

Statistical Natural Language Processing

Tokenization, normalization, segmentation

Çağrı Çöltekin

University of Tübingen
Seminar für Sprachwissenschaft

Summer Semester 2017

Tokenization – a solved problem?

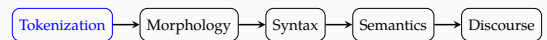
- Typically, we (in NLP/CL/IR/...) process text as a sequence of tokens
- Tokens are word-like units
- A related task is *sentence segmentation*
- Tokenization is a language dependent task, where it becomes more challenging in some languages
- Tokenization is often regarded as trivial, and a mostly solved task

Classical NLP pipeline

- *Tokenization*
Sentences, (normalized) words, stems / lemmas
- *Lexical / morphological processing*
POS tags, morphological features, stems / lemmas, named entities
- *Parsing*
Constituency / dependency trees
- *Semantic processing*
word-senses, logical forms
- *Discourse*
Co-reference resolution, discourse representation

We do not always use a pipeline, not all steps are necessary for all applications

Tokenization in the classical NLP pipeline



- Tokenization is the first in the pipeline
- Even for end-to-end approaches, tokenization is often considered given (needs to be done in advance)
- Errors propagate!

But, can't we just tokenize based on spaces?

...and get rid of the punctuation

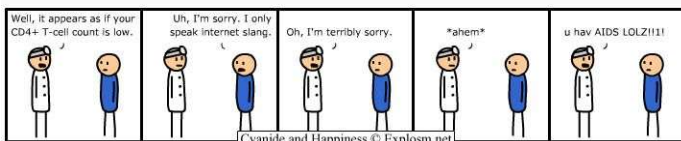
Some examples from English:

- \$10 billion
- rock 'n' roll
- he's
- can't
- O'Reilly
- 5-year-old
- B-52
- C++
- C4.5
- 29.05.2017
- 134.2.129.121
- sfs.uni-tuebingen.de
- New York-based
- wake him up

Gets more interesting in other languages

- Chinese: 猫占领了婴儿床
'The cat occupied the crib'
- German: Lebensversicherungsgesellschaftsangestellter
'life insurance company employee'
- Turkish: İstanbullulaştırılabileceklerimizdenmişsiniz
'You were (evidentially) one of those who we may not be able to convert to an Istanbulite'
- Even more interesting when we need to process 'mixed' text with *code-switching*

Specialized and non-standard text



- Much more difficult for non-standard text
 - Many specialized terms use a mixture of letters, numbers, punctuation
 - Frequent misspelling, omitting space (e.g., after sentence final punctuation)
- The problem is more severe for
 - Specialized domains, e.g., bio-medical texts
 - Informal communication, e.g., social media

Normalization

- For most applications (e.g., IR) we want to treat the following the same
 - Linguistics – linguistics
 - color – colour
 - lower case – lowercase – lower-case
 - Tübingen – Tuebingen – Tubingen
 - see – see
 - flm – film
 - Different date/time formats, phone numbers
- Most downstream tasks require the 'normalized' forms of the words

So, what is a token?

- One token or multiple?
 - John's
 - New York
 - German: *im* (*in + dem*)
 - Turkish: *İstanbullulaştırılamayabileceklerimizdenmişsiniz*
- Answer is language and application dependent
- Tokenization decisions are often arbitrary
- Consistency is important

Rule based tokenization

Regular expressions and finite-state automata

- The 'easy' solution to the tokenization is rule-based
- Using regular expressions,
 - we can define regular expressions for allowed tokens
 - split after match, disregard/discard the remaining parts
- For example,
 - All alphabetic characters, *word*, $[a-z]^+$
 - Capitalization, *John*, $[A-Z]?[a-z]^+$
 - Abbreviations, *Prof.*, $[A-Z]?[a-z]^+[.]?$
 - Numbers too, *123*, $[A-Z]?[a-z]^+[.]?[0-9]^+$
 - Numbers with decimal parts $[A-Z]?[a-z]^+[.]?[0-9.]^+$
 - ...
- Result is typically imprecise, difficult to maintain

Splitting sentences

- Another relevant task is *sentence tokenization*
- For most applications, we need sentence boundaries
- Sentence-final markers, $[. ! ?]$ are useful
- But the dot '.' is ambiguous: can either be end-of- sentence or abbreviation marker, or both
 - The U.N. is the largest intergovernmental organisation.
 - I had the impression he'll be ambassador to U.N.
- Again, heuristics along with a list of abbreviations is possible

Problems with rule-based approaches

- Rule-based approaches are (still) common in practice, however
 - it is difficult to build a rule set that works well in practice
 - it is difficult to maintain
 - it is not domain or language general: needs re-implementation, re-adjustment for every case

Machine learning for word / sentence tokenization

- Another approach is to use machine learning
- Label each character in the text with
 - I inside a token
 - O outside tokens
 - B beginning of a token,
 alternatively to combine word/sentence tokenization
 - T beginning of a token
 - S beginning of a sentence
- How do we create the training data?
- What are the features for the ML?

