

Part I

Natural Language Processing

1 Introduction

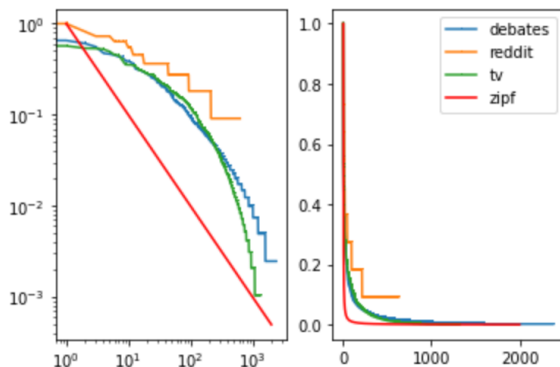
Challenges

- Complicated structure in sentence
- Syntactic ambiguities (“*time flies like an arrow*”, “*get the cat with the gloves*”).
- Metaphores, humor, irony, ...
- Semantics can be very rich and dependent on context, not easy to distinguish
- Language requires knowledge about the world
- Hard to really formalise the notion of “meaning”

Basic approach Gather information on word based on its context. Given a large text corpus, we can assume this to be meaningful statistics.

Zipf's Law The number of elements with a given frequency follows a power law distribution. That is, there is a small number of elements which appear very often and the majority of elements appears rarely. In its most simple form, the probability of the n -th most common word $p(n)$ is

$$p(n) = 1/n$$



Elements of language We can have different points of view on language:

- Phonetics (sound)
- Grammar

- Phonology
- Morphology
- Syntax
- Semantics (meaning)

Morphology How words are built up from smaller meaningful units, for instance *un-lady-like*, *dog-s*

- **Inflection** – variation in the form of a word (usually affix) that expresses a grammatical contrast
 - adds tense, number, person, mood, aspect, etc
 - e.g. *run* → *run* — *running*
 - does not change word class
- **Derivation** – formation of a new word from another
 - e.g. nominalization (*computer* → *computerization*)
 - e.g. formation of adjectives (*computational*, *clueless*)
 - changes word class

Morphemes

- **Root** – equivalence class of a word when all affixes are removed; not further decomposable into meaningful elements.
- **Stem** – part of word that never changes when morphologically inflected, *i.e. without affixes describing tense, number, person, ...*
- **Lemma** – Base form of word
- From *produced*, lemma is *produce* but stem is *produc*

2 Tokenization

Token an individual occurrence of a word (as opposed to a vocabulary/dictionary item)

Challenges

- Keep abbreviations, dates, numbers as single tokens
- distinguish abbreviations from words
- names and phrases (*queen of england*, *TU Wien*)
- compound words
- apostrophes, umlauts, etc and other linguistic characteristics
- encoding issues like RTL/LTR

Maximum Matching algorithm Use a dictionary of known terms. Take the longest prefix of the input string that matches a dictionary item. Does not always make sense (“*Theta bled own there*”).

3 Stemming

Stemming/Lemmatization reduce tokens to equivalence classes. Usually to gather words that are morphologically different but semantically quite similar to the same set, i.e. to improve comparability.

Porter Stemmer Rules for stripping suffixes. Applicability of rules is based on *measure* of a word w , which is the number m s.t. $w = C(VC)\{m\}V$ where C, V are arbitrary sequences of consonants, vowels, resp. — Indeed reduces the words to their *stems*.

WordNet MORPHY

- Has a sophisticated set of rules about inflections
- exception list
- Checks the result of transformation against an extensive dictionary

Note MORPHY reduces to *lemmas*, while PORTER reduces to *stems*.

- **Over-stemming** Two words are reduced to the same root when they should not be
- **Under-stemming** Should be reduced to the same root but are not.

4 POS-Tagging

Given some input text and some tags (usually word types such as *noun*, *verb*, etc.), want to assign tags to tokens. (*Sequence classification problem*).

Tagging can help with other procedures such as stemming, NER, parsing, ...

Def. Can divide words into two different classes

- **Closed class** — can enumerate all members, e.g. determiners, pronouns, prepositions, ...
- **Open class** — don't know all members, e.g. nouns, verbs, adjectives, ...

Note

- A single term (dictionary entry) can have different optimal POS-tags depending on its context.
- Tagging helps to resolve ambiguities that exist on term-level (e.g. *leaves* as NN or as VB)

- Tagging removes unnecessary distinctions e.g. all personal pronouns are PRP, determiners
- Naive method (assigning most frequent tag in training data to term) already has 90% accuracy.

