# WordNet

Marina Sedinkina
- Folien von Desislava Zhekova -

CIS, LMU
marina.sedinkina@campus.lmu.de
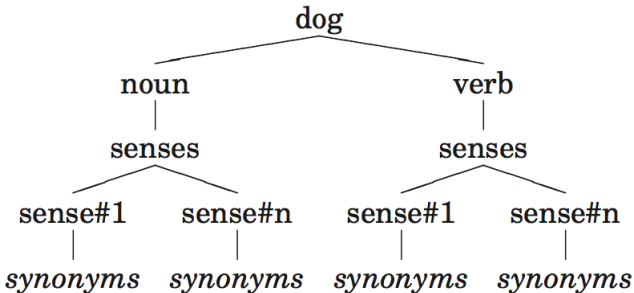
January 9, 2018

# Outline

## WordNet

- WordNet is a large lexical database of English (**semantically-oriented**)
- Nouns, verbs, adjectives and adverbs are grouped into sets of synonyms (**synsets**)
- Basis for grouping the words is their meanings.

# WordNet

English WordNet online: `http://wordnet.princeton.edu`



**WordNet Search - 3.1**
- WordNet home page - Glossary - Help

Word to search for: `motorcar`    Search WordNet

Display Options: (Select option to change) ▾  Change
Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations
Display options for sense: (gloss) "an example sentence"

**Noun**

- S: (n) car, auto, automobile, machine, **motorcar** (a motor vehicle with four wheels; usually propelled by an internal combustion engine) *"he needs a car to get to work"*
    - *direct hyponym* / *full hyponym*
        - S: (n) ambulance (a vehicle that takes people to and from hospitals)
        - S: (n) beach wagon, station wagon, wagon, estate car, beach

# WordNet

http://globalwordnet.org/

**Wordnets in the World**

| Language | Resource name | Developer(s) | Contact | Online Browsing | License | Other Resources |
|---|---|---|---|---|---|---|
| Afrikaans | Afrikaans WordNet ⧉ | North-West University, South Africa ⧉ | Gerhard van Huyssteen 🖾 Ané Bekker 🖾 | NO | OPEN FOR ACADEMIC USE ⧉ | |
| Albanian | AlbaNet ⧉ | Vlora University, Vlora, Albania ⧉ | Ervin Ruci 🖾 | YES ⧉ | OPEN (GPL) ⧉ | |
| Arabic | Arabic WordNet ⧉ | Arabic WordNet 🖾 | Horacio Rodriguez 🖾 | NO | OPEN | |
| Multilingual (Arabic/ English/ Malaysian/ Indonesian/ Finnish/ Hebrew/ Japanese/ Persian/ Thai/ French) | Open Multilingual Wordnet ⧉ | Linguistics and Multilingual Studies, NTU ⧉ | Francis Bond 🖾 | NO | OPEN | |

# WordNet

- NLTK includes the English WordNet (155,287 words and 117,659 synonym sets)
- NLTK graphical WordNet browser: `nltk.app.wordnet()`

Current Word: [          ]  Next Word: [          ]  Search
Help Shutdown

---

**noun**

- S: (noun) **wordnet** (any of the machine-readable lexical databases modeled after the Princeton WordNet)
- S: (noun) **WordNet**, Princeton WordNet (a machine-readable lexical database organized by meanings; developed at Princeton University)

## Senses and Synonyms

Consider the sentence in (1). If we replace the word motorcar in (1) with automobile, to get (2), the meaning of the sentence stays pretty much the same:

1. Benz is credited with the invention of the motorcar.
2. Benz is credited with the invention of the automobile.

⇒ Motorcar and automobile are synonyms.

