

Concepts and Applications in NLP

Word Senses and WordNet

Marion Di Marco

December 17, 2024

Challenges in NLP: Word Senses

SAME WORD · DIFFERENT MEANINGS



seal

SAME WORD · DIFFERENT MEANINGS



plant

SAME WORD · DIFFERENT MEANINGS



sink



nut



bark



nail



crane



letter



note



bow



bat



fan



trunk



table



button

Illustrations from <https://www.englishforkidz.com/2020/01/vocabulary-cards-word-different-meanings.html>

Outline

Word Senses

Relations between Senses

WordNet

Word Sense Disambiguation

Credits and References

Word Senses

- **Word sense:**
discrete representation of one aspect of the meaning of a word
- Context helps to differentiate the meanings:
- *mouse*¹ : ... a mouse controlling a computer system in 1968.
*mouse*² : ... a quiet animal like a mouse.
- *bank*¹ : ... a bank can hold the investments in a custodial account ...
*bank*² : ... as agriculture burgeons on the east bank, the river ...

Defining Word Senses

- Embeddings: define the meaning of a word by its co-occurrences
- Contextual embeddings: embedding that represents the meaning of a word in its textual context
- How to define the meaning of a a word *sense*?
- Textual definitions for each sense (glosses) as given in dictionaries or thesauri
- Glosses for *bank*:
 1. financial institution that accepts deposits and channels the money into lending activities
 2. sloping land (especially the slope beside a body of water)

Glosses

- Glosses are not a formal meaning representation
- Examples from the definitions for *right*, *left*, *red*, *blood*

right *adj.* located nearer the right hand esp. being on the right when facing the same direction as the observer.

left *adj.* located nearer to this side of the body than the right.

red *n.* the color of blood or a ruby.

blood *n.* the red liquid that circulates in the heart, arteries and veins of animals.

- Circularity in definitions (e.g. *right*)
red and *blood* reference each other
- Useful for humans → sufficient grasp of the other terms
- Still useful for computational modeling of senses
 - use sentence to build embeddings
 - make use of *sense relations*

Sense Relations

right *adj.* located nearer the right hand esp. being on the right when facing the same direction as the observer.

left *adj.* located nearer to this side of the body than the right.

red *n.* the color of blood or a ruby.

blood *n.* the red liquid that circulates in the heart, arteries and veins of animals.

- *left – right*: similar words that stand in alternation or contrast
 - *red*: color
 - *blood*: liquid
-
- Sense relations: “*is-a*” relation or antonymy listed in databases like WordNet
 - Given a large database of relations → allows model to perform semantic tasks

How Many Senses Do Words Have?

- Dictionaries and thesauri: discrete lists of senses
Embeddings (static/contextual): high dimensional model of meaning, does not divide up into discrete senses
- How to determine when differing uses of a word should be represented with different sense?
- Example: *serve* from WSJ corpus
 - (1) They rarely *serve* red meat, preferring to prepare seafood.
 - (2) He *served* as U.S. ambassador to Norway in 1976 and 1977.
 - (3) He might have *served* his time, come out and led an upstanding life.
- *serve red meat* and *serve time*: different truth conditions and presuppositions
serve as ambassador: distinct subcategorization structure *serve as NP*

⇒ 3 distinct senses of *serve*

Outline

Word Senses

Relations between Senses

WordNet

Word Sense Disambiguation

Credits and References

Relations: Synonymy

- General idea of synonyms: when two senses of two different words/lemmas are (nearly) identical
 - *couch* – *sofa*
 - *filbert* – *hazelnut*
 - *car* – *automobile*
- Synonymy: relationship between senses rather than words
- Example: *big* – *large*
 - How big is that plane?
Would I be flying on a large or small plane?
 - Miss Nelson, for instance, became a kind of big sister to Benjamin.
Miss Nelson, for instance, became a kind of large sister to Benjamin. (?)
- *big* has a sense of *older/grown up*, while *large* lacks this sense

Relations: Antonymy

- Antonyms: words with an opposite meaning
 - *long – short*
 - *big – little*
 - *fast – slow*
 - *dark – light*
 - *rise – fall*
- Binary opposition, or at opposite ends of a scale
- Reversives: change or movement into in opposite directions
- Antonyms differ with respect to one aspect of their meaning, but otherwise are very similar.

