# Regular Expressions

# Regular expressions

- A formal language for specifying text strings

- string search methods:
  - 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 |

https://regexr.com/

# Regular Expressions: Negation in Disjunction

- Negations `[^Ss]`

  - Carat means negation only when first in []

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

# Regular Expressions: More Disjunction

- Woodchucks is another name for groundhog!

- The pipe | for disjunction

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

# Regular Expressions: ? * + . Kleene *, Kleene +

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

# Regular Expressions: Anchors  ^  $

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

# Example

- Finding "the" in an article

    ○ `the`

    ○ `[tT]he`

    ○ `[^a-zA-Z][tT]he[^a-zA-Z]`

# 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

# Summary

# Shell Script

# The IT specialist who automated his job for a year without getting caught

His script, scanned the on-site drive for any new files, generated hash values for them, and transferred them to the Cloud.

BY MEETA RAMNANI

14

https://analyticsindiamag.com/the-it-specialist-who-automated-his-job-for-a-year-without-getting-caught/

# Shell Script

- Start with **#!/bin/bash**
  - **tr**
  - **sed**
  - **grep**
  - **awk**
  - **cat**
  - **head**
  - **tail**
  - **sort**
  - **………………………..**

# 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

$\quad$ |$V$| 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    |

# Simple Tokenization in UNIX

- word tokens and their frequencies

```
tr -sc 'A-Za-z' '\n' < inputfile
     | sort
     | uniq -c
```

Change all non-alpha to newlines

Sort in alphabetical order

Merge and count each type

```
1945 A
  72 AARON
  19 ABBESS
   5 ABBOT
... ...
```

```
25 Aaron
 6 Abate
 1 Abates
 5 Abbess
 6 Abbey
 3 Abbot
…. …
```

# The first step: tokenizing

```
tr -sc 'A-Za-z' '\n' < inputfile | head
```

THE

SONNETS

by

William

Shakespeare

From

fairest

creatures

We

...

# The second step: sorting

```
tr -sc 'A-Za-z' '\n' < inputfile | sort | head
```

```
A
A
A
A
A
A
A
A
A
...
```

# More counting

- Merging upper and lower case

```
tr 'A-Z' 'a-z' < inputfile | tr -sc 'A-Za-z' '\n' | sort | uniq -c
```

- Sorting the counts

```
tr 'A-Z' 'a-z' < inputfile | 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
            10839 my
            10005 in
            8954  d
```