Let's explore these words with the help of WordNet

## Senses and Synonyms

```
1  >>> from nltk.corpus import wordnet as wn
2  >>> wn.synsets("motorcar")
3  [Synset("car.n.01")]
```

- Motorcar has one meaning **car.n.01** (=the first noun sense of car).
- The entity **car.n.01** is called a **synset**, or "synonym set", a collection of synonymous words (or "lemmas"):

```
1  >>> wn.synset("car.n.01").lemma_names()
2  ["car", "auto", "automobile", "machine", "
       motorcar"]
```

## Senses and Synonyms

Synsets are described with a **gloss** (= definition) and some example sentences

```
1 >>> wn.synset("car.n.01").definition()
2 "a motor vehicle with four wheels; usually propelled
      by an internal combustion engine"
3 >>> wn.synset("car.n.01").examples()
4 ["he needs a car to get to work"]
```
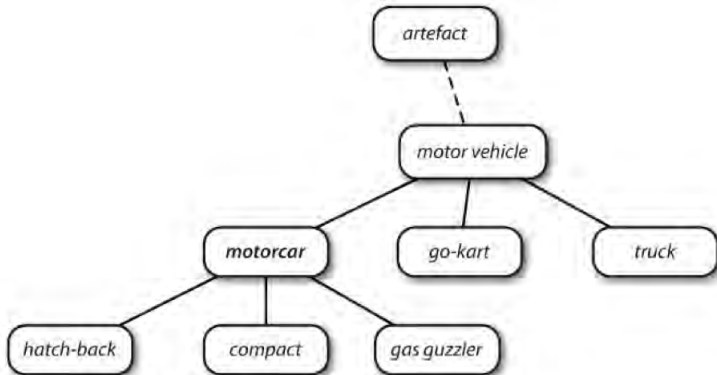
## Senses and Synonyms

Unlike the words automobile and motorcar, which are unambiguous
and have one synset, the word car is ambiguous, having five synsets:

```
1  >>> wn.synsets("car")
2  [Synset("car.n.01"), Synset("car.n.02"), Synset("car.
      n.03"), Synset("car.n.04"), Synset("cable_car.n.
      01")]
3  >>> for synset in wn.synsets("car"):
4  ...   print synset.lemma_names()
5  ...
6  ["car", "auto", "automobile", "machine", "motorcar"]
7  ["car", "railcar", "railway_car", "railroad_car"]
8  ["car", "gondola"]
9  ["car", "elevator_car"]
10 ["cable_car", "car"]
```

## The WordNet Hierarchy

Hypernyms and hyponyms ("is-a relation")

- *motor vehicle* is a hypernym of *motorcar*
- *ambulance* is a hyponym of *motorcar*

## The WordNet Hierarchy

```
1  >>> motorcar = wn.synset("car.n.01")
2  >>> types_of_motorcar = motorcar.hyponyms()
3  >>> types_of_motorcar[26]
4  Synset("ambulance.n.01")
5  >>> sorted([lemma.name() for synset in types_of_motorcar
       for lemma in synset.lemmas()])
6  ["Model_T", "S.U.V.", "SUV", "Stanley_Steamer", "ambulance"
       , "beach_waggon", "beach_wagon", "bus", "cab", "
       compact", "compact_car", "convertible", "coupe", "
       cruiser", "electric", "electric_automobile", "
       electric_car", "estate_car", "gas_guzzler", "hack", "
       hardtop", "hatchback", "heap", "horseless_carriage", "
       hot-rod", "hot_rod", "jalopy", "jeep", "landrover", "
       limo", "limousine", "loaner", "minicar", "minivan", "
       pace_car", "patrol_car", "phaeton", "police_car", "
       police_cruiser", "prowl_car", "race_car", "racer", "
       racing_car" ... ]
```

## The WordNet Hierarchy

```
1  >>> motorcar.hypernyms()
2  [Synset("motor_vehicle.n.01")]
3  >>> paths = motorcar.hypernym_paths()
4  >>> len(paths)
5  2
6  >>> [synset.name() for synset in paths[0]]
7  ["entity.n.01", "physical_entity.n.01", "object.n.01"
       , "whole.n.02", "artifact.n.01", "instrumentality
       .n.03", "container.n.01", "wheeled_vehicle.n.01",
        "self−propelled_vehicle.n.01", "motor_vehicle.n.
       01", "car.n.01"]
8  >>> [synset.name() for synset in paths[1]]
9  ["entity.n.01", "physical_entity.n.01", "object.n.01"
       , "whole.n.02", "artifact.n.01", "instrumentality
       .n.03", "conveyance.n.03", "vehicle.n.01", "
       wheeled_vehicle.n.01", "self−propelled_vehicle.n.
       01", "motor_vehicle.n.01", "car.n.01"]
```

