

# Natural Language Processing

## Lecture III. Preprocessing

- because not all corpora are just words you can find in Merriam-Webster

Forrest Sheng Bao, Ph.D.

Dept. of Computer Science  
Iowa State University  
Ames, IA 50011

Sept. 7, 2021

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Not everyone speaks American English or ASCII

Although in Hollywood movies, aliens all speak English the moment they hit Early (and always in NYC), this is the reality:

```
In [20]: s
Out[20]: 'Universitäten'

In [21]: s.encode()
Out[21]: b'Universit\xc3\xa4ten'

In [22]: s.encode().decode('utf-8')
Out[22]: 'Universitäten'

In [23]: s.encode().decode('latin-1')
Out[23]: 'UniversitÃ¤ten'

In [24]: s.encode().decode('GBK')
Out[24]: 'Universit\u0399ten'
```

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for "A", 1001 for "B", etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).



# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
  - ▶ Local encodings in Europe and Asia
  - ▶ UTF8
  - ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Text encoding for plain text

- ▶ How does computers store text information?
- ▶ Digital computers represent data using a sequence of numbers (e.g., 1s and 0s).
- ▶ Natural languages have basic building blocks: letters, symbols, characters, etc.
- ▶ A piece of text is just a character *string*.
- ▶ Use one-to-one mapping between binary sequences and the building blocks, e.g., 1100 for “A”, 1001 for “B”, etc.
- ▶ ASCII
- ▶ Local encodings in Europe and Asia
- ▶ UTF8
- ▶ Many languages are not yet digitized. Actually many languages still do not have scripts (i.e., writing system).

# Unicode – too big to install



For CJK, Unicode has 48 strokes (U+31C0..U+31EF), 224 radicals (U+2F00..U+2FDF) and 12 ideographic description characters (U+2FF0..U+2FFF). But it still has each character!

2FF0

Ideographic Description Characters

2FFF

	2FF0
0	𠄎
1	𠄏
2	𠄐
3	𠄑
4	𠄒
5	𠄓
6	𠄔
7	𠄕
8	𠄖
9	𠄗

## Ideographic description characters

These are visibly displayed graphic characters, not invisible composition controls.

- 2FF0 𠄎 IDEOGRAPHIC DESCRIPTION CHARACTER LEFT TO RIGHT
- 2FF1 𠄏 IDEOGRAPHIC DESCRIPTION CHARACTER ABOVE TO BELOW
- 2FF2 𠄐 IDEOGRAPHIC DESCRIPTION CHARACTER LEFT TO MIDDLE AND RIGHT
- 2FF3 𠄑 IDEOGRAPHIC DESCRIPTION CHARACTER ABOVE TO MIDDLE AND BELOW
- 2FF4 𠄒 IDEOGRAPHIC DESCRIPTION CHARACTER FULL SURROUND
- 2FF5 𠄓 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM ABOVE
- 2FF6 𠄔 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM UPPER LEFT
- 2FF7 𠄕 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM LEFT
- 2FF8 𠄖 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM UPPER RIGHT
- 2FF9 𠄗 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM LOWER LEFT
- 2FFB 𠄙 IDEOGRAPHIC DESCRIPTION CHARACTER OVERLAP

2F00

Kangxi Radicals

2FDF

	2F00	2F01	2F02	2F03	2F04	2F05	2F06	2F07	2F08	2F09	2FA	2FB	2FC	2FD
0	一	丨	十	乚	支	比	瓜	示	聿	衣	辰	革	髑	鼻
1	冫	刀	夂	巾	支	毛	瓦	肉	西	彡	韋	鬼	齊	
2	勹	力	夂	干	文	氏	甘	禾	臣	見	邑	非	魚	齒
3	勹	勹	夕	彡	斗	气	生	穴	自	角	酉	音	鳥	龍
4	乙	匕	大	广	斤	水	用	立	至	言	采	頁	齒	龜
5	冫	冫	女	彡	方	火	田	竹	白	谷	里	風	鹿	禽
6	二	冫	子	丹	无	爪	疋	米	舌	豆	金	飛	麥	
7	勹	勹	彡	弋	日	父	疋	糸	舛	豕	長	食	麻	
8	人	卜	寸	弓	曰	彡	缶	舟	豕	門	首	黃		
9	儿	卩	小	彡	月	升	白	网	艮	貝	阜	香	黍	
A	入	厂	尤	彡	木	片	皮	羊	色	赤	隶	馬	黑	

# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.
- ▶ The Unicode standard is not designed for NLP tasks. For example, it does not have a way to represent a word that is not in the standard. If you are using a word that is not in the standard, you are using a word that is not in the standard.

# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.
- ▶ [https://www.unicode.org/reports/tr35/#Text\\_Representations](https://www.unicode.org/reports/tr35/#Text_Representations)



# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.

▶ ▶ Bad Characters: Imperceptible NLP Attacks, Boucher et al., 2021 <https://arxiv.org/pdf/2106.09899.pdf>

# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.
- ▶
  - ▶ Bad Characters: Imperceptible NLP Attacks, Boucher et al., 2021 <https://arxiv.org/pdf/2106.09898.pdf>
  - ▶ Hey, AI software developers, you are taking Unicode into account, right ... right? [https://www.theregister.com/2021/08/06/unicode\\_ai\\_bug/](https://www.theregister.com/2021/08/06/unicode_ai_bug/)

# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.
- ▶ ▶ Bad Characters: Imperceptible NLP Attacks, Boucher et al., 2021 <https://arxiv.org/pdf/2106.09898.pdf>
  - ▶ Hey, AI software developers, you are taking Unicode into account, right ... right? [https://www.theregister.com/2021/08/06/unicode\\_ai\\_bug/](https://www.theregister.com/2021/08/06/unicode_ai_bug/)

# Unicode can make NLP thorny

- ▶ Lots of duplicated characters or glyphs, e.g., there is a cross (U+07d9) in Nko, while there many other crosses in Unicode.
- ▶ Invisible characters.
- ▶ Potential adversarial attack by using code in different parts of Unicode.
- ▶ ▶ Bad Characters: Imperceptible NLP Attacks, Boucher et al., 2021 <https://arxiv.org/pdf/2106.09898.pdf>
  - ▶ Hey, AI software developers, you are taking Unicode into account, right ... right? [https://www.theregister.com/2021/08/06/unicode\\_ai\\_bug/](https://www.theregister.com/2021/08/06/unicode_ai_bug/)