Taxonomic Relations

- A word is a **hyponym** of another word if the first sense is more specific
 - *car* is a hyponym of *vehicle*
 - *dog* is a hyponym of *animal*
 - *mango* is a hyponym of *fruit*
- A **hypernym** is a more generic word
 - *vehicle* is a hypernym of *car*
 - *animal* is a hypernym of *dog*
 - *fruit* is a hypernym of *mango*
- **Meronymy**: a part–whole relation
 - a *leg* is a part of a *chair*
 - a *wheel* is a part of a *car*
- Also: *car* is a **holonym** of *wheel*

Outline

Word Senses

Relations between Senses

WordNet

Word Sense Disambiguation

Credits and References

WordNet: A Database of Lexical Relations

- Most commonly used resource for sense relations in English Fellbaum (1998)
- English WordNet consists of three databases:
nouns, verbs, adjectives/adverbs
- Each database contains a set of lemmas annotated with a set of senses
 - 117,798 nouns (on average: 1.23 senses)
 - 11,529 verbs (on average: 2.16 senses)
 - 22,479 adjectives
 - 4,481 adverbs

WordNet: Example

The noun “bass” has 8 senses in WordNet.

1. bass¹ - (the lowest part of the musical range)
2. bass², bass part¹ - (the lowest part in polyphonic music)
3. bass³, basso¹ - (an adult male singer with the lowest voice)
4. sea bass¹, bass⁴ - (the lean flesh of a saltwater fish of the family Serranidae)
5. freshwater bass¹, bass⁵ - (any of various North American freshwater fish with lean flesh (especially of the genus Micropterus))
6. bass⁶, bass voice¹, basso² - (the lowest adult male singing voice)
7. bass⁷ - (the member with the lowest range of a family of musical instruments)
8. bass⁸ - (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

Figure G.1 A portion of the WordNet 3.0 entry for the noun *bass*.

- Set of near synonyms: “synset”
 - {bass¹, deep⁶}
 - {bass⁶, bass voice¹, basso²}
- Lists of word senses that can be used to express the concept
- Glosses are properties of a synset
 - each sense included in the synset has the same gloss

WordNet: Synset Example

- {chump¹, fool², gull¹, mark⁹, patsy¹, fall guy¹, sucker¹, soft touch¹, mug²}
- **Gloss:** a person who is gullible and easy to take advantage of.
- Synsets are the fundamental unit associated with WordNet entries
- Synsets are labeled with a lexicographic category

WordNet: Synsets

- Categories are called **supersenses**: coarse semantic categories or groupings of senses

Category	Example	Category	Example	Category	Example
ACT	<i>service</i>	GROUP	<i>place</i>	PLANT	<i>tree</i>
ANIMAL	<i>dog</i>	LOCATION	<i>area</i>	POSSESSION	<i>price</i>
ARTIFACT	<i>car</i>	MOTIVE	<i>reason</i>	PROCESS	<i>process</i>
ATTRIBUTE	<i>quality</i>	NATURAL EVENT	<i>experience</i>	QUANTITY	<i>amount</i>
BODY	<i>hair</i>	NATURAL OBJECT	<i>flower</i>	RELATION	<i>portion</i>
COGNITION	<i>way</i>	OTHER	<i>stuff</i>	SHAPE	<i>square</i>
COMMUNICATION	<i>review</i>	PERSON	<i>people</i>	STATE	<i>pain</i>
FEELING	<i>discomfort</i>	PHENOMENON	<i>result</i>	SUBSTANCE	<i>oil</i>
FOOD	<i>food</i>			TIME	<i>day</i>

Figure G.2 Supersenses: 26 lexicographic categories for nouns in WordNet.

- Nouns: 26 supersenses
Verbs: 15 supersenses

Sense Relations in WordNet

Relation	Also Called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> ¹ → <i>meal</i> ¹
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> ¹ → <i>lunch</i> ¹
Instance Hypernym	Instance	From instances to their concepts	<i>Austen</i> ¹ → <i>author</i> ¹
Instance Hyponym	Has-Instance	From concepts to their instances	<i>composer</i> ¹ → <i>Bach</i> ¹
Part Meronym	Has-Part	From wholes to parts	<i>table</i> ² → <i>leg</i> ³
Part Holonym	Part-Of	From parts to wholes	<i>course</i> ⁷ → <i>meal</i> ¹
Antonym		Semantic opposition between lemmas	<i>leader</i> ¹ ↔ <i>follower</i> ¹
Derivation		Lemmas w/same morphological root	<i>destruction</i> ¹ ↔ <i>destroy</i> ¹

Figure G.3 Some of the noun relations in WordNet.

