

Natural Language Processing: Introduction and Preliminaries

CSE538 - Spring 2024
Instructor: H. Andrew Schwartz

1. Computers and Natural Language
2. Goal of NLP
3. Course Overview
4. Fundamentals Review
 - a. Regular Expressions
 - b. Probability Theory
5. Words and Corpora

uL8kLyze8kz.F8Yk(.eukuL8k?.zf!

t : u
h : L
e : 8
 : k
h : L
o : y
r : z
s : e
e : 8
 : k
r : z
a : .
c : F
e : 8
d : Y
 : k
p : (
a : .
s : e
t : u
 : k
t : u
h : L
e : 8
 : k
b : ?
a : .
r : z
n : f
 . : !

uL8kLyze8kz.F8Yk(.eukuL8k?.zf!
the horse raced past the barn.

uL8kLyze8kz.F8Yk(.eukuL8k?.zf!

*Most of modern NLP language understanding works by simply analyzing the patterns of language without any external knowledge.
(over massive datasets and very large models)*

uL8kLyze8kz.F8Yk(.eukuL8k?.zf!



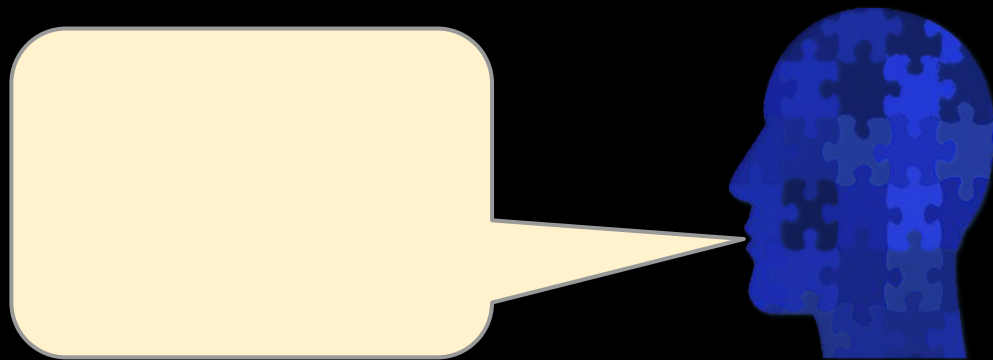
> 7 words? Likely a unique sequence



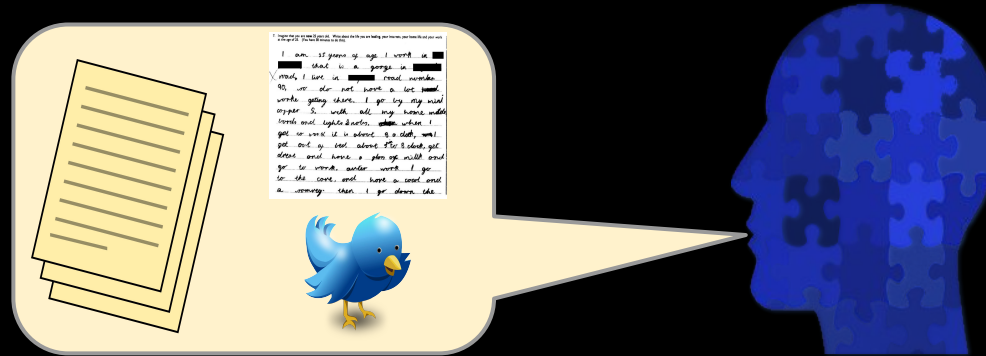
Most of modern NLP language understanding works by simply analyzing the patterns of language without any external knowledge. (over massive datasets and very large models)

Natural language is complicated!

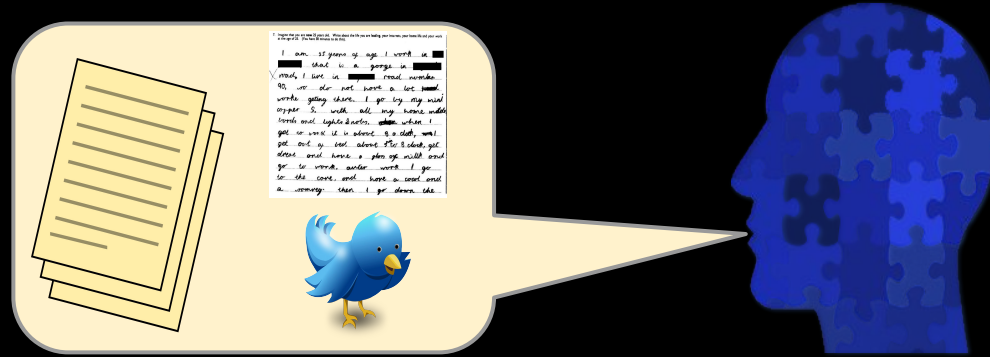
Natural language is complicated!



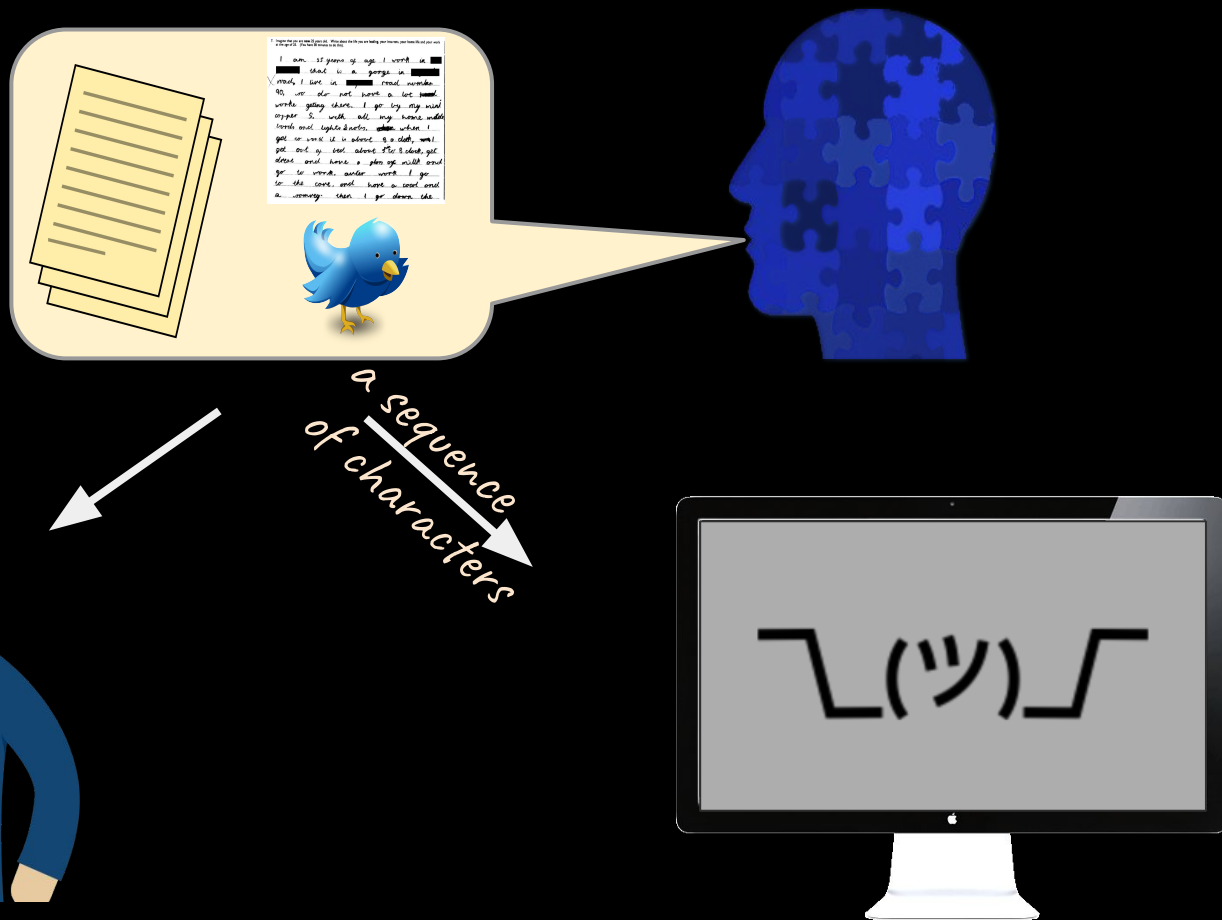
Natural language is complicated!



Natural language is complicated!

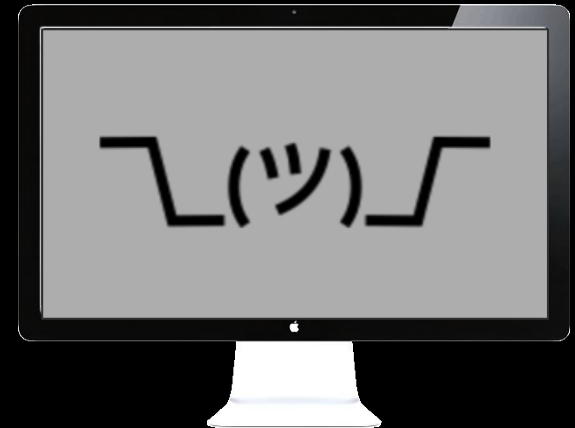


Natural language is complicated!



