Basic Text Processing

Regular Expressions



Regular expressions

- A formal language for specifying text strings
- How can we search for any of these?
 - woodchuck
 - woodchucks
 - Woodchuck
 - Woodchucks





Regular Expressions: Disjunctions

Letters inside square brackets []

Pattern	Matches
[wW]oodchuck	Woodchuck, woodchuck
[1234567890]	Any digit

Ranges [A-Z]

Pattern	Matches	
[A-Z]	An upper case letter	Drenched Blossoms
[a-z]	A lower case letter	my beans were impatient
[0-9]	A single digit	Chapter 1: Down the Rabbit Hole





Regular Expressions: Negation in Disjunction

- **Negations** [^Ss]: not any of the following
 - Carat means negation only when first in []
 - Check: regexpal.com
 - Try it in Arabic also!

Pattern	Matches	
[^A-Z]	Not an upper case letter	O <u>y</u> fn pripetchik
[^Ss]	Neither 'S' nor 's'	\underline{I} have no exquisite reason"
[^e^]	Neither e nor ^	Look <u>h</u> ere
a^b	The pattern a carat b	Look up a^b now [^not first]

STUDENTS-HUB.com



Regular Expressions: More Disjunction

- Woodchucks is another name for groundhog!
- The pipe | for disjunction

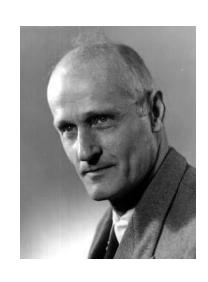
Pattern	Matches
groundhog woodchuck	
yours mine	yours mine
a b c	= [abc]
[gG]roundhog [Ww]oodchuck	





Regular Expressions: ? * +

Pattern	Matches	
colou?r	Optional previous char	<u>color</u> <u>colour</u>
oo*h!	0 or more of previous char	oh! ooh! oooh!
o+h!	1 or more of previous char	oh! ooh! oooh!
baa+		baa baaa baaaa
beg.n	. Any character	begin begun beg3n



Stephen C Kleene

Kleene *, Kleene +



Regular Expressions: Anchors ^ \$

Pattern		Matches	
^[A-Z]	^: Start of line	Palo Alto	
^[^A-Za-z]		1 "Hello"	
\.\$	\$: end of line	The end.	
. \$		The end? The end!	



Example

Find me all instances of the word "the" in a text.

the

Misses capitalized examples

[tT]he

Incorrectly returns other or theology

 $[^a-zA-Z][tT]he[^a-zA-Z]$ try this!



Errors

- The process we just went through was based on fixing two kinds of errors
 - Matching strings that we should not have matched (there, then, other)
 - False positives (Type I)
 - Not matching things that we should have matched (The)
 - False negatives (Type II)



Errors cont.

- In NLP we are always dealing with these kinds of errors.
- Reducing the error rate for an application often involves two antagonistic efforts:
 - Increasing accuracy or precision (minimizing false positives)
 - Increasing coverage or recall (minimizing false negatives).



Summary

- Regular expressions play a surprisingly large role
 - Sophisticated sequences of regular expressions are often the first model for any text processing text
- For many hard tasks, we use machine learning classifiers
 - But regular expressions are used as features in the classifiers
 - Can be very useful in capturing generalizations

Basic Text Processing

Word tokenization





Text Normalization

- 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



How many words?

- I do uh main- mainly business data processing
 - Fragments, filled pauses
- Seuss's cat in the hat is different from other cats!
 - Lemma: same stem, part of speech, rough word sense
 - cat and cats = same lemma مدرس،مدرسة، مدرسون، مدرسات
 - Wordform: the full inflected surface form
 - cat and cats = different wordforms (words)





How many words?

they lay back on the San Francisco grass and looked at the stars and their

- Type: an element of the vocabulary (could be a word or stem or even a root, repitition doesn't count).
- Token: an instance of that type in running text.
- How many?
 - 15 tokens (or 14)
 - 13 types (or 12) (or 11?)
 - Generally more tokens than types (usually much more)!

STUDENTS-HUB.com



How many words?

