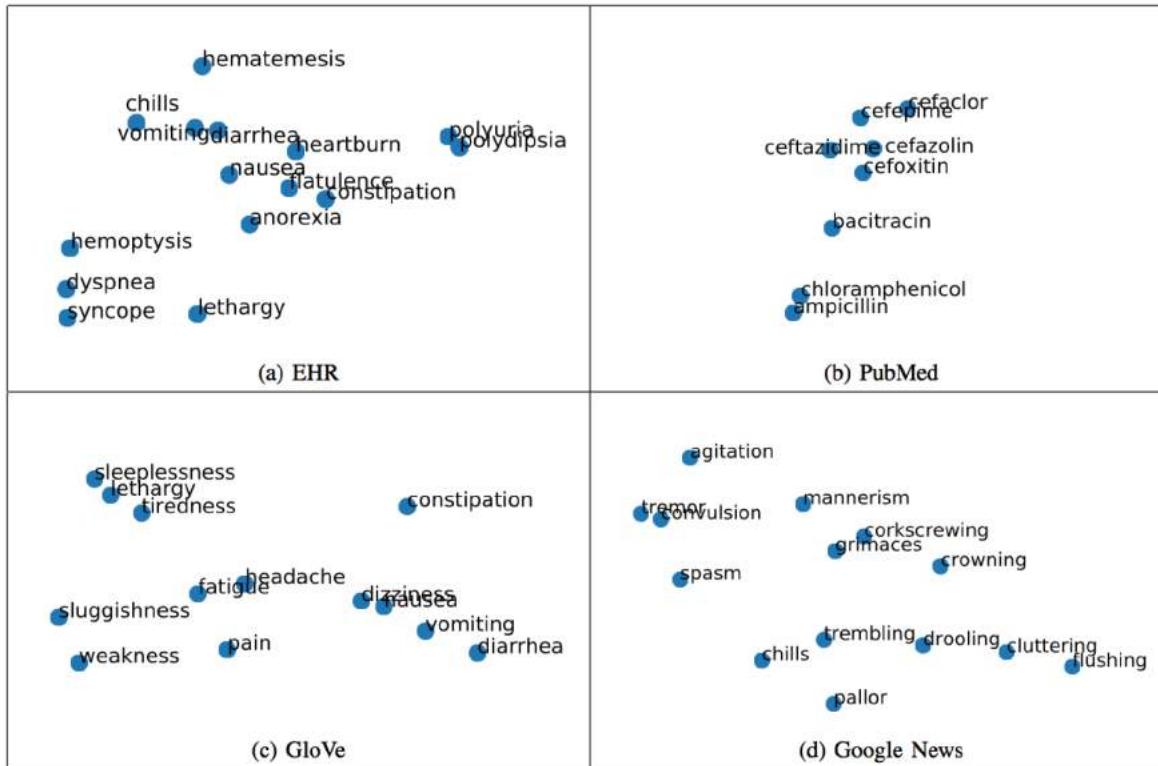
QUALITATIVE EVALUATION

- EHR and PubMed can capture the semantics of medical terms better than GloVe and Google News and find more relevant similar medical terms.
- However, EHR and PubMed find similar medical terms from different perspectives due to their different focuses. EHR contains clinical narratives and thus it is closer to clinical language. Yet, it contains terms with different morphologies and even typos, such as *melitis*, *caner* and *thraot*. Differently, PubMed contains more medical terminology used in medical articles, and finds similar words mostly from a medical research perspective.

QUALITATIVE EVALUATION

■ Visualization – t-SNE

A cluster of symptoms



A cluster of antibiotic medications

A cluster of symptoms

A cluster of symptoms

Fig. 1: Example clusters in different word embeddings.

TAKE-HOME MESSAGE...

