

Natural Language Processing

Portfolio I

Leon F.A. Wetzel
Information Science
Faculty of Arts - University of Groningen
l.f.a.wetzel@student.rug.nl

February 26, 2021

Abstract

In this document, you can find the results and explanations for the assignments of the first part of the portfolio for the course Natural Language Processing, taught at the University of Groningen. The corresponding Python code can be found at <https://github.com/leonwetzel/natural-language-processing>¹. Note that version control of Jupyter notebooks is done via `jupyter`, so do not forget to convert the relevant Python scripts to notebooks yourself!

1 Week 1 - Regular Expressions

We used <https://regex101.com/> to test our regular expressions. In this document, we use \wedge to better display/represent a caret. Spaces cannot be displayed in \LaTeX sadly, so these are not shown below.

1.1 Regular Expressions I

1. Set of all alphabetic strings.
 $\wedge[a-zA-z]^+\$$
2. Set of all lower case alphabetic strings ending in b .
 $\wedge[a-z]^*[b]\$$
3. Set of all strings from the alphabet a, b such that each a is immediately preceded by a b and immediately followed by a b .
 $\wedge(bab)^+\$$

1.2 Regular Expressions II

1. Set of all strings with two consecutive repeated words (e.g., “*Humbert Humbert*” and “*the the*” but not “*the bug*” or “*the big and the bug*”).
 $\wedge b(\wedge w^+)\wedge b\wedge s^+\wedge 1$
2. All strings that start at the beginning of the line with an integer and that end with a word.
 $\wedge \wedge d^+ \wedge . \wedge \wedge b[a-zA-Z]^+ \wedge b\$$
3. All strings that have both the word *grotto* and the word *raven* in them (but not, e.g., words like *grottos* that merely contain the word *grotto*).
 $\wedge bgrotto\wedge b. \wedge \wedge braven\wedge b.*$

¹All code will be published after the course has been completed

1.3 ELIZA

We added extra regular expressions and substitutions to handle more variants of user input. We also implemented a substitution to handle question marks, so that ELIZA's responses feel more like proper questions when relevant. The other implemented substitutions are related to how one feels, and takes into account words such as *encounter*, *experience* and *suffer from*. Snippets from example dialogue can be found below.

ELIZA : How are you?

LEON : I experience coldness.

ELIZA : Why do you experience coldness?

LEON : I feel as if it is snowing today

ELIZA : Why do you feel as if it is snowing today?

LEON : Because it is snowing today

ELIZA : Because it is snowing today?

LEON : Yes, but I am inside

ELIZA : Yes, but Why are you inside?

LEON : Because it is cold outside

ELIZA : Because it is cold outside?

LEON : Yes

1.4 Byte-Pair Encoding

Our example text for the experiment is a snippet from the Wikipedia page on the GameStop short squeeze², which occurred in January 2021. According to <https://wordcounter.net>, our example consists of 137 words and 827 characters.

"In January 2021, a short squeeze of the stock of the American video game retailer GameStop (NYSE: GME) and other securities took place, causing major financial consequences for certain hedge funds and large losses for short sellers. Approximately 140 percent of GameStop's public float had been sold short, and the rush to buy shares to cover those positions as the price rose caused it to rise even further. The short squeeze was initially and primarily triggered by users of the subreddit r/wallstreetbets, an Internet forum on the social news website Reddit. At its height, on January 28, the short squeeze caused the retailer's stock price to reach a pre-market value of over US\$500 per share, nearly 30 times the \$17.25 valuation at the beginning of the month. The price of many other heavily shorted securities increased."

The default value for `vocab_size` (v) is 30000. With this value, we can count 164 words and 190 segments found in our example text by the tokenizer. For $v = 25000$, we count 164 words and 195 segments, and the counts for $v = 20000$ are 164 words and 197 segments. In line with the assignment description, we can observe here that more words are segmented into subwords when the number of segments increase (as a result of v decreasing). An overview of the results can be found in table 1.

v	30000	25000	20000	15000	10000	7500	5000	2500
Words	164	164	164	164	164	164	164	164
Segments	190	195	197	204	222	246	692	692

Table 1: Overview of `vocab_size` values and the word and segment counts

²https://en.wikipedia.org/wiki/GameStop_short_squeeze

When $v = 7500$, the number of segments (246) is exactly 150% of the number of words (164). The longest words in our example text during this setup that were not segmented, are *increased*, *beginning* and *initially*. All these words consist of nine characters each.