# Text encoding in Python3

- ▶ In Python 3, a string is by default in UTF-8.
- ▶ `https://docs.python.org/3/howto/unicode.html`
- ▶ There is a separate type for strings, bytes.
- ▶ When opening a file, ensure the encoding:

```
In [1]: with open("BIG5.example", 'r') as f:
...:     print (list(f))
...:
-----
UnicodeDecodeError                                Traceback (most recent call last)
<ipython-input-1-91b059321606> in <module>()
----> 1 with open("BIG5.example", 'r') as f:
      2     print (list(f))
      3

/usr/lib/python3.6/codecs.py in decode(self, input, final)
   319         # decode input (taking the buffer into account)
   320         data = self.buffer + input
--> 321         (result, consumed) = self._buffer_decode(data, self.errors, final)
   322         # keep undecoded input until the next call
   323         self.buffer = data[consumed:]

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xa4 in position 0: invalid start byte

In [2]: with open("BIG5.example", 'r', encoding="BIG5") as f:
...:     print (list(f))
...:
['反攻大陸\n', '消滅共匪\n', '\n']
```

Plain Text - Tab Width: 4 - Ln 5, Col 1

# Text encoding in Python3

- ▶ In Python 3, a string is by default in UTF-8.
- ▶ `https://docs.python.org/3/howto/unicode.html`
- ▶ There is a separate type for strings, bytes.
- ▶ When opening a file, ensure the encoding:

```
In [1]: with open("BIG5.example", 'r') as f:
...:     print (list(f))
...:
-----
UnicodeDecodeError                                Traceback (most recent call last)
<ipython-input-1-91b059321606> in <module>()
1 with open("BIG5.example", 'r') as f:
----> 2     print (list(f))
3

/usr/lib/python3.6/codecs.py in decode(self, input, final)
319         # decode input (taking the buffer into account)
320         data = self.buffer + input
--> 321         (result, consumed) = self._buffer_decode(data, self.errors, final)
322         # keep undecoded input until the next call
323         self.buffer = data[consumed:]

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xa4 in position 0: invalid start byte

In [2]: with open("BIG5.example", 'r', encoding="BIG5") as f:
...:     print (list(f))
...:
['反攻大陸\n', '消滅共匪\n', '\n']
```

Plain Text   Tab Width: 4   Ln 5, Col 1

# Text encoding in Python3

- ▶ In Python 3, a string is by default in UTF-8.
- ▶ `https://docs.python.org/3/howto/unicode.html`
- ▶ There is a separate type for strings, bytes.
- ▶ When opening a file, ensure the encoding:

```
In [1]: with open("BIG5.example", 'r') as f:
...:     print (list(f))
...:
-----
UnicodeDecodeError                                Traceback (most recent call last)
<ipython-input-1-91b059321606> in <module>()
      1 with open("BIG5.example", 'r') as f:
----> 2     print (list(f))
      3

/usr/lib/python3.6/codecs.py in decode(self, input, final)
    319         # decode input (taking the buffer into account)
    320         data = self.buffer + input
--> 321         (result, consumed) = self._buffer_decode(data, self.errors, final)
    322         # keep undecoded input until the next call
    323         self.buffer = data[consumed:]

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xa4 in position 0: invalid start byte

In [2]: with open("BIG5.example", 'r', encoding="BIG5") as f:
...:     print (list(f))
...:
['反攻大陸\n', '消滅共匪\n', '\n']
```

Plain Text   Tab Width: 4   Ln 5, Col 1

# Text encoding in Python3

- ▶ In Python 3, a string is by default in UTF-8.
- ▶ `https://docs.python.org/3/howto/unicode.html`
- ▶ There is a separate type for strings, bytes.
- ▶ When opening a file, ensure the encoding:

```
In [1]: with open("BIG5.example", 'r') as f:
...:     print (list(f))
...:
-----
UnicodeDecodeError                                Traceback (most recent call last)
<ipython-input-1-91b059321606> in <module>()
      1 with open("BIG5.example", 'r') as f:
----> 2     print (list(f))
      3

/usr/lib/python3.6/codecs.py in decode(self, input, final)
    319         # decode input (taking the buffer into account)
    320         data = self.buffer + input
--> 321         (result, consumed) = self._buffer_decode(data, self.errors, final)
    322         # keep undecoded input until the next call
    323         self.buffer = data[consumed:]

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xa4 in position 0: invalid start byte

In [2]: with open("BIG5.example", 'r', encoding="BIG5") as f:
...:     print (list(f))
...:
['反攻大陸\n', '消滅共匪\n', '\n']
```

Plain Text — Tab Width: 4 — Ln 5, Col 1



# Beyond plain text

Some text data are more than *plain text*, including information about formatting and field/structure.

- ▶ PDF: font and coordinate for each character – no strings
- ▶ (Encapsulated) Postscript
- ▶ Markup languages: HTML, XML (including Microsoft Office XML). Easier but not that easy: e.g., the same type of information is in different tags or tags of different attributes, e.g.,

```
<a href="a.html" id="first_link">Go back</a>  
<a href="b.html" id="second_link">Next</a>
```

- ▶ (The evil) Microsoft formats (prior to Office XML)
- ▶ Dictionary-like data: YAML, JSON

# Outline

Representation of text data in computers

**Pipeline for NLP Preprocessing**

Corpus and Crawling

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.

# Pipeline for NLP preprocessing

You need to have the right data! (Although a lot of research is “garbage in, garbage out”)

1. Download the corpus or Crawl the web or other ways to get the text (e.g., OCR)
2. If the data is (semi-)structured, parse it, e.g., parsing the HTML.
3. Clean up: handling errors (e.g., introduced by OCR or text extraction from Powerpoint slides), and others...
4. Tokenization
5. normalization: turn all tokens to their canonical form (including stemming/lemmatization)
6. Optionally, you wanna index your text so you can locate them easily later. String search will take very long.

Some text data are more difficult to handle than others, e.g., Twitter vs. NYTimes.



# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

**Corpus and Crawling**

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Corpus

- ▶ The body of text to be processed. This class focuses on English.
- ▶ E.g., all articles on Wikipedia form a corpus.
- ▶ E.g., all news from a newspaper form another corpus.
- ▶ There are many corpora, e.g.,

```
nltk.corpus.AlignedCorpusReader
nltk.corpus.AlpinoCorpusReader
nltk.corpus.BNCCorpusReader
nltk.corpus.BracketParseCorpusReader
nltk.corpus.CHILDESCorpusReader
nltk.corpus.CMUDictCorpusReader
nltk.corpus.CategorizedBracketParseCorpusReader
nltk.corpus.CategorizedCorpusReader
nltk.corpus.CategorizedPlaintextCorpusReader
nltk.corpus.CategorizedSentencesCorpusReader
nltk.corpus.CategorizedTaggedCorpusReader
nltk.corpus.ChasenCorpusReader
nltk.corpus.ChunkedCorpusReader
nltk.corpus.ComparativeSentencesCorpusReader
nltk.corpus.ConllChunkCorpusReader
nltk.corpus.ConllCorpusReader
nltk.corpus.CorpusReader
nltk.corpus.CrubadanCorpusReader
nltk.corpus.DependencyCorpusReader
nltk.corpus.EuroparlCorpusReader
nltk.corpus.FramenetCorpusReader
nltk.corpus.IEERCorpusReader
nltk.corpus.IPIPANCorpusReader
nltk.corpus.IndianCorpusReader
nltk.corpus.KNBCorpusReader
```

```
nltk.corpus.PropbankCorpusReader
nltk.corpus.ProsConsCorpusReader
nltk.corpus.RTECorpusReader
nltk.corpus.RegexpTokenizer
nltk.corpus.ReviewsCorpusReader
nltk.corpus.SemcorCorpusReader
nltk.corpus.SensevalCorpusReader
nltk.corpus.SentiSynset
nltk.corpus.SentiWordNetCorpusReader
nltk.corpus.SinicaTreebankCorpusReader
nltk.corpus.StringCategoryCorpusReader
nltk.corpus.SwadeshCorpusReader
nltk.corpus.SwitchboardCorpusReader
nltk.corpus.SyntaxCorpusReader
nltk.corpus.TEICorpusView
nltk.corpus.TaggedCorpusReader
nltk.corpus.TimitCorpusReader
nltk.corpus.TimitTaggedCorpusReader
nltk.corpus.ToolboxCorpusReader
nltk.corpus.TwitterCorpusReader
nltk.corpus.UdhrCorpusReader
nltk.corpus.UnicharsCorpusReader
nltk.corpus.VerbnetsCorpusReader
nltk.corpus.WordListCorpusReader
nltk.corpus.WordNetCorpusReader
```

## Another dataset catalog: Tensorflow datasets

`https:  
//www.tensorflow.org/datasets/catalog/overview`

# A corpus ?= a set of sentences

- ▶ Building a corpus is not that easy when your input data is not plain-text sentences only.
- ▶ Should you include non-sentence text (table content, headings, external links, etc.) from Wikipedia into the corpus when training word embeddings?
- ▶ What about style information? Should “this is not right” as a line in a movie be treated as “this is not right”?
- ▶ What about equations in a paper? (Suppose it is in  $\text{\LaTeX}$ /MathML)

# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)

# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)

# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)

# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)



# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)

# The lack of labeled corpora is limiting NLP research

- ▶ Having the text alone is not enough for many NLP (research) tasks.
- ▶ NLP tasks are far more complicated and diverse than those in other AI areas, such as CV (where most of the tasks are related to object detection/recognition for which anyone can be a human annotator).
- ▶ For example, for years, there is no dataset for using supervised approach to text summarization – no corpus has sentence-level binary labels.
- ▶ As another example, QA research was a chaos before Stanford SQUAD dataset.
- ▶ Today, lots of research has no common ground truth. People publish papers on their own small-scale datasets, e.g., REALSumm and SummEval have only 100 articles annotated.
- ▶ It is also wrong to justify: “because everyone else used this corpus, so do I.” (“IBM does this xxx. We don’t do this xxx.” –Steve Jobs)

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.
- ▶ Example 1:  
<http://jmcauley.ucsd.edu/data/amazon/>
- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>
- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>
- ▶ Example 4:  
<https://catalog.ldc.upenn.edu/ldc2008t19>
- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.
- ▶ Example 1:  
`http://jmcauley.ucsd.edu/data/amazon/`
- ▶ Example 2: `https://webscope.sandbox.yahoo.com/catalog.php?datatype=1`
- ▶ Example 3: `https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets`
- ▶ Example 4:  
`https://catalog.ldc.upenn.edu/ldc2008t19`
- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.



# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.

# Corpora are very diverse

- ▶ We are in a very open domain.
- ▶ But datasets differ in the way that they are stored (e.g., JSON), organized and shared.
- ▶ Corpora are floating in the homepages of researchers in various ways.

- ▶ Example 1:

<http://jmcauley.ucsd.edu/data/amazon/>

- ▶ Example 2: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

- ▶ Example 3: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#datasets>

- ▶ Example 4:

<https://catalog.ldc.upenn.edu/ldc2008t19>

- ▶ Lots of subjectivity in annotated data.

## Crawling data from the Web (Python 2.x)

- ▶ Often times, the data is not available. So you crawl. (So do Google, Bing, and many others.)
- ▶ In Python, you may use `urllib2`:

```
import urllib2
response = urllib2.urlopen('http://gozips.
    uakron.edu/~fbao5/')
html=response.read()
```

- ▶ What if the webpage is dynamic? With the help of another library `urllib`. Example below for POST method.

```
import urllib
url="http://10.24.47.178/search.py"
values = {"single":"landing"}
data=urllib.urlencode(values)
response=urllib2.urlopen(urllib2.Request (
    url, data))
html=response.read()
```

# Make your crawler automatically cool down

An example for using `sysvinit` back in 2015, placed under `/etc/init`. Recent Linux distros use `systemd`.

```
description      "Airbnb review fetcher"

setuid forrest
setgid forrest

start on (net-device-up
          and local-filesystems
          and runlevel [2345])
stop on runlevel [016]

respawn
respawn limit 10 5
limit nofile 500000 500000

env LANG=en_US.UTF-8
env LC_CTYPE=en_US.UTF-8

script
exec /home/forrest/Apps/airbnb/bin/python\
/home/forrest/Apps/airbnb-reviews/tools/app 2>&1 >> \
/home/forrest/airbnb_review_fetcher.log
end script
```

Thanks for Wu Jiang, head of robotics research at Boxed Labs.

# Parse HTML

See iPython notebook notes.

