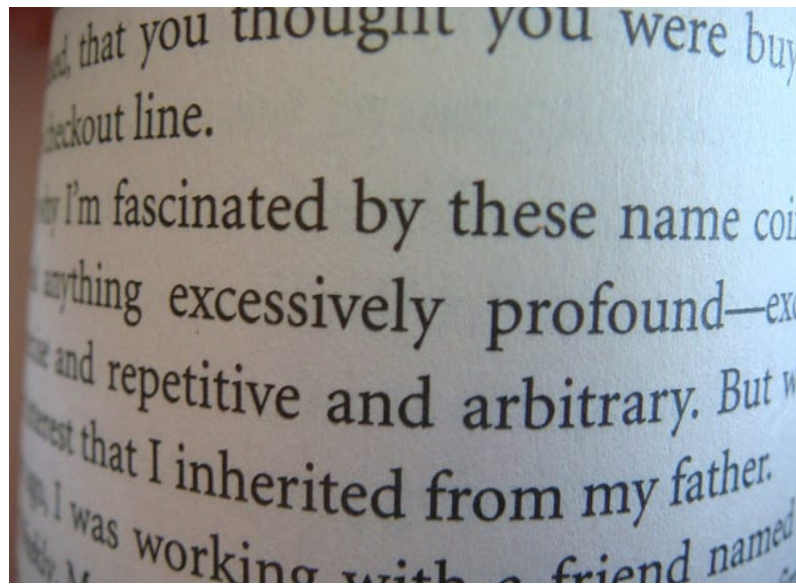


CS 505: Introduction to Natural Language Processing

Wayne Snyder
Boston University

Lecture 3: Basic Notions and Low-Level Text Normalization



Basic Notions: Characters, Words, Tokens, Documents, Corpora

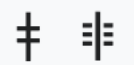
The human race has developed many different writing systems, based on several categories of graphemes (atomic symbols). To *vastly over-simplify*, we have

Alphabets: A set of < 100 symbols, each roughly corresponding to a speech sound:

Example: English: a b c z

Greek: $\alpha \beta \gamma \delta \dots \omega$

Syllabaries: A set of 100s of symbols, each roughly corresponding to spoken syllable:

Linear B (early Greek):  pa-te = πατήρ = pater = "father"

Japanese: さか, *saka*, "hill"

Logographies: A set of 1000s of symbols, each roughly corresponding to a spoken word or concept:

Egyptian Hieroglyphs:



owl



water



mouth



reed
shelter

Chinese: 山 · mountain. 刀 knife 下 down

What about
Emojis?



Basic Notions: Characters, Words, Tokens, Documents, Corpora

For NLP, we want to process text as a **sequence of atomic symbols** and these may be any of the preceeding categories.

Thus: In NLP, **textual data is presented in its most basic form as a sequence of atomic symbols from some finite collection (think Unicode!)**.

In CS 505, our language is English, and this collection will be ASCII characters, and we will generally just call them **characters**. Thus, in its most basic form, a text is simply **one long string**.

| dec | hex | oct | char | dec | hex | oct | char | dec | hex | oct | char | dec | hex | oct | char |
|-----|-----|-----|------|-----|-----|-----|-------|-----|-----|-----|------|-----|-----|-----|------|
| 0 | 0 | 000 | NUL | 32 | 20 | 040 | space | 64 | 40 | 100 | @ | 96 | 60 | 140 | z |
| 1 | 1 | 001 | SOH | 33 | 21 | 041 | ! | 65 | 41 | 101 | A | 97 | 61 | 141 | a |
| 2 | 2 | 002 | STX | 34 | 22 | 042 | " | 66 | 42 | 102 | B | 98 | 62 | 142 | b |
| 3 | 3 | 003 | ETX | 35 | 23 | 043 | # | 67 | 43 | 103 | C | 99 | 63 | 143 | c |
| 4 | 4 | 004 | EOF | 36 | 24 | 044 | \$ | 68 | 44 | 104 | D | 100 | 64 | 144 | d |
| 5 | 5 | 005 | ENQ | 37 | 25 | 045 | % | 69 | 45 | 105 | E | 101 | 65 | 145 | e |
| 6 | 6 | 006 | ACK | 38 | 26 | 046 | & | 70 | 46 | 106 | F | 102 | 66 | 146 | f |
| 7 | 7 | 007 | BEL | 39 | 27 | 047 | ' | 71 | 47 | 107 | G | 103 | 67 | 147 | g |
| 8 | 8 | 010 | BS | 40 | 28 | 050 | (| 72 | 48 | 110 | H | 104 | 68 | 150 | h |
| 9 | 9 | 011 | TAB | 41 | 29 | 051 |) | 73 | 49 | 111 | I | 105 | 69 | 151 | i |
| 10 | a | 012 | LF | 42 | 2a | 052 | * | 74 | 4a | 112 | J | 106 | 6a | 152 | j |
| 11 | b | 013 | VT | 43 | 2b | 053 | + | 75 | 4b | 113 | K | 107 | 6b | 153 | k |
| 12 | c | 014 | FF | 44 | 2c | 054 | , | 76 | 4c | 114 | L | 108 | 6c | 154 | l |
| 13 | d | 015 | CR | 45 | 2d | 055 | - | 77 | 4d | 115 | M | 109 | 6d | 155 | m |
| 14 | e | 016 | SO | 46 | 2e | 056 | . | 78 | 4e | 116 | N | 110 | 6e | 156 | n |
| 15 | f | 017 | SI | 47 | 2f | 057 | / | 79 | 4f | 117 | O | 111 | 6f | 157 | o |
| 16 | 10 | 020 | DLE | 48 | 30 | 060 | 0 | 80 | 50 | 120 | P | 112 | 70 | 160 | p |
| 17 | 11 | 021 | DC1 | 49 | 31 | 061 | 1 | 81 | 51 | 121 | Q | 113 | 71 | 161 | q |
| 18 | 12 | 022 | DC2 | 50 | 32 | 062 | 2 | 82 | 52 | 122 | R | 114 | 72 | 162 | r |
| 19 | 13 | 023 | DC3 | 51 | 33 | 063 | 3 | 83 | 53 | 123 | S | 115 | 73 | 163 | s |
| 20 | 14 | 024 | DC4 | 52 | 34 | 064 | 4 | 84 | 54 | 124 | T | 116 | 74 | 164 | t |
| 21 | 15 | 025 | NAK | 53 | 35 | 065 | 5 | 85 | 55 | 125 | U | 117 | 75 | 165 | u |
| 22 | 16 | 026 | SYN | 54 | 36 | 066 | 6 | 86 | 56 | 126 | V | 118 | 76 | 166 | v |
| 23 | 17 | 027 | ETB | 55 | 37 | 067 | 7 | 87 | 57 | 127 | W | 119 | 77 | 167 | w |
| 24 | 18 | 030 | CAN | 56 | 38 | 070 | 8 | 88 | 58 | 130 | X | 120 | 78 | 170 | x |
| 25 | 19 | 031 | EM | 57 | 39 | 071 | 9 | 89 | 59 | 131 | Y | 121 | 79 | 171 | y |
| 26 | 1a | 032 | SUB | 58 | 3a | 072 | : | 90 | 5a | 132 | Z | 122 | 7a | 172 | z |
| 27 | 1b | 033 | ESC | 59 | 3b | 073 | ; | 91 | 5b | 133 | [| 123 | 7b | 173 | { |
| 28 | 1c | 034 | FS | 60 | 3c | 074 | < | 92 | 5c | 134 | \ | 124 | 7c | 174 | |
| 29 | 1d | 035 | GSK | 61 | 3d | 075 | = | 93 | 5d | 135 |] | 125 | 7d | 175 | } |
| 30 | 1e | 036 | RS | 62 | 3e | 076 | > | 94 | 5e | 136 | ^ | 126 | 7e | 176 | ~ |
| 31 | 1f | 037 | US | 63 | 3f | 077 | ? | 95 | 5f | 137 | _ | 127 | 7f | 177 | DEL |