- ✓ The word embeddings trained on EHR and PubMed can capture the semantics of medical terms better than those trained on GloVe and Google News and find more relevant similar medical terms. However, EHR finds similar terms via clinical language while PubMed contains more medical terminology used in medical articles, and finds similar words mostly from a medical research perspective.
- ✓ The medical semantic similarity captured by the word embeddings trained on EHR and PubMed are closer to human experts' judgments, compared to those trained on GloVe and Google News.
- ✓ There does not exist a consistent global ranking of word embedding quality for downstream biomedical NLP applications. However, adding word embeddings as extra features will improve results on most downstream tasks.
- ✓ Word embeddings trained from biomedical domain corpora do not necessarily have better performance than those trained on other general domain corpora. That is, there might be no significant difference when word embeddings trained from an out-domain corpus were employed for a biomedical NLP application. However, the performance of word embeddings trained from a local corpus is better for local NLP tasks.

AGENDA

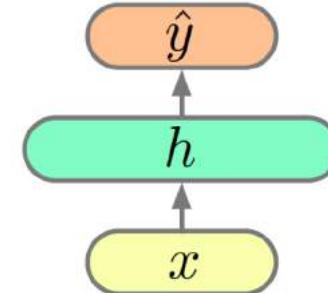
- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

RECURRENT NEURAL NETWORK

■ Feed forward network

$$h = g(Vx + c)$$

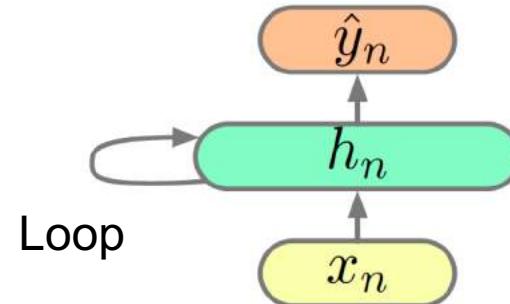
$$\hat{y} = Wh + b$$



■ Recurrent neural network

$$h_n = g(V[x_n; h_{n-1}] + c)$$

$$\hat{y}_n = Wh_n + b$$

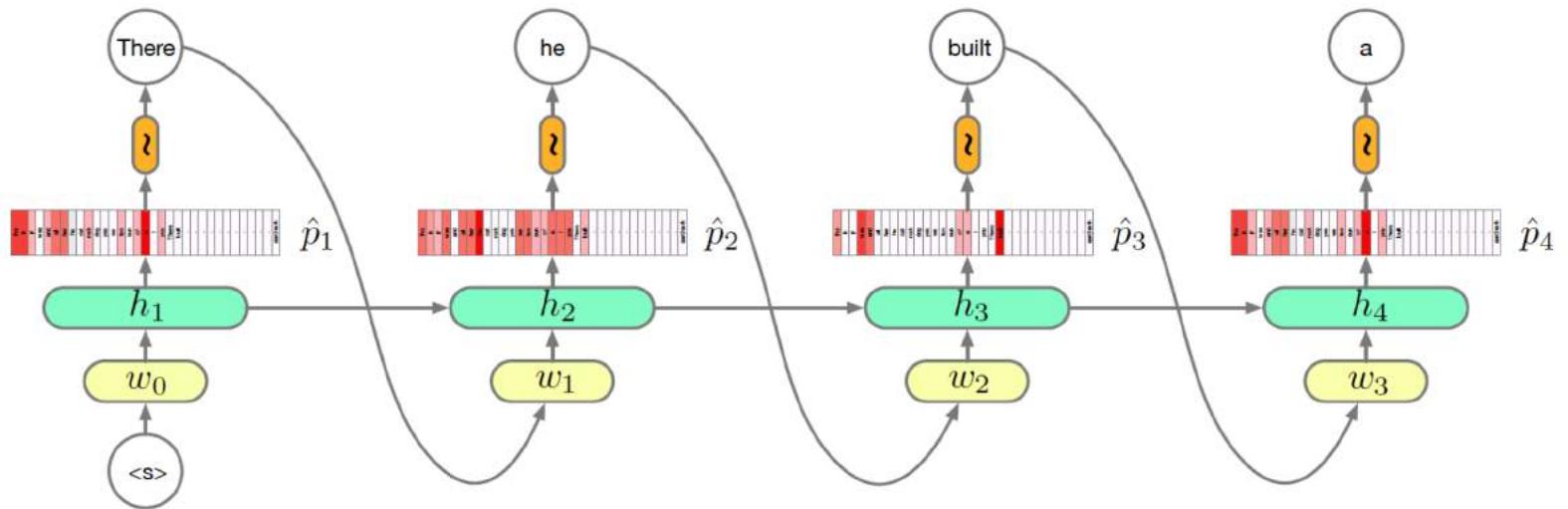


- Humans don't start their thinking from scratch. For example, we want to classify what kind of event is happening at every point in a movie. We want to use reasoning about previous events in the film to inform later ones.
- RNN adds loops, allowing information to persist.

RECURRENT NEURAL NETWORK

- A RNN can be thought of as multiple copies of the same network, each passing a message to a successor.

$$h_n = g(V[x_n; h_{n-1}] + c)$$

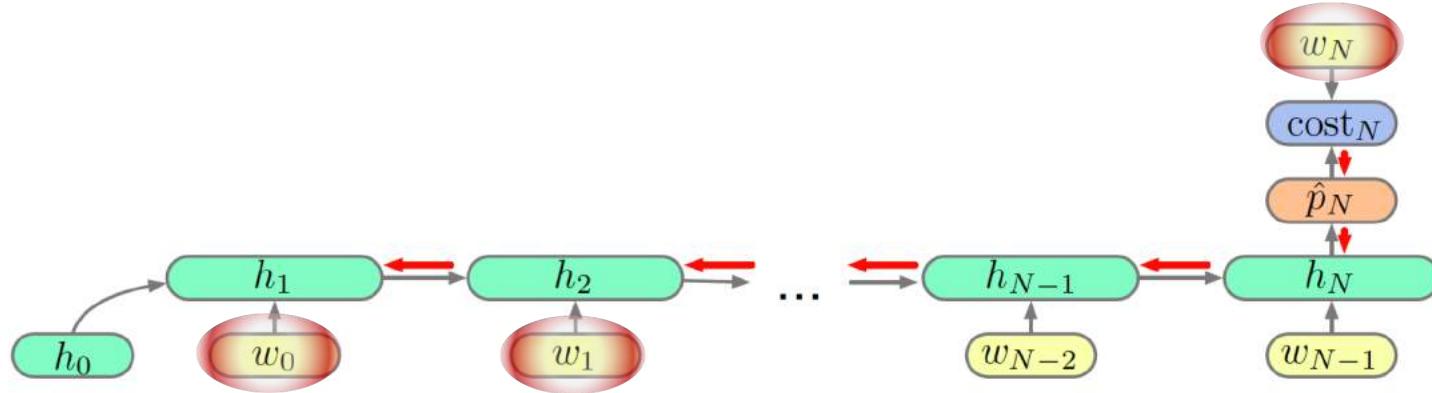


RECURRENT NEURAL NETWORK

■ Problem

- RNNs can learn to use the past information. However, not all the past information is equally important to predict the next word.

$p(\text{sandcastle} | \text{Alice went to the beach. There she built a})$



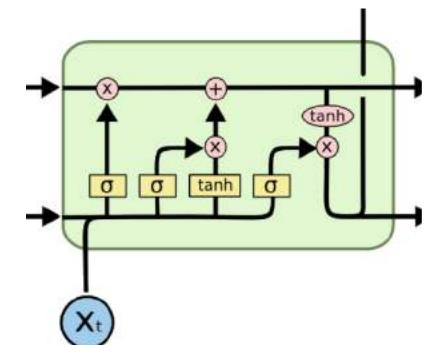
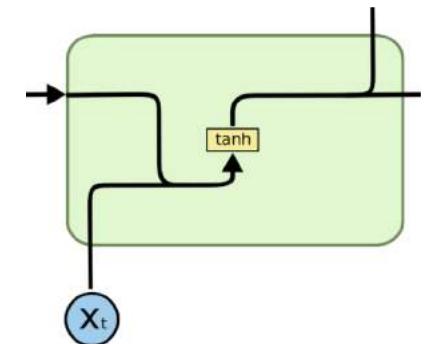
- Unfortunately, as that gap grows, RNNs become unable to learn to connect the information

AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

LSTM

- Long-Short Term Memory networks – usually just called “LSTMs” – are a special kind of RNN, capable of learning long-term dependencies.
- In standard RNNs, this repeating module will have a very simple structure, such as a single tanh layer.
- LSTMs also have this chain like structure, but the repeating module has a different structure. Instead of having a single neural network layer, there are four, interacting in a very special way.

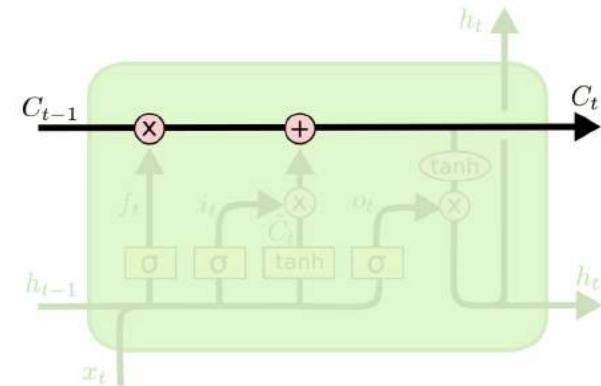


■ Neural Network Layer
Pointwise Operation
Vector Transfer
Concatenate
Copy

LSTM

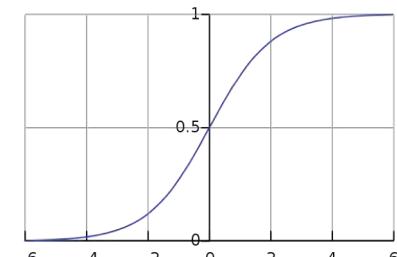
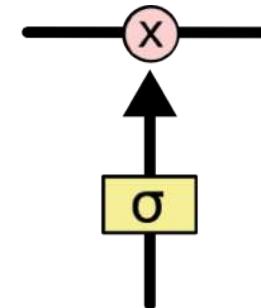
■ Cell state

- The key to LSTMs is the cell state.
- Information flow along the horizontal line unchanged.



■ Gate

- In order to have the ability to remove or add information to the cell state, LSTM uses gates to optionally let information through.
- Gates are composed out of a sigmoid neural net layer and a pointwise multiplication operation.
- An LSTM has three of these gates, to protect and control the cell state.



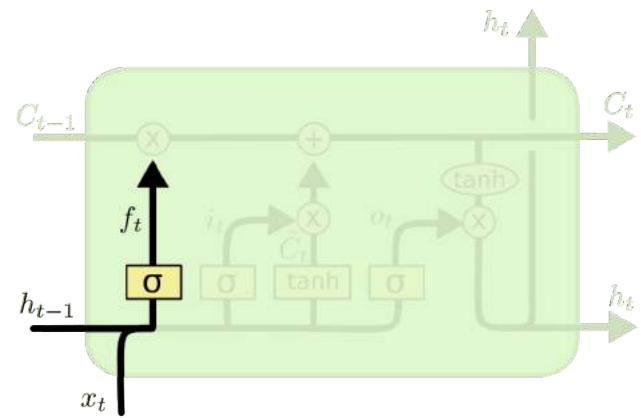
LSTM – STEP 1

■ Forget gate layer

- To decide what information we're going to throw away from the cell state.

$$f_t = \sigma(W_f[x_t; h_{t-1}] + c_f)$$

- It looks at h_{t-1} and x_t , and outputs a number between 0 and 1 for each number in the cell state C_{t-1} . A 1 represents “completely keep this” while a 0 represents “completely get rid of this.”



