

WordNet

Marina Sedinkina

- Folien von Desislava Zhekova -

CIS, LMU

`marina.sedinkina@campus.lmu.de`

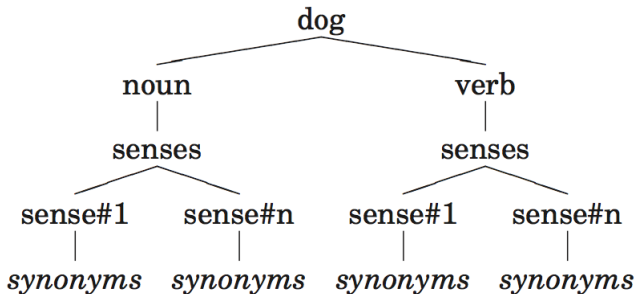
December 18, 2018

Outline

- 1 WordNet
- 2 Lesk Algorithm
- 3 Finding Hypernyms with WordNet
- 4 Relation Extraction with spaCy
- 5 References

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:

Display Options:

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:

[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

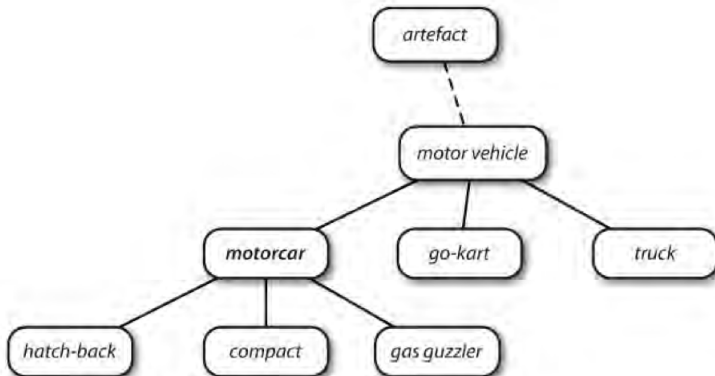
Unlike the words `automobile` and `motorcar`, which are unambiguous and have one synset, the word `car` is ambiguous, having five synsets:

```
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
6             for lemma in synset.lemmas()])
["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 hold between lemmas, e.g., **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.  
7 02.inclined")]  
8 >>> wn.lemma("staccato.r.01.staccato").antonyms()  
9 [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 >>>[ ... , '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_to', '
  substance_holonyms', 'substance_meronyms', 'topic_domains', '
  tree', 'unicode_repr', 'usage_domains', 'verb_groups', '
  wup_similarity']
```

Semantic Similarity

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") = [synset11, synset12, synset13]`
`nlk.corpus.wordnet.synsets("car")`
- `synsets2("automobile") = [synset21, synset22, synset23]`
`nlk.corpus.wordnet.synsets("automobile")`
- consider all combinations of synsets formed by the synsets of the words in the word pair ("car", "automobile")
`[(synset11, synset21), (synset11, synset22), (synset11, synset23), ...]`
- determine score of each combination e.g.:
`synset11.path_similarity(synset21)`
- determine the maximum score → 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 >>>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 >>>
5 [Synset('car.n.01'), Synset('car.n.02'), Synset('car.n.03'),
6  Synset('car.n.04'), Synset('cable_car.n.01')]
```

Polysemy of nouns

```
1 def average_polysemy(part_of_speech):
2
3     lemmas = set(nltk.corpus.wordnet.all_lemma_names(
4         part_of_speech))
5
6     nr_of_synsets = 0
7     for lemma in lemmas:
8         nr_of_synsets += len(nltk.corpus.wordnet.synsets(lemma,
9             pos=part_of_speech))
10
11     return nr_of_synsets / len(lemmas)
```

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

s1: tree

Definition

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

s1: tree

s2: burned stuff

Definition

a tree of the olive family

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

s1 s2

Context

This cigar burns slowly and
creates a stiff ash

Table: Disambiguation of ash with Lesk's algorithm

Lesk Algorithm: Example

Sense

s1: tree

s2: burned stuff

Definition

a tree of the olive family

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

s1 s2

Context

This cigar **burns** slowly and
creates a stiff ash

Table: Disambiguation of ash with Lesk's algorithm

Lesk Algorithm: Example

Sense

s1: tree

s2: burned stuff

Definition

a tree of the olive family

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

s1 s2
0 1

Context

This cigar **burns** slowly and
creates a stiff ash

Table: Disambiguation of ash with Lesk's algorithm

Lesk Algorithm: Example

Sense

s1: tree

s2: burned stuff

Definition

a tree of the olive family

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

s1 s2

???

Context

The ash is one of the last trees
to come into leaf

Table: Disambiguation of ash with Lesk's algorithm

Lesk Algorithm: Example

Sense

s1: tree

s2: burned stuff

Definition

a tree of the olive family

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

s1 s2
1 0

Context

The ash is one of the last trees
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:

```
1 >>> from nltk.corpus import wordnet as wn
2 >>> for ss in wn.synsets("bank"):
3     ...     print(ss, ss.definition())
4 Synset('bank.n.01') sloping land (especially the slope beside a body of water)
5 Synset('depository_financial_institution.n.01') a financial institution that accepts
  deposits and channels the money into lending activities
6 Synset('bank.n.03') a long ridge or pile
7 Synset('bank.n.04') an arrangement of similar objects in a row or in tiers
8 Synset('bank.n.05') a supply or stock held in reserve for future use (especially in
  emergencies)
9 Synset('bank.n.06') the funds held by a gambling house or the dealer in some gambling
  games
10 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
11 Synset('savings_bank.n.02') a container (usually with a slot in the top) for keeping
  money at home
12 Synset('bank.n.09') a building in which the business of banking transacted
13 Synset('bank.n.10') a flight maneuver; aircraft tips laterally about its longitudinal
  axis (especially in turning)
14 Synset('bank.v.01') tip laterally
15 Synset('bank.v.02') enclose with a bank
```

Lesk Algorithm

Check implementation via

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

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

Task

TASK TO SOLVE

In the Wikipedia article on Ada Lovelace,

- how many words refer to a **relative**? (excluding names)
- how many words refer to an **illness**?
- how many words refer to a **science**?

In each case: which words?

Task

TASK TO SOLVE

In the Wikipedia article on Ada Lovelace,

- how many words refer to a **relative**? (excluding names)
- how many words refer to an **illness**?
- how many words refer to a **science**?

In each case: which words?

Let's solve this using WordNet...

Step 1: Read in file

Read `ada_lovelace.txt` as one text string.

```
1 >>> print text
2 "Augusta Ada King, Countess of Lovelace (10 December 1815
3 27 November 1852), born Augusta Ada Byron and
4 now commonly known as Ada Lovelace, was an
5 English mathematician and writer chiefly known
6 for her work on Charles Babbage's early mechanical
7 general-purpose computer, the Analytical Engine. ... "
```

Step 2: Sentence Splitting

Split the text into sentences: `nlk.sent_tokenize(text)`

```
1 >>> print sentences[:3]
2 ["Augusta Ada King, Countess of Lovelace (10 December
3 1815 27 November 1852), born Augusta Ada Byron
4 and now commonly known as Ada Lovelace, was an English
5 mathematician and writer chiefly known for her work on
6 Charles Babbage's early mechanical general-purpose
7 computer, the Analytical Engine.", 'Her notes on
8 the engine include what is recognised as the first
9 algorithm intended to be carried out by a machine.',
10 "Because of this, she is often described as the
11 world's first computer programmer.", ... ]
```

Step 3: Tokenize

Split the sentences into tokens: `nlk.word_tokenize(text)`
Create one list of tokens (containing all tokens of the text).

```
1 >>> print tokens
2 ['Augusta', 'Ada', 'King', ',', 'Countess', 'of',
3 'Lovelace', '(', '10', 'December', '1815', '27',
4 'November', '1852', ')', ',', 'born', 'Augusta',
5 'Ada', 'Byron', 'and', 'now', 'commonly', 'known',
6 'as', 'Ada', 'Lovelace', ',', 'was', 'an',
7 'English', 'mathematician', 'and', 'writer',
8 'chiefly', 'known', 'for', 'her', 'work', 'on',
9 'Charles', 'Babbage', "'s", 'early', 'mechanical',
10 'general-purpose', 'computer', ',', 'the',
11 'Analytical', 'Engine', '.', 'Her', 'notes', 'on',
12 'the', 'engine', 'include', 'what', 'is',
13 'recognised', 'as', 'the', 'first', 'algorithm',
14 'intended', 'to', 'be', 'carried', 'out',
```

Step 4: Part-of-Speech tagging

Find the POS-tag of each token using NLTK's recommended POS tagger.