| dec | hex | oct | char |
|-----|-----|-----|------|
| 64 | 40 | 100 | @ |
| 65 | 41 | 101 | A |
| 66 | 42 | 102 | B |
| 67 | 43 | 103 | C |
| 68 | 44 | 104 | D |
| 69 | 45 | 105 | E |
| 70 | 46 | 106 | F |
| 71 | 47 | 107 | G |
| 72 | 48 | 110 | H |

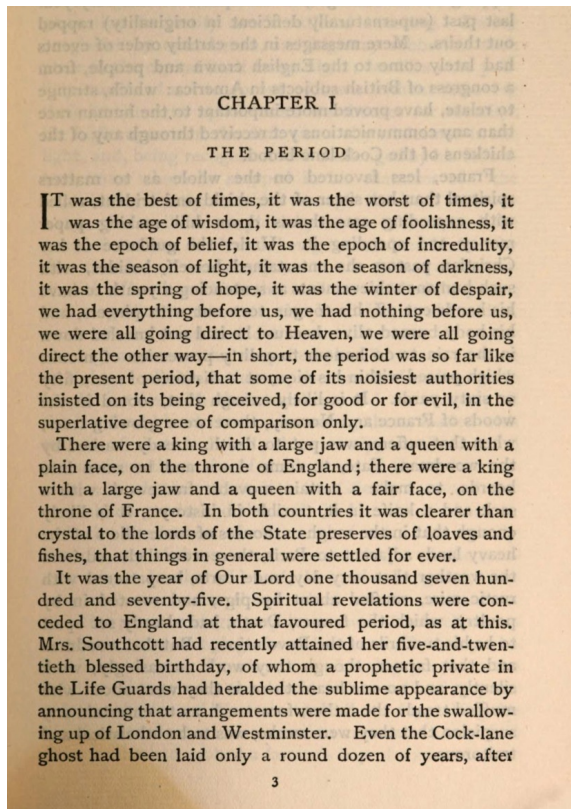
Caveat: We may find it useful if we analyze social media texts to consider Emojis (which are given Unicode numbers!).



Basic Notions: Characters, Words, Tokens, Documents, Corpora

This string form of a text is the minimal representation of the information content of the text (excluding formatting, diagrams, different fonts, illustrations, etc.) and may include some minimal formatting (white space, \n, \t, etc.):

A Tale of Two Cities, Charles Dickens



```
"CHAPTER I\nThe Period\nIt was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way—in short, the period was so far like the present period, that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.\nThere were a king with a large jaw and a queen with a plain face, on the throne of England; there were a king with a large jaw and a queen with a fair face, on the throne of France. In both countries it was clearer than crystal to the lords of the State preserves of loaves and fishes, that things in general were settled for ever. .... "
```


Basic Notions: Characters, Words, Tokens, Documents, Corpora

Although we will have occasion to use the string form when we study character-level machine learning models, almost all NLP uses data which has been grouped into larger units:

Words: Sequence of characters, separated by white space or punctuation;

Tokens: Words possibly preprocessed into some more useful form (the rest of this lecture);

Sentences: Sequences of words/tokens

Paragraphs: Sequences of sentences

Chapters/Sections/Topics: Sequences of paragraphs

Optional

Document: Sequence of paragraphs

Corpus: Set of documents

Brown Corpus:

String form (a list of characters):

```
"The Fulton County Grand Jury said Friday an investigation of Atlanta
```

Words (a list of strings):

```
['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', ...]
```

Sentences (a list of lists of strings):

```
[['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an',  
'primary', 'election', 'produced', 'no', 'evidence', 'the',  
'.'], ['The', 'jury', 'further', 'said', 'in', 'term-end', 'presentme  
ttee', 'which', 'had', 'over-all', 'charge', 'of', 'the', 'elect  
nd', 'thanks', 'of', 'the', 'City', 'of', 'Atlanta', 'for', 'th  
'was', 'conducted', '.'], ...]
```

Paragraphs (a list of lists of lists of strings):

```
[[['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an'  
'primary', 'election', 'produced', 'no', 'evidence', 'the',  
'.'], ['The', 'jury', 'further', 'said', 'in', 'term-end', 'present  
mittee', 'which', 'had', 'over-all', 'charge', 'of', 'the', 'ele  
'and', 'thanks', 'of', 'the', 'City', 'of', 'Atlanta', 'for', 'n',  
'was', 'conducted', '.'], ...]
```

| ID | File | Genre | Description |
|-----|------|-----------------|--|
| A16 | ca16 | news | Chicago Tribune: <i>Society Reportage</i> |
| B02 | cb02 | editorial | Christian Science Monitor: <i>Editorials</i> |
| C17 | cc17 | reviews | Time Magazine: <i>Reviews</i> |
| D12 | cd12 | religion | Underwood: <i>Probing the Ethics of Realtors</i> |
| E36 | ce36 | hobbies | Norling: <i>Renting a Car in Europe</i> |
| F25 | cf25 | lore | Boroff: <i>Jewish Teenage Culture</i> |
| G22 | cg22 | belles_lettres | Reiner: <i>Coping with Runaway Technology</i> |
| H15 | ch15 | government | US Office of Civil and Defence Mobilization: <i>The Family Fallout Shelter</i> |
| J17 | cj19 | learned | Mosteller: <i>Probability with Statistical Applications</i> |
| K04 | ck04 | fiction | W.E.B. Du Bois: <i>Worlds of Color</i> |
| L13 | cl13 | mystery | Hitchens: <i>Footsteps in the Night</i> |
| M01 | cm01 | science_fiction | Heinlein: <i>Stranger in a Strange Land</i> |
| N14 | cn15 | adventure | Field: <i>Rattlesnake Ridge</i> |
| P12 | cp12 | romance | Callaghan: <i>A Passion in Rome</i> |
| R06 | cr06 | humor | Thurber: <i>The Future, If Any, of Comedy</i> |

Corpora for Natural Language Processing

There are many publicly-available corpora for NLP, often categorized (and preprocessed) for specific tasks, in various languages, etc. Dr Google will help you find these....