2 Week 2 - N-gram Language Models

For this week's assignment, we take a look at several exercises from [Jurafsky and Martin \(2020\)](#).

2.1 J&M exercise 3.1

Write out the equation for trigram probability estimation (modifying Eq. 3.11). Now write out all the non-zero trigram probabilities for the I am Sam corpus on page 41.

The original - simplified - formula for bigram probability estimation is as follows:

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

For trigram probability estimation, the formula would be:

$$P(w_n|w_{n-2}w_{n-1}) = \frac{C(w_{n-2}w_{n-1}w_n)}{C(w_{n-2}w_{n-1})}$$

Our corpus contains the following sentences:

<s> I am Sam </s> <s> Sam I am </s> <s> I do not like green eggs and ham </s>

The non-zero trigram probabilities would be:

$$\begin{aligned}
 P(\text{Sam}|\text{I am}) &= \frac{C(\text{I am Sam})}{C(\text{I am})} = \frac{1}{2} = 0.5 \quad (1) & P(\text{like}|\text{do not}) &= \frac{C(\text{do not like})}{C(\text{do not})} = \frac{1}{1} = 1.0 \quad (9) \\
 P(\text{am}|\text{<s> I}) &= \frac{C(\text{<s> I am})}{C(\text{<s> I})} = \frac{1}{2} = 0.5 \quad (2) & P(\text{green}|\text{not like}) &= \frac{C(\text{not like green})}{C(\text{not like})} = \frac{1}{1} = 1.0 \quad (10) \\
 P(\text{</s>}|\text{am Sam}) &= \frac{C(\text{am Sam </s>})}{C(\text{am Sam})} = \frac{1}{1} = 1.0 \quad (3) & P(\text{eggs}|\text{like green}) &= \frac{C(\text{like green eggs})}{C(\text{like green})} = \frac{1}{1} = 1.0 \quad (11) \\
 P(\text{I}|\text{<s> Sam}) &= \frac{C(\text{<s> Sam I})}{C(\text{<s> Sam})} = \frac{1}{1} = 1.0 \quad (4) & P(\text{and}|\text{green eggs}) &= \frac{C(\text{green eggs and})}{C(\text{green eggs})} = \frac{1}{1} = 1.0 \quad (12) \\
 P(\text{am}|\text{Sam I}) &= \frac{C(\text{Sam I am})}{C(\text{Sam I})} = \frac{1}{1} = 1.0 \quad (5) & P(\text{ham}|\text{eggs and}) &= \frac{C(\text{eggs and ham})}{C(\text{eggs and})} = \frac{1}{1} = 1.0 \quad (13) \\
 P(\text{</s>}|\text{I am}) &= \frac{C(\text{I am </s>})}{C(\text{I am})} = \frac{1}{2} = 0.5 \quad (6) & P(\text{</s>}|\text{and ham}) &= \frac{C(\text{and ham </s>})}{C(\text{and ham})} = \frac{1}{1} = 1.0 \quad (14) \\
 P(\text{do}|\text{<s> I}) &= \frac{C(\text{<s> I do})}{C(\text{<s> I})} = \frac{1}{2} = 0.5 \quad (7) \\
 P(\text{not}|\text{I do}) &= \frac{C(\text{I do not})}{C(\text{I do})} = \frac{1}{1} = 1.0 \quad (8)
 \end{aligned}$$

2.2 J&M exercise 3.2

Calculate the probability of the sentence *i want chinese food*. Give two probabilities, one using Fig. 3.2, and another using the add-1 smoothed table in Fig. 3.6.