# Index your document (if applicable)

- ▶ To allow easy localization of string in documents, you wanna build an inversed index.
- ▶ The very basic searching or information retrieval.
- ▶ See a different iPython notebook notes.

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

**Clean up the text**

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Why data clean up is needed

Run the code below and compare the sample text from three summarization datasets.

```
import tensorflow_datasets as tfds

for piece in tfds.load("cnn_dailymail", split="test"):
    print (piece['article'])
    break

for piece in tfds.load("billsum", split="test"):
    print (piece['text'])
    break

for piece in tfds.load("big_patent", split="test"):
    print (piece['description'])
    break
```



# Clean up the text is not that easy

- ▶ ```we are happy''.split()`
- ▶ The ambiguity of non-alphanumeric symbols, e.g., the dash in “low-hanging fruit” vs. “express”, or the spaces between “New York City” (“New York” + “City” or “New” + “York” + “City”?)
- ▶ You cannot change all words to lowercase, e.g., “It’s a CAT 4 hurricane.” Kitty or category?
- ▶ Numbers make things complicated, e.g., CO2.
- ▶ Some words shouldn’t be separated, e.g., in vitro, de facto, et al.
- ▶ A common practice is dropping all non-alphabets – but it won’t work for technical or scientific documents.

# Help the downstream tasks

- ▶ Suppose you are processing a post from StackOverflow:  
“Apply `get` onto the method `cat`.”
- ▶ Downstream tasks like POS tagging may fail on “`get`” and “`cat`”.
- ▶ Some string substitution is needed.

# Help the downstream tasks

- ▶ Suppose you are processing a post from StackOverflow:  
“Apply `get` onto the method `cat`.”
- ▶ Downstream tasks like POS tagging may fail on “get” and “cat”.
- ▶ Some string substitution is needed.

# Help the downstream tasks

- ▶ Suppose you are processing a post from StackOverflow:  
“Apply `get` onto the method `cat`.”
- ▶ Downstream tasks like POS tagging may fail on “get” and “cat”.
- ▶ Some string substitution is needed.

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

**Word and sentence tokenization**

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.

- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?

Should we include the "P" in "The left **P**up" as a token?  
We see 2 "p"s in the text above. Should we regard  
all "p"s as "P"? Are "P" and "p" semantically different here?  
Should we count "The left"?

- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/  
19970614072242if_/http://www.cis.upenn.edu:  
80/~treebank/tokenization.html
```

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like "n't" in "he isn't happy" be split?
  - ▶ How can I remove the words "the" and "and" in a sentence?
  - ▶ What are "stop" words?
  - ▶ Are "of" and "or" semantically different here?
  - ▶ What are "pre" and "post"?
- ▶ Some tokenization is fairly simple.  
[https://web.archive.org/web/19970614072242if\\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html](https://web.archive.org/web/19970614072242if_/http://www.cis.upenn.edu:80/~treebank/tokenization.html)

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like “n’t” in “he isn’t happy” be split?
  - ▶ “We see 2 people in the park after 3pm” Should we reject all digits? Are “2” and “3” semantically connected here?
  - ▶ Shall we break “New York”?
- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/19970614072242if\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html
```



# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like “n’t” in “he isn’t happy” be split?
  - ▶ “We see 2 people in the park after 3pm” Should we reject all digits? Are “2” and “3” semantically connected here?
  - ▶ Shall we break “New York”?

- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/19970614072242if\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html
```

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like “n’t” in “he isn’t happy” be split?
  - ▶ “We see 2 people in the park after 3pm” Should we reject all digits? Are “2” and “3” semantically connected here?
  - ▶ Shall we break “New York”?
- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/19970614072242if\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html
```

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like “n’t” in “he isn’t happy” be split?
  - ▶ “We see 2 people in the park after 3pm” Should we reject all digits? Are “2” and “3” semantically connected here?
  - ▶ Shall we break “New York”?
- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/19970614072242if\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html
```

# Word Tokenization

- ▶ Carrying an concrete meaning, a token is the basic element to construct a sentence – at least linguistics say so.
- ▶ Things are not as easy as

```
"This is a class about NLP".split()
```

- ▶ What are acceptable tokens?
  - ▶ Should contractions like “n’t” in “he isn’t happy” be split?
  - ▶ “We see 2 people in the park after 3pm” Should we reject all digits? Are “2” and “3” semantically connected here?
  - ▶ Shall we break “New York”?
- ▶ Some tokenization is fairly simple.

```
https://web.archive.org/web/19970614072242if\_/http://www.cis.upenn.edu:80/~treebank/tokenization.html
```

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.



# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Regular expression

- ▶ A regular expression defines how to generate a string by drawing characters from alphabets. You can draw characters from different sets at different “steps”, like the example  $(abc) * (123)?$  below.
- ▶  $(a)^*$  means  $\{ \epsilon, a, aa, aaa, aaaa, \dots \}$ , where  $\epsilon$  is the empty string.
- ▶  $(ab)^*$  means  $\{ \epsilon, ab, abab, ababab, abababab, \dots \}$
- ▶  $(a|b)^*$  means  $\{ \epsilon, a, b, ab, ba, aaa, aab, abb, bbb, bba, baa, bab, \dots \}$  where  $|$  means “or”.
- ▶  $(a)^+$  means  $\{ a, aa, aaa, \dots \}$  because  $+$  means repeating at least once.
- ▶  $(ab)?$  means  $\{ ab, \epsilon \}$  because  $?$  means once or none.
- ▶  $(abc) * (123)?$  means  $\{ \epsilon, abc, abc123, abcabc, abcabc123, \dots \}$  any string that begins with the string “abc” and ends with or without the string “123”.
- ▶  $0|(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^+$  means any natural number in common writing format. Note that it does not allow an integer to start with 0.

# Use regular expressions

- ▶ See iPython notebook notes: the syntax of regex for most UNIX-like OSes and C-style languages may be different from the syntax we just saw but the principles are the same.
- ▶ Do not tokenize stupidly.

jst-ph 4pin connector


All Categories

[Refine your search for jst-ph 4pin connector](#) ☐ Include description

Sort:  View:

43 results for **jst-ph 4pin connector** [Save this search](#)

Did you mean: **jst-sh 4pin connector** (33 items)?




**15 SETS JST PH 2.0MM 4 Pin Female Single Connector with Wires 200MM USA Shipping**

**\$7.99**

Buy It Now

Free Shipping

 Top Rated Plus

- ▶ You may need to craft your own tokenizer.

<https://stackoverflow.com/questions/15929233/writing-a-tokenizer-in-python>

# Regex-based feature engineering