```
1 pos_tags = nltk.pos_tag(tokens)
2 print pos_tags
3
4 [('Augusta', 'NNP'), ('Ada', 'NNP'), ('King', 'NNP'),
5  ('', ''), ('Countess', 'NNP'), ('of', 'IN'),
6  ('Lovelace', 'NNP'), ('(', 'NNP'), ('10', 'CD'),
7  ('December', 'NNP'), ('1815', 'CD'), ('27', 'CD'),
8  ('November', 'NNP'), ('1852', 'CD'), (')', 'CD'),
9  ('', ''), ('born', 'NN'), ('Augusta', 'NNP'),
10 ('Ada', 'NNP'), ('Byron', 'NNP'), ('and', 'CC'),
11 ('now', 'RB'), ('commonly', 'RB'), ('known',
12 'VBN'), ('as', 'IN'), ('Ada', 'NNP'), ... ]
```

Print out all the nouns occurring in the text. 

Step 4: Part-of-Speech tagging

```
1 [ ... (',', ', ', ', '), ('born', 'NN'), ('Augusta', 'NNP'),  
2 ('Ada', 'NNP'), ('Byron', 'NNP'), ('and', 'CC'),  
3 ('now', 'RB'), ('commonly', 'RB'), ('known',  
4 'VBN'), ('as', 'IN'), ('Ada', 'NNP'), ... ]
```

- CC – coordinating conjunction
- RB – adverb
- IN – preposition
- NN – noun
- JJ – adjective
- VB – verb

Step 4: Part-of-Speech tagging

NLTK provides documentation for each tag, which can be queried using the tag, e.g:

```
1 >>> nltk.help.upenn_tagset('NN')
2 NN: noun, common, singular or mass
3   common-carrier cabbage knuckle-duster Casino
4     afghan shed thermostat investment slide
5     humour falloff slick wind hyena override
6     subhumanity machinist ...
7 >>> nltk.help.upenn_tagset('CC')
8 CC: conjunction, coordinating
9   & and both but either et for less minus neither
10  nor or plus so therefore times v. versus vs.
11  whether yet
```

Step 4: Part-of-Speech tagging

Note!

Some POS tags denote variation of the same word type, e.g. NN, NNS, NNP, NNPS, such can be looked up via regular expressions.

```
1 >>> nltk.help.upenn_tagset('NN*')
2 NN: noun, common, singular or mass
3     common-carrier cabbage knuckle-duster Casino ...
4 NNP: noun, proper, singular
5     Motown Venneboerger Czestochwa Ranzer Conchita
6     ...
7 NNPS: noun, proper, plural
8     Americans Americas Amharas Amityvilles ...
9 NNS: noun, common, plural
10    undergraduates scotches bric-a-brac ...
```


Step 4: Lemmatize

Now, put the lemma of each noun from the text into one list.

```
1 from nltk.stem.wordnet import WordNetLemmatizer
2 from nltk.corpus import wordnet
3 lemmatizer = WordNetLemmatizer()
4 # your code ...
5 lemmatizer.lemmatize(lemma, wordnet.NOUN)
6 # your code ...
7 >>> print noun_lemmas
8 ['Augusta', 'Ada', 'King', 'Countess', 'Lovelace'
9  ('', 'December', 'November', 'born', 'Augusta',
10 'Ada', 'Byron', 'Ada', 'Lovelace',
11 'mathematician', 'writer', 'work', 'Charles',
12 'Babbage', 'computer', ... ]
```

Ada Lovelace Task: Hypernyms

These are the three hypernyms of interest:
(as there are multiple synsets for a lemma, we pick the first one
in each list returned by `nltk.wordnet`)

```
1 relative = wordnet.synsets("relative", pos='n')[0]
2 science = wordnet.synsets("science", pos='n')[0]
3 illness = wordnet.synsets("illness", pos='n')[0]
```

Ada Lovelace Task: Hypernyms

These are the three hypernyms of interest:
(as there are multiple synsets for a lemma, we pick the first one
in each list returned by `nltk.wordnet`)

