

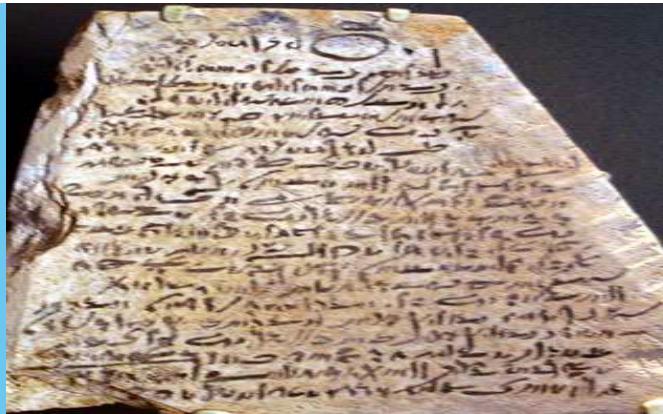
# Natural Language Processing

Week2: Text Processing with Regular Expressions, *Unix tools and Scikit-Learn*

Dr Haithem Aflī

@AflīHaithem

Haithem.aflī@mtu.ie



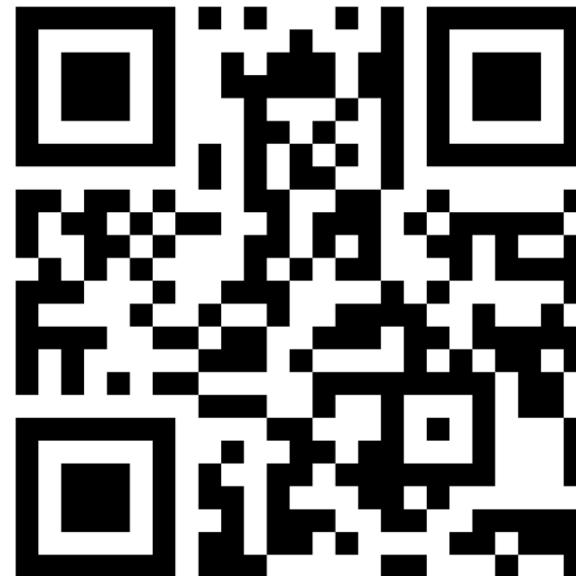
[www.mtu.ie](http://www.mtu.ie)