## More Lexical Relations

Meronyms and holonyms

- *branch* is a meronym (*part meronym*) of *tree*
- *heartwood* is a meronym (*substance meronym*) of *tree*
- *forest* is a holonym (*member holonym*) of *tree*

## More Lexical Relations

```
1  >>> wn.synset("tree.n.01").part_meronyms()
2  [Synset("burl.n.02"), Synset("crown.n.07"), Synset("
      stump.n.01"), Synset("trunk.n.01"), Synset("limb.
      n.02")]
3  >>> wn.synset("tree.n.01").substance_meronyms()
4  [Synset("heartwood.n.01"), Synset("sapwood.n.01")]
5  >>> wn.synset("tree.n.01").member_holonyms()
6  [Synset("forest.n.01")]
```

## More Lexical Relations

Relationships between verbs:

- the act of walking involves the act of stepping, so walking entails stepping
- some verbs have multiple entailments

```
1  >>> wn.synset("walk.v.01").entailments()
2  [Synset("step.v.01")]
3  >>> wn.synset("eat.v.01").entailments()
4  [Synset("swallow.v.01"), Synset("chew.v.01")]
5  >>> wn.synset("tease.v.03").entailments()
6  [Synset("arouse.v.07"), Synset("disappoint.v.01")]
```

## More Lexical Relations

Some lexical relationships can express antonymy:

```
1  >>> wn.lemma("supply.n.02.supply").antonyms()
2  [Lemma("demand.n.02.demand")]
3  >>> wn.lemma("rush.v.01.rush").antonyms()
4  [Lemma("linger.v.04.linger")]
5  >>> wn.lemma("horizontal.a.01.horizontal").antonyms()
6  [Lemma("vertical.a.01.vertical"), Lemma("inclined.a.
       02.inclined")]
7  >>> wn.lemma("staccato.r.01.staccato").antonyms()
8  [Lemma("legato.r.01.legato")]
```

# More Lexical Relations

You can see the lexical relations, and the other methods defined on a synset, using `dir()`. For example:

```
1  import nltk
2  from nltk.corpus import wordnet as wn
3
4  print(wn.synsets("motorcar"))
5  # [Synset( car.n.01 )]
6
7  print(dir(wn.synsets("motorcar")[0]))
8  # [ ... common_hypernyms , definition , entailments , examples
        , frame_ids , hypernym_distances , hypernym_paths ,
     hypernyms , hyponyms , instance_hypernyms ,
     instance_hyponyms , jcn_similarity , lch_similarity ,
     lemma_names , lemmas , lexname , lin_similarity ,
     lowest_common_hypernyms , max_depth , member_holonyms ,
     member_meronyms , min_depth , name , offset ,
     part_holonyms , part_meronyms , path_similarity , pos ,
     region_domains , res_similarity , root_hypernyms ,
     shortest_path_distance , similar_tos , substance_holonyms ,
       substance_meronyms , topic_domains , tree , unicode_repr
       , usage_domains , verb_groups , wup_similarity ]
```

## Semantic Similarity

Two synsets linked to the same root may have several hypernyms in common. If two synsets share a very specific hypernym (low down in the hypernym hierarchy), they must be closely related.

```
1  >>> right = wn.synset("right_whale.n.01")
2  >>> orca = wn.synset("orca.n.01")
3  >>> minke = wn.synset("minke_whale.n.01")
4  >>> tortoise = wn.synset("tortoise.n.01")
5  >>> novel = wn.synset("novel.n.01")
6  >>> right.lowest_common_hypernyms(minke)
7  [Synset("baleen_whale.n.01")]
8  >>> right.lowest_common_hypernyms(orca)
9  [Synset("whale.n.02")]
10 >>> right.lowest_common_hypernyms(tortoise)
11 [Synset("vertebrate.n.01")]
12 >>> right.lowest_common_hypernyms(novel)
13 [Synset("entity.n.01")]
```

## Semantic Similarity

We can quantify this concept of generality by looking up the depth of each synset:

```
1  >>> wn.synset("baleen_whale.n.01").min_depth()
2  14
3  >>> wn.synset("whale.n.02").min_depth()
4  13
5  >>> wn.synset("vertebrate.n.01").min_depth()
6  8
7  >>> wn.synset("entity.n.01").min_depth()
8  0
```

