

Natural Language Processing

Unit 3

Semantic Analysis

Vector Semantics - Words and Vectors - Cosine Similarity - TF-IDF -
Positive PN1 - Word2Vec - Semantic properties of embeddings ;
Lexical Semantics : word senses - relations between senses - WordNet -
Word Sense Disambiguation

* Vector Semantics

- refers to the study of the vector representation of words - and models used to represent the meaning of words
 - Refers to representation of a word as some point in some multi-dimensional space
 - Each word is represented as a vector of numerical feature values.
- These features are called word embeddings. (embed a word into a vectorspace)

* Motivation behind vector representations

- words that occur in similar contexts have similar meanings.
- The meaning of a word is related to the distribution of words around it

fast ≈ rapid

tall ≈ similar

* Word Similarity

- can refer to relations that exist beyond synonyms / antonyms
- Word similarity: words with similar meanings but not synonyms
e.g. cat, dog

→ Words are related by a semantic frame or field

e.g. cat, dog = similar

student, teach = related, but not similar

Relatedness refers to co-participation in a shared event

* Semantic Field

→ set of words covering a particular semantic domain and having structured relations between them

e.g. University

→ teacher, student, project, class, assignment, professor

House

→ room, door, furniture, bedroom

* Generating Embeddings for Vector Semantics

→ The meaning of a word is computed from the distribution of words around it.

→ The word vectors are an array of numbers, related in some way to counts.

→ Vector semantics utilizes:

(i) distributionalist intuition

(ii) vector intuition

→ The meaning of word can be defined simply by how often it occurs near other words.

(3)

→ This method results in very long vectors that are sparse. This can be changed into short, dense vectors that have useful semantic properties.

* Vectors and Documents

- Vector or distributional models of meaning are based on a co-occurrence matrix - it shows how often words co-occur.
- ↳ term-term matrix
- In a term-document matrix, each row represents a word in the vocabulary and each column represents a document.

Term - Document Representation

- An example of a term-document matrix is as follows:

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	86	62	89
fool	36	58	1	4
wit	20	15	2	3

- Each cell in the matrix represents the number of times a particular word occurs in a particular document.
- The vector for a document identifies a point in a $|V|$ -dimensional space.
- The documents are points in a 4-dimensional space.
- ↳ Here it is a
- This is the representation of documents as vectors in vector space

each vector = column

- The term-document representation can also represent the meaning of words, by associating each word with a vector.
- The word ~~vec~~ matrix is a row vector

e.g.

$$\text{foot} = [\begin{array}{cccc} 36 & 58 & 1 & 4 \end{array}]$$

- Similar words have similar vectors because they tend to occur in similar documents.

Term-Term Representation

- also called word-word matrix or term-context matrix.
- The rows & columns of the term-term matrix are labeled by words rather than documents.
- This matrix has a $|V| \times |V|$ dimensionality.
- Each cell records the number of times the row word and the column word co-occur in some context.
- A window can be created around the word, say 4 words to the left and 4 words to the right.
- The cell would thus represent the number of times the column word occurs in a ± 4 word window around the row word.

e.g.

	aardvark	...	computer	data	pinch	result	sugar
apricot	0		0	0	1	0	1
pineapple	0		0	0	1	0	1
digital	0		2	1	0	1	0
information	0		1	6	0	4	0

→ This co-occurrence matrix shows that words like digital & information are more similar to each other.

* Measuring Similarity between Words: Cosine Similarity

→ To find the similarity between 2 word vectors v and w , the cosine of the angle between the vectors can be computed.

→ The cosine similarity between vectors \vec{v} and \vec{w} can be computed as:

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\sum_{i=1}^n v_i w_i}{\sqrt{\sum_{i=1}^n v_i^2} \sqrt{\sum_{i=1}^n w_i^2}}$$

Q. Use cosine similarity to find which of the words apricot or digital is closer in meaning to information.

	large	data	computer
apricot	2	0	0
digital	0	1	2
information	1	6	1

$$\cos(\text{apricot}, \text{information}) = \frac{2+0+0}{\sqrt{4+0+0} \sqrt{1+36+1}} = \frac{2}{\sqrt{38}} = 0.16$$

$$\cos(\overset{\text{digital}}{\cancel{\text{apricot}}}, \text{information}) = \frac{0+6+2}{\sqrt{4+0+0} \sqrt{1+36+1}} = \frac{8}{\sqrt{38}} = 0.58$$