Our sentence is: <s> i want chinese food </s>

The formula and result using figure 3.2 would be:

$$\begin{aligned}
 P(< s> \text{ i want chinese food } </s>) &= P(\text{i}|<s>) * P(\text{want}|\text{i}) * P(\text{chinese}|\text{want}) * P(\text{food}|\text{chinese}) * P(</s>|\text{food}) \\
 &= 0.25 * 0.33 * 0.0065 * 0.00092 * 0.68 \\
 &= \mathbf{0.00000034}
 \end{aligned}
 \tag{15}$$

The formula and result using the add-1 smoothness from figure 3.6 would be:

$$\begin{aligned}
 P_{\text{Laplace}}(< s> \text{ i want chinese food } </s>) &= P_{\text{Laplace}}(\text{i}|<s>) * P_{\text{Laplace}}(\text{want}|\text{i}) * P_{\text{Laplace}}(\text{chinese}|\text{want}) \\
 &\quad * P_{\text{Laplace}}(\text{food}|\text{chinese}) * P_{\text{Laplace}}(</s>|\text{food}) \\
 &= 0.25 * 0.21 * 0.0029 * 0.052 * 0.68 \\
 &= \mathbf{0.0000538356}
 \end{aligned}
 \tag{16}$$

2.3 J&M exercise 3.6

Suppose we train a trigram language model with add-one smoothing on a given corpus. The corpus contains V word types. Express a formula for estimating $P(w_3|w_1, w_2)$, where w_3 is a word which follows the bigram (w_1, w_2) , in terms of various N -gram counts and V . Use the notation $c(w_1, w_2, w_3)$ to denote the number of times that trigram (w_1, w_2, w_3) occurs in the corpus, and so on for bigrams and unigrams.

The original formula with Laplacian smoothing is:

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V} \tag{17}$$

A similar formula supporting trigrams would look like:

$$P_{\text{Laplace}}^*(w_n|w_{n-1}w_{n-2}) = \frac{C(w_{n-2}w_{n-1}w_n) + 1}{\sum_w (C(w_{n-2}w_{n-1}w) + 1)} = \frac{C(w_{n-2}w_{n-1}w_n) + 1}{C(w_{n-2}w_{n-1}) + V} \tag{18}$$

2.4 J&M exercise 3.7

We are given the following corpus, modified from the one in the chapter:

<s> I am Sam </s>
 <s> Sam I am </s>
 <s> I am Sam </s>
 <s> I do not like green eggs and Sam </s>

If we use linear interpolation smoothing between a maximum-likelihood bi-gram model and a maximum-likelihood unigram model with $\lambda_1 = 1/2$ and $\lambda_2 = 1/2$, what is $P(\text{Sam}|\text{am})$? Include <s> and </s> in your counts just like any other token.

By the use of linear interpolation smoothing between the maximum-likelihood n -gram models (with $\lambda_1 = 1/2$ and $\lambda_2 = 1/2$), the formula would look like:

$$\hat{P}(w_n|w_{n-1}) = \lambda_1 P(w_n|w_{n-1}) + \lambda_2 P(w_n) \tag{19}$$

$$\begin{aligned}
 \hat{P}(\text{Sam}|\text{am}) &= 1/2 * P(\text{Sam}|\text{am}) + 1/2 * P(\text{Sam}) \\
 &= 1/2 * \frac{C(\text{am Sam}) + 1}{C(\text{am}) + 1} + 1/2 * \frac{C(\text{am}) + 1}{C(\text{am}) + 1} \\
 &= 1/2 * \frac{2 + 1}{3 + 1} + 1/2 * 1 \\
 &= 0.375 + 0.5 \\
 &= \mathbf{0.875}
 \end{aligned}
 \tag{20}$$

2.5 N-grams in the notebook

What are the different and common n-grams? You may ignore n-grams containing punctuation symbols.