## Semantic Similarity

Similarity measures have been defined over the collection of WordNet synsets that incorporate this insight

- path_similarity() assigns a score in the range 0-1 based on the shortest path that connects the concepts in the hypernym hierarchy
- -1 is returned in those cases where a path cannot be found
- Comparing a synset with itself will return 1

# Semantic Similarity

```
1  >>> right.path_similarity(minke)
2  0.25
3  >>> right.path_similarity(orca)
4  0.16666666666666666
5  >>> right.path_similarity(tortoise)
6  0.076923076923076927
7  >>> right.path_similarity(novel)
8  0.043478260869565216
```

## Similarity between nouns

- ("car", "automobile")
- synsets1("car") = $[synset_{11}, synset_{12}, synset_{13}]$
  nltk.corpus.wordnet.synsets("car")
- synsets2("automobile") = $[synset_{21}, synset_{22}, synset_{23}]$
  nltk.corpus.wordnet.synsets("automobile")
- consider all combinations of synsets formed by the synsets of the words in the word pair ("car", "automobile")
  $[(synset_{11}, synset_{21}), (synset_{11}, synset_{22}), (synset_{11}, synset_{23}), ...]$
- determine score of each combination e.g.:
  $synset_{11}$.path_similarity($synset_{21}$)
- determine the maximum score $\rightarrow$ indicator of similarity

# Semantic Similarity

### ???

Can you think of an NLP application for which semantic similarity will be helpful?

## Semantic Similarity

### ???

Can you think of an NLP application for which semantic similarity will be helpful?

### Suggestion

**Coreference Resolution:**

I saw an **orca**. The **whale** was huge.

## Polysemy

- The **polysemy** of a word is the number of senses it has.
- The noun **dog** has 7 senses in WordNet:

```
1  from nltk.corpus import wordnet as wn
2  num_senses=len(wn.synsets("dog","n"))
3
4  print(num_senses)
5  #prints 7
```

- We can also compute the average polysemy of nouns, verbs, adjectives and adverbs according to WordNet.

## Polysemy of nouns

We can also compute the average polysemy of nouns.

- Fetch all lemmas in WordNet that have a given POS:

  `nltk.corpus.wordnet.all_lemma_names(POS)`

```
1  from nltk.corpus import wordnet as wn
2  all_lemmas=set(wn.all_lemma_names("n"))
3  print(len(all_lemmas))
4  #prints 117798
```

- Determine meanings of each lemma:

  `nltk.corpus.wordnet.synsets(lemma,pos)` returns
  list of senses to a given lemma and POS, e.g. for "car"

```
1  from nltk.corpus import wordnet as wn
2  meanings=wn.synsets("car","n")
3  print(meanings)
4  #[Synset( car.n.01 ), Synset( car.n.02 ), ... )]
```

- Sum up the number of meanings of each lemma (restricted to nouns) and devide this by the total number of lemmas

# Lesk Similarity

### ???

Compute the average polysemy of nouns 'car', 'automobile', 'motorcar'

```
1  all_lemma_nouns = [ car , automobile , motorcar ]
2  senses_car = [Synset( car.n.01 ), Synset( car.n.02 ),
       Synset( car.n.03 ),Synset( cable_car.n.01 )]
3  senses_automobile = [Synset( car.n.01 )]
4  senses_motorcar = [Synset( car.n.01 )]
```

### average polysemy

average_polysemy = ???