→ Digital is closer in meaning to information

Q9 Find out the cosine similarity between the words (cherry, information) and (digital, information)

	pie	data	computer
cherry	442	8	2
digital	5	1683	1676
information	5	3982	3325

$$\cos(\text{cherry, information}) = \frac{(442 * 5) + (8 * 3982) + (2 * 3325)}{\sqrt{442^2 + 8^2 + 2^2} \sqrt{5^2 + 1683^2 + 1676^2}}$$

$$= \underline{\underline{0.18}}$$

$$\cos(\text{digital, information}) = \frac{(5 * 5) + (1683 * 3982) - (1676 * 3325)}{\sqrt{5^2 + 1683^2 + 1676^2} \sqrt{5^2 + 3982^2 + 3325^2}}$$

$$= \underline{\underline{0.996}}$$

\Rightarrow Information is much closer to digital than to cherry

* Term-Frequency - Inverse Document Frequency (TF-IDF)

Disadvantages of using Frequency for word Association

- Simple frequency is not the best measure of association between words.
- Words that occur nearby frequently are more important than words that appear only once or twice.
- On the other hand, words that are too frequent such as 'the' or 'a' are important.

→ These 2 conflicting constraints must be balanced.

* Term Frequency

→ Refers to the frequency of the word in the document

→ Frequency is weighted down by using its log value.

(This is because a word appearing 100 times in a document does not make that word 100 times more likely to be relevant to the meaning of the document)

$$tf_{t,d} = \log_{10} (\text{count}(t,d) + 1)$$

* Document Frequency

The document frequency df_t of a term t is the number of documents it occurs in

Inverse Document Frequency

$$idf_t = \log_{10} \left(\frac{N}{df_t} \right)$$

N = total no. of documents

* TF-IDF

$$w_{t,d} = tf_{t,d} \times idf_t$$

→ The TF-IDF vector model of a target word is a vector with dimensions corresponding to all the words in the vocabulary.

→ The values in each dimension are the frequency with which the target word co-occurs with each neighboring context word, weighted by tf-idf.

→ The model computes the similarity between two words $\times 84$
by taking the cosine of their tf-idf vectors.

high cosine \Rightarrow high similarity

Example Use TF-IDF model and rank the documents

Q : The cat .

D₁ : The cat is on the mat .

D₂ : My dog and cat are the best .

D₃ : The locals are playing .

$$\text{Ans} \quad \text{TF ("the", D)} = \frac{2}{6} = \frac{1}{3} = 0.33$$

$\boxed{\begin{aligned} \text{TF(word, doc)} &= \frac{\text{no. of occurrences}}{\text{no. of words in the document}} \\ &\text{of the word} \end{aligned}}$

$$\text{TF ("the", D}_2\text{)} = \frac{1}{7} = 0.14$$

$$\text{TF ("the", D}_3\text{)} = \frac{1}{4} = 0.25$$

$$\text{TF ("cat", D}_1\text{)} = \frac{1}{6} = 0.17$$

$$\text{TF ("cat", D}_2\text{)} = \frac{1}{7} = 0.14$$

$$\text{TF ("cat", D}_3\text{)} = \frac{0}{4} = 0$$

$$\text{IDF} = \log \left(\frac{\text{no. of docs}}{\text{no. of docs w/ that word}} \right)$$

$$\text{IDF ("the")} = \log \left(\frac{3}{3} \right) = \log(1) = 0$$

$$\text{IDF ("cat")} = \log \left(\frac{3}{2} \right) = 0.18$$

TF-IDF (the, D₁)

$$= 0.33 \times 0$$

$$= 0$$

TF-IDF (word, doc)

(9)

$$= \text{TF}(\text{word}, \text{doc}) \times$$

$$\text{IDF}(\text{word})$$

$$\text{TF-IDF}(\text{the}, \text{D}_2) = 0$$

$$\text{TF-IDF}(\text{the}, \text{D}_3) = 0$$

$$\text{TF-IDF}(\text{cat}, \text{D}_1) = 0.17 \times 0.18 = 0.0306$$

$$\text{TF-IDF}(\text{cat}, \text{D}_2) = 0.14 \times 0.18 = 0.0252$$

$$\text{TF-IDF}(\text{cat}, \text{D}_3) = 0$$

$$\text{TF-IDF score for D}_1 = 0 + 0.0306 = 0.0306$$

$$\text{TF-IDF score for D}_2 = 0 + 0.0252 = 0.0252$$

$$\text{TF-IDF score for D}_3 = 0$$

* Pointwise Mutual Information (PMI)

→ An alternative weighting function to tf-idf is called PMI.