```
https:  
//github.com/forrestbao/autoscholar/commit/  
b732ae5fa14d65558cd0eff7a7a9e103ea7cd937
```

# Get our hands dirty!

- ▶ NLTK: Natural Language Toolkit, classical NLP toolkit in Python
- ▶ SpaCy: Modern, fancy, claimed to be industry-level, friendly APIs, leveraging GPUs (not for all tasks)
- ▶ Stanza: from Stanford NLP group, derived from Core (also model-based), native in Python, models trained in PyTorch, leveraging GPU

# Tokenization and Sentence segmentation in NLTK

Using default ones:

```
In [1]: import nltk
```

```
In [2]: nltk.tokenize.word_tokenize \  
...: ("I am happy. mr. Wang is happy")
```

```
Out[2]: ['I', 'am', 'happy', '.', 'mr.', 'Wang', 'is', '  
        'happy']
```

```
In [3]: nltk.tokenize.sent_tokenize \  
...: ("I am happy. mr. Wang is happy")
```

```
Out[3]: ['I am happy.', 'mr. Wang is also happy']
```

```
In [4]: nltk.tokenize.sent_tokenize \  
...: ("I am happy, mr. wang is happy")
```

```
Out[4]: ['I am happy, mr. wang is also happy']
```

As you can see, sentence segmentation heavily depends on punctuations. The default sentencizer in NLTK is Punkt.

Many other varieties at

<https://www.nltk.org/api/nltk.tokenize.html>



# Tokenization and sentence segmentation in Spacy

## Pipeline-based NLP.

<https://spacy.io/usage/spacy-101#pipelines>

```
In [32]: nlp=spacy.load("en_core_web_sm", \
...: exclude=["tok2vec",'tagger','parser','ner', 'attribute_ruler', 'lemmatizer'])
...: nlp.add_pipe("sentencizer")
Out[32]: <spacy.pipeline.sentencizer.Sentencizer at 0x7f487e80fdc0>

In [33]: [[sent.text for sent in doc.sents ]\
...:       for doc in nlp.pipe(\
...: ['today is monday. it is the first day.', 'soup is yammy. pizza sucks.'])]
Out[33]: [['today is monday.', 'it is the first day.',
          'soup is yammy.', 'pizza sucks.']]

In [34]: [[word.text for word in doc ]
...:       for doc in nlp.pipe(\
...: ['today is monday. it is the first day.', 'soup is yammy. pizza sucks.'])]
Out[34]: [['today', 'is', 'monday', '.', 'it', 'is', 'the', 'first', 'day', '.'],
          ['soup', 'is', 'yammy', '.', 'pizza', 'sucks', '.']]
```

# In Stanza (model-based, slower)

Like SpaCy, Stanza is also pipelined. But it takes one string rather than a collection of strings as the input.

```
In [1]: import stanza
```

```
In [2]: nlp = stanza.Pipeline(lang='en', processors='tokenize')
2021-08-26 14:15:27 INFO: Loading these models for language: en (English):
=====
| Processor | Package |
-----
| tokenize | ewt     |
=====
```

```
In [7]: sentences = 'i am very very happy. the weather is very good.'
```

```
In [8]: [[token.text for token in sentence.tokens] for sentence in nlp(sentences).sentences]
Out[8]:
[['i', 'am', 'very', 'very', 'happy', '.'],
 ['the', 'weather', 'is', 'very', 'good', '.']]
```