1.1 Gutenberg Corpus

NLTK includes a small selection of texts from the Project Gutenberg electronic text archive, which contains the following file identifiers in this corpus:

```
>>> nltk.corpus.gutenberg.fileids()
['austen-emma.txt', 'austen-persuasion.txt', 'austen-sense.txt', 'bible-kjv.txt', 'blake-plutarch.txt', 'bly-denham.txt', 'bly-little-dorchester.txt', 'bly-little-dorchester-letters.txt', 'bly-little-dorchester-letters-2.txt', 'bly-little-dorchester-letters-3.txt', 'bly-little-dorchester-letters-4.txt', 'bly-little-dorchester-letters-5.txt', 'bly-little-dorchester-letters-6.txt', 'bly-little-dorchester-letters-7.txt', 'bly-little-dorchester-letters-8.txt', 'bly-little-dorchester-letters-9.txt', 'bly-little-dorchester-letters-10.txt', 'bly-little-dorchester-letters-11.txt', 'bly-little-dorchester-letters-12.txt', 'bly-little-dorchester-letters-13.txt', 'bly-little-dorchester-letters-14.txt', 'bly-little-dorchester-letters-15.txt', 'bly-little-dorchester-letters-16.txt', 'bly-little-dorchester-letters-17.txt', 'bly-little-dorchester-letters-18.txt', 'bly-little-dorchester-letters-19.txt', 'bly-little-dorchester-letters-20.txt', 'bly-little-dorchester-letters-21.txt', 'bly-little-dorchester-letters-22.txt', 'bly-little-dorchester-letters-23.txt', 'bly-little-dorchester-letters-24.txt', 'bly-little-dorchester-letters-25.txt', 'bly-little-dorchester-letters-26.txt', 'bly-little-dorchester-letters-27.txt', 'bly-little-dorchester-letters-28.txt', 'bly-little-dorchester-letters-29.txt', 'bly-little-dorchester-letters-30.txt', 'bly-little-dorchester-letters-31.txt', 'bly-little-dorchester-letters-32.txt', 'bly-little-dorchester-letters-33.txt', 'bly-little-dorchester-letters-34.txt', 'bly-little-dorchester-letters-35.txt', 'bly-little-dorchester-letters-36.txt', 'bly-little-dorchester-letters-37.txt', 'bly-little-dorchester-letters-38.txt', 'bly-little-dorchester-letters-39.txt', 'bly-little-dorchester-letters-40.txt', 'bly-little-dorchester-letters-41.txt', 'bly-little-dorchester-letters-42.txt', 'bly-little-dorchester-letters-43.txt', 'bly-little-dorchester-letters-44.txt', 'bly-little-dorchester-letters-45.txt', 'bly-little-dorchester-letters-46.txt', 'bly-little-dorchester-letters-47.txt', 'bly-little-dorchester-letters-48.txt', 'bly-little-dorchester-letters-49.txt', 'bly-little-dorchester-letters-50.txt', 'bly-little-dorchester-letters-51.txt', 'bly-little-dorchester-letters-52.txt', 'bly-little-dorchester-letters-53.txt', 'bly-little-dorchester-letters-54.txt', 'bly-little-dorchester-letters-55.txt', 'bly-little-dorchester-letters-56.txt', 'bly-little-dorchester-letters-57.txt', 'bly-little-dorchester-letters-58.txt', 'bly-little-dorchester-letters-59.txt', 'bly-little-dorchester-letters-60.txt', 'bly-little-dorchester-letters-61.txt', 'bly-little-dorchester-letters-62.txt', 'bly-little-dorchester-letters-63.txt', 'bly-little-dorchester-letters-64.txt', 'bly-little-dorchester-letters-65.txt', 'bly-little-dorchester-letters-66.txt', 'bly-little-dorchester-letters-67.txt', 'bly-little-dorchester-letters-68.txt', 'bly-little-dorchester-letters-69.txt', 'bly-little-dorchester-letters-70.txt', 'bly-little-dorchester-letters-71.txt', 'bly-little-dorchester-letters-72.txt', 'bly-little-dorchester-letters-73.txt', 'bly-little-dorchester-letters-74.txt', 'bly-little-dorchester-letters-75.txt', 'bly-little-dorchester-letters-76.txt', 'bly-little-dorchester-letters-77.txt', 'bly-little-dorchester-letters-78.txt', 'bly-little-dorchester-letters-79.txt', 'bly-little-dorchester-letters-80.txt', 'bly-little-dorchester-letters-81.txt', 'bly-little-dorchester-letters-82.txt', 'bly-little-dorchester-letters-83.txt', 'bly-little-dorchester-letters-84.txt', 'bly-little-dorchester-letters-85.txt', 'bly-little-dorchester-letters-86.txt', 'bly-little-dorchester-letters-87.txt', 'bly-little-dorchester-letters-88.txt', 'bly-little-dorchester-letters-89.txt', 'bly-little-dorchester-letters-90.txt', 'bly-little-dorchester-letters-91.txt', 'bly-little-dorchester-letters-92.txt', 'bly-little-dorchester-letters-93.txt', 'bly-little-dorchester-letters-94.txt', 'bly-little-dorchester-letters-95.txt', 'bly-little-dorchester-letters-96.txt', 'bly-little-dorchester-letters-97.txt', 'bly-little-dorchester-letters-98.txt', 'bly-little-dorchester-letters-99.txt', 'bly-little-dorchester-letters-100.txt']
```

1.3 Brown Corpus

The Brown Corpus is a large corpus of American English, consisting of 1 million words, collected from a variety of sources, including newspapers, magazines, and books.

1.4 Reuters Corpus

The Reuters Corpus is a large corpus of news articles, consisting of 1 million words, collected from the Reuters news agency.

1.5 Inaugural Address Corpus

The Inaugural Address Corpus is a small corpus of inaugural addresses, consisting of 100 words, collected from the Library of Congress.

1.7 Corpora in Other Languages

NLTK comes with corpora for many languages, though in some cases you will need to download them separately. For example, the Spanish corpus is located at `nltk.corpus.ess_esp`.

```
>>> nltk.corpus.ess_esp.words()
['El', 'grupo', 'estatal', 'Electricit\xe9 de France', ...]
>>> nltk.corpus.floresta.words()
['Um', 'revivalismo', 'refrescante', 'O', '7_e_Meio', ...]
>>> nltk.corpus.indian.words('hindi.pos')
['पूर्ण', 'प्रतिबंध', 'इटाओ', ':', 'इराक', 'संयुक्त', ...]
```

Common Crawl

The Data ▾ Resources ▾

Common Crawl maintains a free repository of web data that can be used for a variety of purposes.