# Recap

Go to **www.menti.com** and use the code **9816 0270**

Or use the link/QR Code below

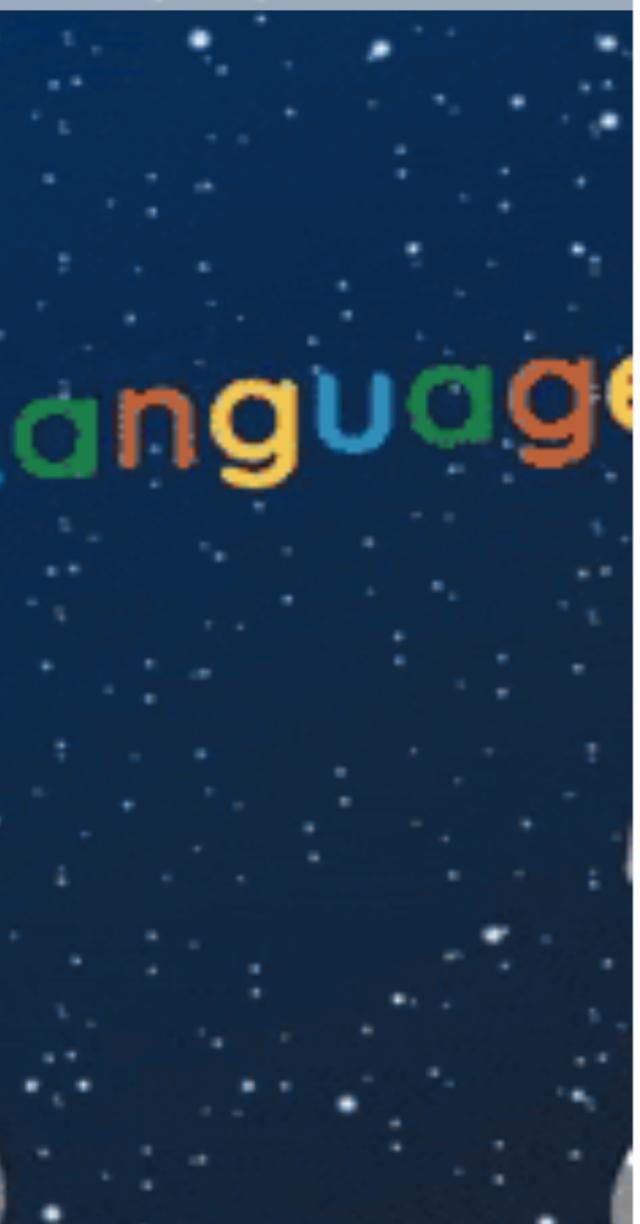
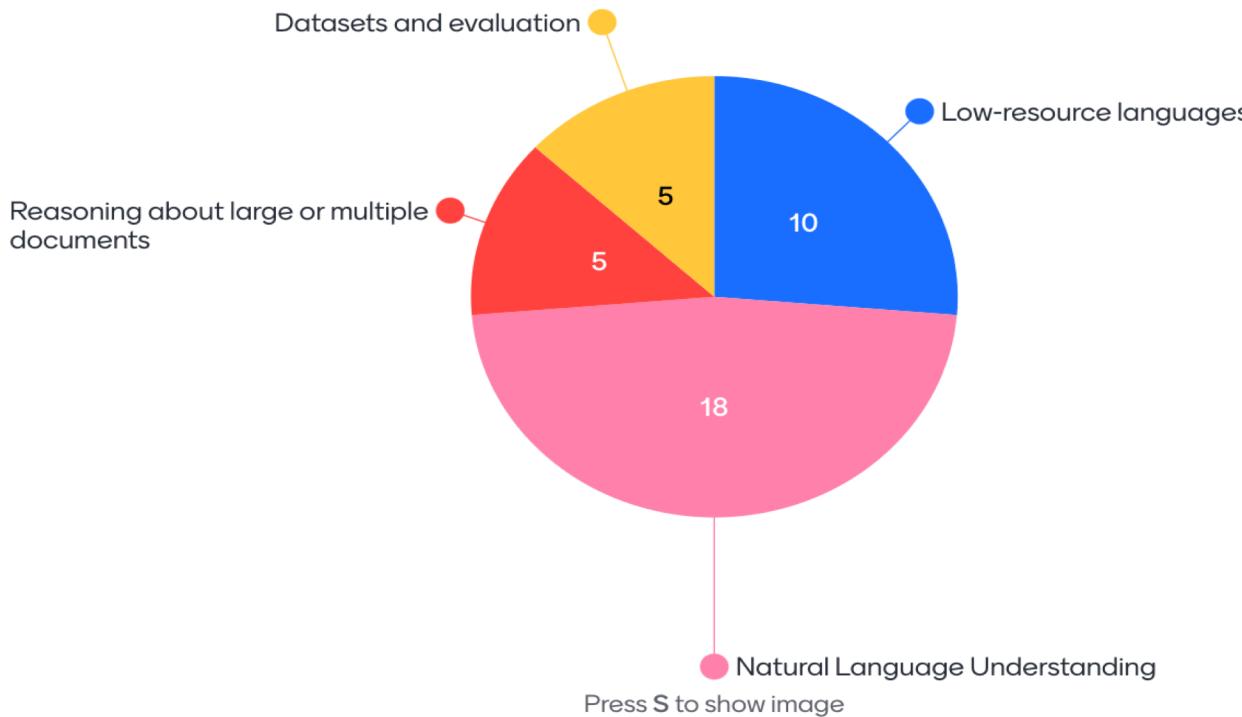
<https://www.menti.com/wxhysyji96>



# Recap

- 1) What are the main Open-Problems in natural language processing?
- 2) Describe NLP technologies with three words!

# What are the main Open-Problems in natural language processing?



# Outline

- 1. Regular Expressions**
- 2. Unix tools*
- 3. Tokenization*
- 4. Text processing with Scikit-Learn*

# Please tell me what's been bothering you.



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

**ELIZA Terminal**

Welcome to

EEEEEE	LL	III	ZZZZZZZ	AAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LL	II	ZZZ	AAAAAAA
EE	LL	II	ZZ	AA AA
EEEEEE	LLLLLL	III	ZZZZZZZ	AA AA

Eliza is a mock Rogerian psychotherapist.  
The original program was described by Joseph Weizenbaum in 1966.  
This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?  
YOU: █

This page is a sample for [mass:werk termlib.js](#) and the ELIZA library [elizabot.js](#).