```
In [9]: [sentence.text for sentence in nlp(sentences).sentences]
Out[9]: ['i am very very happy.', 'the weather is very good.']
```

# Other tokenizers

- ▶ WordPiece, used in BERT, supports subword
- ▶ SentPiece, used in XLnet,

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

Word and sentence tokenization

**Stemming and Lemmatization**

Normalization and spell correction

Building an inverse index

# Stemming and lemmatization

- ▶ The English language has inflexions, *inflectional morphology*. We modify words for different grammatical purposes. “I give him a book” vs. “He gave me a book.”
- ▶ A language of very little inflexion is called an *analytic language* or an *isolating language*.
- ▶ Most mainland southeast Asian and Oceanic languages (Chinese, Vietnamese, Thai, etc.) are isolating language.
- ▶ Nevertheless, English is almost the most isolating language among *synthetic languages*, such as High German or Scottish Gaelic.

# Stemming and lemmatization

- ▶ The English language has inflexions, *inflectional morphology*. We modify words for different grammatical purposes. “I give him a book” vs. “He gave me a book.”
- ▶ A language of very little inflexion is called an *analytic language* or an *isolating language*.
- ▶ Most mainland southeast Asian and Oceanic languages (Chinese, Vietnamese, Thai, etc.) are isolating language.
- ▶ Nevertheless, English is almost the most isolating language among *synthetic languages*, such as High German or Scottish Gaelic.

# Stemming and lemmatization

- ▶ The English language has inflexions, *inflectional morphology*. We modify words for different grammatical purposes. “I give him a book” vs. “He gave me a book.”
- ▶ A language of very little inflexion is called an *analytic language* or an *isolating language*.
- ▶ Most mainland southeast Asian and Oceanic languages (Chinese, Vietnamese, Thai, etc.) are isolating language.
- ▶ Nevertheless, English is almost the most isolating language among *synthetic languages*, such as High German or Scottish Gaelic.

# Stemming and lemmatization

- ▶ The English language has inflexions, *inflectional morphology*. We modify words for different grammatical purposes. “I give him a book” vs. “He gave me a book.”
- ▶ A language of very little inflexion is called an *analytic language* or an *isolating language*.
- ▶ Most mainland southeast Asian and Oceanic languages (Chinese, Vietnamese, Thai, etc.) are isolating language.
- ▶ Nevertheless, English is almost the most isolating language among *synthetic languages*, such as High German or Scottish Gaelic.

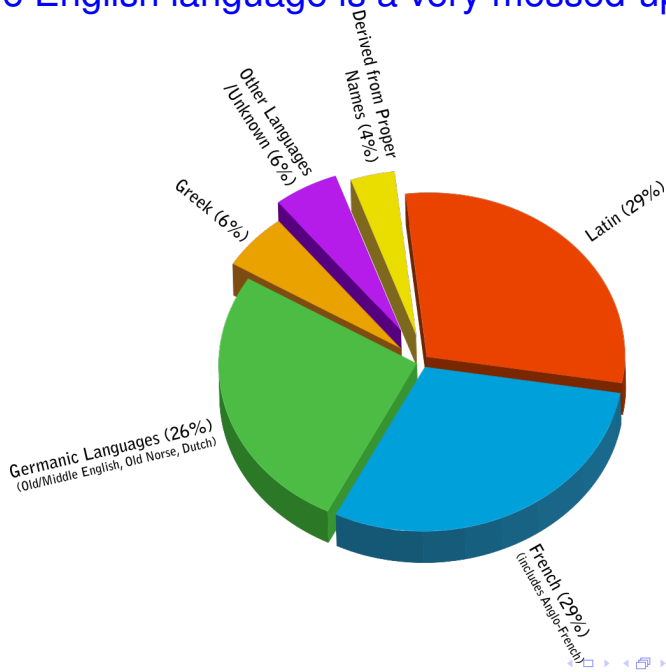


# The English language is a very messed-up one. I

- 1 Ðan ſhe ðiſſe Old Angliſh Tunȝan eode ſpecan
- 2 Than ſhe ȝan to-spaken biſ Middle Englyſſhe Tongue
- 3 Then ſhe wente to ſpake thiſ Early Modern Englyſh Tongue
- 4 Then ſhe went to ſpeak thiſ Late Modern English Tongue

“English irregular verbs” are actually “regular verbs” in Germanic languages.

# The English language is a very messed-up one. II



# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “saw” and “see” are not the same word
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ The WordNet is the dictionary
  - ▶ What words are affixes? E.g., “transform” (Noun) → “transformer” (Noun) → “transforming” (Verb) → “transformed” (Verb) → “transformations” (Noun)
- ▶ Don’t take things for granted: NLTK’s “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present participle: **seeking**

- ▶ How to build a lemmatizer?
- ▶ Don't take things for granted: NLTK's “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
  - ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present participle: **seeking**

- ▶ How to build a lemmatizer?
- ▶ Don't take things for granted: NLTK's “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
- 
- ▶ Don't take things for granted: NLTK's “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ Take advantage of the dictionary
- ▶ Don't take things for granted: NLTK's “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ Take advantage of the dictionary
  - ▶ What about words beyond? E.g., “transformative” (NSF style) vs “transforming” (a real English word in dictionaries.)
- ▶ Don’t take things for granted: NLTK’s “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.



# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ Take advantage of the dictionary
    - ▶ What about words beyond? E.g., “transformative” (NSF style) vs “transforming” (a real English word in dictionaries.)
  - ▶ Don’t take things for granted: NLTK’s “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ Take advantage of the dictionary
  - ▶ What about words beyond? E.g., “transformative” (NSF style) vs “transforming” (a real English word in dictionaries.)
- ▶ Don't take things for granted: NLTK's “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Stemming and lemmatization

- ▶ For many tasks, we want the words “dog” and “dogs” to be considered the same.
- ▶ Stemming: get the stem, e.g., removing the affix.
- ▶ Stemming may not be enough:
  - ▶ “I saw you” should become “I see you”
- ▶ Lemmatization: get the lemma, the item word in a dictionary.

seek

/sēk/ 🔊

*verb*

verb: **seek**; 3rd person present: **seeks**; past tense: **sought**; past participle: **sought**; gerund or present

participle: **seeking**

- ▶ How to build a lemmatizer?
  - ▶ Take advantage of the dictionary
  - ▶ What about words beyond? E.g., “transformative” (NSF style) vs “transforming” (a real English word in dictionaries.)
- ▶ Don’t take things for granted: NLTK’s “The WordNet lemmatizer only removes affixes if the resulting word is in its dictionary.” which means that “lying” will not be restored back to “lie” by it.

# Downside of stemming and lemmatization

- ▶ By converting, e.g., “are”, “is”, “am” and “be” into “be”, we lose information that may help use on other tasks.
- ▶ Some combinations are meaningless when tokenized or lemmatized, e.g., in “I should have quited grad school when Sequoia gave me \$2M”, “should have quited” together mean something.
- ▶ Accurate lemmatization requires understanding the role of each word, e.g., the “saw” in “I saw the table (into half)” should be kept instead of being changed to “see.”

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

**Normalization and spell correction**

Building an inverse index

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing stopwords, words that are not important to the core document topic, e.g., “and”, “the”, etc.
  - ▶ Pluralizing nouns, e.g., “New York” to “New Yorks”
  - ▶ Expanding abbreviations, e.g., “PhD” to “PhD holder” and “Computer” to “computer science”
  - ▶ Normalizing the case, e.g., “New York” to “New York” and “New Yorks” to “New Yorks”
  - ▶ Normalizing the format, e.g., “New York” to “New York” and “New Yorks” to “New Yorks”

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Removing *stop tags*, e.g., “http://”, “www”, “mailto:”
  - ▶ Expanding abbreviations, e.g., “NYC” to “New York City” and “NY” to “New York”
  - ▶ Expanding acronyms, e.g., “FBI” to “Federal Bureau of Investigation”
  - ▶ Expanding contractions, e.g., “can’t” to “cannot”
  - ▶ Expanding numbers, e.g., “100” to “one hundred”
  - ▶ Expanding dates, e.g., “1/1/2010” to “January 1, 2010”
  - ▶ Expanding URLs, e.g., “http://www.example.com” to “www.example.com”
  - ▶ Expanding email addresses, e.g., “mailto:example@example.com” to “example@example.com”
  - ▶ Expanding phone numbers, e.g., “(123) 456-7890” to “123-456-7890”
  - ▶ Expanding social media handles, e.g., “@example” to “example”