Enron Corpus

Article Talk

From Wikipedia, the free encyclopedia

The **Enron Corpus** is a database of over 600,000 emails in the years leading up to the company's collapse in December 2001. It was released by the Enron Commission (FERC) for \$10,000 and computer-processed.

AI Training Datasets : TEXT contents

the OpenAI GPT-3 model has been fairly well documented as having been trained on about 45 TB (terabytes. 1 TB = 1,000 GB, or gigabytes) of pure text data from multiple AI training datasets which include the entirety of our beloved Wikipedia (well, the English-language portion, at least), and books... lots of books (but *not*, perhaps, the timeless classics that you'd think would be required reading for the training of a genius).

The outline of the primary datasets used to train the model are shown below:

| DATASET | RAW SIZE (tokens, ~words, in training mix billions) | WEIGHT 60% | COMPOSITION (actual % of total content mass) | AMPLIFICATION (or Suppression) |
|----------------------------|---|---------------|--|-----------------------------------|
| <i>CommonCrawl</i> | 410 B | 60% | 82% | -0.27 |
| <i>WebText2</i> | 19 B | 20% | 4% | +5.00 |
| <i>Books1 & Books2</i> | 67 B | 15% | 13% | +0.15 |
| <i>Wikipedia</i> | 3 billion | 5% | 1% | +5.00 |
| TOTAL | 500B tokens | | | |

Kaggle Datasets

Search

Create

Home

Hugging Face

Search models, datasets, users...

Datasets: bookcorpus

Tasks: Text Generation Fill-Mask Sub-tasks: language-modeling masked-language-modeling

Annotations Creators: no-annotation Source Datasets: original ArXiv: arxiv:2105.05241 License:

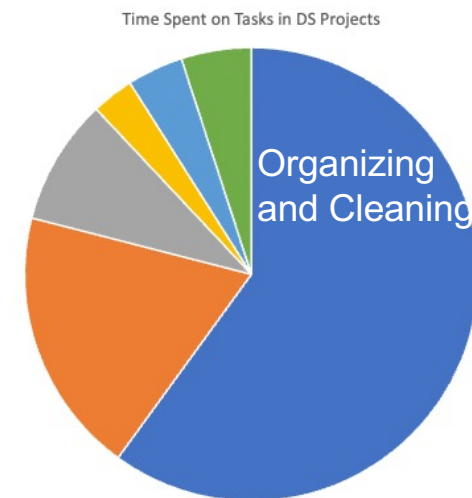
Textual Data Preparation

Why do we need to learn low-level text processing?

Because these corpora are for education, contests, creating general language models (e.g., chatGPT), etc. Most NLP projects involve taking some raw textual data and wrangling it into a corpus of your own.

CrowdFlower, provider of a “data enrichment” [platform for data scientists](#), conducted a survey of about 80 data scientists and found that data scientists spend –

- 60% of the time in organizing and cleaning data.
- 19% of the time is spent in collecting datasets.
- 9% of the time is spent in mining the data to draw patterns.
- 3% of the time is spent in training the datasets.
- 4% of the time is spent in refining the algorithms.
- 5% of the time is spent in other tasks.



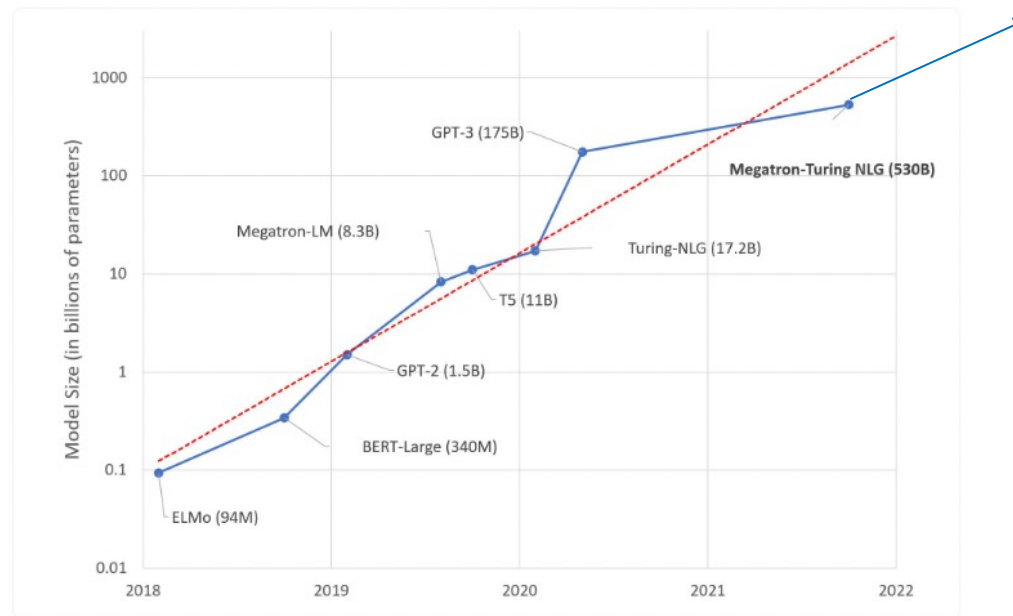
At the same time.....

Textual Data Preparation

ML Algorithms are growing exponentially!

GPT-4: 1.76 trillion parameters!

