

Natural Language Processing using Deep Learning

Chapman Siu

August 2017

Agenda

1. Preliminaries (why?): differing representations
 - ▶ No Free Lunch Theorem
2. Understanding When to use Word2Vec
 - ▶ Term frequency
 - ▶ Topic models
 - ▶ Vectors?
3. Building Word2Vec from scratch
 - ▶ Unsupervised Learning
 - ▶ Extensions to FastText (enriching word2vec and supervised learning)
4. Coming back to CNNs...

1. Motivation

To understand the different parameters in word2vec models:

- ▶ What is skip-gram or CBOW?
- ▶ What is negative sampling?

This session is not. . .

- ▶ Detailed introduction to Neural Networks
- ▶ About deep learning

. . . though hopefully you will learn a bit about these things

(let me know if you want more theoretical sessions)

1. No Free Lunch

(AI theory)

There is no representation/algorithm/model that will outperform all other algorithms for any problem (paraphrased)

Sometimes term frequency/topic models/word2vec models are better, other times they are not.

1. Consequences of No Free Lunch (NFL)

Several Scenarios I recently encountered:

- ▶ I have running in production a model generating scores off 50 features
- ▶ I have new data coming in (50 features) which are shown to be equally predictive

What do I do?

1. Combine all features (100 features) and rebuild the model
2. Build a model using only the new 50 features and do an ensemble
3. ???

2. Understanding Word2Vec

Bag of words/Topic models/word2vec all aim to convert: word(s) to numbers

- ▶ Term frequency: word counts, can be normalised (TFIDF)
- ▶ Topic model: vector represents distribution of words, i.e. association to a particular topic (supervised or unsupervised)
- ▶ Word2Vec: some arbitrary vector in some vector space???

2. Understanding Word2Vec

Vectors allow you to measure things:

- ▶ How close vectors are (how similar)
- ▶ Are vectors orthogonal to each other (dissimilar)
- ▶ Can do “arithmetic”!

Examples from the original paper:

$$\text{vec}(\text{King}) - \text{vec}(\text{man}) + \text{vec}(\text{woman}) = \text{vec}(\text{Queen})$$

...and many other examples.

2. Understanding Word2Vec

Differences with other approaches:

- ▶ Context!

Word2Vec considers context of a word in its construction. The 2 approaches in “converting” the unsupervised problem to a supervised one:

- ▶ skip-gram: $Pr(\text{context}|\text{target word})$
- ▶ continuous bag of words (CBOW): $Pr(\text{target word}|\text{context})$

3. Building Word2Vec from Scratch (Building Training Set)

Skip-gram/CBOW

Training set construction:

1. Pick window size (odd number)
2. Extract all tokens based on this chosen window size
3. Remove the middle word in each window; this becomes your target word, other words are your context

3. Building Word2Vec from Scratch (Building Training Set)

Skip-gram (window size 3)

The cat sat on the mat

window size of 3:

- ▶ the cat sat
- ▶ cat sat on
- ▶ sat on the
- ▶ on the mat

3. Building Word2Vec from Scratch (Building Training Set)

Skip-gram (window size 3)

The cat sat on the mat

window size of 3:

context: the sat, target: cat

context: cat on, target: sat

context: sat them, target: on

context: on mat, target: the

Now we can perform some supervised learning!

3. Building Word2Vec from Scratch

With this frame work we can create a (shallow) neural network

$$Pr(\text{context}|\text{target word})$$

3. From Word2Vec to FastText

Some downsides to word2vec:

- ▶ Does not take into account neighbouring word information
- ▶ Cannot create word vectors for unseen vocabulary

FastText remedies both through use of n -grams, on both character and word level

3. From Word2Vec to FastText

- ▶ Character n -gram, for word where and $n = 3$

wh
whe
her
ere
re

- ▶ Word-level n -gram: cat sat on the mat

cat sat
cat sat on
sat on the
on the mat
the mat

3. From Word2Vec to FastText for Classification

To go from word2vec to classification problem:

1. Take each **word** representation which are *averaged* into a **text** representation
2. Add another layer for prediction of classification problem

4. Convolution Neural Networks for Text Models

What is our “image”?

$$[w_1, w_2, \dots, w_s]$$

Where $w_i \in \mathbb{R}^n$ is a word vector with a vocabulary size n , and s is max length of sentence.

- ▶ Convolutions allow taking neighbourhood information (context)
- ▶ It will generate feature maps of interest