What is natural language like for a computer?

The horse raced past the barn.



What is natural language like for a computer?

The horse raced past the barn.

The horse raced past the barn fell.



What is natural language like for a computer?

The horse raced past the barn. ✓

The horse raced past the barn fell. ✓



What is natural language like for a computer?

The horse raced past the barn.



The horse raced past the barn fell.



The horse **runs** past the barn.



The horse **runs** past the barn fell.



What is natural language like for a computer?

The horse raced past the barn.



The horse **raced** past the barn fell.



that was

The horse **runs** past the barn.



The horse **runs** past the barn fell.

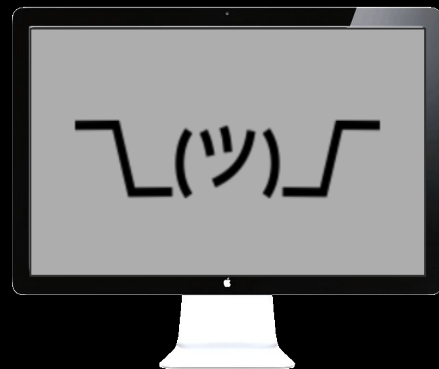


More empathy for the computer...



Colorless purple ideas sleep furiously. (Chomsky, 1956; “purple”=> “green”)

More empathy for the computer...



Colorless purple ideas sleep furiously. (Chomsky, 1956; “purple”=> “green”)

Fruit flies like a banana. Time flies like an arrow.

Daddy what did you bring that book that I don't want to be
read to out of up for?

(Pinker, 1994)

More empathy for the computer...

She ate the cake with the frosting.



More empathy for the computer...

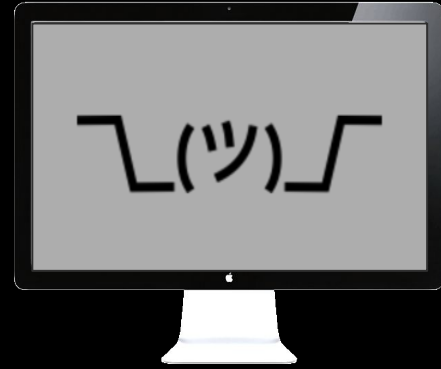
She ate the cake with the frosting.



['She', 'ate', X, 'with', Y, '.']

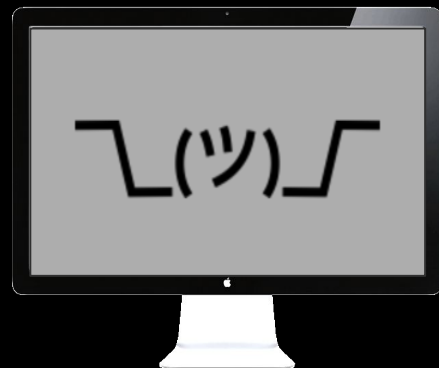
More empathy for the computer...

She ate the cake with the frosting.



['She', 'ate', X, 'with', Y, '.']
=> Y is a part of X

More empathy for the computer...



She ate the cake with the frosting.

She ate the cake with the fork.

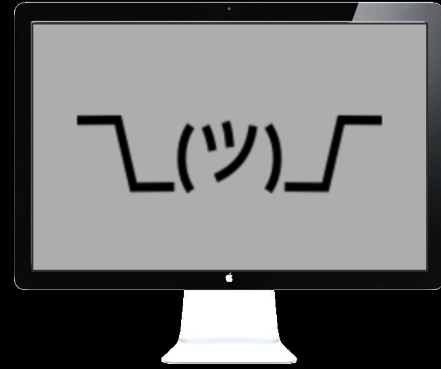
['She', 'ate', X, 'with', Y, '.']

=> Y is a part of X

More empathy for the computer...

She ate the cake with the frosting.

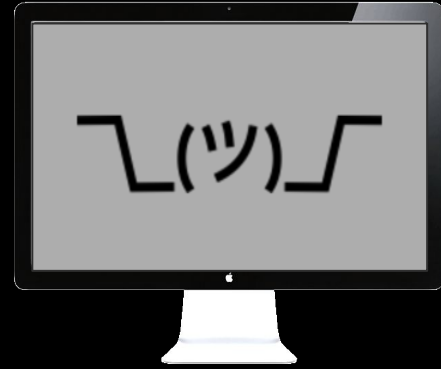
She ate the cake with the fork.



More empathy for the computer...

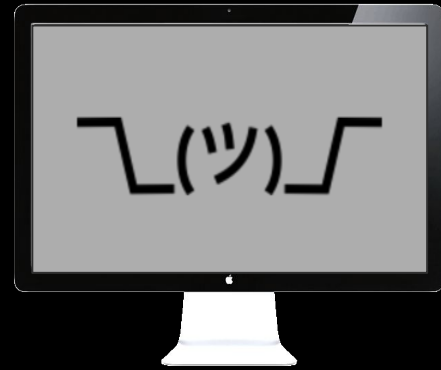
She ate the cake with the frosting.

She ate the cake with the fork.



He walked along the port next to the ship.

More empathy for the computer...



She ate the cake with the frosting.

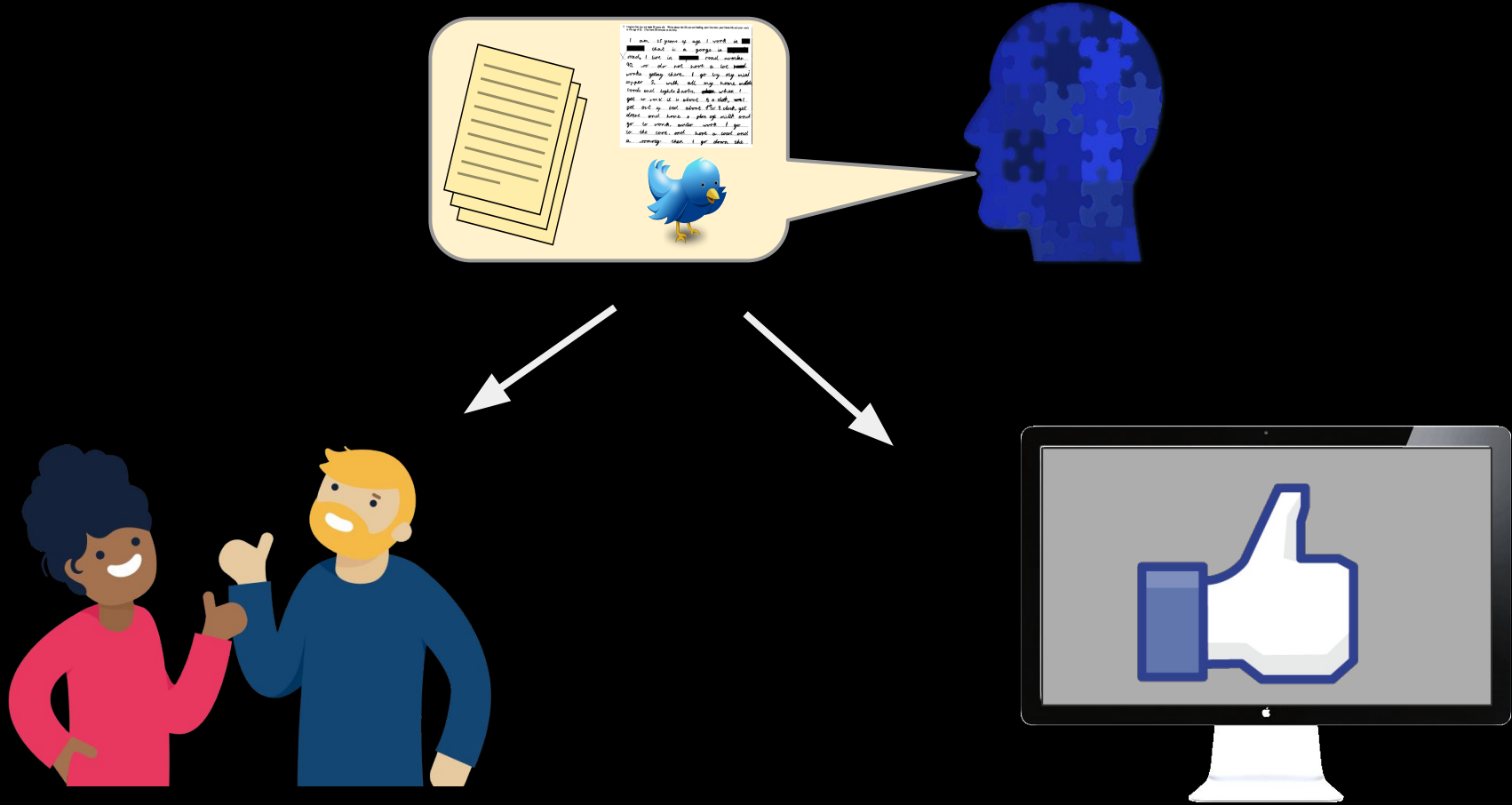
She ate the cake with the fork.