For a chat with two Eliza therapists see the [ELIZA-Chat](#).

ELIZA is a natural language conversation program described by Joseph Weizenbaum in January 1966 [1]. It features the dialog between a human user and a computer program representing a mock Rogerian psychotherapist.

- Eliza is a mock Rogerian psychotherapist.
  - Rogerian psychotherapy, or person-centered therapy, is a form of talk-psychotherapy developed by psychologist Carl Rogers in the 1940s & 1950s.
- The original program was described by Joseph Weizenbaum in 1966.
  - Primitive NLP, using simple scripts and **regular expressions**
  - <http://www.masswerk.at/elizabot/eliza.html>
- Chatterbots or chatbots
- Eliza effect
- Turing test
- Regular expressions



# Please tell me what's been bothering you.

```
Welcome to
      EEEEEE  LL      IIII  ZZZZZZ  AAAAAA
      EE      LL      II     ZZ     AA     AA
      EEEEEE  LL      II     ZZZ    AAAAAAAA
      EE      LL      II     ZZ     AA     AA
      EEEEEE  LLLLLL  IIII  ZZZZZZ  AA     AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.

ELIZA: Please tell me what's been bothering you.
YOU: My boyfriend made me come here.
ELIZA: Your boyfriend made you come here ?
YOU: I am really unhappy.
ELIZA: Do you think coming here will help you not to be unhappy ?
YOU: I hope it will.
ELIZA: Do you say you hope it will for some special reason ?
YOU: 
```



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Regular Expressions

- In computer science, RE (or regex, regexp) is a language used for **specifying text search string**.
- A *regular expression* is a **formula in a special language** that is used for specifying a simple class of *string*.
- Formally, a regular expression is an algebraic notation for characterizing a set of strings.
- RE search requires
  - a *pattern* that we want to search for, and
  - a *corpus* of texts to search through.
- <http://www.regexr.com>

# Regular Expressions

- A RE search function will search through the corpus returning all texts that contain the pattern.
  - In a Web search engine, they might be the entire documents or Web pages.
  - In a word-processor, they might be individual words, or lines of a document. (We take this paradigm.)
    - E.g., the UNIX grep command

# Regular Expressions

## Basic Regular Expression Patterns

- The use of the brackets [ ] to specify a disjunction of characters.

RE	Match	Example Patterns
/ [wW] oodchuck /	Woodchuck or woodchuck	“ <u>Woodchuck</u> ”
/ [abc] /	‘a’, ‘b’, <i>or</i> ‘c’	“In uomini, in soldati”
/ [1234567890] /	any digit	“plenty of <u>7</u> to 5”

- The use of the brackets [ ] plus the dash – to specify a range.

RE	Match	Example Patterns Matched
/ [A-Z] /	an uppercase letter	“we should call it ‘ <u>Drenched Blossoms</u> ’”
/ [a-z] /	a lowercase letter	“ <u>my</u> beans were impatient to be hoed!”
/ [0-9] /	a single digit	“Chapter <u>1</u> : Down the Rabbit Hole”

# Regular Expressions

## Basic Regular Expression Patterns

- Uses of the caret ^ for negation or just to mean ^

RE	Match (single characters)	Example Patterns Matched
[^A-Z]	not an uppercase letter	“Oyfn pripetchik”
[^Ss]	neither ‘S’ nor ‘s’	“I have no exquisite reason for’t”
[^\.]	not a period	“our resident Djinn”
[e^]	either ‘e’ or ‘^’	“look up ^ now”
a^b	the pattern ‘a^b’	“look up a^b now”

- The question-mark ? marks optionality of the previous expression.

