# Basic Text Processing

Dr Muhammad Sarim

# Contents

# Contents

## Text Normalization

- Every NLP task requires text normalization.

## Text Normalization

- Every NLP task requires text normalization.

  1. Segmenting and Tokenizing words in a running text.

## Text Normalization

- Every NLP task requires text normalization.
    1. Segmenting and Tokenizing words in a running text.
    2. Normalizing word formats.

## Text Normalization

- Every NLP task requires text normalization.

  1. Segmenting and Tokenizing words in a running text.
  2. Normalizing word formats.
  3. Segmenting sentences in running text.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Contents

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

I do uh main- mainly business data processing.

**Tokenization**
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

I do uh main- mainly business data processing.

- Fragments, filled pauses.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the 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.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the 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 are different wordforms.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

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

- Type: an element of the vocabulary.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

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

- Type: an element of the vocabulary.
  - 13 types (or 12) (or 11?)

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

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

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

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

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

$N$ = Number of tokens

$V$ = vocabulary = set of types

$|V|$ is the size of the vocabulary

$|V| > O(\sqrt{N})$ (Church and Gale (1990))

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Can you count the words?

N = Number of tokens

V = vocabulary = set of types

$|V|$ is the size of the vocabulary

$|V| > O(\sqrt{N})$ (Church and Gale (1990))

|  | **Tokens = $N$** | **Types = $|V|$** |
|---|---|---|
| Switchboard phone conversations | 2.4 million | 20 thousand |
| Shakespeare | 884,000 | 31 thousand |
| Google N-grams | 1 trillion | 13 million |

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Tokenization using Unix/Linux

For a text file, output the word tokens and their frequencies.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Tokenization using Unix/Linux

For a text file, output the word tokens and their frequencies.

- change all non alphabets to new line
  tr -sc 'A-Za-z' '\n' < mytext.txt

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Tokenization using Unix/Linux

For a text file, output the word tokens and their frequencies.

- change all non alphabets to new line
  tr -sc 'A-Za-z' '\n' < mytext.txt

- sort in alphabetical order
  tr -sc 'A-Za-z' '\n' < mytext.txt | sort

**Tokenization**
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Tokenization using Unix/Linux

For a text file, output the word tokens and their frequencies.

- change all non alphabets to new line
  tr -sc 'A-Za-z' '\n' < mytext.txt

- sort in alphabetical order
  tr -sc 'A-Za-z' '\n' < mytext.txt | sort

- merge and count each type tr -sc 'A-Za-z' '\n' < mytext.txt |
  sort | uniq -c

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Tokenization using Unix/Linux

For a text file, output the word tokens and their frequencies.

- change all non alphabets to new line
  tr -sc 'A-Za-z' '\n' < mytext.txt

- sort in alphabetical order
  tr -sc 'A-Za-z' '\n' < mytext.txt | sort

- merge and count each type tr -sc 'A-Za-z' '\n' < mytext.txt | sort | uniq -c

- sorting the count tr -sc 'A-Za-z' '\n' < mytext.txt | sort | uniq -c | sort -n -r

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Issues in Tokenization

- Finland's capital $\rightarrow$ Finland Finlands Finland's ?

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
**Issues in Tokenization**
Word Segmentation

## Issues in Tokenization

- Finland's capital $\rightarrow$ Finland Finlands Finland's ?
- what're, I'm, isn't $\rightarrow$ What are, I am, is not

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Issues in Tokenization

- Finland's capital $\rightarrow$ Finland Finlands Finland's ?
- what're, I'm, isn't $\rightarrow$ What are, I am, is not
- Hewlett-Packard $\rightarrow$ Hewlett Packard ?

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
**Issues in Tokenization**
Word Segmentation

## Issues in Tokenization

- Finland's capital $\rightarrow$ Finland Finlands Finland's ?
- what're, I'm, isn't $\rightarrow$ What are, I am, is not
- Hewlett-Packard $\rightarrow$ Hewlett Packard ?
- state-of-the-art $\rightarrow$ state of the art ?

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
**Issues in Tokenization**
Word Segmentation

## 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 ?

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## 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?

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## 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. → ??

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

- French

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

- French
  - L'ensemble → one token or two?
    L ? L' ? Le ?
    Want l'ensemble to match with un ensemble
- German noun compounds are not segmented

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

- 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

## Issues in Tokenization: Language

- Chinese and Japanese no spaces between words:

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

- 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Issues in Tokenization: Language

- 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

- Japanese has more complications, with multiple alphabets intermingled

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Word Tokenization in Chinese

Also called Word Segmentation

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Word Tokenization in Chinese

Also called Word Segmentation

- Chinese words are composed of characters

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Word Tokenization in Chinese

Also called Word Segmentation

- Chinese words are composed of characters
  - Average word is 2.4 characters long.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Word Tokenization in Chinese

Also called Word Segmentation

- Chinese words are composed of characters
  - Average word is 2.4 characters long.
- Standard baseline segmentation algorithm:

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Word Tokenization in Chinese