LSTM – STEP 2

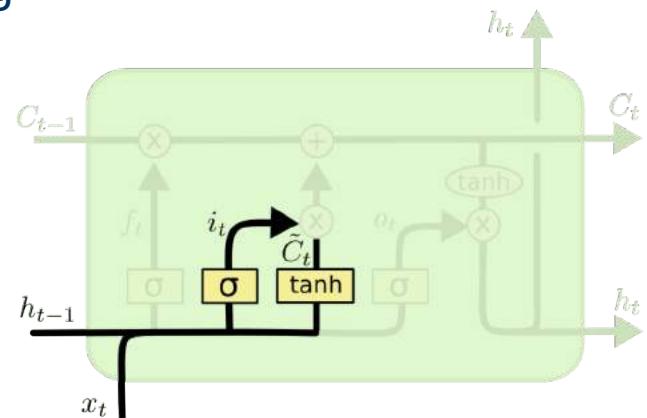
■ Input gate layer

- To decide what new information we're going to store in the cell state.

$$i_t = \sigma(W_i[x_t; h_{t-1}] + b_i)$$

$$\tilde{C}_t = \tanh(W_C[x_t; h_{t-1}] + b_C)$$

- First, a sigmoid layer decides which values we'll update. Next, a tanh layer creates a vector of new candidate values, \tilde{C}_t , that could be added to the state.



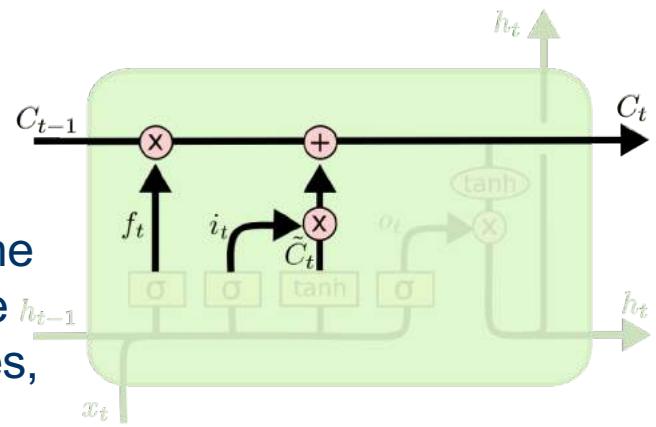
LSTM – STEP 3

■ Update

- To update the old cell state, C_{t-1} , into the new cell state C_t .

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- We multiply the old state by f_t , forgetting the things we decided to forget earlier. Then we add $i_t * \tilde{C}_t$. This is the new candidate values, scaled by how much we decided to update each state value.



LSTM – STEP 4

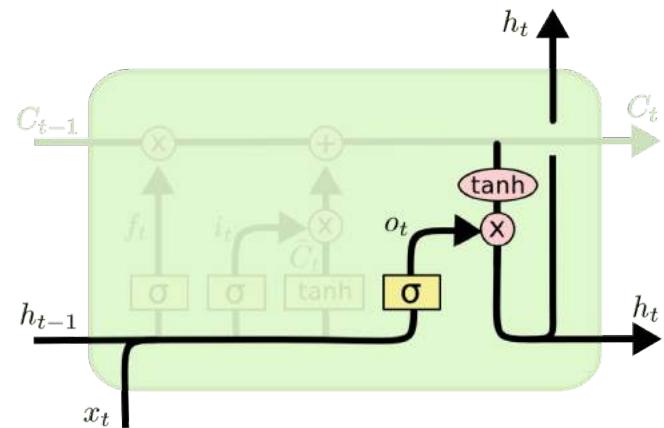
■ Output

- To decide what we're going to output.

$$o_t = \sigma(W_o[x_t; h_{t-1}] + b_o)$$

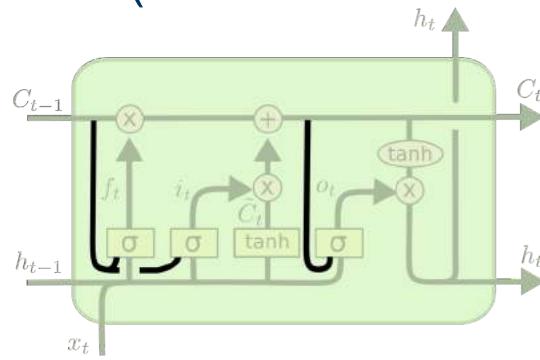
$$h_t = o_t * \tanh(C_t)$$

- This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we're going to output. Then, we put the cell state through tanh (to push the values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.