RE	Match	Example Patterns Matched
woodchucks?	woodchuck or woodchucks	“ <u>woodchuck</u> ”
colou?r	color or colour	“ <u>colour</u> ”

- The use of period . to specify any character

RE	Match	Example Patterns
/beg.n/	any character between beg and n	<u>begin</u> , <u>beg'n</u> , <u>begun</u>

# Regular Expressions: Anchors ^ \$



Pattern	Matches
<code>^ [A-Z]</code> Starting with uppercase letter	<u>Palo Alto</u>
<code>^ [ ^A-Za-z ]</code>	<u>1</u> <u>"Hello"</u>
<code>\. \$</code>	The end <u>.</u>
<code>. \$</code>	The end <u>?</u> The end <u>!</u>

# Regular Expressions

## Disjunction, Grouping, and Precedence



- Disjunction
  - /cat|dog/
- Precedence
  - /gupp(y|ies)/
- Operator precedence hierarchy

()  
+ ? { }  
the ^my end\$  
|

# Regular Expressions

## A Simple Example

- To find the English article ***the***

/the/

/ [tT] he /

/ \b [tT] he \b /

/ [ ^a-zA-Z ] [tT] he [ ^a-zA-Z ] /

/ (^ | [ ^a-zA-Z ]) [tT] he ( [ ^a-zA-Z ] | \$ ) /



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# 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)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# 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).

# Regular Expressions

## Advanced Operators

### Aliases for common sets of characters

RE	Expansion	Match	Example Patterns
\d	[0-9]	any digit	Party_of_5
\D	[^0-9]	any non-digit	Blue_moon
\w	[a-zA-Z0-9_]	any alphanumeric or space	Daiyu
\W	[^\w]	a non-alphanumeric	!!!
\s	[ \r\t\n\f]	whitespace (space, tab)	
\S	[^\s]	Non-whitespace	in_Concord

# Regular Expressions: ? \* + .



Pattern	Matches
colou?r	Optional previous char
oo*h!	0 or more of previous char
o+h!	1 or more of previous char
baa+	
beg.n	any char



Stephen C Kleene

Kleene \*, Kleene +

# Regular Expressions

## Advanced Operators

### Regular expression operators for counting

RE	Match
*	zero or more occurrences of the previous char or expression
+	one or more occurrences of the previous char or expression
?	exactly zero or one occurrence of the previous char or expression
{n}	$n$ occurrences of the previous char or expression
{n, m}	from $n$ to $m$ occurrences of the previous char or expression
{n, }	at least $n$ occurrences of the previous char or expression

# Regular Expressions

## Advanced Operators

Some characters that need to be backslashed

RE	Match	Example Patterns Matched
\*	an asterisk “*”	“K_A_P*L*A*N”
\.	a period “.”	“Dr._Livingston, I presume”
\?	a question mark	“Would you light my candle?”
\n	a newline	
\t	a tab	

# Regular Expressions

## Regular Expression Substitution, Memory, and ELIZA

`s/regexp1/regexp2/`

- E.g. *the 35 boxes* → *the <35> boxes*

`s/ ([0-9]+) /<\1>/`

- The following pattern matches “*The bigger they were, the bigger they will be*”, not “*The bigger they were, the faster they will be*”  
`/the (.*)er they were, the\1er they will be/`
- The following pattern matches “*The bigger they were, the bigger they were*”, not “*The bigger they were, the bigger they will be*”  
`/the (.*)er they (.*) , the\1er they \2/`

registers

# Regular Expressions

## Regular Expressions Substitution, Memory, and ELIZA



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

- Eliza worked by having a cascade of regular expression substitutions that each match some part of the input lines and changed them
  - *my → YOUR, I'm → YOU ARE ...*

```
s/.* YOU ARE (depressed|sad) .*/I AM SORRY TO HEAR YOU ARE \1/
```

```
s/.* YOU ARE (depressed|sad) .*/WHY DO YOU THINK YOU ARE \1/
```

```
s/.* all .* /IN WHAT WAY/
```

```
s/.* always .* /CAN YOU THINK OF A SPECIFIC EXAMPLE/
```

User<sub>1</sub>: Men are all alike.

ELIZA<sub>1</sub>: IN WHAT WAY

User<sub>2</sub>: They're always bugging us about something or other.

ELIZA<sub>2</sub>: CAN YOU THINK OF A SPECIFIC EXAMPLE

User<sub>3</sub>: Well, my boyfriend made me come here.

ELIZA<sub>3</sub>: YOUR BOYFRIEND MADE YOU COME HERE

User<sub>4</sub>: He says I'm depressed much of the time.