```
1 relative = wordnet.synsets("relative", pos='n')[0]
2 science = wordnet.synsets("science", pos='n')[0]
3 illness = wordnet.synsets("illness", pos='n')[0]
```

How can we find out whether one synset is a hyponym of
another?

Ada Lovelace Task: Hypernym Code

```
1 def hypernymOf(synset1, synset2):  
2     """ Returns True if synset2 is a hypernym of  
3     synset1, or if they are the same synset.  
4     Returns False otherwise. """  
5     if synset1 == synset2:  
6         return True  
7     for hypernym in synset1.hypernyms():  
8         if synset2 == hypernym:  
9             return True  
10        if hypernymOf(hypernym, synset2):  
11            return True  
12    return False
```

Ada Lovelace Task: Finding Hypernyms

Reminder:

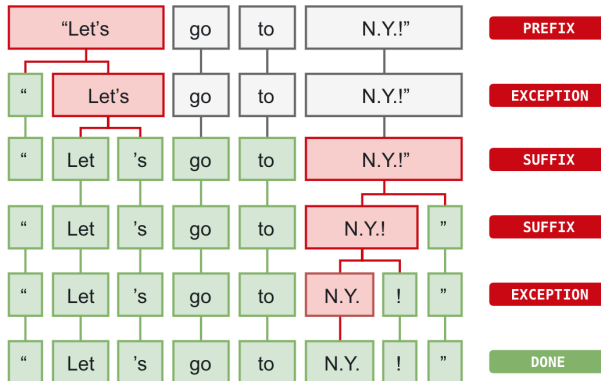
- We have a list of the lemmas of all nouns, `noun_lemmas`.
- Retrieve the synsets for each lemma.
- Check whether it's a hyponym of one of the three synsets of interest.
- Counts the relevant nouns, and collect them.

Tokenization with spaCy

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'Apple is looking at buying U.K. startup for $1 billion
5         ')
6 for token in doc:
7     print(token.text)
8 >>>
9     Apple
10    is
11    looking
12    at
13    buying
14    U.K.
15    ...
```

Tokenization with spaCy

- Does the substring match a tokenizer exception rule? (U.K.)
- Can a prefix, suffix or infix be split off? (e.g. punctuation)



Tokenization with spaCy

- Tokenizer exceptions strongly depend on the specifics of the individual language
- Global and language-specific tokenizer data is supplied via the language data in **spacy/lang**

| LANGUAGE | CODE | LANGUAGE DATA | MODELS |
|------------|------|-------------------------|--------------------------|
| English | en | <code>lang/en</code> <> | 4 models |
| German | de | <code>lang/de</code> <> | 1 model |
| Spanish | es | <code>lang/es</code> <> | 2 models |
| Portuguese | pt | <code>lang/pt</code> <> | 1 model |
| French | fr | <code>lang/fr</code> <> | 2 models |
| Italian | it | <code>lang/it</code> <> | 1 model |

Adding special case tokenization rules

- The tokenizer exceptions define special cases like "don't" in English, which needs to be split into two tokens: {ORTH: do} and {ORTH: n't, LEMMA: not}

```
1 import spacy
2 from spacy.symbols import ORTH, LEMMA, POS, TAG
3
4 nlp = spacy.load('en_core_web_sm')
5 doc = nlp(u'gimme that') # phrase to tokenize
6 print([w.text for w in doc]) # ['gimme', 'that']
7
8 # add special case rule
9 special_case = [{ORTH: u'gim', LEMMA: u'give', POS: u'VERB'},
10                {ORTH: u'me'}]
11 nlp.tokenizer.add_special_case(u'gimme', special_case)
12
13 # check new tokenization
14 print([w.text for w in nlp(u'gimme that')]) # ['gim', 'me', 'that']
```

Adding special case tokenization rules

```
1 doc = nlp(u'I like New York in Autumn.')
```

```
2 span = doc[2:4]
```

```
3 span.merge()
```

```
4 assert len(doc) == 6
```

```
5 assert doc[2].text == 'New York'
```

Relation extraction with spaCy

TASK TO SOLVE

Extract money and currency values (entities labelled as MONEY) and find the noun phrase they are referring to - for example:

“Net income was \$9.4 million compared to the prior year of \$2.7 million.”