→ PMI is a measure of how often two events x and y occur compared with what we expect if they were independent.

$$\text{PMI}(x, y) = \log_2 \frac{P(x, y)}{P(x) P(y)}$$

→ PMI ranges from $-\infty$ to $+\infty$, but the negative values may be problematic.

(i) Things co-occur less than expected by chance

(ii) unreliable without an enormous corpora

Solution - Replace the negative PMI values by 0, thus the

positive PMI between 2 words are:

$$\text{PPMI}(\text{word}_1, \text{word}_2) = \max \left(\log_2 \frac{P(\text{word}_1, \text{word}_2)}{P(\text{word}_1) P(\text{word}_2)}, 0 \right)$$

* Computing PPMI on a Term-Context Matrix

Given a co-occurrence matrix with W rows (target word w_i) and C columns (contexts), and f_{ij} gives the number of times word w_i occurs in context c_j

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$p_{j*} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{pmi}_{ij} = \log_2 \frac{p_{ij}}{p_{i*} p_{j*}}$$

$$\text{ppmi}_{ij} = \begin{cases} \text{pmi}_{ij} & , \text{pmi}_{ij} > 0 \\ 0 & , \text{otherwise} \end{cases}$$

Example Compute the PPMI of information, data from the following corpus,

	computer	data	result	pie	sqai	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

Ans $p_{ij} = \frac{f_{ij}}{\sum_{i=1}^w \sum_{j=1}^c f_{ij}}$ (no. of times the word w_i occurs in the context c_j)

$$p_{ij} (w=\text{information}, c=\text{data}) = \frac{3982}{11716} = 0.3399$$

$$p_{i*} = P(w=\text{information}) = \frac{7703}{11716} = 0.6575$$

$$p_{*j} = P(c=\text{data}) = \frac{5673}{11716} = 0.4842$$

$$\text{PPMI}(\text{information}, \text{data}) = \log_2 \left(\frac{\frac{0.3399}{0.6575 \times 0.4842}}{} \right) \\ = \underline{\underline{0.0944}}$$

Example Consider a small corpus consisting of three documents:

D₁: The quick brown fox jumps over the lazy dog

D₂: The dog barked at the fox, but the fox was too quick.

D₃: The lazy dog sleeps all day.

Calculate the TF-IDF for the term 'fox' in each document.

Ans $TF(\text{the}, D_1) = \frac{2}{9} = 0.22$

$TF(\text{the}, D_2) =$

$TF(\text{Fox}, D_1) = \frac{1}{9} = 0.11$

$TF(\text{Fox}, D_2) = \frac{2}{12} = \frac{1}{6} = 0.17$

$TF(\text{Fox}, D_3) = 0$

$IDF(\text{the}) = \log\left(\frac{3}{3}\right) = 1 = 0$

$IDF(\text{Fox}) = \log\left(\frac{3}{2}\right) \approx 0.176$

$TF-IDF(\text{Fox}, D_1) = 0.11 \times 0.18 = 0.0198$

$TF-IDF(\text{Fox}, D_2) = 0.17 \times 0.18 = 0.0306$

$TF-IDF(\text{Fox}, D_3) = 0$

* Short and Dense Vector Representations

→ short vectors : 50 - 500

dense : most values are non-zero

short

→ Dense vectors work better in every NLP task than sparse vectors

→ This may be because of the following reasons

- ① Short vectors are easier to use as features in NL models, as there would be fewer weights to tune.

(2) Dense vectors may generalize better than storing explicit counts

(13)

(3) Capture synonymy better

(i) car and automobile are synonyms, but are distinct dimensions in sparse vectors.

(ii) a word w/ car as a neighbor and a word with automobile as a neighbor should be similar, but they aren't.

