



Computational Linguistics Introduction

Computational Linguistics: Jordan Boyd-Graber
University of Maryland

WARMUP

Roadmap

My the end of this class you should ...

- Write regular expressions to find (linguistic) patterns
- Do simple counting using nltk
- Play around with Python interpreter
- Access corpora from nltk
- Manipulate conditional probabilities

eliza: a cautionary tale

```
Human: Men are all alike
Eliza: IN WHAT WAY?
Human: They're always bugging us about something or other
Eliza: CAN YOU THINK OF A SPECIFIC EXAMPLE?
Human: Well, my boyfriend made me come here
Eliza: YOUR BOYFRIEND MADE YOU COME HERE
Human: He says I'm depressed much of the time
Eliza: I AM SORRY TO HEAR YOU ARE DEPRESSED
Human: It's true. I am unhappy
Eliza: DO YOU THINK COMING HERE WILL HELP YOU NOT TO BE UNHAPPY
Human: I need some help, that much seems certain
```

- Claim: an electronic psychiatrist
- Is there anything interesting going on?

<http://www.masswerk.at/elizabot/>

What are eliza's tricks?

I feel Y

How often do you feel Y?

I want Y

Suppose you got Y soon ...

If Y

Do you think that it's likely that Y?

Other tricks

- Convert “my” to “your” in reply (and other pronouns)
- Randomly produce a change of subject if no rule matches: “tell me about your mother”

How do they do it?

- eliza is about finding patterns
- But users can type many different things
- We thus need a system for expressing many general patterns

How do they do it?

- eliza is about finding patterns
- But users can type many different things
- We thus need a system for expressing many general patterns
- Regular expressions

Wait a minute!

- Very stupid
- Brute-force

Wait a minute!

- Very elegant
- Low resource

Wait a minute!

- Very elegant
- Low resource
- But still require clever humans to write

Wait a minute!

- Very elegant
- Low resource
- But still require clever humans to write
- Even if you know regexps inside and out, it's important know how to apply them to language

Why in an NLP course?

- Searching for linguistic phenomena (does eat ever take the object “loss”)?
- Creating features for supervised algorithms
- Useful for morphology
- Thinking about regular expressions (nice tool) will help you think about finite state machines (theoretical framework)

Symbols and Operators

Symbol	Meaning
<code>[]</code>	Set of characters
<code>^</code>	Start of line / Negation
<code>\$</code>	End of the line
<code> </code>	Or
<code>-</code>	Range of Characters
<code>+</code>	At least one appearance
<code>*</code>	Any number of appearances
<code>{N}</code>	Exactly <i>N</i> appearances

Sets

<code>\d</code>	digits
<code>\D</code>	non-digits
<code>\s</code>	whitespace
<code>\S</code>	non-whitespace
<code>\w</code>	“words”
<code>\W</code>	non-“words”
<code>\b</code>	empty string at word start
<code>.</code>	any character except for newline

Sets

<code>\d</code>	digits	<code>[0-9]</code>
<code>\D</code>	non-digits	<code>[^0-9]</code>
<code>\s</code>	whitespace	<code>[\t\n\r\f\v]</code>
<code>\S</code>	non-whitespace	<code>[^\t\n\r\f\v]</code>
<code>\w</code>	“words”	<code>[a-zA-Z0-9_]</code>
<code>\W</code>	non-“words”	<code>[^a-zA-Z0-9_]</code>
<code>\b</code>	empty string at word start	<code>\W\b\w</code>
<code>.</code>	any character except for newline	<code>b.d</code>

Backreference

- If you enclose a subexpression in parens (a .)
- You can reference that expression again \1 (for most recent)
- For less recent, the numbers increment \2, etc.

Ranges

What does this RegEx do?

```
\b[a-z]+l
```


Ranges

What does this RegEx do?

`\b[a-z]+l`

```
^I|\.\$
```

```
I am the very model of a modern Major-General,  
I've information vegetable, animal, and mineral,  
I know the kings of England, and I quote the fights historical  
From Marathon to Waterloo, in order categorical;  
I'm very well acquainted, too, with matters mathematical,  
I understand equations, both the simple and quadratical,  
About binomial theorem I'm teeming with a lot o' news, (bothered for a rhyme)  
With many cheerful facts about the square of the hypotenuse.
```

Ranges

What does this RegEx do?

```
[aeiou]{2,}
```

Ranges

What does this RegEx do?

`[aeiou]{2,}`

```
[aeiou]{2,}
```

```
I am I the very model of a modern Major-General,  
I've information on vegetable, animal, and mineral,  
I know the kings of England, and I quote the fights historical  
From Marathon to Waterloo, in order categorical;  
I'm very well acquainted, too, with matters mathematical,  
I understand equations, both the simple and quadratical,  
About binomial theorem I'm teeming with a lot o' news, (bothered for a rhyme)  
With many cheerful facts about the square of the hypotenuse.
```

Ranges

What does this RegEx do?

```
[^aeiou]{2,}
```

Ranges

What does this RegEx do?

`[^aeiou]{2,}`

```
[^aeiou]{2,}
```

I am I the very model of a modern Major-General,
I've information vegetable, animal, and mineral,
I know the kings of England, and I quote the fights historical
From Marathon to Waterloo, in order categorical;
I'm very well acquainted, too, with matters mathematical,
I understand equations, both the simple and quadratical,
About binomial theorem I'm teeming with a lot o' news, (bothered for a rhyme)
With many cheerful facts about the square of the hypotenuse.