We looked for common and different n-grams between the speeches of Theodore Roosevelt and Franklin D. Roosevelt. The common bigrams are:

('can', 'not'), ('is', 'the'), ('the', 'country'), ('have', 'been'), ('there', 'is'), ('of', 'the'), ('to', 'make'), ('the', 'people'), ('that', 'the'), ('and', 'to'), ('is', 'a'), ('the', 'other'), ('of', 'a'), ('and', 'in'), ('in', 'this'), ('the', 'same'), ('they', 'are'), ('is', 'not'), ('has', 'been'), ('it', 'is'), ('with', 'the'), ('and', 'the'), ('united', 'states'), ('the', 'nation'), ('the', 'government'), ('on', 'the'), ('to', 'the'), ('the', 'congress'), ('in', 'the'), ('all', 'the'), ('by', 'the'), ('of', 'our'), ('not', 'be'), ('for', 'the'), ('as', 'a'), ('to', 'be'), ('of', 'all'), ('from', 'the'), ('of', 'this'), ('in', 'a'), ('can', 'be'), ('we', 'have'), ('at', 'the'), ('the', 'united'), ('will', 'be')

The different bigrams are:

('is', 'to'), ('under', 'the'), ('to', 'do'), ('the', 'great'), ('such', 'a'), ('should', 'be'), ('the', 'public'), ('during', 'the'), ('it', 'has'), ('be', 'made'), ('the', 'isthmus'), ('and', 'of'), ('the', 'republic'), ('as', 'to'), ('the', 'present'), ('to', 'secure'), ('not', 'only'), ('as', 'the'), ('the', 'work'), ('of', 'panama'), ('which', 'the'), ('may', 'be'), ('the', 'most'), ('upon', 'the'), ('the', 'national'), ('the', 'law'), ('in', 'which'), ('republic', 'of'), ('of', 'their'), ('must', 'be'), ('would', 'be')

The common trigrams are:

('people', 'of', 'the'), ('of', 'the', 'government'), ('can', 'not', 'be'), ('of', 'the', 'nation'), ('it', 'is', 'a'), ('it', 'is', 'not'), ('of', 'the', 'congress'), ('the', 'government', 'of'), ('as', 'a', 'whole'), ('and', 'in', 'the'), ('the', 'united', 'states'), ('of', 'the', 'united'), ('of', 'the', 'world'), ('to', 'the', 'congress'), ('part', 'of', 'the'), ('of', 'the', 'country'), ('that', 'it', 'is'), ('the', 'fact', 'that'), ('the', 'people', 'of'), ('one', 'of', 'the'), ('as', 'well', 'as'), ('of', 'the', 'people'), ('it', 'will', 'be'), ('in', 'order', 'to')

The different trigrams are:

('the', 'department', 'of'), ('of', 'the', 'republic'), ('some', 'of', 'the'), ('the', 'interests', 'of'), ('interest', 'of', 'the'), ('should', 'be', 'made'), ('to', 'the', 'united'), ('so', 'as', 'to'), ('government', 'of', 'the'), ('a', 'matter', 'of'), ('is', 'to', 'be'), ('of', 'the', 'present'), ('but', 'it', 'is'), ('the', 'interest', 'of'), ('has', 'been', 'made'), ('on', 'the', 'isthmus'), ('of', 'the', 'navy'), ('on', 'the', 'other'), ('as', 'to', 'the'), ('in', 'the', 'interest'), ('there', 'has', 'been'), ('the', 'rights', 'of'), ('and', 'it', 'is'), ('in', 'connection', 'with'), ('the', 'national', 'government'), ('the', 'number', 'of'), ('the', 'work', 'of'), ('the', 'secretary', 'of'), ('of', 'the', 'national'), ('of', 'the', 'canal'), ('of', 'the', 'law'), ('there', 'is', 'no'), ('attention', 'to', 'the'), ('the', 'question', 'of'), ('secretary', 'of', 'the'), ('it', 'would', 'be'), ('of', 'the', 'great'), ('so', 'far', 'as'), ('republic', 'of', 'panama'), ('of', 'the', 'public'), ('it', 'should', 'be'), ('there', 'should', 'be'), ('the', 'republic', 'of'), ('the', 'use', 'of'), ('and', 'of', 'the')

Describe in some detail (but in at most approximately 100 words) how one can implement a generation program based on a trigram language model.

Trigram language models can be used effectively in a text generation program by exploiting the probabilities of words occurring in the same context (i.e. using maximum likelihood estimates). This approach could presumably lead to generated sentences that could directly replicate the training corpus if the highest probabilities are picked for text generation. On the other hand, trigrams are less unique than bigrams and far less unique than unigrams. The use of trigrams can lead to more comprehensible sentences to a certain degree.