I/O/B tokenization: an example

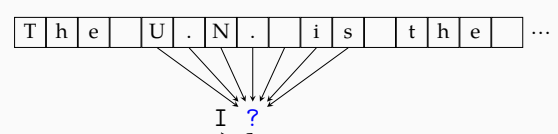
The U.N. is the largest intergovernmental
 BIIIOBIIIOBIIIOBIIIIIOBIIIIIIIIIIIIIIIO
 organisation. I had the impression he'll be
 BIIIIIIIIIIIOOBIIIOBIIIOBIIIIIIIIIOBIBIIIOBIO
 ambassador to U.N.
 BIIIIIIIIIIIOBIOBIIIO

I/O/B tokenization example

with sentence boundary markers

The U.N. is the largest intergovernmental
 SIIOTIIIOBIIIOBIIIIIIIOBIIIIIIIIIIIIIIIO
 organisation. I had the impression he'll be
 TIIIIIIIIIIIOOSOTIIIOBIIIOBIIIIIIIIIOBIIIOBIO
 ambassador to U.N.
 TIIIIIIIIIIIOBIIIOBIIIO

Features for tokenization



- We predict label of each character
- Typical features are the other characters around the target
- Choice of features and the machine learning method vary
- Using the previous prediction is also useful

Segmentation

- Segmentation is a related problem in many areas of computational linguistics
 - In some languages, the word boundaries are not marked
猫占领了婴儿床 → 猫 占领 了 婴儿床
 - We often want to split words into their morphemes
Lebensversicherungsgesellschaftsangestellter →
Leben+s+versicherung+s+gesellschaft+s+angestellter
 - In spoken language there are no reliable word boundaries

Supervised segmentation

- I/O/B tokenization is applicable to segmentation as well
- Often produces good accuracy
- The main drawback is the need for labeled data
- Some unsupervised with reasonable accuracy also exist
- In some cases, unsupervised methods are useful and favorable

A simple ‘unsupervised’ approach

- Using a lexicon, segment at maximum matching lexical item
- Serves as a good baseline, but fails in examples like
theman
where maximum match suggests segmentation ‘them an’
- The out-of-vocabulary words are problematic

Unsupervised segmentation

- Two main approaches
 - Learn a compact lexicon that maximizes the likelihood of the data

$$P(s) = \prod_{i=1}^n P(w_i)$$

$$P(w) = \begin{cases} (1 - \alpha)f(w) & \text{if } w \text{ is known} \\ \alpha \prod_{i=1}^m P(a_i) & \text{if } w \text{ is unknown} \end{cases}$$

- Segment at points where predictability (entropy) is low
The general idea: the predictability within words is high, predictability between words is low

Summary

- Tokenization is an important part of an NLP application
- Tokens are word-like units that are
 - linguistically meaningful
 - useful in NLP applications
- Tokenization is often treated as trivial, has many difficulties of its own
- White spaces help, but does not solve the tokenization problem completely
- Segmentation is tokenization of input where there are no boundary markers
- Solutions include rule-based (regex) or machine learning approaches

Next

Wed More machine learning
Fri First graded assignment

Some extra: modeling segmentation by children

NLP can be ‘sciency’, too

- An interesting application of unsupervised segmentation methods is modeling child language acquisition
- How children learn languages has been one of the central topics in linguistics and cognitive science
- Computational models allow us to
 - test hypotheses
 - create explicit models
 - make predictions

The puzzle to solve

```

1juuzuibutsjhiuljuuz
1juuztbzjubhbjomwfljuuz
xibutuibu
1juuz
epzpvxbounpsfnjmlipofz
1juuzljuuzephjhf
opnjxibuepftbljuuztbz
xibuepftbljuuztbz
ephjhfeph
ephjhf
opnjxibuepftuifephjftbz
xibuepftuifephjftbz
mjuumfcbczcjsejf
cbczcjsejf
zpvpeoumjluiupof
plbzpnanzubluijtpvu
dpx
uifdpxtbztnpppp
xibuepftuifdpxtbzopnj

```

- No clear boundary markers
- No lexical knowledge