## Lesk Algorithm

- classical algorithm for Word Sense Disambiguation (WSD) introduced by Michael E. Lesk in 1986
- idea: word's dictionary definitions are likely to be good indicators for the senses they define

## Lesk Algorithm: Example

| Sense | Definition |
|-------|------------|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

## Lesk Algorithm: Example

| **Sense** | **Definition** |
|---|---|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

Score = number of (stemmed) words that are shared by sense definition and context

| **Scores** | **Context** |
|---|---|
| s1 s2 | This cigar burns slowly and creates a stiff ash |

Table: Disambiguation of ash with Lesk's algorithm

## Lesk Algorithm: Example

| **Sense** | **Definition** |
|---|---|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

Score = number of (stemmed) words that are shared by sense definition and context

| **Scores** | **Context** |
|---|---|
| s1 s2 | This cigar burns slowly and creates a stiff ash |

Table: Disambiguation of ash with Lesk's algorithm

## Lesk Algorithm: Example

| **Sense** | **Definition** |
|-----------|----------------|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

Score = number of (stemmed) words that are shared by sense definition and context

| **Scores** | | **Context** |
|-----------|-----|-------------|
| s1 | s2 | This cigar burns slowly and |
| 0 | 1 | creates a stiff ash |

Table: Disambiguation of ash with Lesk's algorithm

Marina Sedinkina- Folien von Desislava Zhekova -  Language Processing and Python  33/67

## Lesk Algorithm: Example

| **Sense** | **Definition** |
|---|---|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

Score = number of (stemmed) words that are shared by sense definition and context

| **Scores** | **Context** |
|---|---|
| s1 s2 | The ash is one of the last trees |
| ??? | to come into leaf |

Table: Disambiguation of ash with Lesk's algorithm

## Lesk Algorithm: Example

| **Sense** | **Definition** |
|---|---|
| s1: tree | a tree of the olive family |
| s2: burned stuff | the solid residue left when combustible material is burned |

Table: Two senses of **ash**

Score = number of (stemmed) words that are shared by sense definition and context

| **Scores** | | **Context** |
|---|---|---|
| s1 | s2 | The ash is one of the last trees |
| 1 | 0 | to come into leaf |

Table: Disambiguation of ash with Lesk's algorithm

# Lesk Algorithm

```
1  >>> from nltk.wsd import lesk
2  >>> sent = [ I , went , to , the , bank , to ,
       deposit , money , . ]
3
4  >>> print(lesk(sent, bank , n ))
5  Synset( savings_bank.n.02 )
```

# Lesk Algorithm

The definitions for "bank" are:

```
>>> from nltk.corpus import wsordnet as wn
>>> for ss in wn.synsets( bank ):
...       print(ss, ss.definition())
...
```

Synset( bank.n.01 ) sloping land (especially the slope beside a body of water)

Synset( depository_financial_institution.n.01 ) a financial institution that accepts deposits and channels the money into lending activities

Synset( bank.n.03 ) a long ridge or pile

Synset( bank.n.04 ) an arrangement of similar objects in a row or in tiers

Synset( bank.n.05 ) a supply or stock held in reserve for future use (especially in emergencies)

Synset( bank.n.06 ) the funds held by a gambling house or the dealer in some gambling games

Synset( bank.n.07 ) a slope in the turn of a road or track; the outside is higher than the inside in order to reduce the effects of centrifugal force

Synset( savings_bank.n.02 ) a container (usually with a slot in the top) for keeping money at home

Synset( bank.n.09 ) a building in which the business of banking transacted

Synset( bank.n.10 ) a flight maneuver; aircraft tips laterally about its longitudinal axis (especially in turning)

Synset( bank.v.01 ) tip laterally

Synset( bank.v.02 ) enclose with a bank

Synset( bank.v.03 ) do business with a bank or keep an account at a bank

# Lesk Algorithm

Check implementation via

`http://www.nltk.org/_modules/nltk/wsd.html`

```python
1  def lesk(context_sentence, ambiguous_word, pos=None,
       synsets=None):
2      context = set(context_sentence)
3      if synsets is None:
4          synsets = wordnet.synsets(ambiguous_word)
5      if pos:
6          synsets = [ss for ss in synsets if str(ss.pos()) ==
               pos]
7      if not synsets:
8          return None
9
10     _, sense = max(
11         (len(context.intersection(ss.definition().split()))
               , ss) for ss in synsets
12     )
13     return sense
```

# Lesk Algorithm

Check implementation via
`http://www.nltk.org/_modules/nltk/wsd.html`

```python
1  def lesk(context_sentence, ambiguous_word, pos=None,
       synsets=None):
2      ...
3      if not synsets:
4          return None
5
6      inters = []
7      for ss in synsets:
8          defin_words = ss.definition().split()
9          intersec_words = context.intersection(defin_words)
10         len_iter = len(intersec_words)
11         inters.append((len_iter,ss))
12     _,sense = max(inters)
13
14     return sense
```

# Lesk Algorithm

- Information derived from a dictionary is insufficient for high quality **Word Sense Disambiguation (WSD)**.
- Lesk reports accuracies between 50% and 70%.
- Optimizations: to expand each word in the context with a list of synonyms

## Lesk Similarity

- The Lesk similarity of two concepts is defined as the textual **overlap between the corresponding definitions**, as provided by a dictionary.