3 Week 3 - POS Tagging and Hidden Markow Models

For this week's assignment, we use data from the LassySmall corpus in exercises 2 and 4 ([van Noord et al., 2013](#); [Bouma and van Noord, 2017](#)).

3.1 Error detection

Find an error in each of the following sentences, labelled with Universal Dependency Part of Speech tags. Check with the on-line documentation, or search for similar examples in this search engine (any of the English corpora (development version) will do, except for ESL (sign language)).

1. They/PRON promise/~~NOUN~~ **VERB** them/PRON the/DET delights/NOUN ./PUNCT mostly/ADV sexual/ADJ ./PUNCT of/ADP the/DET next/ADJ world/NOUN
2. The/DET incident/NOUN proves/VERB that/SCONJ Sharon/PROPN has/AUX lost/VERB his/~~DET~~ **PRON** patience/NOUN
3. Over/ADP the/DET last/ADJ twenty/~~DET~~ **NUM** years/NOUN ./PUNCT use/NOUN of/ADP the/DET term/NOUN has/AUX been/AUX frequently/ADV criticized/VERB
4. I/PRON do/AUX quickly/~~ADJ~~ **ADV** want/VERB to/PART comment/VERB
5. I/PRON just/ADV wanted/VERB to/PART ask/VERB whether/~~ADV~~ **SCONJ** you/PRON read/VERB the/DET draft/NOUN

3.2 POS tagger

Implement a baseline system for your corpus that assigns each word its most frequent PoS tag. Empty lines and lines starting with a # can be ignored. Also lines starting with a code that is not an integer (i.e. 13.1 or 13-14) can be ignored. Data is tab-separated. The second column is the word, the fourth column the POS-tag. All other columns can be ignored. Collect statistics for the most frequent PoS per word from the *-train.conllu file and compute the performance of the baseline method on the *-dev.conllu. Dealing with unknown words requires special attention. Report your accuracy for the 2 versions of the baseline: one that has the simplest strategy for handling unknowns and scores for an improved version as suggested in the book.

The tagger with the *simple* assignment strategy (assign UNK to unknown words) obtained an accuracy score of 0.8151069050122678. After applying the improvement suggested in [Jurafsky and Martin \(2020\)](#) (assign most frequent POS tag NOUN to unknown words), we obtained an accuracy score of 0.8751314405888538.

3.3 HMM tagger

Consider a bigram HMM tagger. Work out in detail the steps of the Viterbi algorithm in tagging the sentence *zij zag hem slapen* where we have the following probabilities (start and end are used to represent the start-state and final-state respectively). The first columns give probabilities for tag sequences ($P(T_n|T_{n-1})$), the final rows for word emission probabilities ($P(W|T)$). So the first row says $P(\text{DET}|\text{start}) = 0.1$ and $P(\text{zij}|\text{PRON}) = 0.75$.

Our sentence is: <s> zij zag hem slapen </s>

3.4 Testing the baseline

The notebook `hmm-ud.ipynb` (or the python script `hmm-ud.py` if you prefer a command line) can be used to train and test a bigram PoS-tagger using the Hidden Markov Model on Universal Dependencies data. The system takes two arguments: one for the UD file containing the training data, and one for testing. Use the same corpus as you used for assignment 2. Use the dev data for initial testing and development. Report the accuracy of the baseline hmm tagger on your data. Now improve the handling of unknown words: try to improve the predictions made by the function `unknown_word_guesser(word)` in such a way that it no longer blindly assumes each unknown word is a NOUN but instead looks e.g. the form of the word (numbers, suffixes, upper case, ...) to make better predictions. See also the discussion in the book. Note that it is also possible to assume an unknown word is always either a NOUN or a VERB, etc. You can use the flag `-v` to obtain some feedback on tagging errors of the system (also showing which words were unknown). Use only the development corpus to improve the unknown word guesser. Report final accuracy scores for the test corpus (i.e. using `test` instead of `dev` for evaluation).