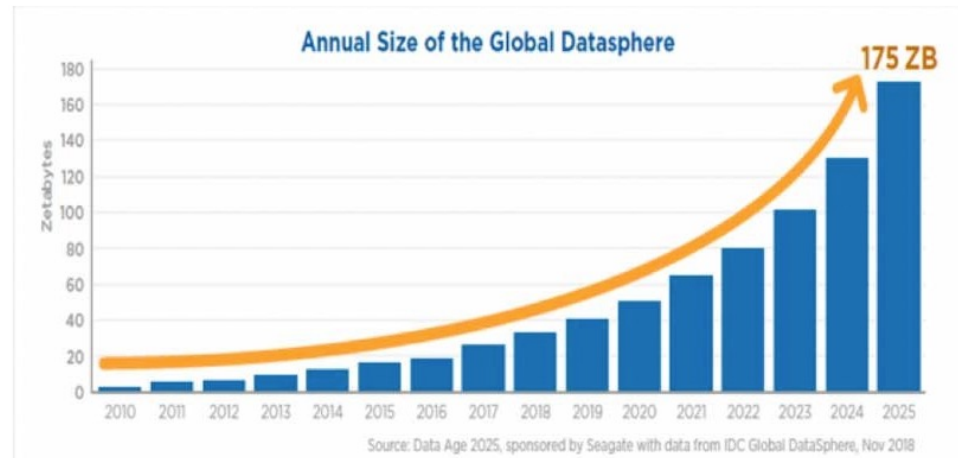
Log Scale



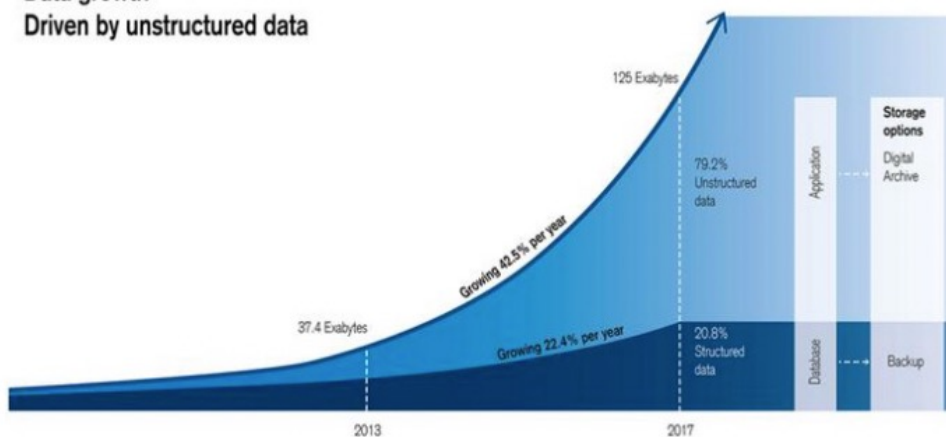
Basic rule: The more parameters, the more data you need... AND...

Textual Data Preparation

The amount of available data is growing exponentially, and it is mostly unstructured. It is simply impossible to process this tsunami of data by (human) hand!



Data growth
Driven by unstructured data



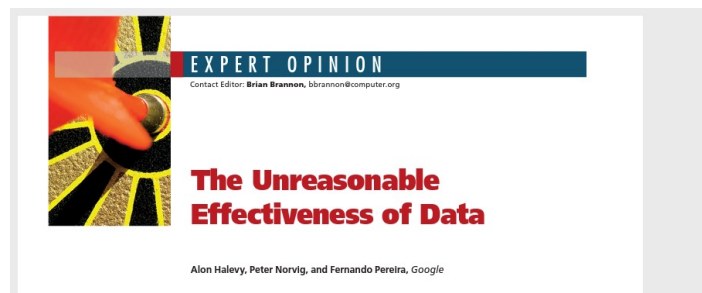
125 Exabytes of enterprise data was stored in 2017; 80% was unstructured data. (Source: Credit Suisse)

Examples of Data Volumes

| | Unit | Value | Example |
|-----------|-----------------|------------------|--|
| 10^3 | Kilobytes (KB) | 1,000 bytes | a paragraph of a text document |
| 10^6 | Megabytes (MB) | 1,000 Kilobytes | a small novel |
| 10^9 | Gigabytes (GB) | 1,000 Megabytes | Beethoven's 5th Symphony |
| 10^{12} | Terabytes (TB) | 1,000 Gigabytes | all the X-rays in a large hospital |
| 10^{15} | Petabytes (PB) | 1,000 Terabytes | half the contents of all US academic research libraries |
| 10^{18} | Exabytes (EB) | 1,000 Petabytes | about one fifth of the words people have ever spoken |
| 10^{21} | Zettabytes (ZB) | 1,000 Exabytes | as much information as there are grains of sand on all the world's beaches |
| 10^{24} | Yottabytes (YB) | 1,000 Zettabytes | as much information as there are atoms in 7,000 human bodies |

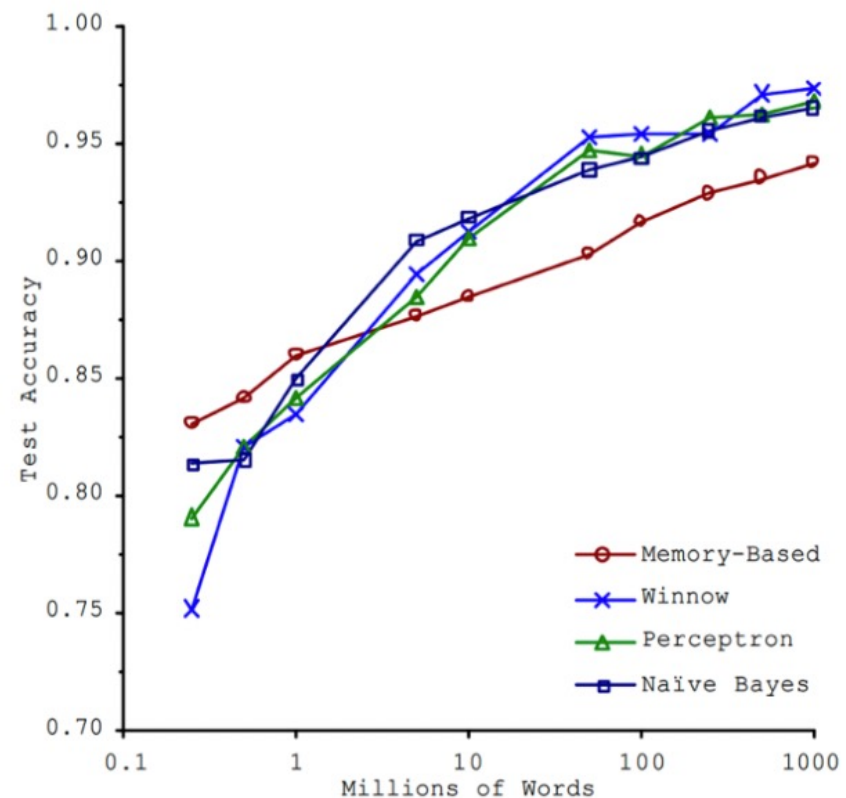
Textual Data Preparation

Plus, **data is the most important resource** – progress in NLP is overwhelmingly dependent on the amount of data, NOT refinements to the algorithms.