Relation	Definition	Example
Hypernym	From events to superordinate events	<i>fly</i> ⁹ → <i>travel</i> ⁵
Troponym	From events to subordinate event	<i>walk</i> ¹ → <i>stroll</i> ¹
Entails	From verbs (events) to the verbs (events) they entail	<i>snore</i> ¹ → <i>sleep</i> ¹
Antonym	Semantic opposition between lemmas	<i>increase</i> ¹ ↔ <i>decrease</i> ¹

Figure G.4 Some verb relations in WordNet.

- WordNet represents hyponymy by relating each synset to its immediately more general and more specific synsets
- Produce longer chains of more general or more specific synsets

Sense Relations in WordNet

```
bass3, basso (an adult male singer with the lowest voice)
=> singer, vocalist, vocalizer, vocaliser
    => musician, instrumentalist, player
        => performer, performing artist
            => entertainer
                => person, individual, someone...
                    => organism, being
                        => living thing, animate thing,
                            => whole, unit
                                => object, physical object
                                    => physical entity
                                        => entity

bass7 (member with the lowest range of a family of instruments)
=> musical instrument, instrument
    => device
        => instrumentality, instrumentation
            => artifact, artefact
                => whole, unit
                    => object, physical object
                        => physical entity
                            => entity
```

Figure G.5 Hyponymy chains for two separate senses of the lemma *bass*. Note that the chains are completely distinct, only converging at the very abstract level *whole, unit*.

Sense Relations in WordNet

- WordNet has two kinds of taxonomic entities: classes and instances
 - instance: individual, a proper noun that is a unique entity
(*San Francisco* is an instance of *city*)
 - *city* is a class, a hyponym of *municipality* and eventually of *location*

Sense Relations in WordNet

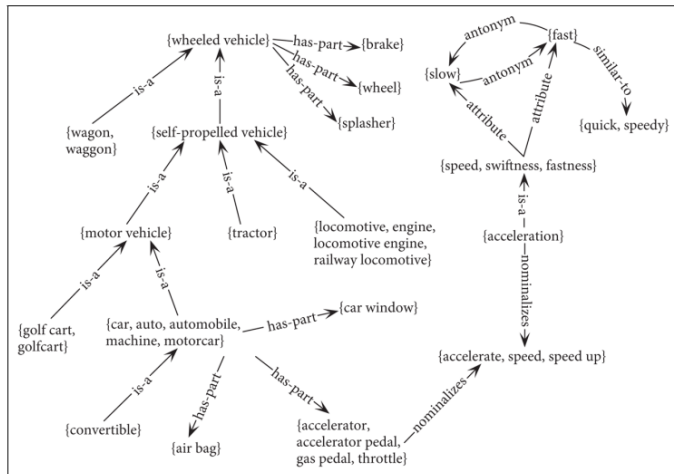


Figure G.6 WordNet viewed as a graph. Figure from Navigli (2016).

Outline

Word Senses

Relations between Senses

WordNet

Word Sense Disambiguation

Credits and References

Word Sense Disambiguation

- **Word sense disambiguation (WSD)**: the task of selecting the correct sense for a word
- Inventory of sense tags depends on task

WordNet Sense	Spanish Translation	WordNet Supersense	Target Word in Context
bass ⁴	lubina	FOOD	... fish as Pacific salmon and striped bass and...
bass ⁷	bajo	ARTIFACT	... play bass because he doesn't have to solo...

Figure G.7 Some possible sense tag inventories for *bass*.