ELIZA<sub>4</sub>: I AM SORRY TO HEAR YOU ARE DEPRESSED



# Discussion

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Outline

1. Regular Expressions
2. ***Unix tools***
3. *Tokenization*
4. *Text processing with Scikit-Learn*

# Unix Tools



- **grep**: search for a pattern (regex)
- **sort**
- **uniq -c** (count duplicates)
- **tr** (translate characters)
- **wc** (word – or line – count)
- **sed** (edit string -- replacement)
- **cat** (send file(s) in stream)
- **man** (manual)
- *Piping commands together can be simple yet powerful in Unix*
- *It gives flexibility.*
- *Traditional Unix philosophy: small tools that can be composed*
- **cut** (columns in tab-separated files)
- **paste** (paste columns)
- **head** (top of the file)
- **tail** (bottom of the file)
- **rev** (reverse lines)
- **comm** (compare)
- **shuf** (shuffle lines of text)
- Input/output redirection:
  - >, <, |



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Simple Tokenization in UNIX

- (Inspired by Ken Church's UNIX for Poets.)
- Given a text file, output the word tokens and their frequencies

```
tr -sc 'A-Za-z' '\n' < shakes.txt
      | sort
      | uniq -c
```

1945 A	25 Aaron
72 AARON	6 Abate
19 ABBESS	1 Abates
5 ABBOT	5 Abbess
	6 Abbey
	3 Abbot

Change all non-alpha to newlines

Sort in alphabetical order

Merge and count each type

# Simple Tokenization in UNIX

- (Inspired by Ken Church's UNIX for Poets.)
- Given a text file, output the word tokens and their frequencies

```
cat shakes.txt
| tr -sc 'A-Za-z' '\n'
| sort
| uniq -c
```

1945 A  
72 AARON  
19 ABBESS  
5 ABBOT

25 Aaron  
6 Abate  
1 Abates  
5 Abbess  
6 Abbey  
3 Abbot

.... ....

Change all non-alpha to newlines

Sort in alphabetical order

Merge and count each type



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Counting

- Merging upper and lower case

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

- Sorting the counts

```
tr 'A-Z' 'a-z' < shakes.txt | 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

What happened here?

# Shakespeare - Sonnet 50



How heavy do I journey on the way,  
When what I seek, my weary travel's end,  
Doth teach that ease and that repose to say,  
'Thus far the miles are measured from thy friend!'  
The beast that bears me, tired with my woe,  
Plods dully on, to bear that weight in me,  
As if by some instinct the wretch did know  
His rider **lov'd** not speed being made from thee.

....

***Not to mention contractions such as had/would:  
he'd said, he'd say...***



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Some of the output

```
• tr -sc 'A-Za-z'  
'\n' <  
nyt_200811.txt  
| sort | uniq -  
c | head -n 5  
  
25476 a  
1271 A  
3 AA  
3 AAA  
1 Aalborg
```

- tr -sc 'A-Za-z'  
'\n' <  
nyt\_200811.txt  
| sort | uniq -  
c | head
- Gives you the first 10 lines
- tail does the same with the end of the input
- (You can omit the “-n” but it’s discouraged.)

# Sorting and reversing lines of text



- sort
- sort -f              Ignore case
- sort -n              Numeric order
- sort -r              Reverse sort
- sort -nr             Reverse numeric sort
  
- echo "Hello" | rev

# grep

- Grep finds patterns specified as regular expressions
  - **globally** search for **regular expression** and **print**
- Finding words ending in –ing:
- `grep 'ing$' nyt.words | sort | uniq -c`

# grep

- grep is a filter – you keep only some lines of the input
- grep gh keep lines containing “gh”
- grep '^con' keep lines beginning with “con”
- grep 'ing\$' keep lines ending with “ing”
- grep -v gh keep lines NOT containing “gh”
  
- grep -P Perl regular expressions (extended syntax)
- grep -P '^ [A-Z ]+ \$' nyt.words | sort | uniq -c ALL UPPERCASE

# Perl Regex



MTU

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

```
cat text|head -4
```

To Sherlock Holmes she is always the woman. I have seldom heard him mention her under any other name. In his eyes she eclipses and predominates the whole of her sex. It was not that he felt any emotion akin to love for Irene Adler. All emotions, and that one particularly, were abhorrent to his cold, precise but admirably balanced mind. He was, I take it, the most perfect reasoning and observing machine that the world has seen, but as a lover he would have placed himself in a false position. He never spoke of the softer passions, save with a gibe and sneer. They were admirable things for the observer--excellent for drawing the veil from men's motives and actions. But for the trained reasoner to admit such intrusions into his own delicate and finely adjusted temperament was to introduce a distracting factor which might throw a doubt upon all his mental results. Grit in a sensitive instrument, or a crack in one of his own high-power lenses, would not be more disturbing than a strong emotion in a nature such as his. And yet there was but one woman to him, and that woman was the late Irene Adler, of dubious and questionable memory.