Also called Word Segmentation

- Chinese words are composed of characters
  - Average word is 2.4 characters long.
- Standard baseline segmentation algorithm:
  - Maximum Matching (also known as Greedy Matching)

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Maximum Matching: Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Maximum Matching: Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
    1. Start a pointer at the beginning of the string

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Maximum Matching: Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
    1. Start a pointer at the beginning of the string
    2. Find the longest word in dictionary that matches the string starting at pointer

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
Word Segmentation

# Maximum Matching: Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
  1. Start a pointer at the beginning of the string
  2. Find the longest word in dictionary that matches the string starting at pointer
  3. Move the pointer over the word in string

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Maximum Matching: Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
  1. Start a pointer at the beginning of the string
  2. 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Maximum Matching: Example

- Thecatinthehat          $\rightarrow$ the cat in the hat

**Tokenization**
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

## Maximum Matching: Example

- Thecatinthehat            $\rightarrow$ the cat in the hat
- Thetabledownthere         $\rightarrow$ the table down there ???

**Tokenization**
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Maximum Matching: Example

- Thecatinthehat        $\rightarrow$ the cat in the hat
- Thetabledownthere     $\rightarrow$ the table down there ???

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Maximum Matching: Example

- Thecatinthehat $\rightarrow$ the cat in the hat
- Thetabledownthere $\rightarrow$ the table down there ???
  theta bled own there

- Doesn't generally work in English!

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

# Maximum Matching: Example

- Thecatinthehat $\rightarrow$ the cat in the hat
- Thetabledownthere $\rightarrow$ the table down there ???
  theta bled own there

- Doesn't generally work in English!
- But works really well in Chinese

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

In Unix/Linux
Issues in Tokenization
**Word Segmentation**

## Maximum Matching: Example

- Thecatinthehat          → the cat in the hat
- Thetabledownthere       → the table down there ???
                            theta bled own there
- Doesn't generally work in English!
- But works really well in Chinese
- Modern probabilistic segmentation algorithms even better

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Contents

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Normalization

- Need to "normalize" terms

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Normalization

- Need to "normalize" terms
  - Information Retrieval: indexed text & query terms must have same form.

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Normalization

- Need to "normalize" terms
  - Information Retrieval: indexed text & query terms must have same form.
    - We want to match U.S.A. and USA

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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:

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Case folding

- Applications like IR: reduce all letters to lower case

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Case folding

- Applications like IR: reduce all letters to lower case
  - Since users tend to use lower case

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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?

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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?
    - General Motors

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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?
        - General Motors
        - Fed vs. fed

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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?
    - General Motors
    - Fed vs. fed
- For sentiment analysis, MT, Information extraction

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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?
    - General Motors
    - Fed vs. fed
- For sentiment analysis, MT, Information extraction
  - Case is helpful
    US versus us is important

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
**Lemmatization**

## Lemmatization

- Reduce inflections or variant forms to base form

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# Lemmatization

- Reduce inflections or variant forms to base form
  - am, are, is → be

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
**Lemmatization**

# Lemmatization

- Reduce inflections or variant forms to base form
  - am, are, is → be
  - car, cars, car's, cars' → car

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
**Lemmatization**

# Lemmatization

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Case folding
Lemmatization

# 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

Tokenization
**Normalization**
Morphology
Stemming
Sentence Segmentation

Case folding
**Lemmatization**

# 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'

# Contents

# Morphology

- Morphemes:

# Morphology

- Morphemes:
  - The small meaningful units that make up words

## Morphology

- Morphemes:
    - The small meaningful units that make up words
    - Stems: The core meaning-bearing units

## 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

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Contents

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Stemming

- Reduce terms to their stems in information retrieval

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Stemming

- Reduce terms to their stems in information retrieval
- Stemming is crude chopping of affixes

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Stemming

- Reduce terms to their stems in information retrieval
- Stemming is crude chopping of affixes
  - language dependent

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# 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.

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# 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.

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# 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 |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Porters 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 |
| s | $\rightarrow$ | $\emptyset$ | cats | $\rightarrow$ | cat |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Porters 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 |
| s | $\rightarrow$ | $\emptyset$ | cats | $\rightarrow$ | cat |

- Step 1b:

| (*v*) ing | $\rightarrow$ | $\emptyset$ | walking | $\rightarrow$ | walk |
|-----------|------|------|---------|------|------|
| | | | sing | $\rightarrow$ | sing |
| (*v*) ed | $\rightarrow$ | $\emptyset$ | plastered | $\rightarrow$ | plaster |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Porters algorithm:
# The most common English stemmer

- Step 2: (long stems)

| ational | $\rightarrow$ | ate | relational | $\rightarrow$ | relate |
|---------|---------------|-----|------------|---------------|--------|
| izer | $\rightarrow$ | ize | digitizer | $\rightarrow$ | digitize |
| ator | $\rightarrow$ | ate | operator | $\rightarrow$ | operate |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Porters algorithm:
# The most common English stemmer

- Step 2: (long stems)

  | ational | $\rightarrow$ | ate | relational | $\rightarrow$ | relate |
  | izer | $\rightarrow$ | ize | digitizer | $\rightarrow$ | digitize |
  | ator | $\rightarrow$ | ate | operator | $\rightarrow$ | operate |