Def.

- **Informativeness** — Assignment of tag adds information, reduces ambiguity
- **Specifiability** — Ease of mapping a term to a tag
- **Example:** Collapsing multiple related tags into one decreases informativeness, decreases specifiability.

Feature selection Can look at word-local features (term, pre-, suffixes, capitalization); but very often the tag of a word depends on its context in the sentence.

Main techniques

- **Probabilistic tagging** — consider lexical frequencies of tag in training data – good when large training corpora are available
- **Rule-based tagging** — use rules based on linguistic understanding – good to tailor solution to very specific problems

4.1 Probabilistic tagging

Consider the definition of *conditional probability*

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

This gives rise to the *chain rule*

$$\Rightarrow P(A, B) = P(A|B) \cdot P(B)$$

or, more generally

$$P(w_1, \dots, w_n) = \prod P(w_i | w_1, \dots, w_{i-1})$$

The problem here is that we cannot realistically obtain all the components of the product because there are way too many possible sequences of words. Instead, we employ the *Markov assumption* that says that we can estimate the probabilities by only considering only the k preceding terms

$$P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-k}, \dots, w_{i-1})$$

For $k = 1$, this yields the *unigram model*, for $k = 2$ the *bigram model* (i.e. $P(w_i | w_{i-1})$).

n -gram modelling is insufficient because language has *long-distance dependencies*.

Unigram Tagger Assume that a unigram model generates the current tagging.

Assign a token w its most frequent tag, i.e.

$$t(w) := \operatorname{argmax}_t P(t \mid w)$$

Improvement: Use Bayes' formula, i.e. $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$, omitting the quotient:

$$t(w) := P(t \mid w) = \operatorname{argmax}_t P(t) \cdot P(w \mid t)$$

TODO example

n -gram tagger Use information about the previous n tokens in addition to information about current token. Can have **word-based** and **tag-based** (tags are more common, training data covers more ground).

Assume a bigram language model (generating the sequence of POS-tags).

Pick the tag t_i for word w_i that maximises

$$P(t_i \mid t_{i-1}) \cdot P(w_i \mid t_{i-1})$$

For finding $P(t_i \mid t_{i-1})$, use the **Maximum Likelihood Estimate** (where c is the count of observations)

$$P(t_i \mid t_{i-1}) \approx \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

(Use start and end symbols to be able to calculate probs for first and last words)

TODO example

4.2 Rule-based tagging

Try to incorporate linguistic insight.

Brill tagger

1. Tag each word using a **baseline tagger** (e.g. unigram tagger, i.e. most common tag)
2. Apply patches that improve the result
 - e.g. if one of the two preceding words is a determiner, change the tag from verb to noun
 - Based on training data, compute the error between any two should-be/is-assigned: (t_a, t_b, freq) .
 - For each error triple, apply the patch that results in the greatest improvement, apply it.
 - Repeat until no further improvement is possible.

5 Parsing

We try to determine the grammatical structure of a sentence.

Parsing problem Identify the parse tree (w.r.t some grammar) for a given sentence. *(There may be multiple parse trees, or none.)*

Parse tree depicts derivation of sentence beginning from start symbol. Obviously requires context-free grammar *(each set of children has exactly one parent)*

Motivation

- Grammar checking
- Question answering
- Machine translation
- ...

5.1 Constituency parsing

Basic assumption Language is made up of **constituents**, i.e. basic, nested building blocks *(terminals and nonterminals of a formal grammar.)*

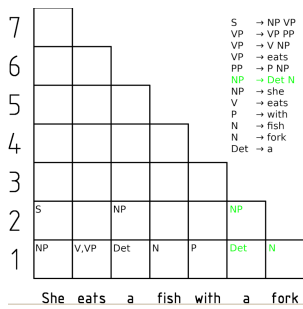
Leftmost derivation Always expand the leftmost expandable nonterminal.

Shift-reduce parser *(bottom-up approach)* Push tokens onto stack until top of stack matches the *rhs* of a rule, then reduce (on stack). Accepts the word if start symbol alone is left on the stack.