* Word2Vec

→ Word2Vec is a popular embedding method, that is based on the intuition that instead of counting how often each word occurs near another word, a classifier is trained on a binary prediction task.

→ The learned classifier weights are taken as embeddings.

→ The running text is used as implicitly supervised training data. - there is no need for hand-labeled supervision.

* Word2Vec - Skipgram with negative sampling

1. Treat the target word and a neighboring context word as positive examples
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish the two cases
4. Use the weights as the embeddings.

* Training Word2Vec

Goal : Train the classifier such that a given tuple (t, c)

t = target word

c = context word

will return the probability that c is a real context word

e.g. $(\text{apricot}, \text{jam}) = \text{True}$

$(\text{apricot}, \text{aardvark}) = \text{False}$

→ This is denoted by $P(+ | t, c)$

→ The probability that the word c is not a real context word for

t is :

$$P(- | t, c) = 1 - P(+ | t, c)$$

Step1 : To compute the probability P , consider the similarity

similarity $(t, c) \approx t \cdot c$

→ The dot product $t \cdot c$ is not a probability, it is just a number ranging from 0-60

Step2 : Convert the dot product into a probability (the logistic or sigmoid function can be used)

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

→ The probability that word c is a real context word for the target word t is :

$$P(+ | t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

→ The probability that the word c is not a real context word is : (15)

$$P(-|t,c) = 1 - P(+|t,c)$$

$$= \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

→ $P(+|t,c)$ and $P(-|t,c)$ give us the probability for one word, but it is necessary to take account of the multiple context words in the window.

* Skip-gram Modelling

→ Skipgram makes the assumption that all context words are independent.

$$P(+|t, c_1:L) = \prod_{i=1}^L \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+|t, c_1:L) = \sum_{i=1}^L \log \frac{1}{1 + e^{-t \cdot c_i}}$$

→ The skipgram trains a probabilistic classifier that given a test target word t and a context window of L words $c_1:L$, a probability is assigned based on how similar the context window is similar to the target window.

→ Word2vec learns embeddings by starting w/ an initial set of embedding vectors and then iteratively shifting the embeddings of each word to be more like the embeddings of words that occur nearby in texts, and less like the embeddings of words that don't occur nearby.

* Skipgram Training

Step 1 : Creation of Positive and Negative Training Instances

Positive Training Instances

→ For a target word, if there is a $L=2$ window, then there would be 4 positive training instances

eq.	t	c
	apricot	tablespoon
	apricot	of
	apricot	preserves
	apricot	or.

Negative Training Instances

→ For training the binary classifier, negative examples are also needed and skip-gram uses more negative examples than positive examples

- For each (t,c) instance, create k negative samples, each consisting of the target and a noise word.
- A ~~random~~ noise word is a random word from the lexicon.
- When $k=2$, there would be 2 negative examples in the training set for each positive example (t,c).

eq.	t	c
	apricot	cardvark
	apricot	puddle
	apricot	where
	apricot	coaxial

	t	c
	apricot	twelve
	apricot	hello
	apricot	dear
	apricot	forever.

→ The noise words are chosen according to their unigram frequency $p(w)$

→ If the sampling is done with unweighted frequency $p(w)$, then by the unigram probability $p('the')$, 'the' would be chosen as the noise word.

→ In practice, a weighted unigram frequency $p_\alpha(w)$ is used, where α is a weight

$$p_\alpha(w) = \frac{\text{count}(w)^\alpha}{\sum_{w'} \text{count}(w')^\alpha}$$

usually $\alpha = 0.75$

Step 2 : Initialization and Optimization

→ Represent the words as vectors of some length, randomly initialized (say 300)

There is thus, $300 * V$ random parameters

→ Over the training set, adjust the word vectors such that:

(i) Maximize the similarity of the target word, context word pairs (t,c) from the positive data

(ii) Minimize the similarity of the (t,c) pairs drawn from the negative data.

i.e maximize

$$\sum_{(t,c) \in +} \log P(+ | t, c) + \sum_{(t,c) \in -} \log P(- | t, c)$$

→ Focusing on one word/context pair (t, c) with k noise words

n_1, n_2, \dots, n_k , the learning objective is:

$$L(\theta) = \log P(+1|t, c) + \sum_{i=1}^k \log P(-1|t, n_i)$$

$$= \log \sigma(c \cdot t) + \sum_{i=1}^k \log \sigma(-n_i \cdot t)$$

$$= \log \left(\frac{1}{1 + e^{-ct}} \right) + \sum_{i=1}^k \log \left(\frac{1}{1 + e^{n_it}} \right)$$

→ Use stochastic gradient descent to train to this objective, iteratively modifying the parameters (embeddings for t and c)

→ The skip-gram model thus has two separate embeddings for each word w : the target embedding t , and the context embedding c .