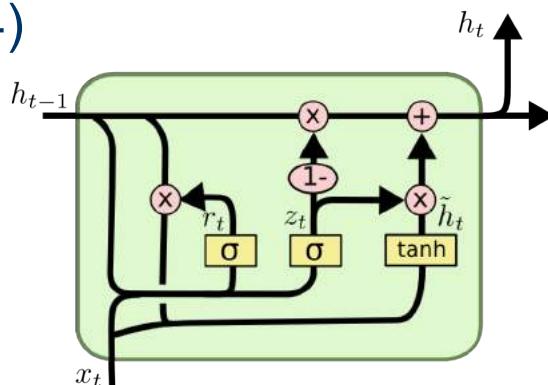
LSTM VARIANTS

- Adding "peephole connections" to let the gate layers look at the cell state. (Gers and Schmidhuber, 2000.)



$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$
$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

- Gated Recurrent Unit, or GRU, combines the forget and input gates into a single “update gate.” It also merges the cell state and hidden state, and makes some other changes. (Cho, et al. 2014)



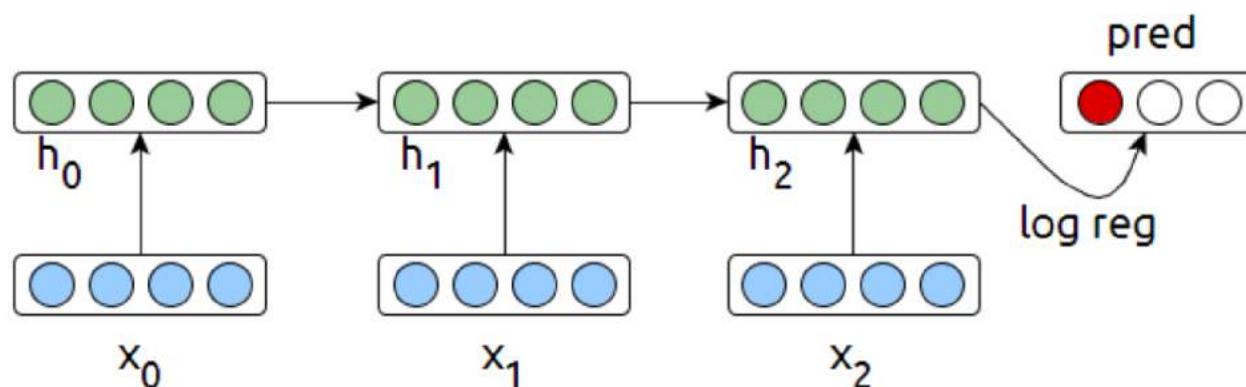
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

APPLICATIONS: TEXT CLASSIFICATION

- Binary classification (true, false)
 - Multi-class classification (politics, sports, gossip)
 - Multi-label classification (#party #FRIDAY #fail)
-
- Clustering (labels unknown)
-
- The diagram illustrates the classification of text applications. It features two main categories: 'Supervised Learning' and 'Unsupervised Learning'. 'Supervised Learning' covers binary, multi-class, and multi-label classification. 'Unsupervised Learning' covers clustering. A large green bracket groups the first three items together under 'Supervised Learning', while another green bracket groups the last item under 'Unsupervised Learning'.

TEXT CLASSIFICATION WITH RNN/LSTM

- So in order to classify text we can simply take a trained language model (RNN/LSTM), extract text representations from the final hidden state C_n , and apply softmax function.