The baseline HMM tagger obtained an accuracy score of 0.8554756054756055 on the test set. After this run, we made adjustments to the handling of unknown words by adding probabilities for PROP, PUNCT and VERB ($P(W|T) = 0.0001$). This increased the accuracy score to 0.868989118989119.

The next step in optimization was checking suffixes, prefixes and capitalisation of letters. After enabling the `-verbose` parameter, we observed that ADJ often end in *er*. We added a check for this suffix and set the probability for ADJ to 0.01. We also checked if the word is numeric, and set the probability for NUM to 0.3. Lastly, we checked if a word starts with a capital letter and set the probability for PROP to 0.5. These improvements - along with $P(W|T) = 0.0001$ for both NOUN and VERB - increased the accuracy score to 0.8852228852228852. An extra check for the suffix *ge* for ADJ raises this number to 0.887065637065637.

For our final adjustments, we added the suffix *ees* for AJD and adjusted the probabilities for PROP (0.1), NUM (0.3) and ADJ (0.001). This lead to the highest accuracy score of 0.8896103896103896.

3.5 Viterbi complexity

What is the complexity of the Viterbi algorithm for tagging with bigrams, in terms of the length of the sentence N , and the number of tags T . Explain your answer.

$$O(N^T)$$

4 Week 4: Distributional Semantics and Word Embeddings

4.1 Cosine similarity

Given the word co-occurrence matrix above, compute the cosine similarity of Apple and Snow, and of Snow and Ski. Is Snow more similar to Apple or Ski?

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

$$\text{cosine}(\text{Apple}, \text{Snow}) = \frac{2 * 2 + 1 * 3 + 5 * 2}{\sqrt{2^2 + 1^2 + 5^2} \sqrt{2^2 + 3^2 + 2^2}} = \frac{17}{\sqrt{30} \sqrt{17}} \approx 0.752773$$

$$\text{cosine}(\text{Snow}, \text{Ski}) = \frac{2 * 0 + 3 * 4 + 2 * 5}{\sqrt{2^2 + 3^2 + 2^2} \sqrt{0^2 + 4^2 + 5^2}} = \frac{22}{\sqrt{17} \sqrt{41}} \approx 0.833309$$

Snow is more similar to Apple than Ski in this context.

4.2 Pointwise mutual information

Given the co-occurrence matrix above, compute the pointwise mutual information (using eq 6.17) of Mountain and Ski and of Apple and Fresh. Assume that the corpus size is 100.

$$\text{PMI}(w_1, w_2) = \log_2 \frac{P(w_1, w_2)}{P(w_1)P(w_2)} \quad N = 100$$

Mountain and Ski:

$$P(\text{Mountain}, \text{Ski}) = \frac{5}{100} = 0.05$$

$$P(\text{Mountain}) = \frac{12}{100} = 0.12 \quad P(\text{Ski}) = \frac{9}{100} = 0.09$$

$$\text{PMI}(\text{Mountain}, \text{Ski}) = \log_2 \frac{0.05}{0.12 * 0.09} \approx \log_2(4.629636) \approx \mathbf{0.665546}$$

Apple and Fresh:

$$P(\text{Apple}, \text{Fresh}) = \frac{2}{100} = 0.02$$

$$P(\text{Apple}) = \frac{8}{100} = 0.08 \quad P(\text{Fresh}) = \frac{4}{100} = 0.04$$

$$\text{PMI}(\text{Apple}, \text{Fresh}) = \log_2 \frac{0.02}{0.08 * 0.04} = \log_2(6.25) \approx \mathbf{0.795880}$$

4.3 Twitter word clusters

Word2vec represents word meaning as a vector based on distributions of these words in a corpus. One way to visualise the result is to create clusters of similar words. The Tweet NLP project of CMU provides clusters for embeddings trained English Tweets, and this is a similar demo for Dutch.

In the English clusters, can you find two examples of emphatic lengthening? (Emphatic lengthening is the phenomenon where some characters in a word are repeated to give extra stress to the word, Ha, haha, hahaha, etc.). Give the cluster ids and some examples of words in this cluster.