He put the **port** on the ship.

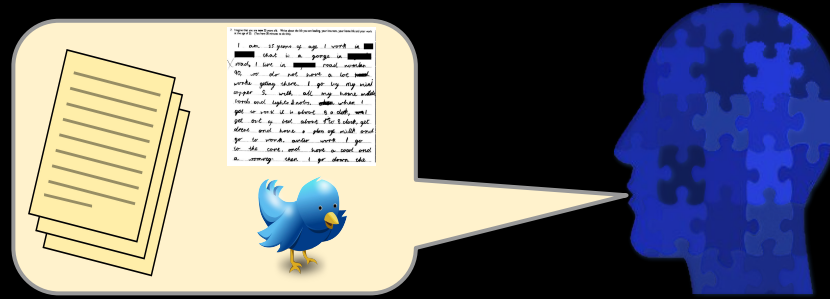
He walked along the **port** of the ship.

He walked along the **port** next to the ship.

NLP's Old grand goal: completely understand natural language.

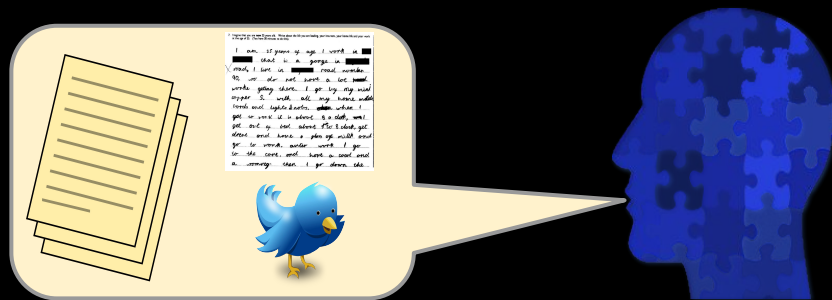


NLP's practical applications <circa 2021>



- Machine translation

NLP's practical applications



- Machine translation

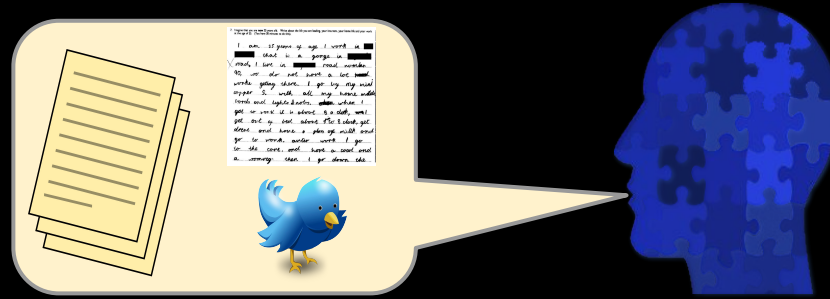
The spirit is willing, but the flesh is weak.

English -> Russian -> English

The vodka is good, but the meat is rotten.

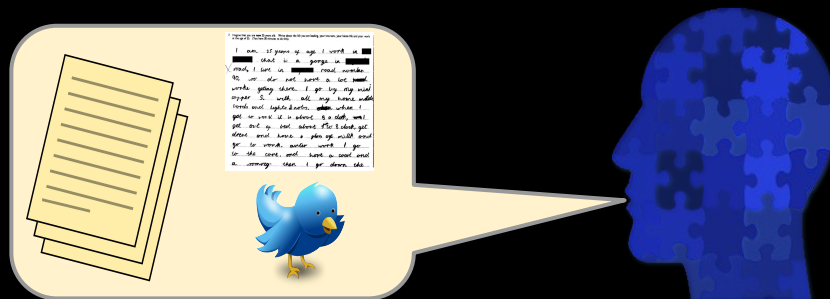
(Garbade, 2018)

NLP's practical applications



- Machine translation
- Sentiment Analysis

NLP's practical applications



- Machine translation
- Sentiment Analysis

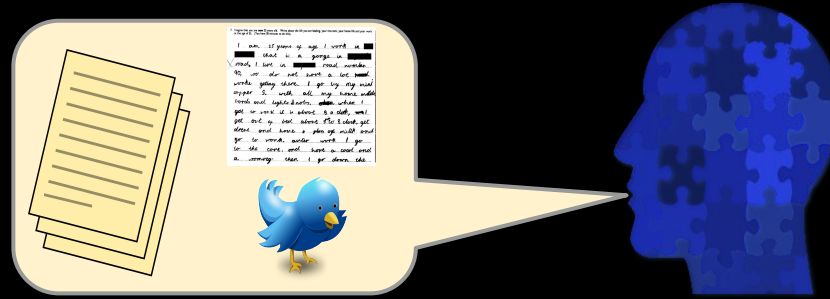
I like the the movie.



The movie is like terrible.



NLP's practical applications



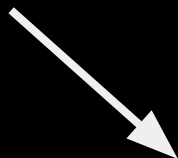
- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service

NLP's practical applications

The author of our
book is Jurafsky!



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service

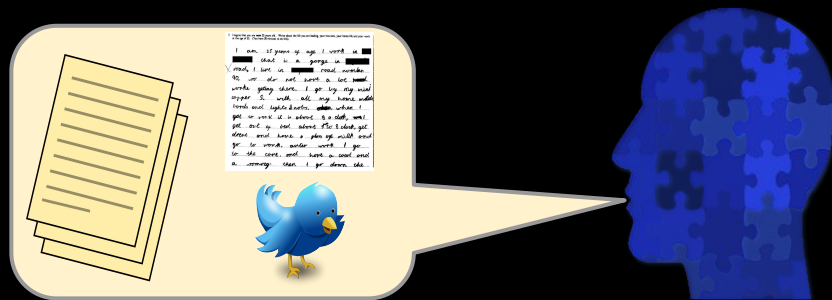


the author of our
book is giraffe
ski

ㄟ(ツ)ㄟ

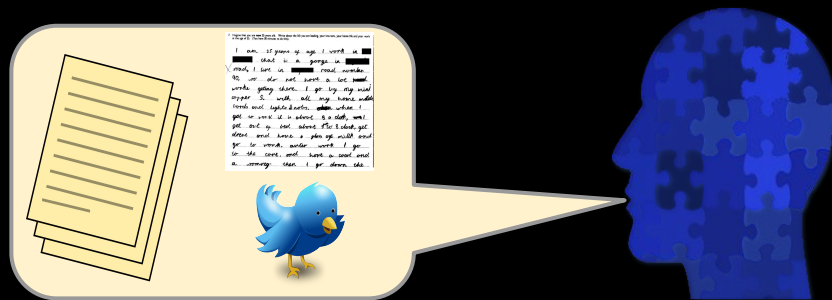


NLP's practical applications



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering

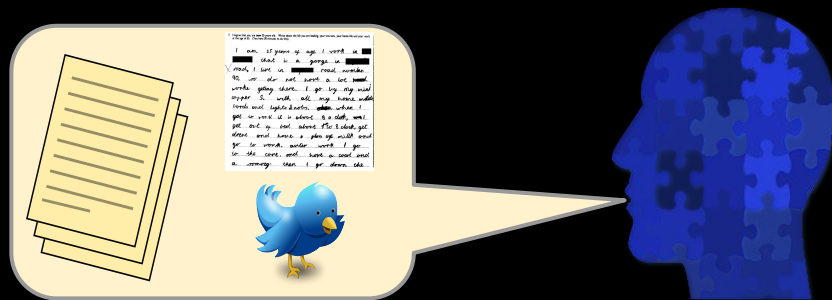
NLP's practical applications



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering

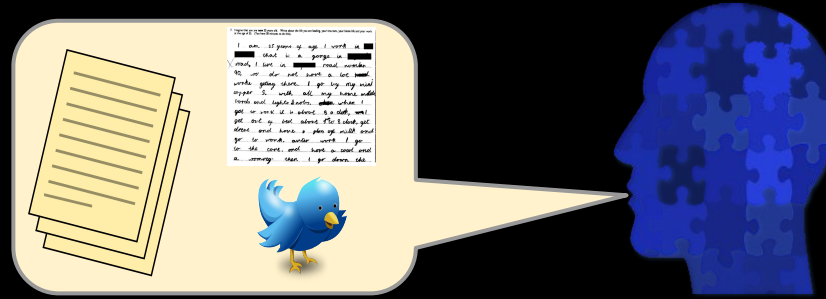
Google™

NLP's practical applications



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering
- Computational Social Science