- Step 3: (longer stems)

  | al | $\rightarrow$ | $\emptyset$ | revival | $\rightarrow$ | reviv |
  | able | $\rightarrow$ | $\emptyset$ | adjustable | $\rightarrow$ | adjust |
  | ate | $\rightarrow$ | $\emptyset$ | activate | $\rightarrow$ | activ |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

## Morphology in a Corpus

| | | | | | |
|---|---|---|---|---|---|
| (\*v\*) ing | $\rightarrow$ | $\emptyset$ | walking | $\rightarrow$ | walk |
| | | | sing | $\rightarrow$ | sing |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
**Morphology in a Corpus**

# Morphology in a Corpus

(*v*) ing  →  ∅   walking  →  walk
                  sing     →  sing

tr -sc 'A-Za-z' '\n' < mytext.txt | grep 'ing$' |
sort | uniq -c | sort -n -r

| | | | | |
|---|---|---|---|---|
| 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 |

Tokenization
Normalization
Morphology
**Stemming**
Sentence Segmentation

Porter's Algorithm
Morphology in a Corpus

# Morphology in a Corpus

| (*v*) ing | $\rightarrow$ | $\emptyset$ | walking | $\rightarrow$ | walk |
|---|---|---|---|---|---|
| | | | sing | $\rightarrow$ | sing |

tr -sc 'A-Za-z' '\n' < mytext.txt | grep 'ing$' |
sort | uniq -c | sort -n -r

| 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' < mytext.txt | grep '[aeiou].*ing$' |
sort | uniq -c | sort -n -r

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Contents

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Sentence Segmentation

- !, ? are relatively unambiguous

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous
  - Sentence boundary

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous
  - Sentence boundary
  - Abbreviations like Inc. or Dr.

Tokenization
Normalization
Morphology
Stemming
**Sentence Segmentation**

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous
  - Sentence boundary
  - Abbreviations like Inc. or Dr.
  - Numbers like .02% or 4.3

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## 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 "."

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# 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 End-Of-Sentence/Not-End-Of-Sentence

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
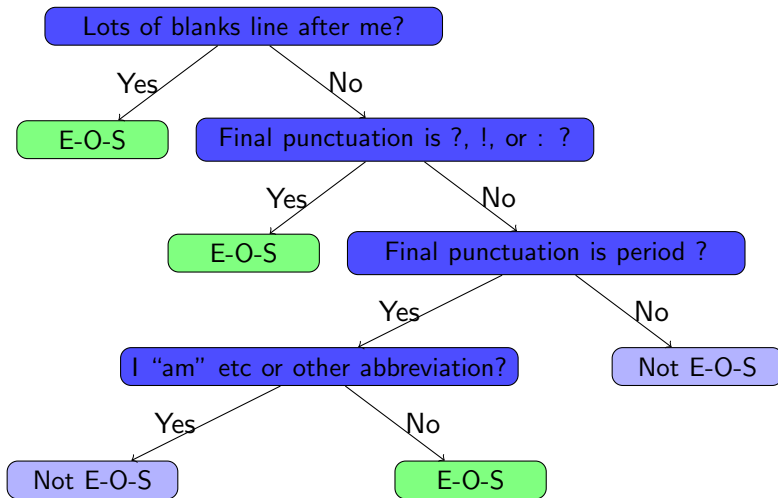Implementing DT
Other Classifiers

# 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 End-Of-Sentence/Not-End-Of-Sentence
    - Classifiers: hand-written rules, regular expressions, or machine-learning

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# A Decision Tree for E-O-S

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

- Non-Numeric features

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

- Non-Numeric features
  - Case of word with ".": Upper, Lower, Cap, Number

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

- Non-Numeric features
  - Case of word with ".": Upper, Lower, Cap, Number
  - Case of word after ".": Upper, Lower, Cap, Number
- Numeric features
  - Length of word with "."

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

- Non-Numeric 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 E-O-S)

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# More Sophisticated Features for DT

- Non-Numeric 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 E-O-S)
    - Probability (word after . occurs at B-O-S)

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Implementing Decision Tree

- A decision tree is just an if-then-else statement

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Implementing Decision Tree

- A decision tree is just an if-then-else statement
- The interesting research is choosing the features

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Implementing Decision Tree

- 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Implementing Decision Tree

- 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Implementing Decision Tree

- 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, its too hard to pick each threshold

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

# Implementing Decision Tree

- 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, its too hard to pick each threshold
  - Instead, structure usually learned by machine learning from a training corpus

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Other Classifiers

- We can think of the questions in a decision tree

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## Other Classifiers

- We can think of the questions in a decision tree
- As features that could be exploited by any kind of classifier

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## 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

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## 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
  - Support Vector Machine

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## 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
  - Support Vector Machine
  - Artificial Neural Network

Tokenization
Normalization
Morphology
Stemming
Sentence Segmentation

Decision Tree for E-O-S
Sophisticated Features
Implementing DT
Other Classifiers

## 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
  - Support Vector Machine
  - Artificial Neural Network
  - etc.