AGENDA

- Introduction
- Background
 - Natural Language Processing
 - Language Models
 - Deep Learning
- Deep Learning for NLP
 - Neural Language Models
 - Recurrent Neural Network (RNN)
 - Long-Short Term Memory (LSTM)
- RNN/LSTM in Python

RNN/LSTM IN PRACTICE

■ Environment

- Python=2.7
- Anaconda 2018.12 OR
- Keras=2.2.4, TensorFlow=1.13.1

```
import numpy as np
import pandas as pd
import re
from bs4 import BeautifulSoup
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM
from keras.layers.embeddings import Embedding
from keras.preprocessing.sequence import pad_sequences
from keras.preprocessing.text import Tokenizer
from keras.utils.np_utils import to_categorical
```

RNN/LSTM IN PYTHON

```
data_train = pd.read_csv('data/labeledTrainData.tsv', sep='\t')

# read text data to sequences
texts = []
labels = []

for idx in range(data_train.sentiment.shape[0]):
    text = BeautifulSoup(data_train.review[idx], "lxml")
    texts.append(clean_str(text.get_text().encode('ascii', 'ignore')))
    labels.append(data_train.sentiment[idx])
```

▼ □ texts = {list} <type 'list'>: <Too big to print. Len: 25000>

01 Unable to handle: = {str} 'Too large to show contents. Max items to show: 300'
01 00000 = {str} 'with all this stuff going down at the moment with mj ive started listening to his music, watching the odd documentary here and there, ... View
01 00001 = {str} 'the classic war of the worlds by timothy hines is a very entertaining film that obviously goes to great effort and lengths to faithfully re... View
01 00002 = {str} 'the film starts with a manager (nicholas bell) giving welcome investors (robert carradine) to primal park . a secret project mutating a pr... View
01 00003 = {str} 'it must be assumed that those who praised this film (the greatest filmed opera ever, didnt i read somewhere?) either dont care for ope... View
01 00004 = {str} 'superbly trashy and wondrously unpretentious 80s exploitation, hooray! the pre-credits opening sequences somewhat give the false i... View
01 00005 = {str} 'i dont know why people think this is such a bad movie. its got a pretty good plot, some good action, and the change of location for ha... View
01 00006 = {str} 'this movie could have been very good, but comes up way short. cheesy special effects and so-so acting. i could have looked past that. View

▼ □ labels = {list} <type 'list'>: <Too big to print. Len: 25000>

01 Unable to handle: = {str} 'Too large to show contents. Max items to show: 300'
01 00000 = {int64} 1
01 00001 = {int64} 1
01 00002 = {int64} 0
01 00003 = {int64} 0
01 00004 = {int64} 1
01 00005 = {int64} 1
01 00006 = {int64} 0

RNN/LSTM IN PYTHON

```
# maximum number of words to keep, based on word frequency
MAX_NB_WORDS = 20000

# Tokenization
tokenizer = Tokenizer(num_words=MAX_NB_WORDS)
tokenizer.fit_on_texts(texts)
sequences = tokenizer.texts_to_sequences(texts)

# maximal length of sequence
MAX_SEQUENCE_LENGTH = 1000

# pad input sequences
data = pad_sequences(sequences, maxlen=MAX_SEQUENCE_LENGTH)
```

RNN/LSTM IN PYTHON

```
# converts a class vector (integers) to binary class matrix.
labels = to_categorical(np.asarray(labels))

# shuffle data
indices = np.arange(data.shape[0])
np.random.shuffle(indices)
data = data[indices]
labels = labels[indices]
print('Shape of data tensor:', data.shape)
print('Shape of label tensor:', labels.shape)

# training/testing data split
VALIDATION_SPLIT = 0.1
nb_validation_samples = int(VALIDATION_SPLIT * data.shape[0])
x_train = data[:data.shape[0]-nb_validation_samples]
y_train = labels[:data.shape[0]-nb_validation_samples]
x_test = data[data.shape[0]-nb_validation_samples:]
y_test = labels[data.shape[0]-nb_validation_samples:]

print('Number of positive and negative reviews in traing and validation set ')
print y_train.sum(axis=0)
print y_test.sum(axis=0)
```