\$9.4 million → Net income.

\$2.7 million → the prior year

???

How can we solve this task?

Relation extraction with spaCy

TASK TO SOLVE

Extract money and currency values (entities labelled as MONEY) and find the noun phrase they are referring to - for example:

“Net income was \$9.4 million compared to the prior year of \$2.7 million.”

\$9.4 million → Net income.

\$2.7 million → the prior year

- Step 1: use spaCy's **named entity recognizer** to extract money and currency values (entities labelled as MONEY)
- Step2: use spaCy's **dependency parser** to find the noun phrase they are referring to.

Relation extraction with spaCy

Step 1: use spaCy's **named entity recognizer** to extract money and currency values (entities labelled as MONEY)

```
1 import spacy
2 model = spacy.load('en_core_web_sm')
3
4 doc = nlp(u'Net income was $9.4 million compared to
5         the prior year of $2.7 million.')
6 print(doc.ents)
7
8 >>> $9.4 million , the prior year , $2.7 million
9
10 print([token.ent_type_ for token in doc])
11 [' ', ' ', ' ', 'MONEY', 'MONEY', 'MONEY', ' ', ' ', 'DATE',
12  ' ', 'DATE', 'DATE', ' ', 'MONEY', 'MONEY', 'MONEY',
13  '']
```

Relation extraction with spaCy

Step2: use spaCy's **dependency parser** to find the noun phrases

```
1 import spacy
2 model = spacy.load('en_core_web_sm')
3
4 doc = nlp(u'Net income was $9.4 million compared to
5 the prior year of $2.7 million.')
6 for noun_phrase in doc.noun_chunks:
7     print(noun_phrase)
8 Net income
9 the prior year
```

Relation extraction with spaCy

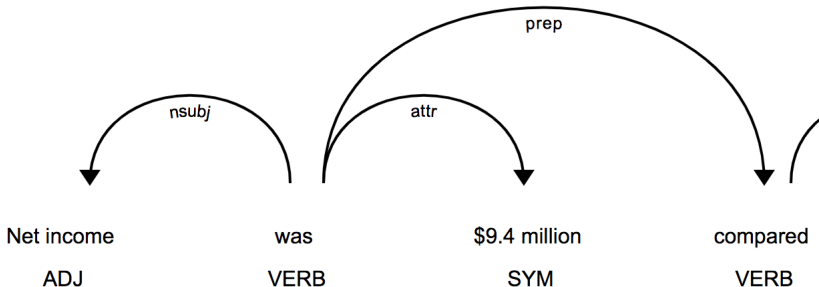
Step 3: convert MONEY phrases and noun phrases to one token

```
1 import spacy
2 model = spacy.load('en_core_web_sm')
3
4 doc = nlp(u'Net income was $9.4 million compared to
5         the prior year of $2.7 million.')
6 #your code
7 for token in doc:
8     print(token.text)
9
10 Net income
11 was
12 $9.4 million
13 ...
```

Relation extraction with spaCy

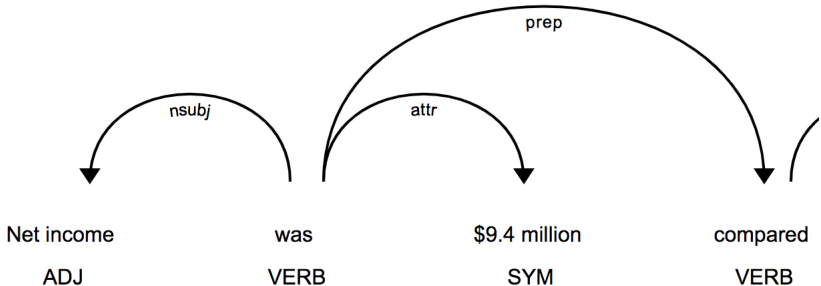
Step4: link named entities (MONEY) to the noun phrases they are referring to: use **dependency labels**

- 1 `from spacy import displacy`
- 2 `displacy.serve(doc, style='dep')`



Relation extraction with spaCy

- An attribute (**attr**) is a noun phrase that is a non-VP (verbal phrase) predicate usually following a copula verb such as “to be”
- A nominal subject (**nsubj**) is a noun phrase which is the syntactic subject of a clause.