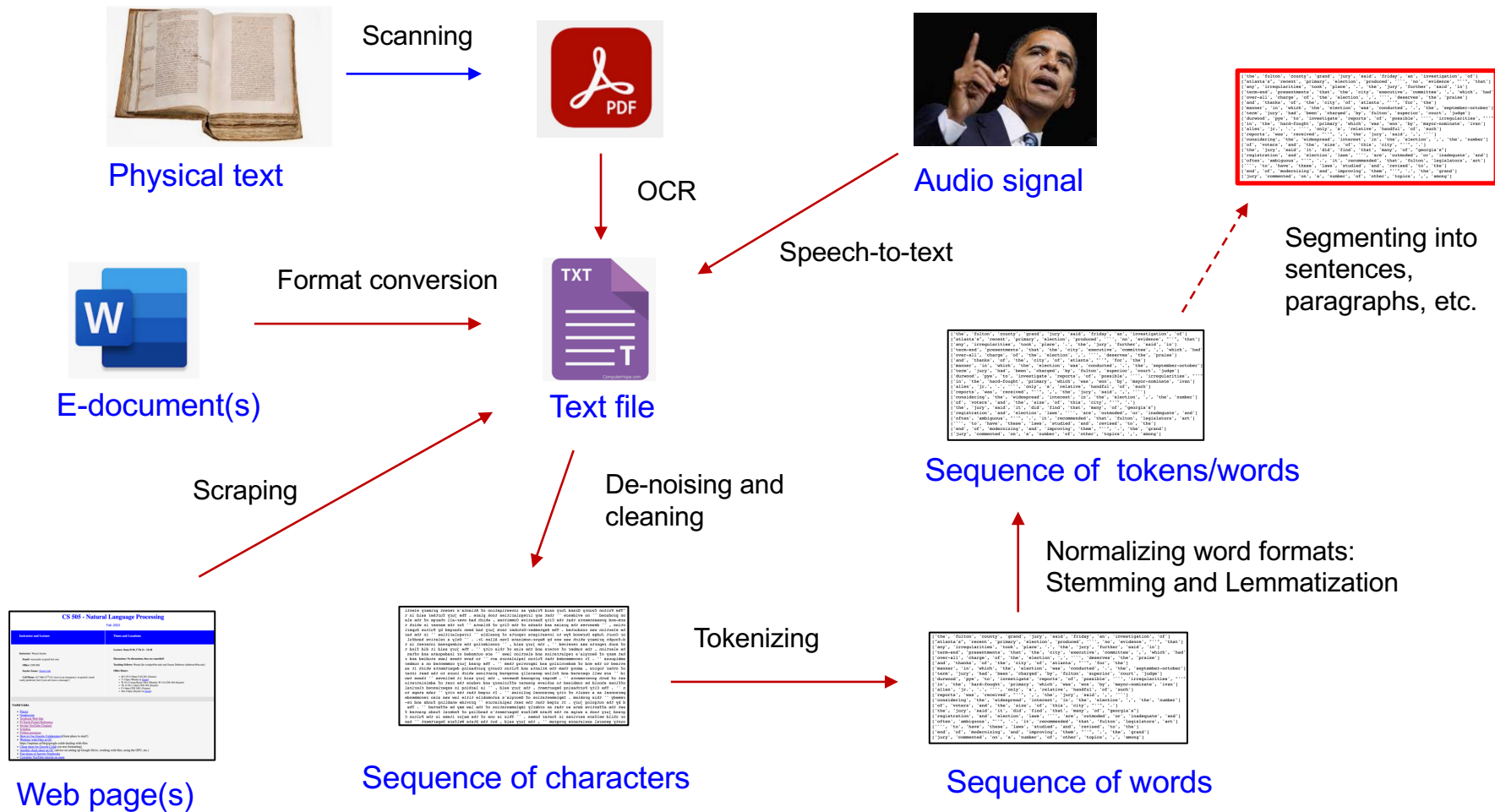
“...consider an experiment done by Microsoft in 2001. Researchers ran a side-by-side test to evaluate the merits of 4 of different approaches to ML translation. They trained each model from scratch with the same input data, running a series of trials with varying data set sizes from 100k to 1 billion words.”

“... simple models and a lot of data trump more elaborate models based on less data.....”



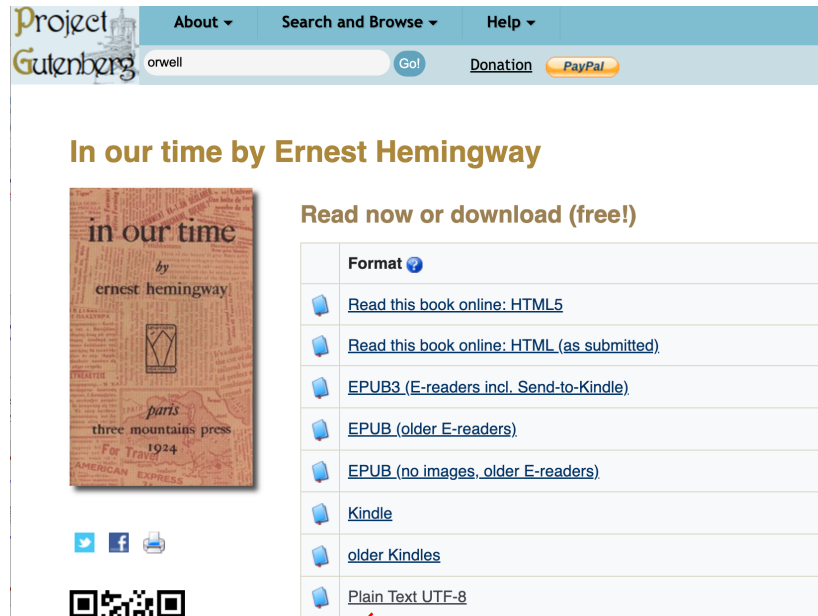
Textual Data Preparation

Most textual data preparation will follow some part of this flowchart:



Tokenizing: Separating a string into words

Example



Project Gutenberg

About Search and Browse Help

orwell Go! Donation PayPal

In our time by Ernest Hemingway

Read now or download (free!)

Format

- [Read this book online: HTML5](#)
- [Read this book online: HTML \(as submitted\)](#)
- [EPUB3 \(E-readers incl. Send-to-Kindle\)](#)
- [EPUB \(older E-readers\)](#)
- [EPUB \(no images, older E-readers\)](#)
- [Kindle](#)
- [older Kindles](#)
- [Plain Text UTF-8](#)

in our time

chapter 1

Everybody was drunk. The whole battery was drunk going along the road in the dark. We were going to the Champagne. The lieutenant kept riding his horse out into the fields and saying to him, "I'm drunk, I tell you, mon vieux. Oh, I am so soused." We went along the road all night in the dark and the adjutant kept riding up alongside my kitchen and saying, "You must put it out. It is dangerous. It will be observed." We were fifty kilometers from the front but the adjutant worried about the fire in my kitchen. It was funny going along that road. That was when I was a kitchen corporal.