- **Punctuation in definitions should be eliminated**, because they do not have a meaning. If two definitions contain punctuation, the score increases.

- The larger a text, the higher can its score be. It should be normalized to allow a fair comparison → **devide overlap by maximum matching number**

# Lesk Similarity

```
1  def lesk_similarity(synset1, synset2):
2      #TODO find tokens of wordnet definition of synset1, ignore
             punctuation
3      definition_words1 = ...
4
5      #TODO find tokens of wordnet definition of synset2, ignore
             punctuation
6      definition_words2 = ...
7
8      #TODO calculate maximum matching number (length of shortest
             definition)
9      max_match = ...
10
11     #TODO find overlap in definitions, consider words occuring
             twice
12     overlap = ...
13
14     return overlap/max_match
15
16 print(lesk_similarity(wn.synset( car.n.01 ),wn.synset( wheel.n.01
       )))
```

# Lesk Similarity

### ???

Find overlap in definitions, consider word occurring twice?

```
1  def1 = [ a , motor , vehicle , propelled , by , a
          , combustion , engine ]
2  def2 = [ a , vehicle , that , takes , people ,
          to , a , hospital ]
```

### overlap_number

overlap_number = ???

## Preprocessing

| | |
|---|---|
| Original | The boy's cars are different colors. |
| Tokenized | ["The", "boy's", "cars", "are", "different", "colors."] |
| Punctuation removal | ["The", "boy's", "cars", "are", "different", "colors"] |
| Lowecased | ["the", "boy's", "cars", "are", "different", "colors"] |
| Stemmed | ["the", "boy's", "car", "are", "differ", "color"] |
| Lemmatized | ["the", "boy's", "car", "are", "different", "color"] |
| Stopword removal | ["boy's", "car", "different", "color"] |

## Tokenization

- Tokenization is the process of breaking raw text into its building parts: words, phrases, symbols, or other meaningful elements called tokens.

- A list of tokens is almost always the first step to any other NLP task, such as part-of-speech tagging and named entity recognition.

## Tokenization

- **token** – is an instance of a sequence of characters in some particular document that are grouped together as a useful semantic unit for processing
- **type** – is the class of all tokens containing the same character sequence

## Tokenization

- What is Token?
- Fairly trivial: chop on whitespace and throw away punctuation characters.
- Tricky cases: various uses of the apostrophe for possession and contractions?

## Tokenization

Mrs. O'Reilly said that the girls' bags from
H&M's shop in New York aren't cheap.

| | |
|---|---|
| Mrs. | "Mrs."; "Mrs", "." |
| O'Reilly | "O'Reilly"; "OReilly"; "O'", "Reilly"; "O", "'", "Reilly"; |
| aren't | "aren't"; "arent"; "are", "n't"; "aren", "t" |
| ... | ... |

## Tokenization

### ???

Tokenize manually the following sentence. How many tokens do you get?

`Mrs.  O'Reilly said that the girls' bags from`
`H&M's shop in New York aren't cheap.`

# Tokenization

## ???

Tokenize manually the following sentence:
```
Mrs.  O'Reilly said that the girls' bags from
H&M's shop in New York aren't cheap.
```

## Answer

NLTK returns the following 20 tokens:
```
["Mrs.", "O'Reilly", "said", "that", "the",
"girls", "'", "bags", "from", "H", "&", "M",
"'s", "shop", "in", "New", "York", "are",
"n't", "cheap", "."]
```

## Tokenization

Most decisions need to be met depending on the language at hand.
Some problematic cases for English include:

- hyphenation – *ex-wife*, *Cooper-Hofstadter*, *the bazinga-him-again maneuver*
- internal white spaces – *New York*, *+49 89 21809719*, *January 1, 1995*, *San Francisco-Los Angeles*
- apostrophe – *O'Reilly*, *aren't*
- other cases – *H&M's*

## Sentence Segmentation

Tokenization can be approached at any level:

- word segmentation
- sentence segmentation
- paragraph segmentation
- other elements of the text

## Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1  >>> from nltk import word_tokenize,
       wordpunct_tokenize
2  >>> s = "Good muffins cost $3.88\nin New York.
       Please buy me\n two of them.\n\nThanks."
3  >>> word_tokenize(s)
4  [ Good , muffins , cost , $ , 3.88 , in , New ,
       York , . , Please , buy , me , two , of ,
       them , . , Thanks , . ]
5  >>> wordpunct_tokenize(s)
6  [ Good , muffins , cost , $ , 3 , . , 88 , in
       , New , York , . , Please , buy , me , two
       , of , them , . , Thanks , . ]
```

## Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1  >>> from nltk.tokenize import *
2  >>> # same as s.split():
3  >>> WhitespaceTokenizer().tokenize(s)
4  [Good, muffins, cost, $3.88, in, New,
     York., Please, buy, me, two, of,
     them., Thanks.]
5  >>> # same as s.split(" "):
6  >>> SpaceTokenizer().tokenize(s)
7  [Good, muffins, cost, $3.88\nin, New, York
     ., , Please, buy, me\ntwo, of, them
     .\n\nThanks.]
```

## Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1  >>> # same as s.split(\n):
2  >>> LineTokenizer(blanklines=keep).tokenize(s)
3  [Good muffins cost $3.88,  in New York.  Please buy
        me,    two of them.,   ,  Thanks.]
4  >>> # same as [l for l in s.split(\n) if l.strip()
        ]:
5  >>> LineTokenizer(blanklines=discard).tokenize(s)
6  [Good muffins cost $3.88,  in New York.  Please buy
        me,    two of them.,  Thanks.]
7  >>> # same as s.split(\t):
8  >>> TabTokenizer().tokenize(a\tb c\n\t d)
9  [a,  b c\n,   d]
```

# Segmentation

NLTK PunktSentenceTokenizer: divides a text into a list of sentences

```
1  >>> import nltk.data
2  >>> text = "Punkt knows that the periods in Mr. Smith and
       Johann S. Bach do not mark sentence boundaries.  And
       sometimes sentences can start with non-capitalized
       words.  i is a good variable name."
3  >>> sent_detector = nltk.data.load( tokenizers/punkt/
       english.pickle )
4  >>> print \n———\n .join(sent_detector.tokenize(text.
       strip()))
5  # Punkt knows that the periods in Mr. Smith and Johann S.
       Bach do not mark sentence boundaries.
6  # ———
7  # And sometimes sentences can start with non-capitalized
       words.
8  # ———
9  # i is a good variable name.
```

## Normalization

Once the text has been segmented into its tokens (paragraphs, sentences, words), most NLP pipelines do a number of other basic procedures for text normalization, e.g.:

- lowercasing
- stemming
- lemmatization
- stopword removal

# Lowercasing

Lowercasing:

```python
import nltk

string = "The boy´s cars are different colors."
tokens = nltk.word_tokenize(string)
lower = [x.lower() for x in tokens]
print(" ".join(lower))

# prints
# the boy ´s cars are different colors .
```

- Often, however, instead of working with all word forms, we would like to extract and work with their base forms (e.g. lemmas or stems)
- Thus with **stemming** and **lemmatization** we aim to reduce inflectional (and sometimes derivational) forms to their base forms.

# Stemming

**Stemming**: removing morphological affixes from words, leaving only the word stem.