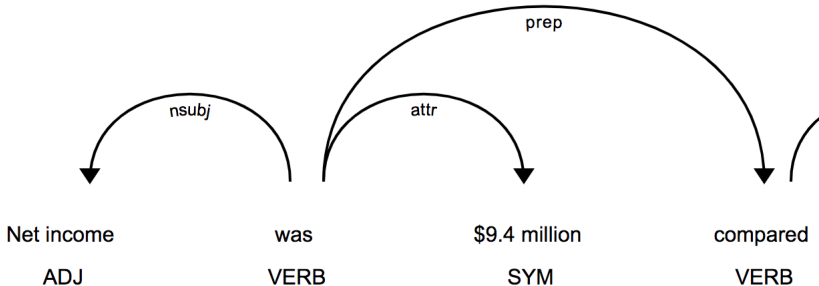
Relation extraction with spaCy

Step4: link named entities (MONEY) to the noun phrases they are referring to: use **dependency labels**

```
1 import spacy
2 model = spacy.load('en_core_web_sm')
3
4 doc = nlp(u'Net income was $9.4 million compared to
5 the prior year of $2.7 million.')
6 for token in doc:
7     print(token.text, token.dep_, token.head.text,
8           [el for el in token.head.lefts])
9 Net income    nsubj was [Net income]
10 was          ROOT  was [Net income]
11 $9.4 million attr was [Net income]
12 ...
```

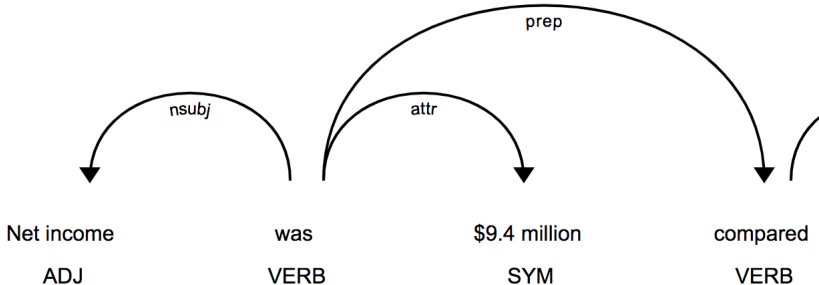
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------------|------------|------------|------------------|
| Net income | nsubj | was | [Net income] |
| was | ROOT | was | [Net income] |
| \$9.4 million | attr | was | [Net income] |



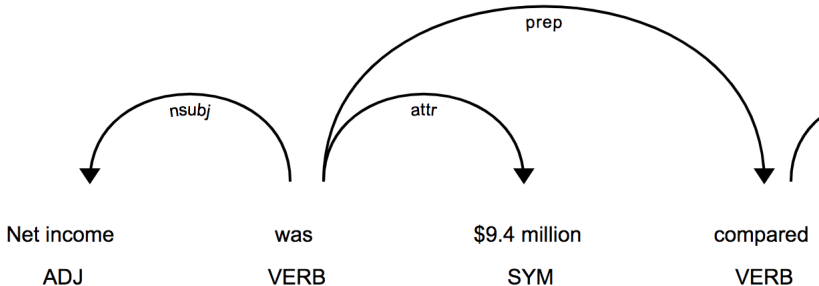
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------------|-------------|------------|-------------------------|
| Net income | nsubj | was | [Net income] |
| was | ROOT | was | [Net income] |
| \$9.4 million | attr | was | [Net income] |



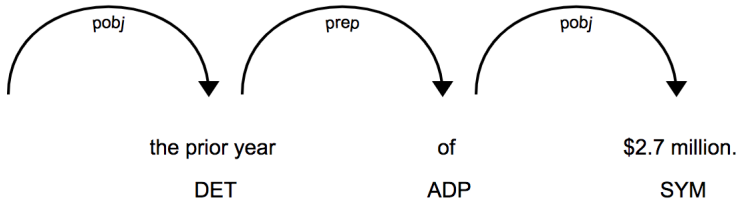
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------------|--------------|------------|---------------------|
| Net income | nsubj | was | [Net income] |
| was | ROOT | was | [Net income] |
| \$9.4 million | attr | was | [Net income] |