Tokenizing: Separating a string into words

First try: separate on white space

in our time

chapter 1

Everybody was drunk. The whole battery was drunk going along the road in the dark. We were going to the Champagne. The lieutenant kept riding his horse out into the fields and saying to him, "I'm drunk, I tell you, mon vieux. Oh, I am so soused." We went along the road all night in the dark and the adjutant kept riding up alongside my kitchen and saying, "You must put it out. It is dangerous. It will be observed." We were fifty kilometers from the front but the adjutant worried about the fire in my kitchen. It was funny going along that road. That was when I was a kitchen corporal.

```
1 import re
2
3 separator_1 = '\s+'          # one or more white space characters
4
5 print(re.split(separator_1,text),'\n')
6
```

```
['in', 'our', 'time', 'chapter', '1', 'Everybody', 'was', 'drunk.', 'The', 'whole', 'battery', 'was', 'drunk', 'going', 'along', 'the', 'road', 'in', 'the', 'dark.', 'We', 'were', 'going', 'to', 'the', 'Champagne.', 'The', 'lieutenant', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him,', '"I'm', 'drunk,', 'I', 'tell', 'you,', 'mon', 'vieux.', 'Oh,', 'I', 'am', 'so', 'soused."', 'We', 'went', 'along', 'the', 'road', 'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'and', 'saying,', '"You', 'must', 'put', 'it', 'out.', 'It', 'is', 'dangerous.', 'It', 'will', 'be', 'observed."', 'We', 'were', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire', 'in', 'my', 'kitchen.', 'It', 'was', 'funny', 'going', 'along', 'that', 'road.', 'That', 'was', 'when', 'I', 'was', 'a', 'kitchen', 'corporal.']
```


Tokenizing: Separating a string into words

First try: separate on white space

in our time

chapter 1

Everybody was drunk. The whole battery was drunk going along the road in the dark. We were going to the Champagne. The lieutenant kept riding his horse out into the fields and saying to him, "I'm drunk, I tell you, mon vieux. Oh, I am so soused." We went along the road all night in the dark and the adjutant kept riding up alongside my kitchen and saying, "You must put it out. It is dangerous. It will be observed." We were fifty kilometers from the front but the adjutant worried about the fire in my kitchen. It was funny going along that road. That was when I was a kitchen corporal.

```
1 import re
2
3 separator_1 = '\s+'          # one or more white space characters
4
5 print(re.split(separator_1, text), '\n')
6
```

```
['in', 'our', 'time', 'chapter', '1', 'Everybody', 'was', 'drunk.', 'The', 'whole', 'battery', 'was', 'drunk', 'going', 'along', 'the', 'road', 'in', 'the', 'dark.', 'We', 'were', 'going', 'to', 'the', 'Champagne.', 'The', 'lieutenant', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him,', '"I'm', 'drunk', 'I', 'tell', 'you,', 'mon', 'vieux.', 'Oh,', 'I', 'am', 'so', 'soused."', 'We', 'went', 'along', 'the', 'road', 'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'and', 'saying,', '"You', 'must', 'put', 'it', 'out.', 'It', 'is', 'dangerous.', 'It', 'will', 'be', 'observed."', 'We', 'were', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire', 'in', 'my', 'kitchen.', 'It', 'was', 'funny', 'going', 'along', 'that', 'road.', 'That', 'was', 'when', 'I', 'was', 'a', 'kitchen', 'corporal.']
```


Tokenizing: Separating a string into words

Second try: separate on white space and punctuation:

```
in our time
```

```
chapter 1
```

```
Everybody was drunk. The whole battery was drunk going along the road
in the dark. We were going to the Champagne. The lieutenant kept
riding his horse out into the fields and saying to him, "I'm drunk, I
tell you, mon vieux. Oh, I am so soused." We went along the road all
night in the dark and the adjutant kept riding up alongside my
kitchen and saying, "You must put it out. It is dangerous. It will be
observed." We were fifty kilometers from the front but the adjutant
worried about the fire in my kitchen. It was funny going along that
road. That was when I was a kitchen corporal.
```

```
1 separator_2 = '[\s,;:!.?]+'      # one or more white space or punctuation characters
2
3 print(re.split(separator_2,text),'\n')
```

```
['in', 'our', 'time', 'chapter', '1', 'Everybody', 'was', 'drunk', 'The', 'whole', 'battery', 'was', 'drunk', 'goin
g', 'along', 'the', 'road', 'in', 'the', 'dark', 'We', 'were', 'going', 'to', 'the', 'Champagne', 'The', 'lieutenan
t', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him', '"I'm', 'drunk',
'I', 'tell', 'you', 'mon', 'vieux', 'Oh', 'I', 'am', 'so', 'soused', '"', 'We', 'went', 'along', 'the', 'road', 'al
l', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'an
d', 'saying', '"You', 'must', 'put', 'it', 'out', 'It', 'is', 'dangerous', 'It', 'will', 'be', 'observed', '"', 'We',
'were', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire',
'in', 'my', 'kitchen', 'It', 'was', 'funny', 'going', 'along', 'that', 'road', 'That', 'was', 'when', 'I', 'was',
'a', 'kitchen', 'corporal', '']
```


Tokenizing: Separating a string into words

Third try: separate on anything other than a word character

in our time

chapter 1

Everybody was drunk. The whole battery was drunk going along the road in the dark. We were going to the Champagne. The lieutenant kept riding his horse out into the fields and saying to him, "I'm drunk, I tell you, mon vieux. Oh, I am so soused." We went along the road all night in the dark and the adjutant kept riding up alongside my kitchen and saying, "You must put it out. It is dangerous. It will be observed." We were fifty kilometers from the front but the adjutant worried about the fire in my kitchen. It was funny going along that road. That was when I was a kitchen corporal.

```
1 separator_3 = '\W+'          # one or more non-word characters
2
3 print(re.split(separator_3,text),'\n')
```

```
['in', 'our', 'time', 'chapter', '1', 'Everybody', 'was', 'drunk', 'The', 'whole', 'battery', 'was', 'drunk', 'goin', 'g', 'along', 'the', 'road', 'in', 'the', 'dark', 'We', 'were', 'going', 'to', 'the', 'Champagne', 'The', 'lieutenan', 't', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him', 'I', 'm', 'drun', 'k', 'I', 'tell', 'you', 'mon', 'vieux', 'Oh', 'I', 'am', 'so', 'soused', 'We', 'went', 'along', 'the', 'road', 'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'and', 'saying', 'You', 'must', 'put', 'it', 'out', 'It', 'is', 'dangerous', 'It', 'will', 'be', 'observed', 'We', 'were', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire', 'in', 'm', 'y', 'kitchen', 'It', 'was', 'funny', 'going', 'along', 'that', 'road', 'That', 'was', 'when', 'I', 'was', 'a', 'kitch', 'en', 'corporal', '']
```


Tokenizing: Separating a string into words

But it gets worse!

- Can't just blindly remove punctuation or non-word characters:
 - m.p.h., Ph.D., AT&T, cap'n
 - prices (\$45.55)
 - dates (01/02/06)
 - URLs (<http://www.stanford.edu>)
 - hashtags (#nlproc)
 - email addresses (someone@cs.colorado.edu)
- Clitic: a word that doesn't stand on its own
 - "are" in we're, French "je" in j'ai, "le" in l'honneur
- When should multiword expressions (MWE) be words?
 - New York, rock 'n' roll

Tokenizing: Separating a string into words

Even worse, many languages do not use consistently use spaces or punctuation to separate words, and the tokenization is quite complicated!

