

# Natural Language Processing

My Note



What is language?

Fixed vocabulary of words shared by humans

How do we express language as a mathematical model?

speech & text  
↑

NLP is broadly defined as the automatic manipulation of natural language through software

Tokenization, Stemming, Lemmatization

Tokenization: breaking down text into words  
→ commonly at spaces.

Stemming → getting base words and removing  
derivational affixes.

Lemmatization: Remove inflection endings added to a word that  
change its meaning

N-Grams

combining  $n$  words for representation purposes

1-gram / unigram model tokenizes into one word combinations

bigram - tokenizes into combinations of 2 words each token

Natural Language processing is essential

Natural Language processing is essential

up to n-gram model

## Transformation methods:

tf-idf — Score vocabulary to provide adequate weight of a word in proportion to impact of its meaning in a sentence.

product of 2 independent scores

→ tf (term frequency)

→ idf (inverse document frequency)  $\Rightarrow$  common / rare across documents

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$i$  in document  $j$ ,

$tf_{i,j}$  = # of occurrences of  $i$  in  $j$

$df_i$  = # of documents containing  $i$

$N$  = total # of documents

## One hot encodings:

represent words in numeric form

row length = length of vocabulary

column length = length of observation

The cat sat on the mat

the	1	0	0	0	1	0
Cat	0	1	0	0	0	0
on	0	0	0	1	0	0
:						

## Word embeddings:

words & phrases are mapped to vectors of real numbers

want similar words to be projected closer to each other

## Bag of words:

Tabular format  $\Rightarrow$  columns representing the total vocabulary

and row representing a single observation

Intersections are the count of the word at the column, in a particular observation.

$\Rightarrow$  Useful for spam & sentiment analysis

observation ID	word 1	word 2	word 3
1111	4	6	3
0011	0	1	8
0222	0	2	3

Disregards order / grammar, losing the context that a word is used

The importance of a word is often inversely related to its use  
in a sentence (limitation / problem)

Embedding form of a word  $\Rightarrow$  each word in a bag encoded form is multiplied by  
the embedding matrix to give word embeddings for the sample. (More on next page)

After one-hot encoding: cat  $\rightarrow [0, 1, 0]$  (One hot encoding)

Embedding matrix: Each row corresponds to the embedding vector for a word  $\rightarrow$  numbers to represent the "meaning" of a word. # of features  $\rightarrow$  # of columns

$$\begin{pmatrix} 2.3 & 1.5 \\ 0.7 & -0.9 \\ 0.2 & 0.5 \end{pmatrix} \rightarrow \begin{array}{l} \text{dog} \\ \text{cat} \\ \text{bird} \end{array} \quad (\text{Learned through training})$$

The result should be word vectors that capture the semantic meaning/context of each word in a continuous vector space, with each word being inputted w.r.t its context.

RNNs  $\Rightarrow$  after a word as input

$\Rightarrow$  can vary input length from words to sentences

Each input occurs at time  $t$ , and time  $t-1$  is also inputted in addition to  $t=t$  input.

Many to many architecture

Many to one architecture

One to many architecture

Many to many architecture

Lexical tokenization: conversion of text into meaningful tokens which are then categorized

⇒ Identifiers, operators, grouping symbols and data types

The lexer is responsible for tokenizing words

Infix: Affix inserted into word stem

↓  
attached to word stem, that forms a new word.

Stop words ⇒ Common words that are filtered out before NLP due to their insignificance

Sentiment analysis: understand the tone of a piece of writing

Tokenize the sentences → break down into words

remove the stop words

Condense all forms of words into their single form

Vectorize text: Turn text into numerical representation for the classifier

① Tokenization: Breaking down phrases into smaller pieces.

⇒ Sentence  $\nrightarrow$  word tokenization

② Remove stop words

Normalize words

Stemming: A word is cut off at its stem

③ Lemmatization (stemming but keeps words in their simplest form rather than truncating)  
only one offered by Spacy

This happens automatically with speech tagging and named entity recognition when nlp() is called

④ Vectorization of a text: Transforms a token into a vector, that uniquely represents the features of the token.