NLP's practical applications

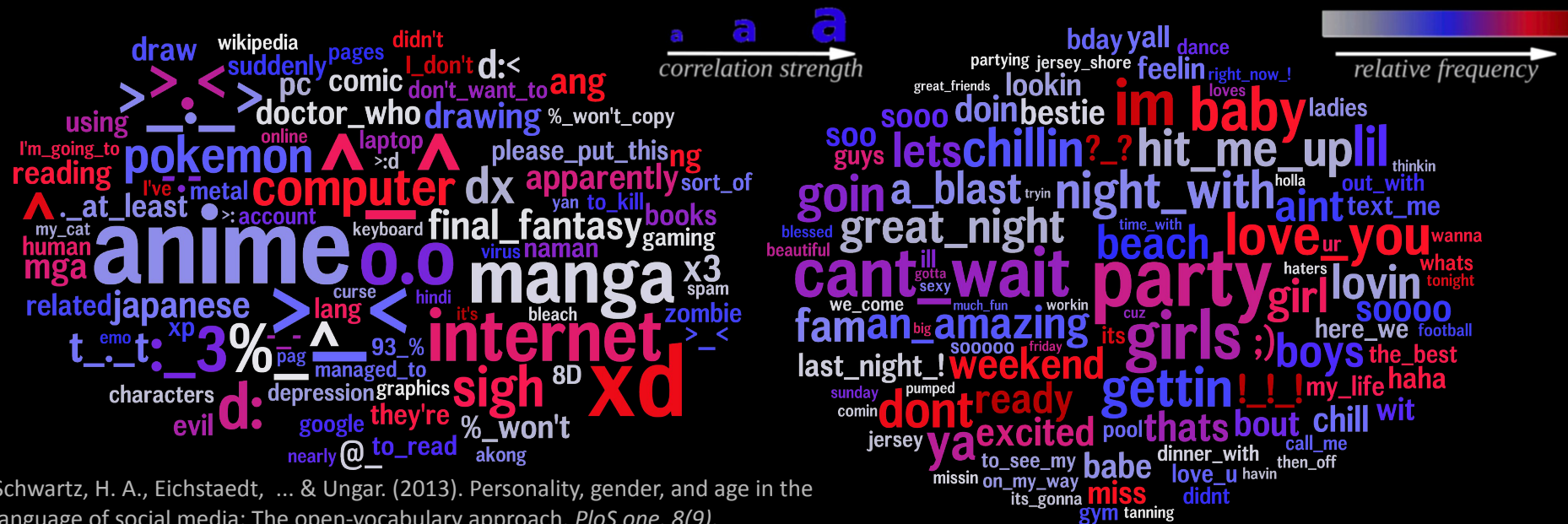
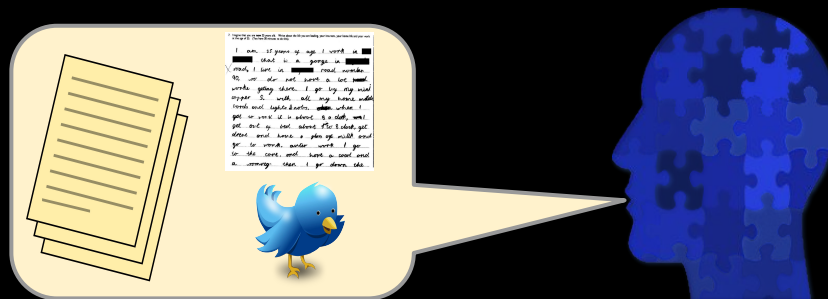


- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering
- Computational Social Science

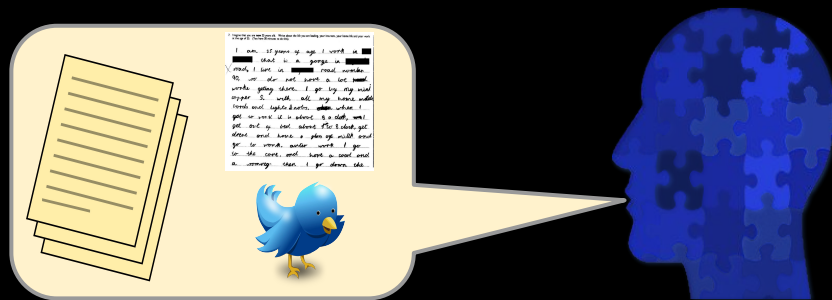


Schwartz, H. A., Eichstaedt, ... & Ungar. (2013). Personality, gender, and age in the language of social media: The open-vocabulary approach. *PLoS one*, 8(9).

NLP's practical applications

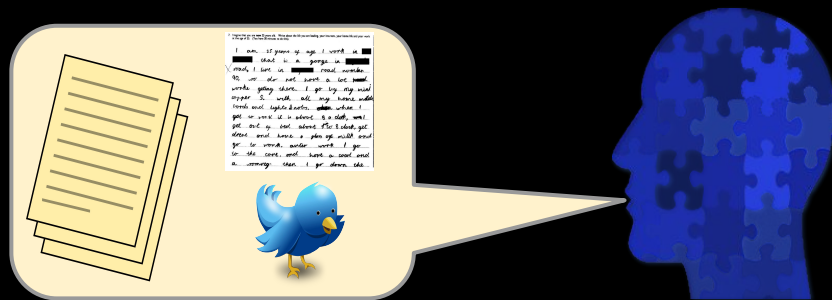


NLP's practical applications



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering
- Computational Social Science

NLP's practical applications



- Machine translation
- Sentiment Analysis
- Automatic speech recognition
 - Personalized assistants
 - Auto customer service
- Information Retrieval
 - Web Search
 - Question Answering
- Computational Social Science

LLMs have enabled:

- Open-ended information tasks. e.g.

Editing emails

Summarizing areas of work

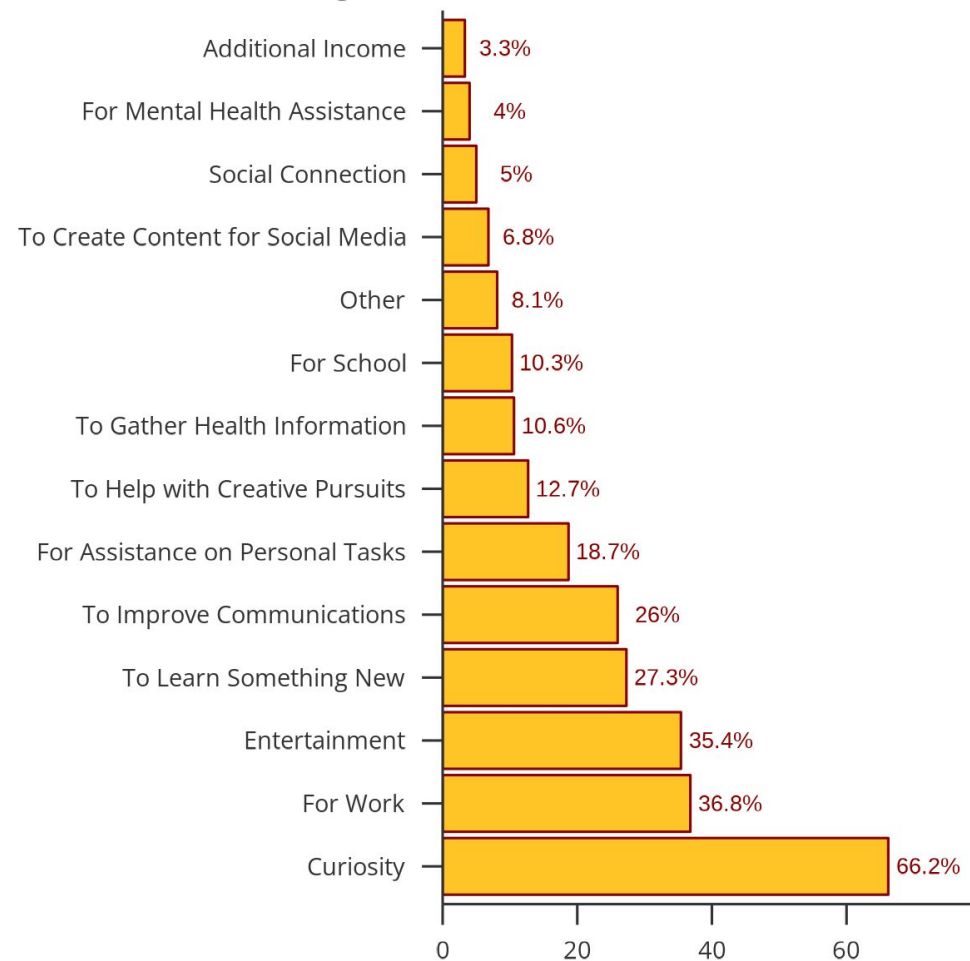
Question Answering

Counseling (not well validated)

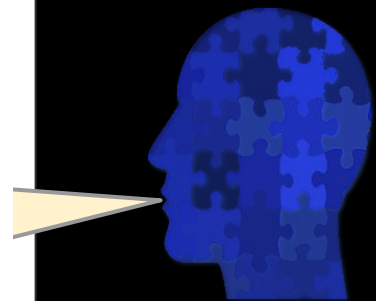
...

What do people use generative text AI tools to do?

% of US adults selecting each reason



Source: Neely Artificial Intelligence Index survey of US adults conducted September 10 - October 29, 2023.



AI has enabled:
Open-ended information tasks. e.g.
Editing emails
Summarizing areas of work
Question Answering
Counseling (not well validated)
...

NLP: The Coarse

Speech and Language Processing

An Introduction to Natural Language Processing,
Computational Linguistics, and Speech Recognition

Third Edition draft

Daniel Jurafsky
Stanford University

James H. Martin
University of Colorado at Boulder

Copyright ©2023. All rights reserved.