It appears that clusters 001010111 and 00101110 include samples with emphatic lengthening. Cluster 001010111³ contains words related to **ever**, and includes words such as *eva*, *evan*, *everrr*, *evah* and *evur*. Cluster 00101110⁴ contains words related to **really**, and includes words such as *rly*, *realy*, *rlly*, *reali* and *reaaaly*.

³<https://www.cs.cmu.edu/~ark/TweetNLP/paths/001010111.html>

⁴<https://www.cs.cmu.edu/~ark/TweetNLP/paths/00101110.html>

In the Dutch clusters, can you find 2 clusters consisting of predominantly English words? Give the cluster id and 5 representative words per cluster.

Clusters 178 and 228 contain a fair amount of predominantly English words. Cluster 178 contains words such as *livestream*, *songfestival*, *popstars*, *friendzone* and *xfactor*. Cluster 228 contains words such as *track*, *clip*, *tune*, *freestyle* and *underground*.

In the Dutch data, can you find a cluster consisting predominantly of words from a language other than Dutch or English?

Cluster 839⁵ seems to include Swedish words, such as *kommer*, *vi*, *til*, *han*, *var* and *jeg*.

4.4 Analogies in Word2Vec

Word Embeddings are trained on huge text corpora. They are sometimes evaluated by checking top N lists of similar words, or by solving analogy problems. It has been observed that word embeddings may contain (undesirable) bias. (see section 6.11 of J&M, and Bolukbasi et al (2016), Man is to computer programmer as woman is to homemaker? and Manzini et al, 2019, Black is to criminal as caucasian is to police).

Try to find two more examples using the on-line demo (choose the corpus trained on English Google news). You can either use analogies (the logic in the analogy demo is word1 is to word2 as word3 is to word4, or in terms of vectors: the most similar word for vector word2-word1+word3 is word4 ie most-similar(king man+woman)=queen, most-similar(Amsterdam-Netherlands Amsterdam + Germany)= Berlin). Note that in many analogy cases, word4 is also the most similar word of word2 if you just check for the most similar words to word2 (you can check this also in the demo), which suggests there is no bias in this case. So try to find cases where the outcome for word4 really depends on the presence of word3. An alternative method to detect bias is to check directly for the similarity of pairs of words, i.e. is similarity(man,programmer) different from similarity(woman,programmer). You may use this method as well.

The similarity score for {man, doctor} is 0.3144896, whereas the similarity score for {woman, doctor} is 0.37945858. The similarity score for {man, engineer} is 0.15128928, whereas the similarity score for {woman, engineer} is 0.09435379. If we take a look at the analogy demo and feed it with the triple {man, doctor, woman}, then the demo's top 3 words are {gynecologist, nurse, doctors}.

Explain the problem and why you think this might be a problem in some settings/applications.

A model is as good as its (training) data, and the bias(es) from its author(s) can have severe impact on the display of the world in corpora. If systems rely on these analogies, they can potentially discriminate certain groups in a systematic way.

Discuss possible solutions. (max approx 100 words.)

Possible solutions to combat this challenge could involve the use of corpora from various sources, instead of a single website or news outlet, and proper evaluations of similarities and analogies. Based on these evaluations, the developers can act on possible cases of discrimination and adjust the training data. Another solution could include adjusting the training parameters, and comparing the new word embeddings with the older word embeddings when it comes to similarity scores and analogies.

⁵<http://www.let.rug.nl/gosse/twitter-clusters/paths/839.html>

References

- Gosse Bouma and Gertjan van Noord. Increasing return on annotation investment: The automatic construction of a Universal Dependency treebank for Dutch. In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 19–26, Gothenburg, Sweden, May 2017. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W17-0403>.
- Daniel Jurafsky and James H Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, December 2020.
- Gertjan van Noord, Gosse Bouma, Frank Van Eynde, Daniël de Kok, Jelmer van der Linde, Ineke Schuurman, Erik Tjong Kim Sang, and Vincent Vandeghinste. *Large Scale Syntactic Annotation of Written Dutch: Lassy*, pages 147–164. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-3-642-30910-6. doi: 10.1007/978-3-642-30910-6_9. URL https://doi.org/10.1007/978-3-642-30910-6_9.