# Basic Text Processing

- Standardization
- Case-folding
- Lemmatization
- Stemming
- Sentence Segmentation

## Word Normalization Sentence Segmentation and other issues

# Word Normalization

Putting words/tokens in a standard format

- ◦ U.S.A. or USA
- ◦ uhhuh or uh-huh
- ◦ Fed or fed
- ◦ am, is, be, are

# 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)

# Lemmatization

Represent all words as their lemma, their shared root
= dictionary headword form:

◦ *am, are, is* $\rightarrow$ *be*

◦ *car, cars, car's, cars'* $\rightarrow$ *car*

◦ *He is reading detective stories*

$\rightarrow$ *He be read detective story*

# Lemmatization is done by Morphological Parsing

## Morphemes:

◦ The small meaningful units that make up words
◦ **Stems**: The core meaning-bearing units
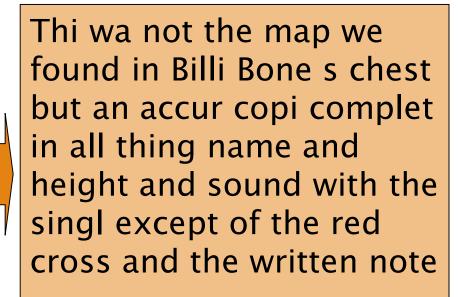◦ **Affixes**: Parts that adhere to stems, often with grammatical functions

## Morphological Parsers:

◦ Parse *cats* into two morphemes *cat* and *s*

# Stemming

Reduce terms to stems, chopping off affixes crudely

This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things-names and heights and soundings-with the single exception of the red crosses and the written notes.

Thi wa not the map we found in Billi Bone s chest but an accur copi complet in all thing name and height and sound with the singl except of the red cross and the written note .

# Porter Stemmer

Based on a series of rewrite rules run in series

◦ A cascade, in which output of each pass fed to next pass

Some sample rules:

$$\text{ATIONAL} \rightarrow \text{ATE} \quad (\text{e.g., relational} \rightarrow \text{relate})$$

$$\text{ING} \rightarrow \epsilon \quad \text{if stem contains vowel (e.g., motoring} \rightarrow \text{motor)}$$

$$\text{SSES} \rightarrow \text{SS} \quad (\text{e.g., grasses} \rightarrow \text{grass})$$

# Normalization

Case conversion

Standardization
◦ Date, currency
◦ Orthographic
◦ Spelling variation

Stemming
◦ Affix elimination

Lemmatization
◦ Geese → Goose

Pang and Lee Movie Reviews (English)

4000

# Sentence Segmentation

!, ? mostly unambiguous but **period** "." is very ambiguous

- ◦ Sentence boundary
- ◦ Abbreviations like Inc. or Dr.
- ◦ Numbers like .02% or 4.3