Chomsky normal form The *rhs* of a rule is of one of these shapes

- $A \rightarrow BC$
- $A \rightarrow a$
- $S \rightarrow \varepsilon$

CYK Parser *(dynamic programming, bottom-up)* Assumes a grammar in CNF, i.e. *rhs* of rules is two nonterminals or a terminal. From the bottom up, in a dynamic programming fashion, remember if two subcells are the *rhs* of a production rule. The topmost cell produces the entire sentence. Note that this finds all possible parse trees.



Part II

Similarity Measures

As there are many different expressions for the same semantic concept, we want to derive concepts of similarity between words. Further, this is useful in spelling correction.

5.2 Statistical Parsing

Motivation: Sentences can have many different parse trees but some are clearly more likely than others.

Basic idea

- Associate rules of a grammar with probabilities
- The probability of a parse tree is the product of all used productions/derivations.
- Probability of a sentence is the sum of all possible parse tree probabilities.
- Probabilities can be learned from a training corpus.
- Can use probabilities in parsing algorithms such as the CYK parser to find the most likely parse tree.

Lexicalised Parsing Extend production rules to be specific to terms (e.g. $VP(ate) \rightarrow VP(ate) PP(with)$). Then, again, assign/learn probabilities. **TODO** Head of sentence

5.3 Dependency Parsing

Basic concept Dependencies (in the linguistic sense) between tokens (such as *determiner*, *subject*, ...), these form a tree. Again, want to find the most likely parse/dependency structure

Constraint-based Parsing Come up with rules describing what a dependency can possibly look like. Begin with a complete graph of pairwise dependencies, then iteratively eliminate dependencies according to rules.

Global linear models Assign a *local feature vector* $g(x, h, m)$ to a specific dependency (h, m) in a specific sentence x . Try to maximise the (weighted) sum of all dependencies in the sentence, i.e. the best parse $y = F(x)$ is given by

$$F(X) = \operatorname{argmax}_y \vec{w} \cdot \sum_{(h,m) \in y} g(x, h, m)$$

6 Word similarity

Phonological similarity How much do two words sound alike?

Soundex Reduce a given term/word into a phonetic representation based on rules.

Morphological similarity basically comes down to stemming.

Spelling similarity Mainly want to identify spelling errors

Levenshtein distance / edit distance (dynamic programming) We already know this. Extensions: Cost matrices, based for example on distance of letters on keyboard.

Semantic similarity

- **Synonym** – same meaning. A *synset* is the equivalence class of synonymity.
- **Antonym** – opposite meaning
- **hypernym** – more general meaning
- **hyponym** – more specific
- **meronym** – part of collective

WordNet provides a tree (forest) of semantic relationships between words. Similarity between v, w can then be defined e.g.

TODO

7 Document similarity

Vector space model See IR lecture notes.

tf-idf weight

- **term frequency** $tf(t, d)$ — How often does term t appear in document d ? Commonly dampened by logarithm.

- *document frequency* $df(t)$ — In how many documents does t appear?
- *inverse document frequency* $idf(t)$ — The basic idea: Rare terms are more informative than frequent terms and should thus receive a high score. Thus we do the inverse fraction and log-dampening:

$$idf(t) := \log \frac{n}{df(t)}$$

The tf-idf weight is given by the product of term frequency and inverse document frequency:

$$tf.idf(t, d) = \log(1 + tf(td)) \cdot \log\left(\frac{n}{df(t)}\right)$$

- term frequency increases with the number of occurrences in the document
- inverse document frequency increases with the rarity of the term in the collection

The score for a multi-term query and a document is

$$score(q, d) := \sum_{t \in q \cap d} tf.idf(t, d)$$

Problems:

- Bag-of-words model, does not consider order
- Score is a sum, longer documents receive higher score

Vector space model Each document is a vector consisting of tf-idf weights of terms. *Distance measure*:

- euclidean distance is not suited because it considers the lengths (norms) of vectors, but we only really care about the distribution of terms
- hence use *cosine similarity*, which measures the angle between vectors.

Part III

Language Modelling

Basic pipeline

1. Corpus
2. high-dimensional vector space
3. latent space (word embeddings)
4. word relationships

Word embedding Encode some notion of similarity to other words in a vector. Motivation: Approaches like WordNet are static and thus insufficient. Possible information

- syntactic similarity (*run, ran*)
- semantic similarity (*large, big*)
- relatedness (*coffee, cup*)

Word context matrix (co-occurrence matrix)

- Each row represents a word
- Each column represents some “context” (specific entity or other word)
- cell represents the strength of association, for example *point-wise mutual information*

This matrix is very sparse. We’d like to find a low-dimensional representation of our word vector (word embedding). We project into a lower-dimensional space, the so-called *latent space*.

