

# Basic Text Processing

## Regular Expressions

SLP3 slides  
(Jurafsky & Martin)

# 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              |
|---------------------------|----------------------|
| <code>[wW]oodchuck</code> | Woodchuck, woodchuck |
| <code>[1234567890]</code> | Any digit            |

- Ranges `[A-Z]`

| Pattern            | Matches              |   |
|--------------------|----------------------|---|
| <code>[A-Z]</code> | An upper case letter | <u>D</u> renched Blossoms               |
| <code>[a-z]</code> | A lower case letter  | <u>m</u> y beans were impatient         |
| <code>[0-9]</code> | A single digit       | Chapter <u>1</u> : Down the Rabbit Hole |

# Regular Expressions: Negation in Disjunction

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

| Pattern             | Matches                  |                                    |
|---------------------|--------------------------|------------------------------------|
| <code>[^A-Z]</code> | Not an upper case letter | O <u>y</u> fn pripetchik           |
| <code>[^Ss]</code>  | Neither 'S' nor 's'      | <u>I</u> have no exquisite reason" |
| <code>[^e^]</code>  | Neither e nor ^          | <u>L</u> ook here                  |
| <code>a^b</code>    | The pattern a carat b    | Look up <u>a^b</u> now             |

# Regular Expressions: More Disjunction

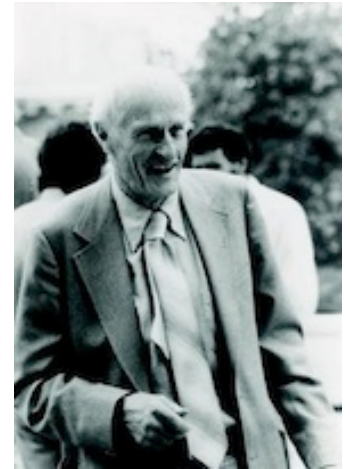
- Woodchucks is another name for groundhog!
- The pipe | for disjunction

| Pattern                                  | Matches              |
|--|----------------------|
| <code>groundhog   woodchuck</code>       |                      |
| <code>yours   mine</code>                | yours<br>mine        |
| <code>a   b   c</code>                   | = <code>[abc]</code> |
| <code>[gG]roundhog   [Ww]oodchuck</code> |                      |



# Regular Expressions: ? \* + .

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



Stephen C Kleene

Kleene \*, Kleene +

# Regular Expressions: Anchors <sup>^</sup> <sup>\$</sup>

| Pattern                 | Matches                             |
|-------------------------|-------------------------------------|
| <code>^[A-Z]</code>     | <u>P</u> alo Alto                   |
| <code>^[^A-Za-z]</code> | <u>1</u> <u>"Hello"</u>             |
| <code>\.\$</code>       | The end <u>.</u>                    |
| <code>.\$</code>        | The end <u>?    The end<u>!</u></u> |

# 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]



Refer to [http://people.cs.georgetown.edu/nshneid/cosc572/s20/02\\_py-notes.html](http://people.cs.georgetown.edu/nshneid/cosc572/s20/02_py-notes.html) and links on that page for further regex notation, and advice for using regexes in Python 3.

# 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 Normalization and  
Stemming

# 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)

# 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'



# 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 equal to compress