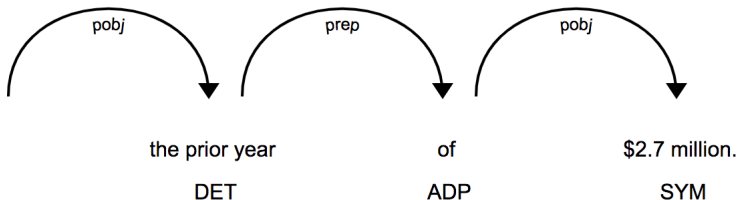
Relation extraction with spaCy

- An object of a preposition (**pobj**) is a noun phrase that modifies the head of a prepositional phrase, which is usually a preposition.
- A prepositional modifier (**prep**) is any prepositional phrase that modifies the meaning of its head.



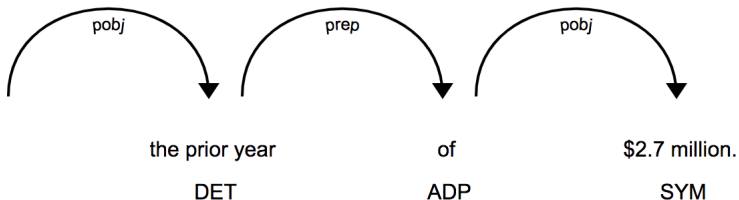
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------------|-------------|----------------|------------------|
| the prior year | pobj | to | [] |
| of | prep | the prior year | [] |
| \$2.7 million | pobj | of | [] |



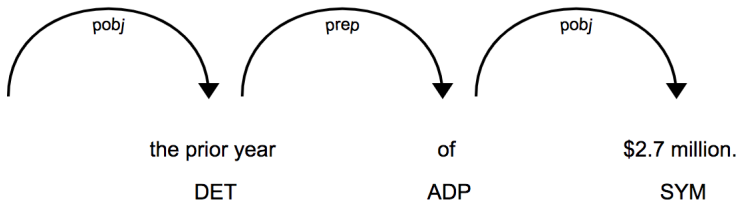
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------------|-------------|----------------|------------------|
| the prior year | pobj | to | [] |
| of | prep | the prior year | [] |
| \$2.7 million | pobj | of | [] |



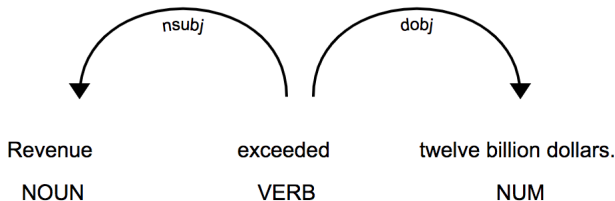
Relation extraction with spaCy

| token | token.dep_ | token.head | token.head.lefts |
|----------------|-------------|-----------------------|-------------------------|
| the prior year | pobj | to | [] |
| of | prep | the prior year | [] |
| \$2.7 million | pobj | of | [] |



Relation extraction with spaCy

- From which sentences the information will be extracted?
 - Research and product development expenses were \$6 million.
 - Net loss for the year ended December 31, 2017 was \$11 million.
 - an increase of \$0.4 million
 - greater by \$2.9 million
- What about a direct object (**dobj**)? It is a noun phrase that is the accusative object of the verb.
 - Revenue exceeded twelve billion dollars.



Conclusion

- **WordNet** is a large lexical database where nouns, verbs, adjectives and adverbs are grouped into sets of synonyms:
 - word sense disambiguation - **Lesk Algorithm** (also implemented in **NLTK**)
 - find hypernyms and hyponyms
- **spaCy** is open-source library for advanced Natural Language Processing (NLP) in Python
 - use pre-trained models (e.g. **en_core_web_sm**)
 - use the models to preprocess the text: e.g. tokenization, pos-tagging and lemmatization
 - customize tokenizer
 - use the models for information extraction: named entities, dependency labels (use both for relation extraction)



References

- <http://www.nltk.org/book/>
- <https://github.com/nltk/nltk>
- <https://spacy.io/>