Draft of January 7, 2023. Comments and typos welcome!

web.stanford.edu/~jurafsky/slp3/

Course Website - Syllabus

www3.cs.stonybrook.edu/~has/CSE538/

Ingredients for success

The following covers the major components of the course and the estimated amount of time one might put into each if they are aiming to fully learn the material.

- **Review Quizzes:** 20 minutes, once a week (start second week)
- **Readings:** 2.5 hours/wk; 12 - 25 pages/wk (best before each class)
- **Study:** 1 - 2 hours/wk to review notes and look up extra content
- **Assignments (3):** 8 to 15 hours each
- **Get help early and be honest:** For anything you struggle to understand, seek office hours and extra learning suggestions.

Course Website - Syllabus

Example grade distribution;
Big Data Analytics 2023

www3.cs.stonybrook.edu/~has/CSE538/

Grade	% of class
A	34%
A-	18%
B+	7%
B	18%
B-	4%
C+	5%
C	5%
C-	9%
D	0%
F	0%

CSE538 - Preliminaries

Regular Expressions - a means for efficiently processing strings or sequences.

Use case: A basic tokenizer

Probability - a measurement of how likely an event is to occur.

Use case: How likely is “force” to be a noun?

Tokenizing Words:

tokens - an individual word instance.

types - distinct words.

CSE538 - Preliminaries

Regular Expressions - a means for efficiently processing strings or sequences.

Use case: A basic tokenizer

Probability - a measurement of how likely an event is to occur.

Use case: How likely is “force” to be a noun?

Tokenizing Words:

tokens - an individual word instance.

types - distinct words.

How many word tokens and word types?

Will, will Will will Will Will's will?

*Rose rose to put rose roes on her rows
of roses.*

Regular Expressions

The unsung hero of NLP



Regular Expressions

Patterns to match in a string.



Example:

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>in</u> gles', 'class'X

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>in</u> gles', 'class'X
[sS]bu	'sbu', 'I like Sbu a lot', 'SBU'	

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

pattern	example strings	matches
<code>ing</code>	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>in</u> gles', 'class'X
<code>[sS]bu</code>	'sbu', 'I like Sbu a lot', 'SBU'	' <u>s</u> bu', 'I like <u>S</u> bu a lot', 'SBU'X

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

character ranges: `[-]` -- matches a range of characters according to ascii order

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>ing</u> les', 'class'X
[sS]bu	'sbu', 'I like Sbu a lot', 'SBU'	' <u>sbu</u> ', 'I like <u>Sbu</u> a lot', 'SBU'X
[A-Z][a-z]	'sbu', 'Sbu' #capital followed by lowercase	
[0-9][MmKk]	'5m', '50m', '2k', '2b'	

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

character ranges: `[-]` -- matches a range of characters according to ascii order

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>ing</u> les', 'class'X
[sS]bu	'sbu', 'I like Sbu a lot', 'SBU'	' <u>sbu</u> ', 'I like <u>Sbu</u> a lot', 'SBU'X
[A-Z][a-z]	'sbu', 'Sbu' #capital followed by lowercase	'sbu'X, ' <u>Sbu</u> '
[0-9][MmKk]	'5m', '50m', '2k', '2b'	' <u>5m</u> ', ' <u>50m</u> ', ' <u>2k</u> ', '2b'X

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

character ranges: `[-]` -- matches a range of characters according to ascii order

not characters: `[^]` -- matches any character except this

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>ing</u> les', 'class'X
[sS]bu	'sbu', 'I like Sbu a lot', 'SBU'	' <u>sbu</u> ', 'I like <u>Sbu</u> a lot', 'SBU'X
[A-Z][a-z]	'sbu', 'Sbu' #capital followed by lowercase	'sbu'X, ' <u>Sbu</u> '
[0-9][MmKk]	'5m', '50m', '2k', '2b'	' <u>5m</u> ', '50m'X, ' <u>2k</u> ', '2b'X
ing[^s]	'kicking ', 'holdings ', 'ingles ', 'kicking'	

Regular Expressions

Patterns to match in a string.

character class: `[]` --matches any single character inside brackets

character ranges: `[-]` -- matches a range of characters according to ascii order

not characters: `[^]` -- matches any character except this

pattern	example strings	matches
ing	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>ing</u> les', 'class'X
[sS]bu	'sbu', 'I like Sbu a lot', 'SBU'	' <u>sbu</u> ', 'I like <u>Sbu</u> a lot', 'SBU'X
[A-Z][a-z]	'sbu', 'Sbu' #capital followed by lowercase	'sbu'X, ' <u>Sbu</u> '
[0-9][MmKk]	'5m', '50m', '2k', '2b'	' <u>5m</u> ', '50m'X, ' <u>2k</u> ', '2b'X
ing[^s]	'kicking ', 'holdings ', 'ingles ', 'kicking'	'kick <u>ing</u> ', 'holdings 'X, ' <u>ing</u> les', 'kicking'X

Regular Expressions

Pattern

In python we denote regular expressions with:

`r'PATTERN'`

character

character range

not characters

according to ascii order

matches any character except this

pattern	example strings	matches
<code>r'ing'</code>	'kicking', 'ingles', 'class'	'kick <u>ing</u> ', ' <u>ing</u> les', 'class' X
<code>r'[sS]bu'</code>	'sbu', 'I like Sbu a lot', 'SBU'	' <u>sbu</u> ', 'I like <u>Sbu</u> a lot', 'SBU' X
<code>r'[A-Z][a-z]'</code>	'sbu', 'Sbu' #capital followed by lowercase	'sbu' X , ' <u>Sbu</u> '
<code>r'[0-9][MmKk]'</code>	'5m', '50m', '2k', '2b'	' <u>5m</u> ', ' <u>50m</u> ', ' <u>2k</u> ', '2b' X
<code>r'ing[^s]'</code>	'kicking ', 'holdings ', 'ingles '	'kick <u>ing</u> ', 'holdings ' X , ' <u>ing</u> les '

Regular Expressions

Matching recurring patterns:

* : match 0 or more

+ : match 1 or more

pattern	example strings	matches
r'ing!*	'swing', 'swing!' 'swing!!!' '!!!'	
r'[sS][oO]+'	'so', 'sooo', 'SOOoo', 'so!', 'soso'	

Regular Expressions

Matching recurring patterns:

* : match 0 or more

+ : match 1 or more

pattern	example strings	matches
r'ing!*	'swing', 'swing!' 'swing!!!' '!!!'	'sw <u>ing</u> ', 'sw <u>ing</u> !', 'sw <u>ing</u> !!!' '!!!'X
r'[sS][oO]+'	'so', 'sooo', 'SOOoo', 'so!', 'soso'	' <u>so</u> ', ' <u>sooo</u> ', ' <u>SOOoo</u> ', ' <u>so</u> !', ' <u>so</u> ' <u>so</u> ' #would match twice

Regular Expressions

Matching recurring patterns:

* : match 0 or more

+ : match 1 or more

? : 0 or 1

pattern	example strings	matches
r'ing!*	'swing', 'swing!' 'swing!!!' '!!!'	'sw <u>ing</u> ', 'sw <u>ing</u> !', 'sw <u>ing</u> !!!' '!!!'X
r'[sS][oO]+'	'so', 'sooo', 'SOOoo', 'so!', 'soso'	' <u>so</u> ', ' <u>sooo</u> ', ' <u>SOOoo</u> ', ' <u>so</u> !', ' <u>so</u> ' <u>so</u> ' #would match twice
r'oranges?	'orange', 'oranges', 'orangess'	

Regular Expressions

Matching recurring patterns:

* : match 0 or more

+ : match 1 or more

? : 0 or 1

pattern	example strings	matches
r'ing!*	'swing', 'swing!' 'swing!!!' '!!!'	'sw <u>ing</u> ', 'sw <u>ing</u> !', 'sw <u>ing</u> !!!' '!!!'X
r'[sS][oO]+'	'so', 'sooo', 'SOOoo', 'so!', 'soso'	' <u>so</u> ', ' <u>sooo</u> ', ' <u>SOOoo</u> ', ' <u>so</u> !', ' <u>so</u> ' <u>so</u> ' #would match twice
r'oranges?	'orange', 'oranges', 'orangess'	' <u>orange</u> ', ' <u>oranges</u> ', ' <u>orangess</u> ' #matches all it can

Regular Expressions

Patterns applied to groups of characters

AA|BB : matches group AA or group BB

pattern	example strings	matches
r'hers his theirs"	'this is hers', 'this is his!'	'this is <u>hers</u> ', 'this is <u>his</u> !'

Regular Expressions

Patterns applied to groups of characters

AA|BB : matches group AA or group BB

(AA) : apply any following operations to group

pattern	example strings	matches
r'hers his'	'this is hers', 'this is his!'	'this is <u>hers</u> ', 'this is <u>his</u> !'
r'([A-Z][a-z]+)+'	'This matches Cap Words followed By a Space.'	