→ The ~~two~~ embeddings are stored in 2 matrices: target matrix W and context matrix C .

Target Matrix = $\begin{bmatrix} & \\ & \\ & \end{bmatrix}$ → Each row is a $1 \times d$ vector $\not\in$ embedding t_i

Context Matrix $\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$ → each column is a $d \times 1$ vector embedding c_i for word $i \in V$

* Visualizing Embeddings

- The simplest way to visualize the meaning of a word w 's embeddings is to list the most similar words to w sorting all the words in the vocabulary by their cosines.
- Another method is to use a clustering algorithm to show a hierarchical representation of which words are similar to others in the embedding space.

* Semantic Properties of Embeddings

Similarity - similarity depends on the window size C

e.g. $C=5$ The nearest words to Hogwarts:

- Dumbledore
- MarfoY
- halfblood

Co-occurrence

First Order Co-occurrence

→ If 2 words are typically nearby one another

e.g. wrote is a first-order associate of book or poem.

Second-Order Co-occurrence

→ Two words have second-order cooccurrence if they have similar neighbors

e.g. wrote is a second-order associate of said / remarked

Relational meanings - Embeddings have the ability to capture relational meanings.

Ex: By projecting vectors onto two dimensions:

'king' - 'man' are close

'queen' - 'woman' are close

→ Similarity offsets in the 2D projections can capture comparative and superlative morphology.

* Lexical Semantics

→ Lexical semantics is the study of meaning-related connections among lexemes, and their internal meaning-related structure of individual lexemes

→ Lexeme = individual entry in the lexicon

Lexicon = finite list of lexemes

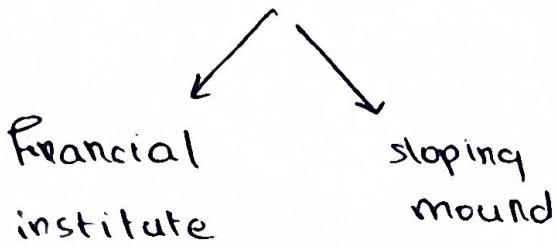
→ Semantics of individual lexemes can be captured by analyzing and labelling their relations to other lexemes

* Word Senses

A: Homonymy

→ Relation that holds between words that have the same form with unrelated meanings

e.g. bank has 2 senses



B. | Homophone

(21)