Ranges

What does this RegEx do?

```
[^aeiou\W]{2,}
```

```
[^aeiou\W]{2,}
```

I am I the very model of a modern Major-General,
I've information vegetable, animal, and mineral,
I know the kings of England, and I quote the fights historical
From Marathon to Waterloo, in order categorical;
I'm very well acquainted, too, with matters mathematical,
I understand equations, both the simple and quadratical,
About binomial theorem I'm teeming with a lot o' news, (bothered for a rhyme)
With many cheerful facts about the square of the hypotenuse.

Backreference

What does this RegEx do?

```
\b\w*(.)\1\w*\b
```

Backreference

What does this RegEx do?

```
\b\w*(.)\1\w*\b
```

```
\b\w*(.)\1\w*\b
```

I am I the very model of a modern Major-General,
I've information vegetable, animal, and mineral,
I know the kings of England, and I quote the fights historical
From Marathon to **Waterloo**, in order categorical;
I'm very **well** acquainted, **too**, with **matters** mathematical,
I understand equations, both the simple and quadratical,
About binomial theorem I'm **teeming** with a lot o' news, (bothered for a rhyme)
With many **cheerful** facts about the square of the hypotenuse.

Thou Must

Challenge

Find all examples of “thou ___t” in the bible; what are the most frequent?

- `nltk.corpus.gutenberg`
- `import re`
- `FreqDist` or `Counter`

Thou Must

Thou Must

```
thou_regexp = re.compile(r"[Tt]hou\s[\w]*t\s")
thou_count = FreqDist()
for ii in thou_regexp.findall(gutenberg.raw('bible-kjv.txt')):
    thou_count[ii] += 1
thou_count.tabulate(5)
```

Find a Street

Challenge

Find all examples of “Capital Word” Street in all of the Gutenberg text.



Find a Street

Find a Street

```
street_regexp = re.compile(r"[A-Z]\w*\s[S]treet")  
for fileid in gutenbergs.fileids():  
    print(fileid, street_regexp.findall(gutenberg.raw(f
```

Repeated Words

Challenge

1. Find all examples of repeated words in all of Gutenberg.
2. Find all examples of repeated words separated by some other word in Gutenberg.

- `finditer`
- `group`
- Back references

Repeated Words

Repeated Words

```
repeat_regexp = re.compile(r'\b(\w+)\s(\1\b)+')  
for fileid in gutenbergs.fileids():  
    matches = list(repeat_regexp.finditer(gutenberg.raw(fileid)))  
    print(fileid, [x.group(0) for x in matches])
```

Repeated Words (with something in between)

Repeated Words (with something in between)

```
repeat_regexp = re.compile(r"\b(\w+)\s\w+\s(\1\b)+")  
for fileid in gutenbergs.fileids():  
    matches = list(repeat_regexp.finditer(gutenberg.raw(fileid)))  
    print(fileid, [x.group(0) for x in matches])
```

Regex Golf



Regex Golf

Regexp	Matches	Doesn't Match
	afoot	Atlas
	tick	trickingly
	abac	beam
	undergrounder	hypergoddess
	civic	cinnabar
	unintelligibility	unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>f o o</code>	afoot tick abac undergrounder civic unintelligibility	Atlas trickingly beam hypergoddess cinnabar unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>f oo</code>	afoot	Atlas
<code>k \$</code>	tick	trickingly
	abac	beam
	undergrounder	hypergoddess
	civic	cinnabar
	unintelligibility	unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>foo</code>	afoot	Atlas
<code>k\$</code>	tick	trickingly
<code>^[a-f]+\$</code>	abac	beam
	undergrounder	hypergoddess
	civic	cinnabar
	unintelligibility	unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>f oo</code>	afoot	Atlas
<code>k\$</code>	tick	trickingly
<code>^[a-f]+\$</code>	abac	beam
<code>(\w3) . * \1</code>	undergrounder	hypergoddess
	civic	cinnabar
	unintelligibility	unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>f o o</code>	afoot	Atlas
<code>k\$</code>	tick	trickingly
<code>^[a-f]+\$</code>	abac	beam
<code>(\w3) .*\1</code>	undergrounder	hypergoddess
<code>(.)(.)?.*\2\1</code>	civic	cinnabar
	unintelligibility	unregainable

Regex Golf

Regexp	Matches	Doesn't Match
<code>foo</code>	<code>afoot</code>	<code>Atlas</code>
<code>k\$</code>	<code>tick</code>	<code>trickingly</code>
<code>^[a-f]+\$</code>	<code>abac</code>	<code>beam</code>
<code>(\w3) .*\1</code>	<code>undergrounder</code>	<code>hypergoddess</code>
<code>(.)(.)?.?\2\1</code>	<code>civic</code>	<code>cinnabar</code>
<code>(.)(.\1){3}</code>	<code>unintelligibility</code>	<code>unregainable</code>

Changin Gears: Bayes Rule

There's a test for Boogie Woogie Fever (BWF). The probability of getting a positive test result given that you have BWF is 0.8, and the probability of getting a positive result given that you do not have BWF is 0.01. The overall incidence of BWF is 0.01.

1. What is the marginal probability of getting a positive test result?
2. What is the probability of having BWF given that you got a positive test result?

Conditional Probabilities

One coin in a collection of 65 has two heads. The rest are fair. If a coin, chosen at random from the lot and then tossed, turns up heads 6 times in a row, what is the probability that it is the two-headed coin?