For sentiment analysis  $\Rightarrow$  use a dense array

vectorization is automatically done with nlp() call.

For classification, a basic workflow looks like:

Split data into training and evaluation sets.

Select a model architecture

use training data to train model

use test data to evaluate performance of model

use trained model to generate predictions (-1, 1)

Spacy provides a pipelining functionality in a JSON file

under en\_core\_web\_sm

Built-in pipeline components includes `textcat` enables you to assign categories / labels to a dataset, for training purposes

Add the `textcat` component to the existing pipeline

Add labels to the `textcat` component

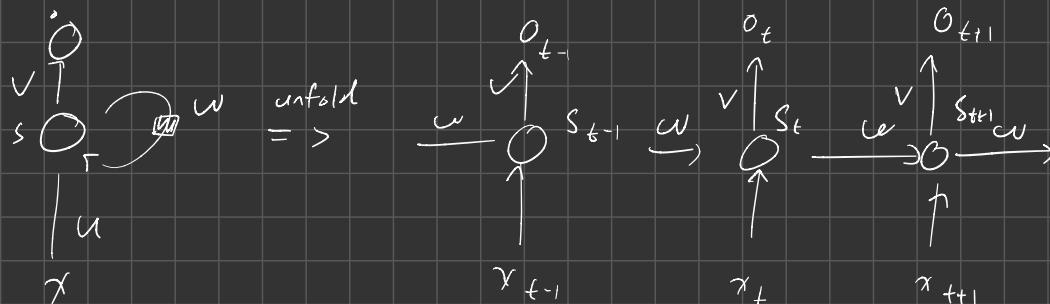
Load & shuttle data

Train model and evaluate each loop

## Introduction to RNNs

RNNs make use of sequential information

Perform the same task for every element of a sequence  
with the output depending on previous computations



$x_t$  is the input at time step  $t$ .

$s_t$  is the hidden state at time  $t$ .

Calculated based on the previous hidden state and input at current time step

$$s_t = f(u x_t + W s_{t-1})$$

↓  
usually non-linearity layer

$o_t$  is the output step at time  $t$

if we wanted

$o_t$  is the output at step  $t$ . If we wanted to predict the

next word in a sentence, it is a vector of probabilities  
across the vocabulary  $o_t = \text{softmax}(v_{st})$

$s_t$  can be thought of as the hidden memory of the network

$\Rightarrow$  Captures what happened in all the previous time steps.

$o_t$  is calculated solely based on memory at time  $t$

A RNN Shares the same parameters ( $M, V, W$ ) across all steps

Some times the output  $o_t$  is not needed

Language modelling  $\Leftrightarrow$  Generating text

predict the next word, given the previous words

We set  $O_t = x_{t+1} \Rightarrow$  output at step  $t$  to be the next actual word

Generating Image Descriptions

paired with CNNs

Similar to training a neural network

$\Rightarrow$  Backpropagation becomes Backpropagation through time

Vanilla RNNs have difficulty learning long-term dependencies

$\Rightarrow$  Vanishing / exploding gradient problem

## LSTM Networks

memory are called cells

Block boxes that take as input the previous state  $h_{t-1}$  and current input  $x_t$ .

They decide what to keep and erase from memory.

Combine previous state, current memory, and input

# Language Models from Scratch

Goal: Assign a probability for all sentences

$10^{50}$  (10 words per sentence)

Graph:

Early -> one -> morning -> the -> sun was -> laying in bed  
shining -> I

Can get crazy samples

$P(z_n | z_{n-1})$  (not very insightful)

Need to cover long range dependencies

Fourier / Taylor Series

Neural networks

$$\sin(x) = \frac{x^2}{\sqrt{2}}, \text{ just need input/output pairs}$$

Word2Vec, Glove

Attention  $\times$  words = probabilities

Context for every word + original words

Stacking Attention allows for complex relationships