  - ▶ Expanding hashtags, e.g., “#example” to “example”
  - ▶ Expanding emojis, e.g., “👍” to “thumbs up”
  - ▶ Expanding slang, e.g., “LOL” to “laugh out loud”
  - ▶ Expanding abbreviations, e.g., “NYC” to “New York City”
  - ▶ Expanding acronyms, e.g., “FBI” to “Federal Bureau of Investigation”
  - ▶ Expanding contractions, e.g., “can’t” to “cannot”
  - ▶ Expanding numbers, e.g., “100” to “one hundred”
  - ▶ Expanding dates, e.g., “1/1/2010” to “January 1, 2010”
  - ▶ Expanding URLs, e.g., “http://www.example.com” to “www.example.com”
  - ▶ Expanding email addresses, e.g., “mailto:example@example.com” to “example@example.com”
  - ▶ Expanding phone numbers, e.g., “(123) 456-7890” to “123-456-7890”
  - ▶ Expanding social media handles, e.g., “@example” to “example”
  - ▶ Expanding hashtags, e.g., “#example” to “example”
  - ▶ Expanding emojis, e.g., “👍” to “thumbs up”
  - ▶ Expanding slang, e.g., “LOL” to “laugh out loud”

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?



# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?

# Normalization

- ▶ Normalization means restoring a word to its canonical form, e.g., “saw” to “see”.
- ▶ Stemming and lemmatization are just one side of normalization.
- ▶ Other common tasks of normalization:
  - ▶ Removing *stop words*, words that are too frequency to become background noise, e.g., “the”, “of”, etc.
  - ▶ Preserving entities, e.g., “New York”.
  - ▶ Expanding abbreviations, e.g., “PhD” as “Piled Higher and Deeper” not “Philosophiae Doctor”.
  - ▶ Changing characters into lowercase. Not for all!
  - ▶ Question: Do you wanna treat “lower case” and “lowercase” the same? How about “E-mail” and “email”?

# PDF is not a mark-up language

Increasing demands for petroleum have stimulated sustainable ways to produce chemicals and biofuels. Specifically, fatty acids of varying chain lengths ( $C_6$ – $C_{16}$ ) naturally synthesized in many organisms are promising starting points for the catalytic production of industrial chemicals and diesel-like biofuels. However, bio-production of fatty acids from plants and other microbial production hosts relies heavily on manipulating tightly regulated fatty acid biosynthetic pathways. In addition, precursors for fatty acids are used along other central metabolic pathways for the production of amino acids and biomass, which further complicates the engineering of microbial hosts for higher yields. Here, we demonstrate an iterative metabolic engineering effort that integrates computationally driven predictions and metabolic flux analysis techniques to meet this challenge. The OptForce procedure was used for suggesting and prioritizing genetic manipulations that overproduce fatty acids of different chain lengths from  $C_6$  to  $C_{16}$  starting with wild-type *E. coli*. We identified some common but mostly chain-specific genetic interventions alluding to the possibility of fine-tuning overproduction for specific fatty acid chain lengths. In accordance with the OptForce prioritization of interventions, *fabZ* and acyl-ACP thioesterase were upregulated and *fadD* was deleted to arrive at a strain that produces 1.70 g/L and 0.14 g fatty acid/g glucose (~39% maximum theoretical yield) of  $C_{14-16}$  fatty acids in minimal M9 medium. These results highlight the benefit of using computational strain design and flux analysis tools in the design of recombinant strains of *E. coli* to produce free fatty acids.

© 2012 Elsevier Inc. All rights reserved.

Increasing demands for petroleum have stimulated sustainable ways to produce chemicals and biofuels.

Specifically, fatty acids of varying chain lengths ( $C_6$ – $C_{16}$ ) naturally synthesized in many organisms are promising starting points for the catalytic production of industrial chemicals and diesel-like biofuels. However, bio-production of fatty acids from plants and other microbial production hosts relies heavily on manipulating tightly regulated fatty acid biosynthetic pathways. In addition, precursors for fatty acids are used along other central metabolic pathways for the production of amino acids and biomass, which further complicates the engineering of microbial hosts for higher yields. Here, we demonstrate an iterative metabolic engineering effort that integrates computationally driven predictions and metabolic flux analysis techniques to meet this challenge. The OptForce procedure was used for suggesting and prioritizing genetic manipulations that overproduce fatty acids of different chain lengths from  $C_6$  to

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".



# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?  
<http://norvig.com/spell-correct.html>
- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
  - ▶ Do not think at token level.
  - ▶ Language model to be discussed. How to train the model?
  - ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?  
<http://norvig.com/spell-correct.html>
- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Tackling spelling errors or OCR errors

- ▶ How to build a spell checker?

<http://norvig.com/spell-correct.html>

- ▶ Using a dictionary, computing editing distance, and ranking candidates.
- ▶ Peter's solution assumes at most 2 edits from original string.
- ▶ String alignment: Smith-waterman algorithm, etc.
- ▶ What else?
- ▶ Do not think at token level.
- ▶ Language model to be discussed. How to train the model?
- ▶ What about symbols used in technical documents? E.g., "10g/mL".

# Outline

Representation of text data in computers

Pipeline for NLP Preprocessing

Corpus and Crawling

Clean up the text

Word and sentence tokenization

Stemming and Lemmatization

Normalization and spell correction

Building an inverse index

# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in `whoosh`.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID ("the entire value of the field as a single unit (that is, it doesn't break it up into individual words)."), text and numeric, respectively.
- ▶ "A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)"



# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in whoosh.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID ("the entire value of the field as a single unit (that is, it doesn't break it up into individual words)."), text and numeric, respectively.
- ▶ "A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)"

# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in `whoosh`.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID ("the entire value of the field as a single unit (that is, it doesn't break it up into individual words)."), text and numeric, respectively.
- ▶ "A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)"

# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in `whoosh`.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID (“the entire value of the field as a single unit (that is, it doesn’t break it up into individual words).”), text and numeric, respectively.
- ▶ “A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)”

# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in `whoosh`.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID (“the entire value of the field as a single unit (that is, it doesn’t break it up into individual words).”), text and numeric, respectively.
- ▶ “A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)”

# Building an inverse index

- ▶ An inverse index allows us to quickly locate a string (all occurrences or just one) in the corpus.
- ▶ Suppose I am building a file scanner for all files on my computer in `whoosh`.
- ▶ Begin with creating a schema, the structure of the index.