```
cat text|perl -pe 'tr/[a-z]/[A-Z]/;' |head -4
```

# Extended Counting Exercises



How common are different sequences of vowels (e.g., ieu)?

# Extended Counting Exercises



How common are different sequences of vowels (e.g., ieu)?

```
cat text|perl -pe 'tr/[A-Z]/[a-z]/; s/[^aieou]/\n/g;'|sort |uniq -c |sort -nr
```



# Discussion

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Recap

Go to **www.menti.com** and use the code **9816 0270**

Or use the link/QR Code below

<https://www.menti.com/wxhysyji96>





**MTU**

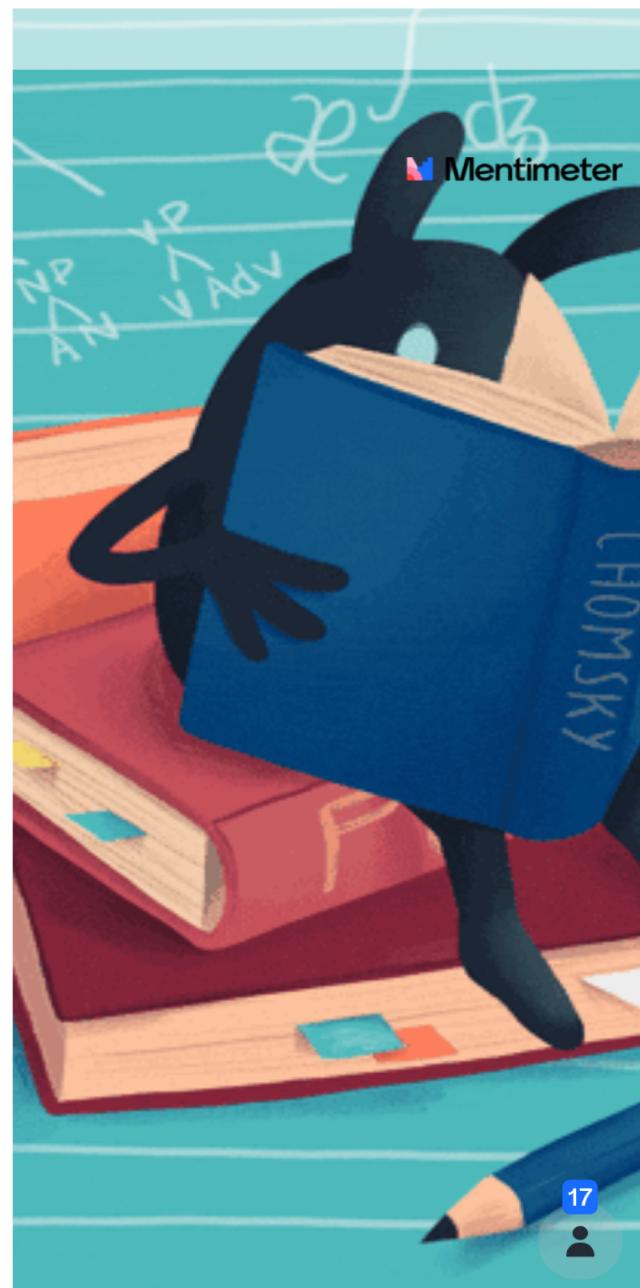
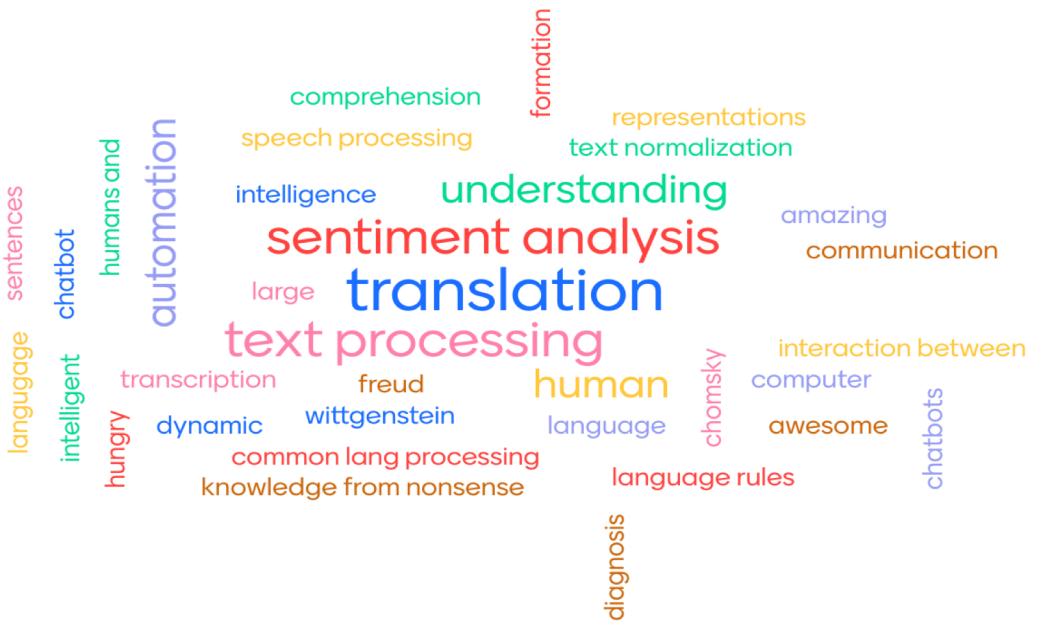
Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Recap

1) What are the main Open-Problems in natural language processing?

2) Describe NLP technologies with three words!

# Describe NLP technologies with three words!



Mentimeter



# Discussion

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Outline

1. Regular Expressions
2. *Unix tools*
3. ***Tokenization***
4. *Text processing with Scikit-Learn*



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# 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: language issues