Regular Expressions

Patterns applied to groups of characters

AA|BB : matches group AA or group BB

(AA) : apply any following operations to group

pattern	example strings	matches
r'hers his'	'this is hers', 'this is his!'	'this is <u>hers</u> ', 'this is <u>his</u> !'
r'([A-Z][a-z]+)+'	'This matches Cap Words followed By a Space.'	' <u>This</u> matches <u>Cap Words</u> followed <u>By</u> a Space.'

Regular Expressions

. : any single character

pattern	example strings	matches
.	'kicking'	<u>'k'</u> <u>'i'</u> <u>'c'</u> <u>'k'</u> ...

Regular Expressions

. : any single character

\$: end of string

pattern	example strings	matches
.	'kicking'	<u>'k'</u> <u>'i'</u> <u>'c'</u> <u>'k'</u>
.\$	'great', 'great!', '50'	

Regular Expressions

. : any single character

\$: end of string

pattern	example strings	matches
.	'kicking'	' <u>k</u> ' ' <u>i</u> ' ' <u>c</u> ' ' <u>k</u> '
.\$	'great', 'great!', '50'	'great <u>t</u> ', 'great! <u>!</u> ', '50 <u>0</u> '

Regular Expressions

. : any single character

\$: end of string

^: beginning of string

pattern	example strings	matches
.	'kicking'	' <u>k</u> ' ' <u>i</u> ' ' <u>c</u> ' ' <u>k</u> '
.\$	'great', 'great!', '50'	'great <u>t</u> ', 'great! <u>!</u> ', '50 <u>0</u> '
^.a	'Happy', 'slate', 'a', 'kick a door'	

Regular Expressions

. : any single character

\$: end of string

^: beginning of string

pattern	example strings	matches
.	'kicking'	' <u>k</u> ' ' <u>i</u> ' ' <u>c</u> ' ' <u>k</u> '
.\$	'great', 'great!', '50'	'great <u>t</u> ', 'great! <u>!</u> ', '50 <u>0</u> '
^a	'Happy', 'slate', 'a', 'kick a door'	' <u>H</u> appy', 'slate', 'a'X, 'kick a door'
a	'Happy', 'slate', 'a', 'kick a door'	' <u>H</u> appy', 's <u>l</u> ate', 'a'X, 'kick <u>a</u> door'

Regular Expressions

`\s` : matches any whitespace (space, tab, newline)

`\b` : matches a word boundary

Tokenizing -- breaking a sentence into simple lexical units (basically words).

Here are a couple simple regular expressions for tokenizing:

pattern	example strings	matches
<code>r'(\s ^)[A-z]+...</code>	'Kick a door.'	

Regular Expressions

`\s` : matches any whitespace (space, tab, newline)

`\b` : matches a word boundary

Tokenizing -- breaking a sentence into simple lexical units (basically words).

Here are a couple simple regular expressions for tokenizing:

pattern	example strings	matches
<code>r'(\s ^)[A-z]+([!?\.\,] \$)?'</code>	'Kick a door.'	

Regular Expressions

`\s` : matches any whitespace (space, tab, newline)

`\b` : matches a word boundary

Tokenizing -- breaking a sentence into simple lexical units (basically words).

Here are a couple simple regular expressions for tokenizing:

pattern	example strings	matches
<code>r'(\s ^)[A-z]+([!?\.\.])\$)?'</code>	'Kick a door.'	'Kick' ' a' ' door.'

Regular Expressions

`\s` : matches any whitespace (space, tab, newline)

`\b` : matches a word boundary

Tokenizing -- breaking a sentence into simple lexical units (basically words).

Here are a couple simple regular expressions for tokenizing:

pattern	example strings	matches
<code>r'(\s ^)[A-z]+([!?\.\,] \\$)?'</code>	'Kick a door.'	' <u>Kick</u> ' ' <u>a</u> ' ' <u>door.</u> '
<code>r'\b[A-z]+\b'</code>	'Kick a door.'	' <u>Kick</u> ', ' <u>a</u> ', ' <u>door</u> .' #3 matches, no whitespace

Regular Expressions

```
import re

words = re.findall(r'\b[A-z]+\b', sentence)

for word in words:
    print(word)
```

pattern	example strings	matches
<code>r'(\s ^)[A-z]+([!?\.\,] \$)?'</code>	'Kick a door.'	' <u>Kick</u> ' ' <u>a</u> ' ' <u>door.</u> '
<code>r'\b[A-z]+\b'</code>	'Kick a door.'	' <u>Kick</u> <u>a</u> <u>door.</u> ' #3 matches, no whitespace

Regular Expressions

```
import re

words = re.split(r'\s', sentence)

for word in words:
    print(word)
```

pattern	example strings	matches
<code>r'(\s ^)[A-z]+([!?\.\] \$)?'</code>	'Kick a door.'	' <u>Kick</u> ' ' <u>a</u> ' ' <u>door.</u> '
<code>r'\b[A-z]+\b'</code>	'Kick a door.'	' <u>Kick</u> <u>a</u> <u>door.</u> ' #3 matches, no whitespace

Probability

Probability

Definitely not the unsung hero...

Probability



Probability

1970

1980s

1990s

2000s

2010s

2020s

Rule-based and Logic Systems

Statistical NLP

Machine Learning

Deep Learning

LLMs



What is Probability?

Examples

1. outcome of flipping a coin
2. side of a die
3. mentioning a word
4. mentioning a word “a lot”

What is Probability?

The chance that something will happen.

Given infinite observations of an event, the proportion of observations where a given outcome happens.

Strength of belief that something is true.

“Mathematical language for quantifying uncertainty” - Wasserman

What is Probability?

The chance that something will happen.

Given infinite observations of an event, the proportion of observations where a given outcome happens. -- *probably describes frequency in data*

Strength of belief that something is true.

--probability describes amount of conviction toward a hypothesis

“Mathematical language for quantifying uncertainty” - Wasserman

Probability

Ω : Sample Space, set of all outcomes of a random experiment

A : Event ($A \subseteq \Omega$), collection of possible outcomes of an experiment

$P(A)$: Probability of event A , P is a function: events $\rightarrow \mathbb{R}$

Probability

Ω : Sample Space, set of all outcomes of a random experiment

A : Event ($A \subseteq \Omega$), collection of possible outcomes of an experiment

$P(A)$: Probability of event A , P is a function: events $\rightarrow \mathbb{R}$

1. $P(\Omega) = 1$

2. $P(A) \geq 0$, for all A

If A_1, A_2, \dots are disjoint events then:

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i)$$

Probability

Ω : Sample Space, set of all outcomes of a random experiment

A : Event ($A \subseteq \Omega$), collection of possible outcomes of an experiment

$P(A)$: Probability of event A , P is a function: events $\rightarrow \mathbb{R}$

P is a *probability measure*, if and only if

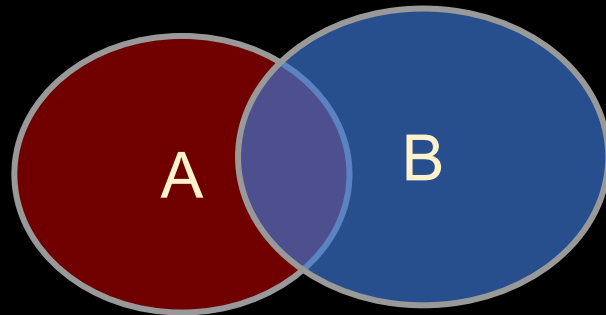
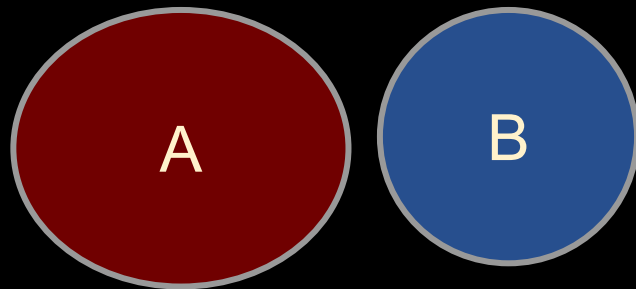
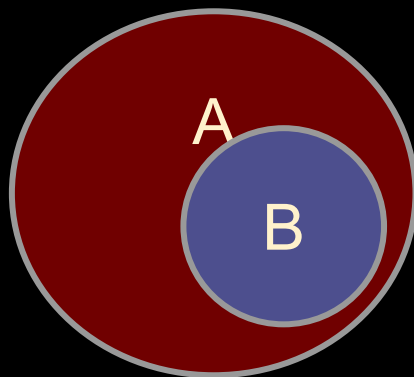
1. $P(\Omega) = 1$
2. $P(A) \geq 0$, for all A

If A_1, A_2, \dots are disjoint events then:

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i)$$

Probability

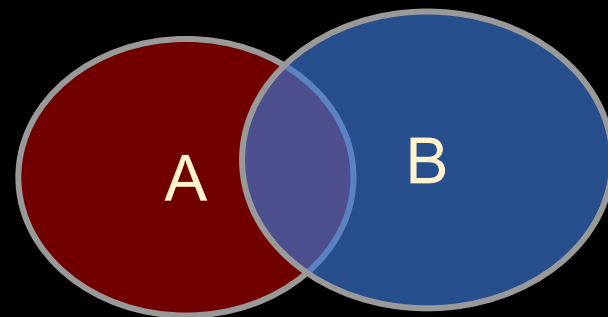
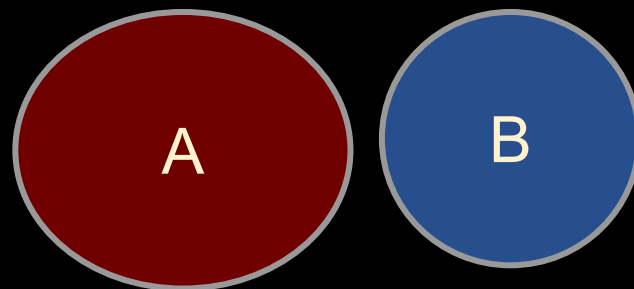
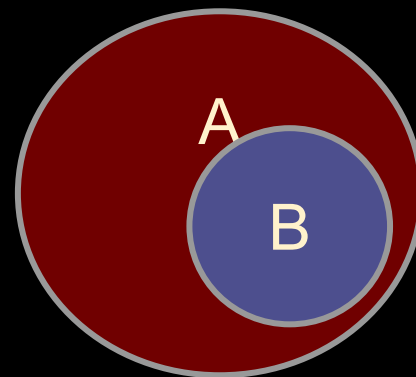
Some Properties:



Probability

Some Properties:

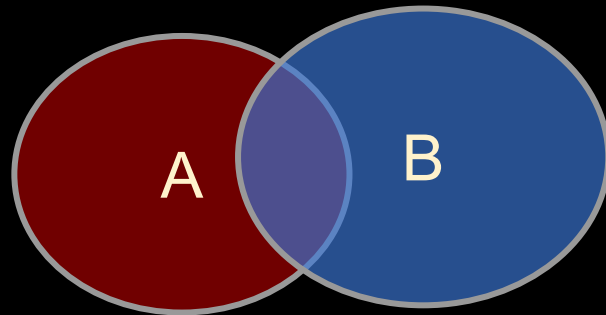
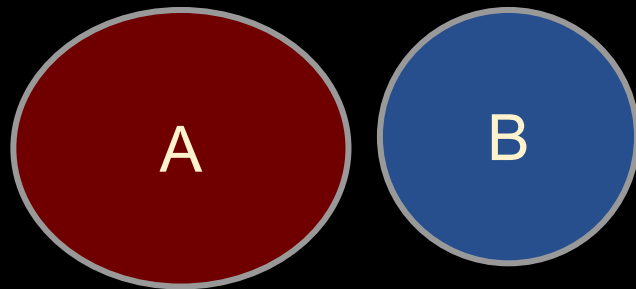
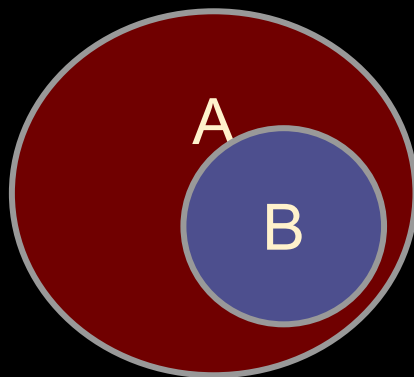
1. If $B \subseteq A$ then $P(A) \geq P(B)$



Probability

Some Properties:

1. If $B \subseteq A$ then $P(A) \geq P(B)$
2. $P(A \cup B) \leq P(A) + P(B)$

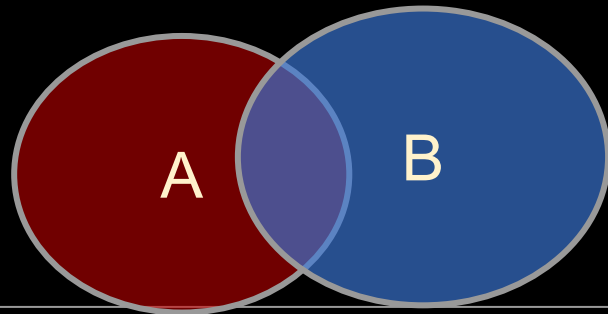
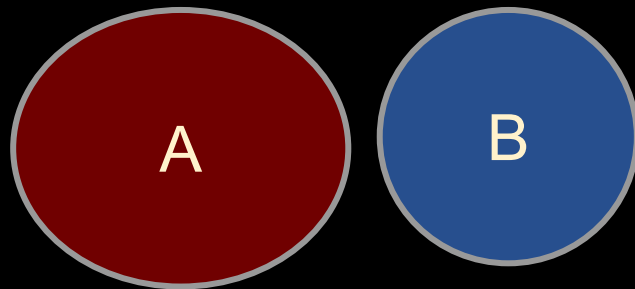
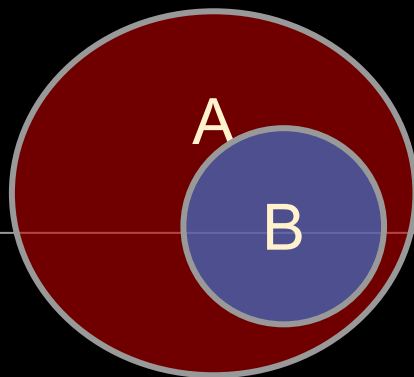


Probability

Some Properties:

1. If $B \subseteq A$ then $P(A) \geq P(B)$
2. $P(A \cup B) \leq P(A) + P(B)$
3. $P(A \cap B) \leq \min(P(A), P(B))$
4. $P(\neg A) = P(\Omega / A) = 1 - P(A)$

/ is set difference



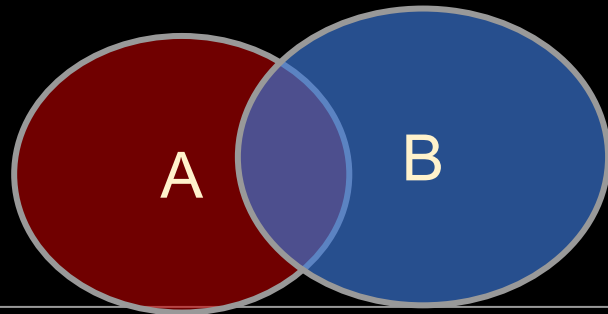
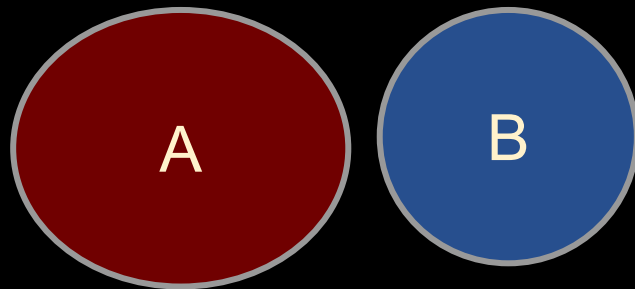
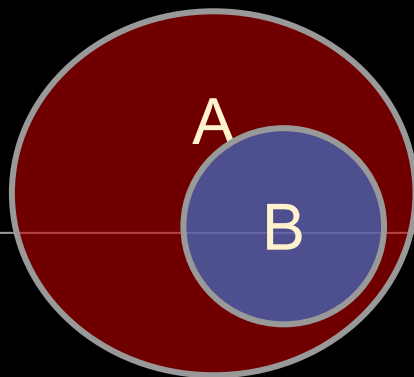
Probability

Some Properties:

1. If $B \subseteq A$ then $P(A) \geq P(B)$
2. $P(A \cup B) \leq P(A) + P(B)$
3. $P(A \cap B) \leq \min(P(A), P(B))$
4. $P(\neg A) = P(\Omega / A) = 1 - P(A)$

$/$ is set difference

$P(A \cap B)$ will be notated as $P(A, B)$



Probability

Independence

Two Events: A and B

Does knowing something about A tell us whether B happens (and vice versa)?

Probability

Independence

Two Events: A and B

Does knowing something about A tell us whether B happens (and vice versa)?

1. A : first flip of a fair coin; B : second flip of the same fair coin
2. A : sentence mentions (or not) the word “happy”
 B : sentence mentions (or not) the word “birthday”

Probability

Independence

Two Events: A and B

Does knowing something about A tell us whether B happens (and vice versa)?

1. A : first flip of a fair coin; B : second flip of the same fair coin
2. A : sentence mentions (or not) the word “happy”
 B : sentence mentions (or not) the word “birthday”

Two events, A and B , are *independent* iff: $P(A, B) = P(A)P(B)$

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

“|” is often referred to as “given”:

*“The probability of A **given** B is ...”*

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Two events, A and B, are *independent* iff: $P(A, B) = P(A)P(B)$

$$P(A, B) = P(A)P(B) \text{ iff } P(B|A) = P(B)$$

Interpretation of Independence:

Observing *A* has no effect on probability of *B*.

(Disjoint events, typically, are not independent!)