N= number of tokens

Church and Gale (1990): $|V| > O(N^{\frac{1}{2}})$

V = vocabulary = set of types

| V | is the size of the vocabulary

	Tokens = N	Types = V
Switchboard phone conversations	2.4 million	20 thousand
Shakespeare	884,000	31 thousand
Google N-grams	1 trillion	13 million
Quran (in Arabic)	77413	18994





How many words? Quran Statistics

- Tokens = 77,430
- Number of *unique* surface forms= 18994
- Token to Term Ration(final)=4.1
- Number of unique words by *stem* = 12183
- Number of unique words by *root* = 1685
- Number of unique words by *lemma* = 3382 (excluding verbs, and other words where lemma is not annotated)

http://quickestwaytoquran.blogspot.com/2013/08/number-of-unique-words-in-guran.html



Simple Tokenization in UNIX

- (Inspired by Ken Church's UNIX for Poets.)
- Given a text file, output the word tokens and their frequencies

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

STUDENTS-HUB.com

....



The first step: tokenizing

```
tr -sc 'A-Za-z' '\n' < shakes.txt | head
THE
SONNETS
by
William
Shakespeare
From
fairest
creatures
We
```

. . .



The second step: sorting

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



More 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
                                      What happened here?
            10839 my
            10005 in
            8954 d
```



Issues in Tokenization

 \rightarrow $\dot{}$ $\dot{}$

```
    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.



Tokenization: language issues

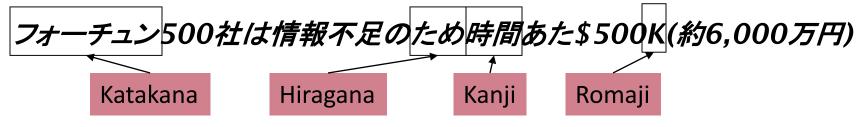
- French
 - *L'ensemble* → one token or two?
 - L? L'? Le?
 - Want *l'ensemble* to match with *un ensemble*

- German noun compounds are not segmented
 - Lebensversicherungsgesellschaftsangestellter [TTR for German?]
 - 'life insurance company employee'
 - German information retrieval needs compound splitter



Tokenization: language issues

- Chinese and Japanese no spaces between words:
 - 莎拉波娃现在居住在美国东南部的佛罗里达。
 - 莎拉波娃 现在 居住 在 美国 东南部 的 佛罗里达
 - Sharapova now lives in US southeastern Florida
- Further complicated in Japanese, with multiple alphabets intermingled
 - Dates/amounts in multiple formats



End-user can express query entirely in hiragana!

STUDENTS-HUB.com



Word Tokenization in Chinese

- Also called Word Segmentation
- Chinese words are composed of characters
 - Characters are generally 1 syllable and 1 morpheme.
 - Average word is 2.4 characters long.
- Standard baseline segmentation algorithm:
 - Maximum Matching (also called Greedy)



Maximum Matching Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
- 1) Start a pointer at the beginning of the string
- Find the longest word in dictionary that matches the string starting at pointer
- 3) Move the pointer over the word in string
- 4) Go to 2 [does it work for Arabic error correction?]



Max-match segmentation illustration

Thecatinthehat

the cat in the hat

Thetabledownthere

the table down there

theta bled own there

Doesn't generally work in English!

- But works astonishingly well in Chinese
 - 莎拉波娃现在居住在美国东南部的佛罗里达。
 - 莎拉波娃 现在 居住 在 美国 东南部 的 佛罗里达
- Modern probabilistic segmentation algorithms even better

Basic Text Processing

Word Normalization and Stemming



Normalization (English)

- Need to "normalize" terms
 - Information Retrieval: indexed text & query terms must have same form.
 - We want to match *U.S.A.* and *USA*,
- We implicitly define equivalence classes of terms
 - e.g., deleting periods in a term
- Alternative: asymmetric expansion:
 - Enter: window Search: window, windows
 - Enter: *windows* Search: *Windows, windows, window*
 - Enter: *Windows* Search: *Windows*
- Potentially more powerful, but less efficient



Normalization (Arabic)

- Need to "normalize" terms
 - Information Retrieval: indexed text & query terms must have same form.
 - Hamza, Alef Mqsoora and Ta Marboota: confusion letters
 - على، على الى، الى فدية، فديه We want to match

- We implicitly define equivalence classes of terms
- By matching different manifestations to a common (most frequent?) term
- Index only one form and convert query terms to the representative of the class:

فءة بناء كفء :are written as فئة بناء كفؤ then فئة بناء كفو



Case folding

- Applications like IR: reduce all letters to lower case
 - Since users tend to use lower case
 - Possible exception: upper case in mid-sentence?
 - e.g., *General Motors*
 - Fed vs. fed
 - SAIL vs. sail
- But: For sentiment analysis, MT, Information extraction
 - Case is helpful (*US* versus *us* is important)