RNN/LSTM IN PYTHON

```
# LSTM model.
```

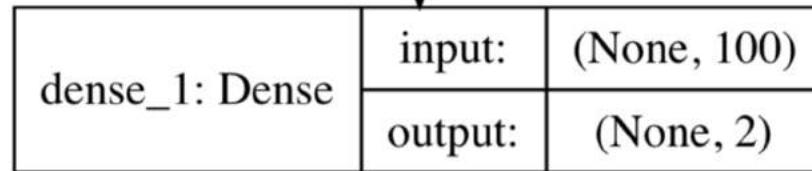
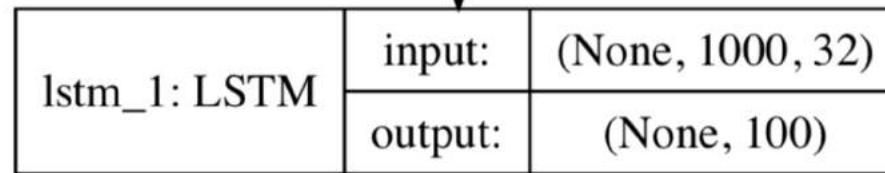
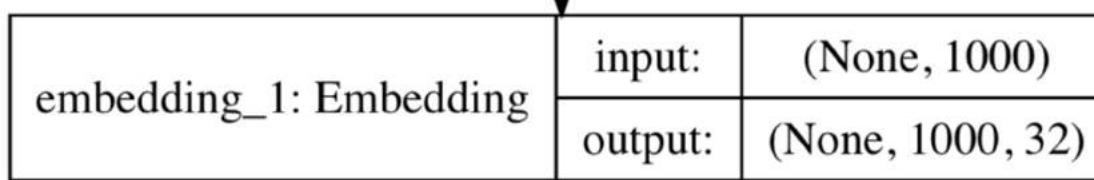
```
# The first layer is the Embedded layer that uses 32 length vectors to represent each word.  
embedding_vecor_length = 32  
model = Sequential()  
model.add(Em
```

```
# The next l  
model.add(LS
```

```
# Finally, b  
# to make 0  
model.add(De
```

```
# plot neur  
from keras.u  
plot_model(m
```

```
# Because it  
# The effici  
model.compil  
print(model.
```



RNN/LSTM IN PYTHON

```
# A large batch size of 64 reviews is used to space out weight updates.  
# The model is fit for 2 epochs because it quickly overfits the problem.  
model.fit(x_train, y_train, validation_data=(x_test, y_test), epochs=2, batch_size=64)  
  
# Final evaluation of the model  
scores = model.evaluate(x_test, y_test, verbose=0)  
print("Accuracy: %.2f%%" % (scores[1]*100))
```

Epoch 1/2

```
64/22500 [........................] - ETA: 11:28 - loss: 0.6930 - acc: 0.5000  
128/22500 [........................] - ETA: 8:57 - loss: 0.6930 - acc: 0.5078  
192/22500 [........................] - ETA: 8:15 - loss: 0.6933 - acc: 0.4896  
256/22500 [........................] - ETA: 7:42 - loss: 0.6933 - acc: 0.4922  
320/22500 [........................] - ETA: 7:19 - loss: 0.6936 - acc: 0.4781  
384/22500 [........................] - ETA: 7:10 - loss: 0.6932 - acc: 0.4948  
448/22500 [........................] - ETA: 7:02 - loss: 0.6933 - acc: 0.4799  
512/22500 [........................] - ETA: 6:55 - loss: 0.6933 - acc: 0.4821
```

Accuracy: 88.00%

RNN/LSTM IN PYTHON

- [Demo](#)
- Hyper-parameters fine-tuning
 - Number of layers
 - Number of hidden units per layer
 - Activation function
 - Optimizer
 - Learning rate
 - Batch size
 - Number of epochs
 - Initialization

■ Thank you!

■ Contact

- Email: wang.yanshan@mayo.edu
- LinkedIn, Twitter (yanshan_wang)