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Independence example:

F1=H: first flip of a fair coin is heads

F2=H: second flip of the same coin is heads

$$P(F1=H) = 0.5 \quad P(F2=H) = 0.5$$

$$P(F2=H, F1=H) = 0.25$$

Two events, A and B, are *independent* iff: $P(A, B) = P(A)P(B)$

$$P(A, B) = P(A)P(B) \text{ iff } P(B|A) = P(B)$$

Interpretation of Independence:

Observing A has no effect on probability of B. (and vice-versa)

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Independence example:

F1=H: first flip of a fair coin is heads

F2=H: second flip of the same coin is heads

$$P(F1=H) = 0.5 \quad P(F2=H) = 0.5$$

$$P(F2=H, F1=H) = 0.25$$

$$P(F2=H|F1=H) = 0.5 = P(H2)$$

Two events, A and B, are *independent* iff: $P(A, B) = P(A)P(B)$

$$P(A, B) = P(A)P(B) \text{ iff } P(B|A) = P(B)$$

Interpretation of Independence:

Observing A has no effect on probability of B. (and vice-versa)

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Dependence example:

W1=happy: first word is “happy”

W2=birthday: second word is “birthday”

from observing language data, we find:

$$P(W1=happy) = 0.1, P(W2=birthday) = 0.05$$

$$P(W1=happy, W2=birthday) = 0.025$$

Two events, A and B, are *independent* iff: $P(A, B) = P(A)P(B)$

$$P(A, B) = P(A)P(B) \text{ iff } P(B|A) = P(B)$$

Interpretation of Independence:

Observing A has no effect on probability of B. (and vice-versa)

Probability

Conditional Probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Dependence example:

$W1=$ happy: first word is “happy”

$W2=$ birthday: second word is “birthday”

from observing language data, we find:

$$P(W1=happy) = 0.1, P(W2=birthday) = 0.05$$

$$P(W1=happy, W2=birthday) = 0.025$$

thus $P(A, B) \neq P(A)P(B)$

also $P(B|A) \neq P(B)$:

$$P(W2=birthday|W1=happy) = .025 / .1 = .25 \neq 0.05 = P(W2=birthday)$$

Two events, A and B, are **independent** iff: $P(A, B) = P(A)P(B)$

$$P(A, B) = P(A)P(B) \text{ iff } P(B|A) = P(B)$$

Interpretation of Independence:

Observing A has no effect on probability of B . (and vice-versa)

Why Probability?

A formality to make sense of the world.

1. To quantify uncertainty in language data.
Should we believe something or not? Is it a meaningful difference?
2. To be able to generalize from one situation to another.
Can we rely on some information? What is the chance Y happens?
3. To create structured data.
Where does X belong? What words are similar to X?
(necessary no matter what approaches take place)

Why Probability?

A formality to make sense of the world.

1. To quantify uncertainty in language data.
Should we believe something or not? Is it a meaningful difference?
2. To be able to generalize from one situation to another.
Can we rely on some information? What is the chance Y happens?
3. To create structured data.
*Where does X belong? What words are similar to X?
(necessary no matter what approaches take place)*

Why Probability?



Words: Tokens and Types

word tokens - an individual word instance. (a list)

word types - distinct words. (a set)

V - "vocabulary" $|V|$ - vocabulary size (number of types)

N - number of tokens

Words: Tokens and Types

word tokens - an individual word instance. (a list)

word types - distinct words. (a set)

V - "vocabulary" $|V|$ - vocabulary size (number of types)

N - number of tokens

Corpus - a natural language dataset

(i.e. observational data of word sequence in the wild!)

Corpus	Tokens = N	Types = $ V $
Shakespeare	884 thousand	31 thousand
Brown corpus	1 million	38 thousand
Switchboard telephone conversations	2.4 million	20 thousand
COCA	440 million	2 million
Google n-grams	1 trillion	13 million

Figure 2.11 Rough numbers of types and tokens for some English language corpora. The largest, the Google n-grams corpus, contains 13 million types, but this count only includes types appearing 40 or more times, so the true number would be much larger. (SLP3, 2023)

V - "vocabulary" $|V|$ - vocabulary size (number of types)

N - number of tokens

Corpus - a natural language dataset

(i.e. observational data of word sequence in the wild!)

Corpus	Tokens = N	Types = $ V $
Shakespeare	884 thousand	31 thousand
Brown corpus	1 million	38 thousand
Switchboard telephone conversations	2.4 million	20 thousand
COCA	440 million	2 million
Google n-grams	1 trillion	13 million

Figure 2.11 Rough numbers of types and tokens for some English language corpora. The largest, the Google n-grams corpus, contains 13 million types, but this count only includes types appearing 40 or more times, so the true number would be much larger. (SLP3, 2023)

V - "vocabulary" $|V|$ - vocabulary size (number of types)

N - number of tokens

Corpus - a natural language dataset

(i.e. observational data of word sequence in the wild!)

Herndon or
Heap's Law:

$$|V| = kN^\beta$$

Tokenizers

1.

2.

3.

Tokenizers

1. Word Tokenizers

- 2.

- 3.

Tokenizers

1. Word Tokenizers

```
import re  
  
def tokenize(sentence):  
    tokens = re.split(r'\s', sentence)  
    return tokens
```

2.

3.

Tokenizers

1. Word Tokenizers

```
import re

def tokenize(sentence):
    tokens = re.split(r'\s', sentence)
    return tokens
```

- a. [nlTK's TreebankWordTokenizer](#)
- b. [DLATK's happierfuntokenizing.py \(latest version\)](#)

2.

3.

Tokenizers

1. Word Tokenizers
2. Byte-Pair Encoding
- 3.

Tokenizers

1. Word Tokenizers

2. Byte-Pair Encoding

Motivations:

- more data-driven; no predefined words or rules
- allow for *subwords* (e.g. "unlikeliest" -> "un", "like", "liest") – better for unseen words or capturing semantics of parts of words.

3.

Tokenizers

1. Word Tokenizers

2. Byte-Pair Encoding

Motivations:

- more data-driven; no pre
- allow for *subwords* (e.g.)
- unseen words or capturing

3.

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure BPE( $D, k$ )
3:    $V \leftarrow$  all unique characters in  $D$ 
4:   (about 4,000 in English Wikipedia)
5:   while  $|V| < k$  do                                ▷ Merge tokens
6:      $t_L, t_R \leftarrow$  Most frequent bigram in  $D$ 
7:      $t_{\text{NEW}} \leftarrow t_L + t_R$                     ▷ Make new token
8:      $V \leftarrow V + [t_{\text{NEW}}]$ 
9:     Replace each occurrence of  $t_L, t_R$  in
10:       $D$  with  $t_{\text{NEW}}$ 
11:   end while
12:   return  $V$ 
13: end procedure
```

([Bostrum & Durrett, 2020](#))

Tokenizers

1. Word Tokenizers

2. Byte-Pair Encoding

Motivations:

- more data-driven; no predefined words
- allow for *subwords* (e.g. "unlikely")
- capture rare or unseen words or capturing semantic

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure BPE( $D, k$ )
3:    $V \leftarrow$  all unique characters in  $D$ 
4:   (about 4,000 in English Wikipedia)
5:   while  $|V| < k$  do           ▷ Merge tokens
6:      $t_L, t_R \leftarrow$  Most frequent bigram in  $D$ 
7:      $t_{\text{NEW}} \leftarrow t_L + t_R$   ▷ Make new token
8:      $V \leftarrow V + [t_{\text{NEW}}]$ 
9:     Replace each occurrence of  $t_L, t_R$  in
10:       $D$  with  $t_{\text{NEW}}$ 
11:   end while
12:   return  $V$ 
13: end procedure
```

([Bostrum & Durrett, 2020](#))

corpus

5 l o w _
2 l o w e s t _
6 n e w e r _
3 w i d e r _
2 n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w

(SLP3, p.18)

Tokenizers

1. Word Tokenizers

2. Byte-Pair Encoding

Motivations:

- more data-driven; no predefined words
- allow for *subwords* (e.g. "unlikely")
- capture unknown words or capturing semantic

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure BPE( $D, k$ )
3:    $V \leftarrow$  all unique characters in  $D$ 
4:   (about 4,000 in English Wikipedia)
5:   while  $|V| < k$  do           ▷ Merge tokens
6:      $t_L, t_R \leftarrow$  Most frequent bigram in  $D$ 
7:      $t_{\text{NEW}} \leftarrow t_L + t_R$   ▷ Make new token
8:      $V \leftarrow V + [t_{\text{NEW}}]$ 
9:     Replace each occurrence of  $t_L, t_R$  in
10:       $D$  with  $t_{\text{NEW}}$ 
11:   end while
12:   return  $V$ 
13: end procedure
```

corpus

```
5  l o w _
2  l o w e s t _
6  n e w e r _
3  w i d e r _
2  n e w _
```

vocabulary

, d, e, i, l, n, o, r, s, t, w, er, er, ne

([Bostrum & Durrett, 2020](#))

(SLP3, p.19)

Tokenizers

1. Word Tokenizers
2. Byte-Pair Encoding
3. Wordpiece

chosen based on what increases likelihood of data.

does putting "a" and "b" together increase ability to model the corpus?

can be quantified by: $p('a','b') / (p('a')p('b'))$

More here: ([Shuster and Nakajima, 2012](#); [Kudo and Richardson, 2018](#))