Word Sense Disambiguation

- WSD task:
given an entire texts and a lexicon with an inventory of senses for each entry → disambiguate every word in the text
- Supervised WSD tasks are typically trained on sense-annotated corpora (semantic concordance)
- For example the *SemCor corpus*: subset of the Brown Corpus (226,036 words manually tagged with WordNet senses)

You will find_v⁹ that avocado_n¹ is_v¹ unlike_j¹ other_j¹ fruit_n¹
you have ever_r¹ tasted_v²

- For *fruit*, choose the correct sense:
 - fruit_n¹: the ripened reproductive body of a seed plant
 - fruit_n²: yield; an amount of a product
 - fruit_n³: the consequence of some effort or action

Word Sense Disambiguation

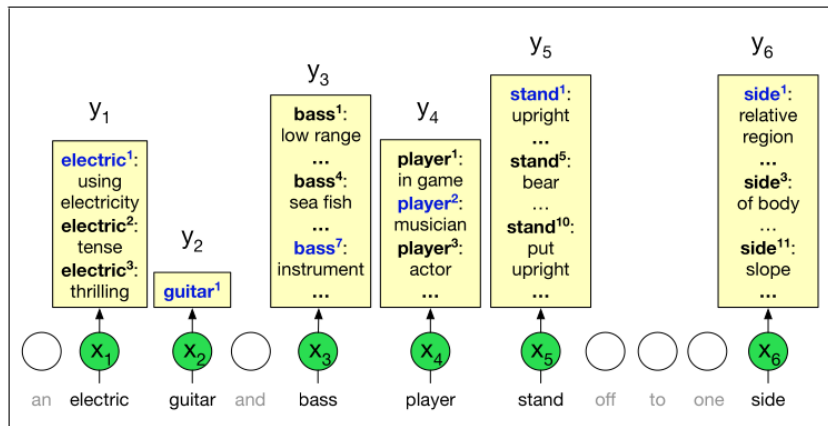


Figure G.8 The all-words WSD task, mapping from input words (x) to WordNet senses (y). Only nouns, verbs, adjectives, and adverbs are mapped, and note that some words (like *guitar* in the example) only have one sense in WordNet. Figure inspired by [Chaplot and Salakhutdinov \(2018\)](#).

Word Sense Disambiguation: Simple Strategies

- **Most frequent sense:** surprisingly strong baseline
- Senses in WordNet are generally ordered from most frequent to least frequent based on their counts
- Quite accurate, and is therefore often used as a default

- **One sense per discourse:** a word appearing multiple times in a text often appears with the same sense
- Better for coarse-grained senses and particularly when a word's senses are unrelated

Word Sense Disambiguation: Contextual Embeddings

- Simple 1-nearest-neighbor algorithm using contextual word embeddings
- Training time:
 - pass each sentence in SemCore labeled dataset through contextual embedding (e.g., BERT)
 - Produce a contextual sense embedding \mathbf{v}_s for s :

$$\mathbf{v}_s = \frac{1}{n} \sum_i \mathbf{v}_i \quad \forall \mathbf{v}_i \in \text{token}(s)$$

for each of the n tokens of sense s : average n contextual representations

- Test time:
 - given target word t in context: compute its contextual embedding \mathbf{t}
 - choose its nearest neighbor sense from the training set

$$\text{sense}(t) = \underset{s \in \text{senses}(t)}{\operatorname{argmax}} \operatorname{cosine}(\mathbf{t}, \mathbf{v}_s)$$

Word Sense Disambiguation: Nearest Neighbor

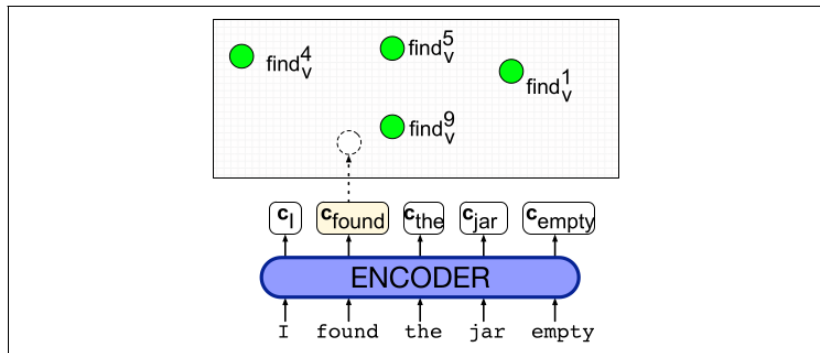


Figure G.9 The nearest-neighbor algorithm for WSD. In green are the contextual embeddings precomputed for each sense of each word; here we just show a few of the senses for *find*. A contextual embedding is computed for the target word *found*, and then the nearest neighbor sense (in this case $find_v^9$) is chosen. Figure inspired by Loureiro and Jorge (2019).