- 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**



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Tokenization: language issues

وسيكتبها في الكتاب

wsyktbhA fy AlktAb

and he will write it in the book

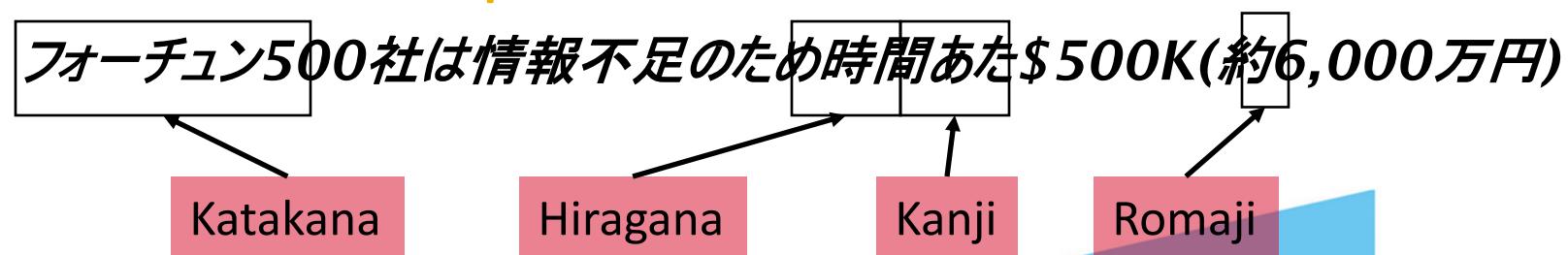
و + س + يكتب + ها + في + الـ + كتاب

w+ s+ yktb +hA fy Al+ ktAb

and he will write it in the book

# Tokenization: language issues

- 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
- Further complicated in Japanese, with multiple alphabets intermingled
  - Dates/amounts in multiple formats



End-user can express query entirely in hiragana!

# Word Tokenization in Chinese

- Also called **Word Segmentation**
- Chinese words are composed of characters
- Characters are generally 1 syllable and 1 morpheme.
- Average word is 2.4 characters long.
- Standard baseline segmentation algorithm:

➤ Maximum Matching (also called Greedy)

# 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

# Max-match segmentation illustration



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



- But works astonishingly well in Chinese
  - 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