→ words with the same pronunciation but different spellings

eg. would / wood

sun / son

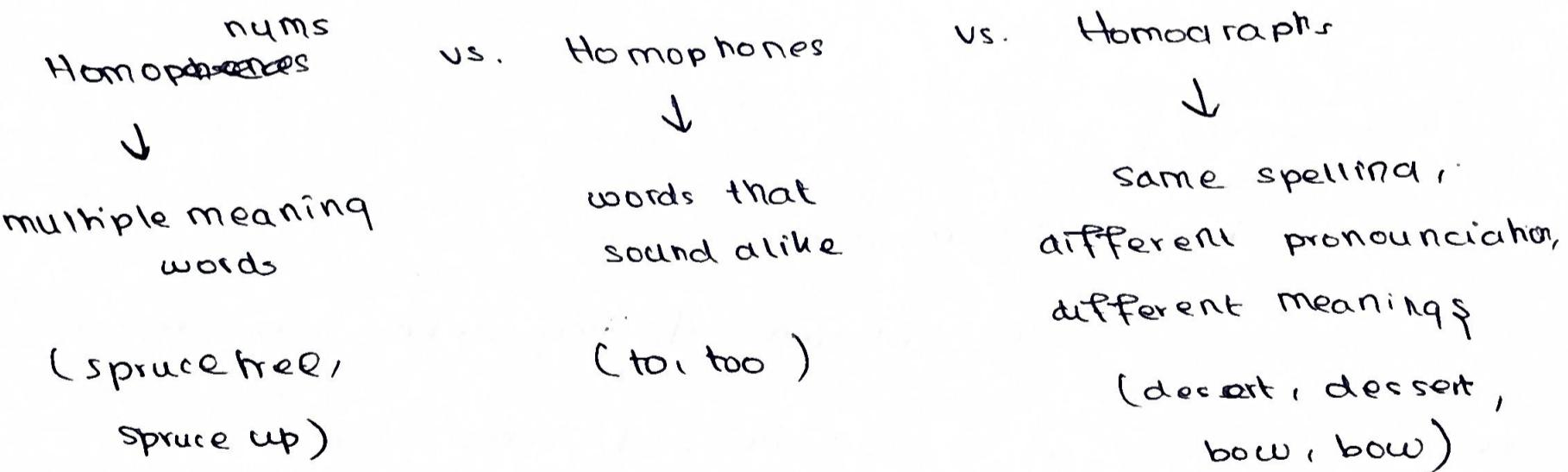
C. | Homographs

→ words with the same orthographic form but different pronunciations

eg. bass → fish
→ music

Text-to-speech & TIR systems are vulnerable to homographs.

Speech recognition, spelling correction systems are vulnerable to
homophones



D. | Polysemy?

→ Polysemy relates to multiple related meanings within a single lexeme

eg. bank = repository of entities

eg. blood bank, egg bank, eye bank

→ This sense of bank is derived from the financial institution sense.

+ Relations between Senses

A. Synonymy

→ Different lexicons with the same meaning

e.g. couch / sofa

vomit / throw up

car / automobile

→ Two words are synonymous if they are substitutable one for another in any sentence without changing the truth conditions of the sentence

→ However, no two words are exactly identical in meaning - e.g. big & large sometimes cannot be used interchangeably

→ In practice synonym is used to describe a relationship of rough synonymy.

B. Antonymy

→ words with opposite meaning

→ Two senses can be antonyms if they define a binary opposition or are at opposite ends of 'some scale'

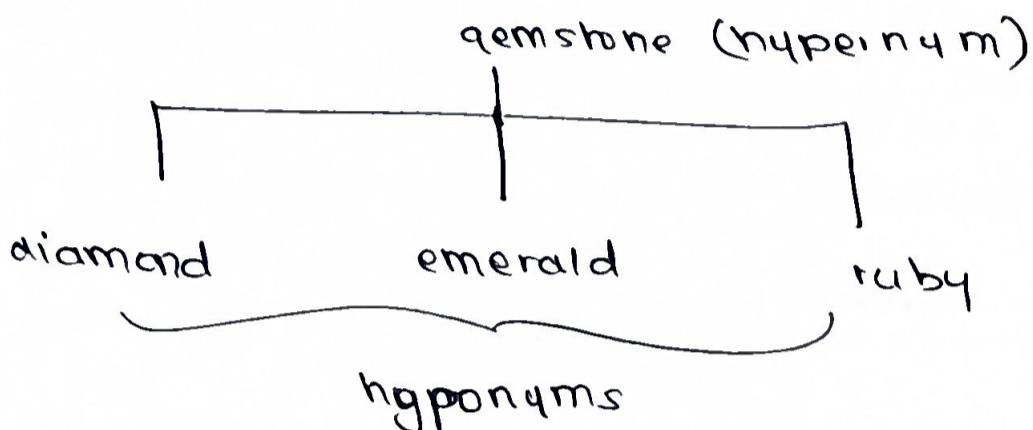
e.g. long / short

fast / slow

→ Though antonyms have very different meanings - they also share all aspects of their meaning except their position on a scale or their direction.

c. Hyponymy

- One lexeme denotes a subclass of the other
- e.g. car is a hyponym of vehicle
mango is a hyponym of fruit
- Conversely, vehicle is a hypernym of car.



* WordNet → a large electronic Lexical database, that can be used to find synonyms, antonyms and other semantic relationships between words

→ It is a hierarchically organized Lexical database.

* Word Sense Disambiguation

* Disambiguation

- Many words have several meanings or senses, and thus there is an ambiguity about how they are to be interpreted.
- The task of disambiguation is to determine which of the senses of an ambiguous word involved in a particular use of the word.
- This is done by looking at the context of the word's use.
- The task of WSD is to make a choice between senses of an ambiguous word, based on the context of use.