```
1   import nltk
2
3   string = "The boy´s cars are different colors."
4   tokens = nltk.word_tokenize(string)
5   lower = [x.lower() for x in tokens]
6   stemmed = [stem(x) for x in lower]
7   print(" ".join(stemmed))
8
9   def stem(word):
10      for suffix in ["ing", "ly", "ed", "ious", "ies", "ive",
                "es", "s", "ment"]:
11          if word.endswith(suffix):
12              return word[:-len(suffix)]
13      return word
14  # prints
15  # the boy ´s car are different color .
```

# Stemming

Stemming:

```
1   import nltk
2   import re
3
4   string = "The boy s cars are different colors."
5   tokens = nltk.word_tokenize(string)
6   lower = [x.lower() for x in tokens]
7   stemmed = [stem(x) for x in lower]
8   print(" ".join(stemmed))
9
10  def stem(word):
11      regexp = r"^(.*?)(ing|ly|ed|ious|ies|ive|es|s|ment)?$"
12      stem, suffix = re.findall(regexp, word)[0]
13      return stem
14
15  # prints
16  # the boy ´s car are different color .
```

## Stemming

NLTK's stemmers:

- **Porter Stemmer** is the oldest stemming algorithm supported in NLTK, originally published in 1979.
  http:
  //www.tartarus.org/~martin/PorterStemmer/

- **Lancaster Stemmer** is much newer, published in 1990, and is more aggressive than the Porter stemming algorithm.

- **Snowball stemmer** currently supports several languages: Danish, Dutch, English, Finnish, French, German, Hungarian, Italian, Norwegian, Porter, Portuguese, Romanian, Russian, Spanish, Swedish.

- **Snowball stemmer**: slightly faster computation time than porter.

# Stemming

NLTK's stemmers:

```
1  import nltk
2
3  string = "The boy´s cars are different colors."
4  tokens = nltk.word_tokenize(string)
5  lower = [x.lower() for x in tokens]
6
7  porter = nltk.PorterStemmer()
8  stemmed = [porter.stem(t) for t in lower]
9  print(" ".join(stemmed))
10 # prints
11 # the boy ´s car are differ color .
12
13 lancaster = nltk.LancasterStemmer()
14 stemmed = [lancaster.stem(t) for t in lower]
15 print(" ".join(stemmed))
16 # prints
17 # the boy ´s car ar diff col .
```

# Stemming

NLTK's stemmers:

```
1   import nltk
2
3   string = "The boy´s cars are different colors."
4   tokens = nltk.word_tokenize(string)
5   lower = [x.lower() for x in tokens]
6
7   snowball = nltk.SnowballStemmer("english")
8   stemmed = [snowball.stem(t) for t in lower]
9   print(" ".join(stemmed))
10  # prints
11  # the boy ´s car are differ color .
```

## Lemmatization

- stemming can often create non-existent words, whereas lemmas are actual words
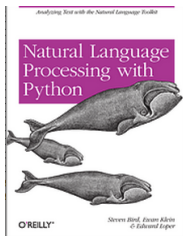- **NLTK WordNet Lemmatizer** uses the WordNet Database to lookup lemmas

```
1  import nltk
2  string = "The boy´s cars are different colors."
3  tokens = nltk.word_tokenize(string)
4  lower = [x.lower() for x in tokens]
5  porter = nltk.PorterStemmer()
6  stemmed = [porter.stem(t) for t in lower]
7  print(" ".join(lemmatized))
8  # prints the boy ´s car are differ color .
9  wnl = nltk.WordNetLemmatizer()
10 lemmatized = [wnl.lemmatize(t) for t in lower]
11 print(" ".join(lemmatized))
12 # prints the boy ´s car are different color .
```

# Stopword removal:

Stopword removal:

```
 1  import nltk
 2
 3  string = "The boy s cars are different colors."
 4  tokens = nltk.word_tokenize(string)
 5  lower = [x.lower() for x in tokens]
 6  wnl = nltk.WordNetLemmatizer()
 7  lemmatized = [wnl.lemmatize(t) for t in lower]
 8
 9  content = [x for x in lemmatized if x not in nltk.
        corpus.stopwords.words("english")]
10  print(" ".join(content))
11  # prints
12  # boy ´s car different color .
```

## References

http://www.nltk.org/book/

- https://github.com/nltk/nltk
- Christopher D. Manning, Hinrich Schütze 2000. Foundations of Statistical Natural Language Processing. *The MIT Press Cambridge, Massachusetts London, England*. http://ics.upjs.sk/~pero/web/documents/ pillar/Manning_Schuetze_StatisticalNLP.pdf