➤ Modern probabilistic segmentation algorithms even better

# Sentence Segmentation

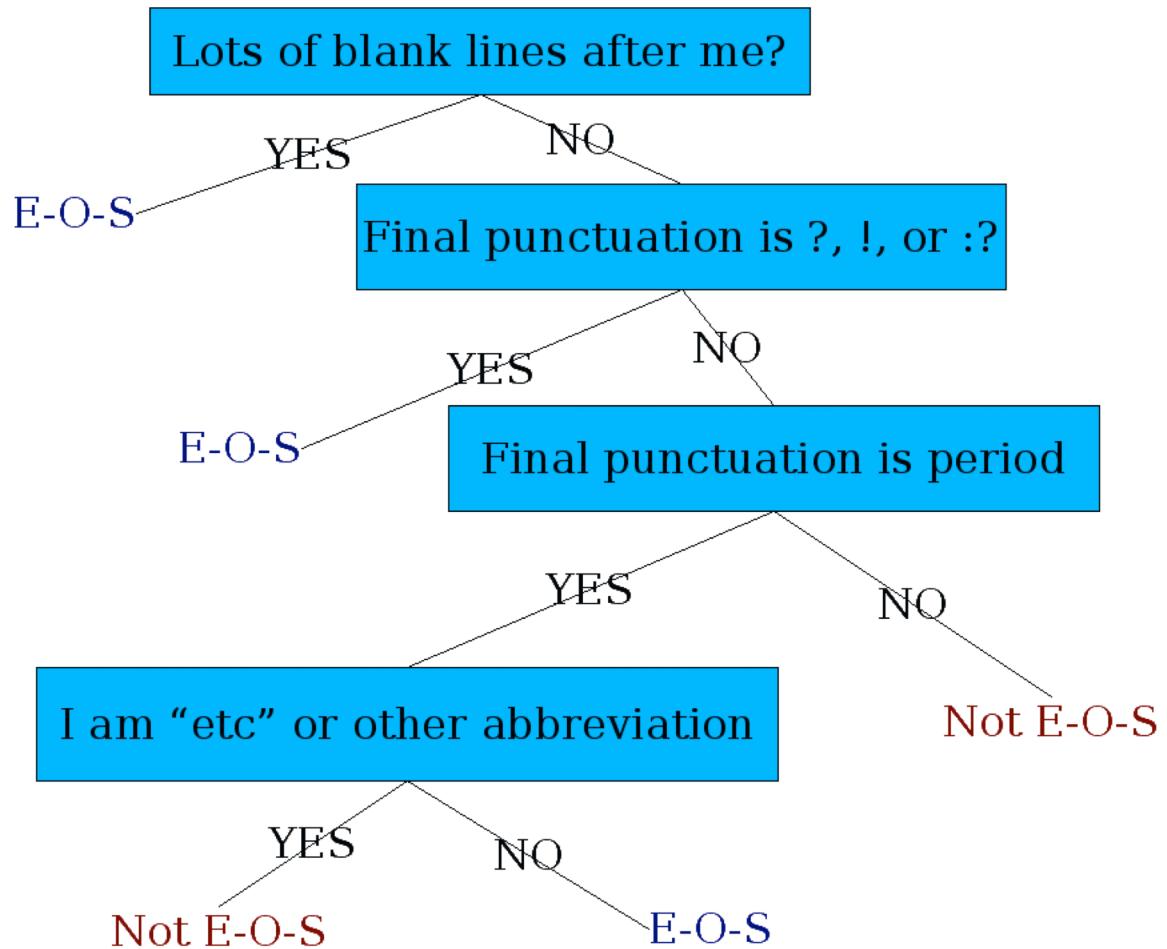
- !, ? are relatively unambiguous
- Period “.” is quite ambiguous
- Sentence boundary
- Abbreviations like Inc. or Dr.
- Numbers like .02% or 4.3
- Build a binary classifier
- Looks at a “.”
- Decides EndOfSentence/NotEndOfSentence
- Classifiers: hand-written rules, regular expressions, or machine-learning

# Determining if a word is end-of-sentence: a Decision Tree



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



# More sophisticated decision tree 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 end-of-s)
- Probability(word after “.” occurs at beginning-of-s)

# Implementing Decision Trees



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

# Decision Trees and 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
- SVM
- Neural Nets
- etc.



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



# Discussion

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Next Week

1. *Text processing with Scikit-Learn*
2. Language Modeling

# Thank You!

Haithem.afli@mtu.ie

Haithem Afli

@AfliHaithem

[www.mtu.ie](http://www.mtu.ie)