## Basic Text Processing

### NLP Text Processing Pipeline

---

### NLP Text Processing Pipeline

**nltk** provides implementations for most operations

- Document → Sections and Paragraphs
- Paragraphs → Sentences (sentence segmentation / extraction)
- Sentences → Tokens
- Tokens → Lemmas or Morphological Variants / Stems
- Tokens → Part-of-speech (POS) Tags
- Tokens, POS Tags → Phrase Chunks (Noun & Verb Phrases)
- Tokens, POS Tags → Parse Trees
  - Augment above with coreference, entailment, sentiment, …

---

## 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
  - **Wordform**: the full inflected surface form
    - cat and cats = different wordforms

---

### How many words?

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

- **Type**: an element of the vocabulary.
- **Token**: an instance of that type in running text.
- How many?
  - 15 tokens (or 14)
  - 13 types (or 12) (or 11?)

## How many words?

$N$ = number of tokens

$V$ = vocabulary = set of types

Church and Gale (1990): $|V| > O(N^{½})$

|$V$| 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 |

7

---

## Issues in Tokenization

- 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. → ??

8

---

## 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*
  - 'life insurance company employee'
  - German information retrieval needs **compound splitter**

9

---

## 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

フォーチュン500社は情報不足のため時間あた$500K(約6,000万円)

| Katakana | Hiragana | Kanji | Romaji |
|---|---|---|---|

10    End-user can express query entirely in hiragana!

---

## 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)

11

---

# Basic Text Processing

## Word Normalization and Stemming

12

---

## Normalization

- 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

## 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*
- For sentiment analysis, MT, Information extraction
  - Case is helpful (*US* versus *us* is important)

14

## Lemmatization

- Reduce inflections or variant forms to base form
  - *am, are, is → be*
  - *car, cars, car's, cars' → car*
- *the boy's cars are different colors → the boy car be different color*
- Lemmatization: have to find correct dictionary headword form
- Machine translation
  - Spanish quiero ('I want'), quieres ('you want') same lemma as querer 'want'

15

## 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

16

## 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 |

17

## Porter's algorithm
## The most common English stemmer

Step 1a
```
sses → ss   caresses → caress
ies  → i    ponies   → poni
ss   → ss   caress   → caress
s    → ø    cats     → cat
```
Step 1b
```
(*v*)ing → ø   walking   → walk
               sing      → sing
(*v*)ed  → ø   plastered → plaster
…
```

Step 2 (for long stems)
```
ational→ ate  relational→ relate
izer→ ize    digitizer → digitize
ator→ ate    operator  → operate
…
```
Step 3 (for longer stems)
```
al   → ø   revival    → reviv
able → ø   adjustable → adjust
ate  → ø   activate   → activ
…
```

18

3

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

```
(*v*)ing → ø   walking   → walk
               sing      → sing
```

19

---

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

```
(*v*)ing → ø   walking   → walk
               sing      → sing
```

```
tr -sc 'A-Za-z' '\n' < shakes.txt | grep 'ing$' | sort | uniq -c | sort -nr
```

| | |
|---|---|
| 1312 King | 548 being |
| 548 being | 541 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
```

20

---

## 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'

21

---

# Basic Text Processing

## Sentence Segmentation and Decision Trees

22

---

## 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

23

---

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

Lots of blank lines after me?
- YES → E-O-S
- NO → Final punctuation is ?, !, or :?
  - YES → E-O-S
  - NO → Final punctuation is period
    - YES → I am "etc" or other abbreviation
      - YES → Not E-O-S
      - NO → E-O-S
    - NO → Not E-O-S

24

4

## 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)

25

---

# Basic Text Processing

Regular Expressions: Detecting word pattern variations

26

---

## Regular expressions

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

27

---

## 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 |

28

---

## Regular Expressions: Negation in Disjunction

- Negations [^Ss]
  - Carat means negation only when first in []

| Pattern | Matches | |
|---------|---------|---|
| [^A-Z] | Not an upper case letter | Oyfn pripetchik |
| [^Ss] | Neither 'S' nor 's' | I have no exquisite reason" |
| [^e^] | Neither e nor ^ | Look here |
| a^b | The pattern a carat b | Look up a^b now |

29

---

## 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 | |

Photo D. Fletcher

30

---

5

## Regular Expressions: ?  *  +  .

| Pattern | Matches | |
|---------|---------|---|
| colou?r | Optional previous char | color    colour |
| oo*h! | 0 or more of previous char | oh! ooh!  oooh! ooooh! |
| o+h! | 1 or more of previous char | oh! ooh!  oooh! ooooh! |
| baa+ | | baa baaa baaaa baaaaa |
| beg.n | | begin begun begun beg3n |

Stephen C Kleene

Kleene *,  Kleene +

31

## Regular Expressions: Anchors  ^  $

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

32

## 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]

33

## 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)

34

## 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).

35

## Exercise

- Write a regular expression to match dates
  - November 9, 1989
  - 17 December 1967
  - 11-09-1989
  - 12/17/67

  Where might you use these matchers?

- Write a regular expression to match time expressions
  - Next Wednesday at noon
  - Tomorrow morning

36

6

# Regex 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

37