Introduction, Basic Text Processing, Edit Distance

Language Technology

making good progress

mostly solved



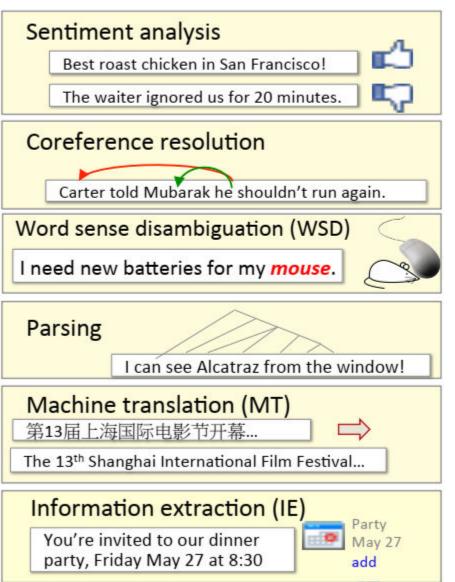
Part-of-speech (POS) tagging

ADJ ADJ NOUN VERB ADV

Colorless green ideas sleep furiously.

Named entity recognition (NER)

___ PERSON ORG LOC __ Einstein met with UN officials in Princeton



still really hard

Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

Summarization

The Dow Jones is up
The S&P500 jumped
Housing prices rose

Economy is good

Dialog

Where is Citizen Kane playing in SF?



Castro Theatre at 7:30. Do you want a ticket?

Text Processing 1. Regular expression

- A formal language for specifying text strings How can we search for any of these?
- Apple
- apple
- Apples
- apples

Regular expression: Disjunction

• Letters inside square brackets []

Patterns	Matches
[Aa]pple	Apple, apple
[1234567890]	Any digit

Ranges

Patterns	Matches	
[A-Z]	Any upper case letter	<u>I</u> am going.
[a-z]	Any lower case letter	what are you doing
[0-9]	Any digit	Chapter <u>1</u> :

Regular expression: Negation in Disjunction

• Negations [^Ss]: carat means negation if used as 1st char inside []

Patterns	Matches	
[^A-Z]	Not an upper case letter	I <u>n</u> ot am going.
[^aA]	Not a or A	what are you doing
[^e^]	Neither e nor ^	e <u>l</u> even. ^go
a^b	Pattern a carat b	Find <u>a^b</u>

Regular expression: More Disjunction

• Pipe | is used for disjunction.

Patterns	Matches
apple banana	apple, banana
1 2 3	=[1-3]
[aA]pple [bB]anana	

Regular expression: * +?.

Patterns	Matches	
oo*h!	Zero or more of the previous character	oh!, ooh!,
oo+h!	one or more of the previous character	ooh!, oooh!,; does not match oh!
colou?r	Zero or one of the previous character	colour, color
col.r	Any character	color, coler, colar,colwr

Regular expression: Anchor ^ \$

• ^ matches the beginning of a sentence, \$ matches the end of a text.

Patterns	Matches	
^[A-Z]	<u>G</u> pple	
^[^a-zA-Z]	1 "	<u>'</u> hi"
\.\$	end <u>.</u>	
.\$	end <u>.</u>	end <u>!</u>

Regular expression: examples

• Find all occurrences of word 'the' in the text.

Patterns	Matches
the	misses capitalized examples
[tT]he	Incorrectly captures other
[^a-zA-Z][tT]he[^a-zA-Z]	correct.

Visit https://regexr.com/ to test all these regular expressions type.

Text Processing 2. Word Tokenization

- · Every NLP task needs to do text normalization:
 - 1. Segmenting/tokenizing words in running text
 - 2. Normalizing word formats
 - 3. Segmenting sentences in running text

Words

- · Lemma: same stem, part of speech, rough word sense
 - · cat and cats = same lemma
- · Wordform: the full inflected surface form
 - cat and cats = different wordforms
- Type: an element of the vocabulary.
- · Token: an instance of that type in running text

they lay back on the San Francisco grass and looked at the stars and their 15 tokens (or 14), 13 types

Tokenization in Linux

Given a text file, output the word tokens and their frequencies

```
tr -sc 'A-Za-z' '\n' < shakes.txt Change all non-alpha to newlines
| sort | Sort in alphabetical order | uniq -c | Merge and count each type
```

```
1945 A 25 Aaron
72 AARON 6 Abate
19 ABBESS 5 Abbess
5 ABBOT 6 Abbey
... 3 Abbot
```

The first step: tokenizing

```
tr -sc 'A-Za-z' '\n' < shakes.txt | head

THE
SONNETS
by
William
Shakespeare
From
fairest
creatures
We</pre>
```

. . .