Chinese

莎拉波娃现在居住在美国东南部的佛罗里达。

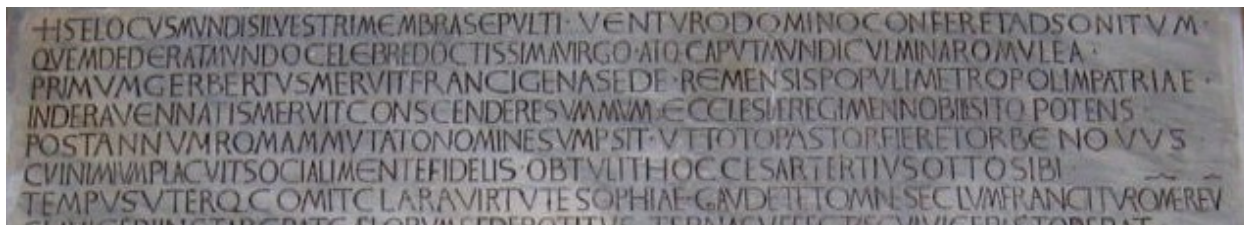
莎拉波娃 现在 居住 在 美国 东南部 的 佛罗里达

Sharapova now lives in US southeastern Florida

Classical Sanskrit

पश्यैतां पाण्डुपुत्राणामाचार्य महतीं चमूम् ।

Latin Inscriptions:



Tokenizing: Separating a string into words

In the case of English and many western languages, there are a variety of useful approaches:

- Rule-based methods, perhaps with list of exceptions
- Sequence-to-Sequence methods:
 - Hidden Markov Models;
 - Neural Networks

Tokenizing: Separating a string into words

Example: Tokenization in the Natural Language Toolkit (NLTK) uses rules based on regular expressions.

```
>>> text = 'That U.S.A. poster-print costs $12.40...'
>>> pattern = r'''(?x)      # set flag to allow verbose regexps
...     ([A-Z]\.)+          # abbreviations, e.g. U.S.A.
...     | \w+(-\w+)*        # words with optional internal hyphens
...     | \$?\d+(\.\d+)?%?  # currency and percentages, e.g. $12.40, 82%
...     | \.\.\.           # ellipsis
...     | [][.,;'"'?()[:-_'] # these are separate tokens; includes ], [
...     '''
>>> nltk.regexp_tokenize(text, pattern)
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...']
```


Tokenizing: Separating a string into words

Example: SpaCy uses a multi-phase approach based on rules and exceptions:



- First, the tokenizer split the text on whitespace similar to the `split()` function.
- Then the tokenizer checks whether the substring matches the tokenizer exception rules.
For example, "don't" does not contain whitespace, but should be split into two tokens, "do" and "n't", while "U.K." should always remain one token.
- Next, it checks for a prefix, suffix, or infix in a substring, these include commas, periods, hyphens, or quotes. If it matches, the substring is split into two tokens.

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp(text)
tst = [ token.text for token in doc]

for k in range(0,len(tst),10):
    print(tst[k:k+10])
```



```
[ 'in', 'our', 'time', '\n\n\n\n\n', 'chapter', '1', '\n\n\n\n', 'Everybody', 'was', 'drunk' ]
[ '.', 'The', 'whole', 'battery', 'was', 'drunk', 'going', 'along', 'the', 'road' ]
[ 'in', 'the', 'dark', '.', 'We', 'were', 'going', 'to', 'the', 'Champagne' ]
[ '.', 'The', 'lieutenant', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the' ]
[ 'fields', 'and', 'saying', 'to', 'him', ' ', 'I', 'm', 'drunk' ]
[ ' ', 'I', 'tell', 'you', ' ', 'mon', 'vieux', '.', 'Oh', ' ', ' ' ]
[ 'I', 'am', 'so', 'soused', ' ', ' ', 'We', 'went', 'along', 'the' ]
[ 'road', 'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept' ]
[ 'riding', 'up', 'alongside', 'my', 'kitchen', 'and', 'saying', ' ', ' ', ' ', 'You' ]
[ 'must', 'put', 'it', 'out', ' ', 'It', 'is', 'dangerous', ' ', 'It' ]
[ 'will', 'be', 'observed', ' ', ' ', ' ', 'We', 'were', 'fifty', 'kilometers', 'from' ]
[ 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire', 'in' ]
[ 'my', 'kitchen', ' ', 'It', 'was', 'funny', 'going', 'along', 'that', 'road' ]
[ '.', 'That', 'was', 'when', 'I', 'was', 'a', 'kitchen', 'corporal', ' ' ]
```

+ Code

+ Text

Word Normalization; Stemming and Lemmatization

Normalization is putting words into a standard format

Simple: Make text case-insensitive by converting all to lower case

More complex:

- Misspellings: “ In tge beginning.... “

- Abbreviations:

PHD PhD Ph.D Ph.D.

etc. &c

US U.S. U.S.A.

- Hyphenation:

lowercase lower-case

- Miscellaneous:

uhhuh uh-huh

Word Normalization; Stemming and Lemmatization

Stemming: Remove suffixes

likes, liked, likely, liking, likable

Stem: **like**

Naïve method: Chop off last part of word (based on list of cases)!

This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things-names and heights and soundings-with the single exception of the red crosses and the written notes.



Thi wa not the map we found in Billi Bone s chest but an accur copi complet in all thing name and height and sound with the singl except of the red cross and the written note

.

Word Normalization; Stemming and Lemmatization

Stemming: Remove suffixes

likes, liked, likely, liking, likable

Stem: like

Better: Multi-phase rule-based systems such as the Porter Stemmer (available in NLTK):

ATIONAL → ATE (e.g., relational → relate)

ING → ϵ if stem contains vowel (e.g., motoring → motor)

SSES → SS (e.g., grasses → grass)

Word Normalization; Stemming and Lemmatization

Lemmatization: Represent all words by their lemma, the shared root word (dictionary headword):

- *am, are, is* → *be*
- *car, cars, car's, cars'* → *car*
- Spanish **quiero** ('I want'), **quieres** ('you want')
→ **querer** 'want'
- *He is reading detective stories*
→ *He be read detective story*

Sentence Segmentation

Segmenting into sentences is based on punctuation.

!, ? Are mostly unambiguous but period “.” is very ambiguous

- Usually “.” is a sentence boundary. Especially followed by capital letter. But:
- Abbreviations like Inc. or Dr. Snyder
- Numbers like .02% or 4.3

Common algorithm: Tokenize first: use rules or ML to classify a period as either (a) part of the word or (b) a sentence-boundary. Capitalization can help too!

Sentence segmentation can then often be done by rules based on this tokenization.