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 → the boy car be different color
- اجتهاد الطالب الطريق الى النجاح اجتهد طالب طريق الى نجاح
- Lemmatization: have to find correct dictionary headword form





Morphology

- 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

Dan Jurafsky



Porter's algorithm The most common English stemmer

```
Step 1a
```

```
sses \rightarrow ss caresses \rightarrow caress
ies \rightarrow i ponies \rightarrow poni
ss \rightarrow ss
                  caress \rightarrow caress
  \rightarrow Ø
                  cats \rightarrow cat
```

Step 2 (for long stems)

```
ational \rightarrow ate relational \rightarrow relate
izer→ ize digitizer → digitize
ator\rightarrow ate operator \rightarrow operate
```

Step 1b

```
(*v*)ing \rightarrow \emptyset walking \rightarrow walk
                       sing \rightarrow sing
(*v*)ed \rightarrow \emptyset plastered \rightarrow plaster
```

Step 3 (for longer stems)

```
al
         \rightarrow \emptyset revival \rightarrow reviv
able \rightarrow \phi adjustable \rightarrow adjust
ate \rightarrow \emptyset activate \rightarrow activ
```





Viewing morphology in a corpus Why only strip –ing if there is a vowel?



Viewing morphology in a corpus Why only strip –ing if there is a vowel?

```
sing \rightarrow sing
tr -sc 'A-Za-z' '\n' < shakes.txt | grep 'ing$' | sort | uniq -c | sort -nr
                1312 King548 being548 being541 nothing
                541 nothing 152 something
                388 king 145 coming
                375 bring 130 morning
                358 thing 122 having
                307 ring 120 living
152 something 117 loving
                145 coming 116 Being
                130 morning 102 going
tr -sc 'A-Za-z' '\n' < shakes.txt | grep '[aeiou].*ing$' | sort | uniq -c | sort -nr
```

(*v*)ing $\rightarrow \emptyset$ walking \rightarrow walk



Dealing with complex morphology is sometimes necessary

- Some languages requires complex morpheme segmentation
 - Turkish
 - Uygarlastiramadiklarimizdanmissinizcasina
 - `(behaving) as if you are among those whom we could not civilize'
 - Uygar `civilized' + las `become'
 - + tir `cause' + ama `not able'
 - + dik `past' + lar 'plural'
 - + imiz 'p1pl' + dan 'abl'
 - + mis 'past' + siniz '2pl' + casina 'as if'

Basic Text Processing

Sentence Segmentation and Decision Trees



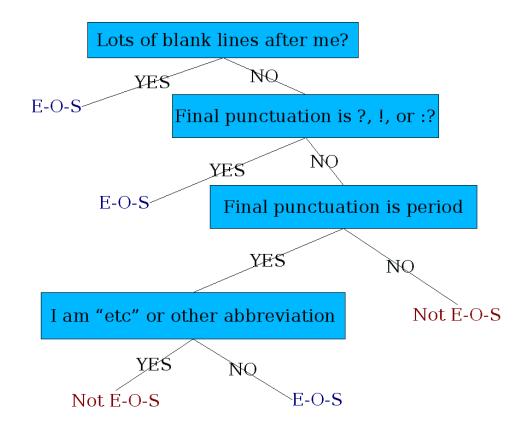
Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous
 - Sentence boundary
 - Abbreviations like Inc. or Dr.
 - Numbers like .02% or 4.3
- Build a binary classifier
 - Looks at a "."
 - Decides EndOfSentence/NotEndOfSentence
 - Classifiers: hand-written rules, regular expressions, or machine-learning





Determining if a word is end-of-sentence: a Decision Tree





More sophisticated decision tree features

- Case of word with ".": Upper, Lower, Cap, Number
- Case of word after ".": Upper, Lower, Cap, Number

- Numeric features
 - Length of word with "."
 - Probability(word with "." occurs at end-of-s)
 - Probability(word after "." occurs at beginning-of-s)



Implementing Decision Trees

- A decision tree is just an if-then-else statement
- The interesting research is choosing the features
- Setting up the structure is often too hard to do by hand
 - Hand-building only possible for very simple features, domains
 - For numeric features, it's too hard to pick each threshold
 - Instead, structure usually learned by machine learning from a training corpus



Decision Trees and other classifiers

- We can think of the questions in a decision tree
- As features that could be exploited by any kind of classifier: Logistic regression, SVM, Neural Nets,....

Here are some videos

https://www.youtube.com/watch?v=hwDhO1GLb_4

https://www.youtube.com/watch?v=RGLldper5II

https://www.youtube.com/watch?v=jBk24DI8kg0