The second step: sorting

```
tr -sc 'A-Za-z' '\n' < shakes.txt | sort | head

A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A</pre>
```

The third step: counting

Merging upper and lower case

```
tr 'A-Z' 'a-z' < shakes.txt | tr -sc 'A-Za-z' 'n' | sort | uniq -c
```

Sorting the counts

```
tr 'A-Z' 'a-z' < shakes.txt | tr -sc 'A-Za-z' '\n' | sort | uniq -c | sort -n -r

23243 the

22225 i
18618 and
16339 to
15687 of
12780 a
12163 you
10839 my
10005 in
8954 d
```

Tokenization issues

```
Finland's capital → Finland Finlands Finland's ?
what're, I'm, isn't → What are, I am, is not
Hewlett-Packard → Hewlett Packard ?
state-of-the-art → state of the art ?
Lowercase → lower-case lowercase lower case ?
San Francisco → one token or two?
m.p.h., PhD. → ??
```

Language issues:

- Chinese and Japanses words are not separated by spaces
- Use maximum matching to find out tokens
- Doesn't generally work in English! Thetabledownthere

Modern Tokenizers: BPE, Wordpiece

- Word Tokenizers
- Character Tokenizers
- Sub-word Tokenizers
 - Developed for machine translation by Sennrich et al., ACL 2016

"The main motivation behind this paper is that the translation of some words is transparent in that they are translatable by a competent translator even if they are novel to him or her, based on a translation of known subword units such as morphemes or phonemes."

- Later used in BERT, T5, RoBERTa, GPT, etc.
- Relies on a simple algorithm called byte pair encoding (Gage, 1994)

Form base vocabulary (all characters that occur in the training data

word	frequency
hug	10
pug	5
pun	12
bun	4
hugs	5

Base vocab: b, g, h, n, p, s, u

 Now, count up the frequency of each character pair in the data, and choose the one that occurs most frequently

word	frequency	character pair	frequency
h+u+g	10	ug	20
p+u+g	5	pu	17
p+u+n	12	un	16
b+u+n	4	hu	15
h+u+g+s	5	gs	5

•

 Now, choose the most common pair (ug) and then merge the characters together into one symbol. Then, retokenize the data

word	frequency	character pair	frequency
h+ <i>ug</i>	10	un	16
p+ug	5	h+ug	15
p+u+n	12	pu	12
b+u+n	4	p+ug	5
h+ug+s	5	ug+s	5

...

Modern Tokenizers: BPE, Wordpiece

Byte pair encoding

 Keep repeating this process! This time we choose un to merge, next time we choose h+ug, etc.

word	frequency	character pair	frequency
h+ <i>ug</i>	10	un	16
p+ <i>ug</i>	5	h+ug	15
p+u+n	12	pu	12
b+u+n	4	p+ug	5
h+ <i>ug</i> +s	5	ug+s	5

...

Eventually, after a fixed number of merge steps, we stop

frequency
10
5
12
4
5

new vocab: b, g, h, n, p, s, u, ug, un, hug

- To avoid <UNK>, all possible characters / symbols need to be included in the base vocab. This can be a lot if including all unicode characters!
- GPT-2 uses bytes as the base vocabulary (size 256) and then applies BPE on top of this sequence (with some rules to prevent certain types of merges).
- Commonly have vocabulary sizes of 32K to 64K

Modern Tokenizers: BPE, Wordpiece

- WordPiece (Schuster et al., ICASSP 2012): merge by likelihood as measured by language model, not by frequency
- SentencePiece (Kudo et al., 2018): can do subword tokenization without pretokenization (good for languages that don't always separate words w/ spaces), although pretokenization usually improves performance

Text Processing 3. Normalization & Stemming

- · We implicitly define equivalence classes of terms
 - · deleting periods in a term: U.S.A and USA
 - asymmetric expansion: window → window, windows, Windows
 Windows → Windows
- Case folding: convert every letter to lowercase
 - Exception: uppercase in the middle of sentence General Motors, SAIL
 - For MT, Sentiment analysis etc, case may be useful (US vs us)

Lemmatization

- Reduce inflections or variant forms to base form
 - am, are, is \rightarrow be
 - car, cars, car's, cars' → car
- the boy's cars are different colors \rightarrow the boy car be different color
- · have to find correct dictionary headword form

Morphemes

- · The small meaningful units that make up words
 - · Stems: The core meaning--bearing units
 - · Affixes: Bits and pieces that adhere to stems
 - · Often with grammatical functions

Stemming

- Reduce terms to their stems in information retrieval
- Stemming is crude chopping of affixes
 - language dependent
 - e.g., automate(s), automatic, automation all reduced to automat.

for example compressed and compression are both accepted as equivalent to compress.



for exampl compress and compress ar both accept as equival to compress

Porters Stemmer

```
Step 1a
                                          Step 2 (for long stems)
   sses → ss caresses → caress
                                             ational → ate relational → relate
   ies → i ponies → poni
                                             izer→ ize digitizer → digitize
   ss → ss caress → caress
                                             ator→ ate operator → operate
   s \rightarrow \phi cats \rightarrow cat
Step 1b
                                          Step 3 (for longer stems)
   (*v*)ing \rightarrow \emptyset walking \rightarrow walk
                                             al \rightarrow \emptyset revival \rightarrow reviv
                   sing \rightarrow sing
                                             able \rightarrow \emptyset adjustable \rightarrow adjust
   (*v*)ed \rightarrow \emptyset plastered \rightarrow plaster
                                             ate → ø activate → activ
           (*v*)ing \rightarrow \emptyset walking \rightarrow walk
                                       sing \rightarrow sing
```