The Simplified Lesk Algorithm as WSD Baseline

- Generating sense labeled corpora like SemCor: difficult and expensive
- Knowledge-based algorithms rely solely on WordNet or and don't require labeled data
 - supervised algorithms typically better
- Lesk algorithm: choose the sense whose dictionary gloss or definition shares the most words with the target word's neighborhood

The Simplified Lesk Algorithm as WSD Baseline

- Disambiguate *bank* in the following context

The bank can guarantee deposits will eventually cover future tuition costs because it invests in adjustable-rate mortgage securities.

<i>bank</i> ¹	Gloss:	a financial institution that accepts deposits and channels the money into lending activities
	Examples:	“he cashed a check at the bank”, “that bank holds the mortgage on my home”
<i>bank</i> ²	Gloss:	sloping land (especially the slope beside a body of water)
	Examples:	“they pulled the canoe up on the bank”, “he sat on the bank of the river and watched the currents”

- sense *bank*¹ has two overlaps with the context: *deposits*, *mortgage*
- sense *bank*²: no overlap
- Variations:
 - weighting overlapping words by inverse document frequency (IDF)
 - use word embedding cosine instead of word overlap

Word-in-Context Evaluation

- Context-free word similarity task: how similar is *cup* to *mug*?
- WSD as contextualized similarity task: distinguish the meaning of a word in one context from another context
- Word-in-Context task:
two sentences with the same target word → used in the same sense?

F There's a lot of trash on the **bed** of the river —
I keep a glass of water next to my **bed** when I sleep

F **Justify** the margins — The end **justifies** the means

T **Air** pollution — Open a window and let in some **air**

T The expanded **window** will give us time to catch the thieves —
You have a two-hour **window** of clear weather to finish working on the lawn

Figure G.11 Positive (T) and negative (F) pairs from the WiC dataset (Pilehvar and Camacho-Collados, 2019).

Wikipedia as Source of Training Data

- Wikipedia as a source of sense-labeled data
- Articles in Wikipedia contain explicit links to concepts → link as sense annotation

In 1834, Sumner was admitted to the **[[bar (law)|bar]]** at the age of twenty-three, and entered private practice in Boston.

It is danced in 3/4 time (like most waltzes), with the couple turning approx. 180 degrees every **[[bar (music)|bar]]**.

Jenga is a popular beer in the **[[bar (establishment)|bar]]**s of Thailand.

- Add senses to training data of a supervised system
- Map Wikipedia concepts to relevant sense inventory (e.g. WordNet)
 - find WordNet sense with lexical overlap with the Wikipedia concept
 - vector of words in WordNet synset, gloss and related senses with
 - vector of words in Wikipedia page title, outgoing links and page category

Thesauri to Improve Embeddings

- Thesauri have been used to improve word embeddings
- Static word embeddings have a problem with antonyms
 - *expensive* is often similar in embedding cosine to *cheap*
- Two strategies to include information from thesauri
- Retraining:
modify the embedding training to incorporate thesaurus relations
- Retro-fitting or counterfitting:
after embedding training, learn a second mapping using thesaurus information: move synonyms closer and antonyms further apart

	Before counterfitting			After counterfitting		
east	west	north	south	eastward	eastern	easterly
expensive	pricey	cheaper	costly	costly	pricy	overpriced
British	American	Australian	Britain	Brits	London	BBC

Figure G.12 The nearest neighbors in GloVe to *east*, *expensive*, and *British* include antonyms like *west*. The right side showing the improvement in GloVe nearest neighbors after the counterfitting method (Mrkšić et al., 2016).

Outline

Word Senses

Relations between Senses

WordNet

Word Sense Disambiguation

Credits and References

Content based on:

- Dan Jurafsky and James H. Martin (2024)
Speech and Language Processing: Appendix Chapter G
<https://web.stanford.edu/~jurafsky/slp3/>

References

- Fellbaum, C., editor. 1998.
WordNet: An Electronic Lexical Database. MIT Press.