Word2Vec (neural network for learning word vectors, context-independent model) Given a training corpus, every word in a fixed vocabulary is represented by a vector.

1. For each center word c and a context window of fixed size o ...
2. use the word vector similarity of c and o to calculate $P(o|c)$ (or $P(c|o)$, see below)
3. Adjust word vector to maximise predictive accuracy
 - *Continuous Bag-of-Words* (CBOW) learns $P(c|o)$, i.e. focus word given the context
 - *Skip-gram* learns $P(o|c)$, i.e. context given some focus word.

Limitations:

- Any given word is represented by a vector which has to be stored
- Polysemy (different meanings of same word) is not addressed at all
- Dependence of meaning on context is not considered (gives rise to *context-dependent models*)

Given some word embeddings, how do we determine their similarity?

- similarity measures like cosine similarity
- visualise by dimension-reduction techniques like t-SNE, PCA, MDS, UMAP, ...
- *context-independent models* output a single word vector for a word
- *context-dependent models* generate multiple word embeddings for a word that capture different contexts.

Sentence embedding

- Basic approach with Bag-of-words assumption, vector of tf-idf weights – Problems:
 - problem: vectors are large and very sparse
 - cosine similarity of sentences with distinct words is zero but there could still be semantic similarity.
- Average word embeddings
- Deep learning approaches

Applications

- Assistants like Google, Siri, Alexa, ... for understanding and generating language, answering questions, ...
- Opinion mining
- Sentiment analysis
- Named Entity Recognition

Limitations Language models potentially contain biases (ethnic, gender) induced by the training data.

Part IV

Text Data Mining

8 Text classification

Applications

- categorise web pages
- spam vs. non-spam
- categorisation of correspondence
- categorisation of news articles
- classify to author
- classify reviews
- easy/hard to read, suitability for age groups
- writing style

As always, have a rule-based and a machine learning approach

Naive Bayes (for text classification) Assumes bag-of-words (independence between words). For a class c and words w_i , find the predicted class y by

$$y = \operatorname{argmax}_c P(c) \prod_i P(w_i | c)$$

For individual probabilities use maximum likelihood estimate, i.e.

- $P(c) = \frac{\text{\# of training docs with class } c}{\text{\# all docs}}$
- $P(w_i, c) = \frac{\text{count}(w_i, c_j)}{\sum_w \text{count}(w, c_j)}$

and the fact that $P(A|B) = P(A, B) / P(B)$

Choice of classification algorithm is dependent in characteristics and amount of training data.

Possible features

- tf-idf weights of terms
- phrases, characters, n -grams
- POS-tags
- metadata about the document
- non-text parts of document

9 Text clustering

10 Topic modelling

Opinion Mining & Information Extraction

11 Information Extraction

to extract • the entities and • relationships between such entities (i.e. clear, factual information)

11.1 Named Entity Extraction

used for • summarizing text • answering questions • integrating into knowledge bases • associating information (e.g. sentiments) to sentiments (e.g. of parts of printer in question)

Possible types of entities • location • time • person • ...

Supervised learning models Based on labelled training sequences (of tokens), train a classifier to predict labels (*Sequence Labelling Problem*)

- new data must fit training data
- time-consuming

To make this easier, use features that go beyond single tokens, e.g. context window of k words.

Sequence Labelling • reminiscent to POS-tagging • assuming that label is dependent on context. Typical models are • Markov models • Conditional Random Fields • Bidirectional LSTMs

Once we have identified the entities, we'd like to find relationships between them (e.g. triples of operators *is-a*, *daughter-of*, ...) (*Can save these triples in a knowledge base e.g. for question-answering; cf RDF-triples*).

Relationship Extraction Try to find type of relationship between two entities. Possibilities

- Extract RDF triples from large corpora like Wikipedia
- Use (specialised) ontologies / knowledge bases (for ex. medical applications)

Methods to extract information:

- handwritten rules – e.g. “*Y such as X*”, “*X, especially Y*”, “*X, including Y*” all express an *is-a-relationship*. – there can be more specific relations that only make sense between certain types of entities (e.g. *cures(drug, disease)*) – pros: • precise •

can be tailored to specific domains – cons: • low recall • high effort

- supervised,
- unsupervised machine learning.

semi-supervised learning: extract less common patterns based on training corpus?

Question Answering & Text Summarization