```
import whoosh, whoosh.fields, whoosh.analysis
schema = whoosh.fields.Schema(\
    path=whoosh.fields.ID(stored=True), \
    body=whoosh.fields.TEXT(\
        analyzer=whoosh.analysis.StemmingAnalyzer(), \
        stored=True), \
    page=whoosh.fields.NUMERIC(int, 64, signed=False, stored=True))
```

- ▶ You may imagine the schema as the headers of a table.
- ▶ This index contains 3 fields, path, body and page, which are ID (“the entire value of the field as a single unit (that is, it doesn’t break it up into individual words).”), text and numeric, respectively.
- ▶ “A field can be indexed (meaning it can be searched) and/or stored (meaning the value that gets indexed is returned with the results)”

# Building an inverse index

- ▶ Then you create an index object.

```
ix = whoosh.index.create_in("/home/forrest/index_dir",  
                             schema)  
writer = ix.writer()
```

- ▶ The index will be saved in the directory.
- ▶ Then you add documents (not necessarily a real file but a piece of text) into it. Just like filling rows in a table.

```
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I was born"), page=1)  
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I went to school"), page=2)
```

- ▶ Always remember to commit after you are done with adding documents each time.

```
writer.commit()
```

# Building an inverse index

- ▶ Then you create an index object.

```
ix = whoosh.index.create_in("/home/forrest/index_dir",  
                             schema)  
writer = ix.writer()
```

- ▶ The index will be saved in the directory.
- ▶ Then you add documents (not necessarily a real file but a piece of text) into it. Just like filling rows in a table.

```
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I was born"), page=1)  
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I went to school"), page=2)
```

- ▶ Always remember to commit after you are done with adding documents each time.

```
writer.commit()
```

# Building an inverse index

- ▶ Then you create an index object.

```
ix = whoosh.index.create_in("/home/forrest/index_dir",  
                             schema)  
writer = ix.writer()
```

- ▶ The index will be saved in the directory.
- ▶ Then you add documents (**not necessarily a real file but a piece of text**) into it. Just like filling rows in a table.

```
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I was born"), page=1)  
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I went to school"), page=2)
```

- ▶ Always remember to commit after you are done with adding documents each time.

```
writer.commit()
```



# Building an inverse index

- ▶ Then you create an index object.

```
ix = whoosh.index.create_in("/home/forrest/index_dir",  
                             schema)  
writer = ix.writer()
```

- ▶ The index will be saved in the directory.
- ▶ Then you add documents (**not necessarily a real file but a piece of text**) into it. Just like filling rows in a table.

```
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I was born"), page=1)  
writer.add_document(path=unicode("/mylife/preschool.log"),  
                    body=unicode("I went to school"), page=2)
```

- ▶ Always remember to commit after you are done with adding documents each time.

```
writer.commit()
```

# Building an inverse index

- ▶ Now, let's search by creating a query.

```
import whoosh.qparser
qp = whoosh.qparser.QueryParser("body", schema = ix.schema)
    # the query is on body field
q = qp.parse(u"born")          # the query value is "born"
```

- ▶ then pass the query to a searcher

```
searcher = ix.searcher()
results = searcher.search(q)
```

- ▶ Here is the results:

```
for result in results: # each result is for one document
    # which could contain multiple matches
    print "found match(es) in file %s".format(result["
        path"])
```

# Building an inverse index

- ▶ Now, let's search by creating a query.

```
import whoosh.qparser
qp = whoosh.qparser.QueryParser("body", schema = ix.schema)
    # the query is on body field
q = qp.parse(u"born")          # the query value is "born"
```

- ▶ then pass the query to a searcher

```
searcher = ix.searcher()
results = searcher.search(q)
```

- ▶ Here is the results:

```
for result in results: # each result is for one document
    # which could contain multiple matches
    print "found match(es) in file %s".format(result["
        path"])
```

# Building an inverse index

- ▶ Now, let's search by creating a query.

```
import whoosh.qparser
qp = whoosh.qparser.QueryParser("body", schema = ix.schema)
    # the query is on body field
q = qp.parse(u"born")          # the query value is "born"
```

- ▶ then pass the query to a searcher

```
searcher = ix.searcher()
results = searcher.search(q)
```

- ▶ Here is the results:

```
for result in results: # each result is for one document
    # which could contain multiple matches
    print "found match(es) in file %s".format(result["
        path"])
```

# Cloud APIs for syntactic analysis

- ▶ **Google:**

<https://cloud.google.com/natural-language/>

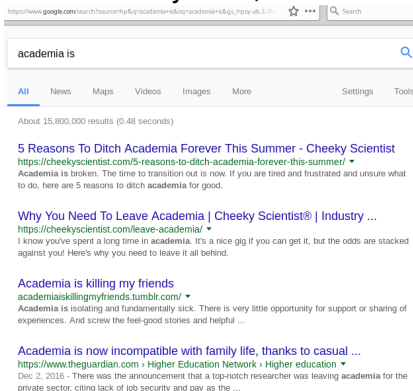
- ▶ **Microsoft:** <https://azure.microsoft.com/en-us/services/cognitive-services/text-analytics/>

# Other libraries

- ▶ TextBlob: built on top of NLTK
- ▶ SpaCy: much faster than NLTK

# Homework 3

- Write a Python program that passes a query (which may contain non-alphanumerics) to Google using GET method and extract the first sentence of each sample result on the first page. In the example below, you should extract “academia is broken”, “academia is isolating and fundamentally sick”, and others.



The screenshot shows a Google search interface. The search bar contains the text "academia is". Below the search bar, there are tabs for "All", "News", "Maps", "Videos", "Images", "More", "Settings", and "Tools". The "All" tab is selected. Below the tabs, it says "About 15,800,000 results (0.48 seconds)". The search results are listed below, with the first three results visible:

- 5 Reasons To Ditch Academia Forever This Summer - Cheeky Scientist**  
<https://cheekyscientist.com/5-reasons-to-ditch-academia-forever-this-summer/> ▼  
Academia is broken. The time to transition out is now. If you are tired and frustrated and unsure what to do, here are 5 reasons to ditch academia for good.
- Why You Need To Leave Academia | Cheeky Scientist® | Industry ...**  
<https://cheekyscientist.com/leave-academia/> ▼  
I know you've spent a long time in academia. It's a nice gig if you can get it, but the odds are stacked against you! Here's why you need to leave it all behind.
- Academia is killing my friends**  
[academiaiskillingmyfriends.tumblr.com/](https://academiaiskillingmyfriends.tumblr.com/) ▼  
Academia is isolating and fundamentally sick. There is very little opportunity for support or sharing of experiences. And screw the feel-good stories and helpful ...

Below these, there is another result:

- Academia is now incompatible with family life, thanks to casual ...**  
<https://www.theguardian.com> > Higher Education Network > Higher education ▼  
Dec 2, 2016 - There was the announcement that a top-notch researcher was leaving academia for the private sector, citing lack of job security and pay as the ...