* Supervised Word Sense Disambiguation

- Extract features from the text and then train a classifier to assign the correct sense given these features.
- A feature vector consisting of numeric values is used to encode this linguistic information as an input to ML algorithms
- Two classes of features are:
 - (i) collocational features
 - (ii) bag-of-words features

A. Co-locational Features

- A collocation is a word or a phrase in a position-specific relationship to a target word
- Collocational features encode information about specific positions located to the left or right of the target word
e.g. exactly 1 word to the right, or exactly 4 words to the left.
- Typical features include the word itself, the root form of the word, and the word's pos.

For e.g. An electric guitar and bass player stand off to one side

The collocational feature-vector to disambiguate bass could be:

$[w_{i-2}, POS_{i-2}, POS_{i-1}, w_{i+1}, POS_i, w_{i+2}, POS_{i+3}]$
[guitar, NN, and, CC, player, NN, stand, VB]

B. Bag-of-words features

95

- A bag of words means an unordered set of words, ignoring their exact position.
- The bag of words approach aims to find the context of a target word by a vector of features. Each binary feature indicates whether a vocabulary word w does or does not occur in the context.
 - e.g. An electric guitar and bass player stand off to one side.
 - = [Fishing, big, sound, player, fly, rod, pound, double, runs, playing
guitar, band]
 - = [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0]

* Methods for WSD

- ① Supervised ML - Naive Bayes Approach
- ② Dictionary-based - & Lesk Algorithm
- ③ Bootstrapping - Yarowsky Algorithm

A. Naive Bayes Approach

- Choosing the best sense, for a feature vector f involves choosing the most probable sense given that vector.

$$\begin{aligned}\hat{s} &= \arg \max P(s|f) \\ &= \arg \max P(f|s) P(s) / P(f) \\ &= \arg \max P(f|s) P(s) \\ &= \arg \max P(s) \prod P(f_i|s)\end{aligned}$$

$P(f)$ = constant for all senses

$P(s)$ = prior probability of each sense
 $P(f_i|s)$ = individual feature probability

Ex1 Identify the word sense of the test data from the given training data.

Training

Doc	Words	Class
1	fish smoked fish	f
2	fish line	f
3	Fish haul smoked	f
4	guitar jazz line	g

Test

line guitar jazz jazz : ?

Ans

Prior Probabilities

$$P(f) = \frac{3}{4}$$

$$P(g) = \frac{1}{4}$$

$V = \{ \text{fish, smoked, line, haul, guitar, jazz} \}$

Conditional probabilities

$$P(\text{line}|f) = \frac{1+1}{8+6} = \frac{2}{14}$$

$$P(\text{guitar}|f) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{jazz}|f) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{line}|g) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{guitar}|g) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{jazz}|g) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

↓

no. of words

that are there in
the training set of

that class

Choosing a class

(27)

$$P(\text{f1ds}) = \frac{3}{14} \times \frac{2}{14} \times \frac{1}{14} \times \left(\frac{1}{14}\right)^2 \\ = 0.00003$$

$$P(\text{q1df}) = \frac{1}{14} \times \frac{2}{9} \times \left(\frac{2}{9}\right)^2 \times \frac{8}{9} \\ = 0.0006$$

The test sentence belongs to class q1.

B. Dictionary-based Approach - Lesk Algorithm

→ Choose the sense whose dictionary gloss or definition shares the most words with the target word's neighborhood.

Algorithm

Function LESK (word, sentence) returns best sense of word

best-sense ← most frequent sense for word

max-overlap ← 0

context ← set of words in sentence

for each sense in senses of word do

signature ← set of words in the gloss

overlap ← COMPUTE OVERLAP (signature, context)

if overlap > max-overlap then

max-overlap ← overlap

end ~~end~~ best-sense ← sense
return (best-sense)

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

bank¹ → financial institution - has more overlapping words than bank² (sloping land), so bank 1 is chosen.

Limitations of Resk Algorithm

- The dictionary entries for the target words are short, and may not provide enough chance of overlap with the context

Solution Apply a weight to each overlapping word. The weight is the inverse document frequency (IDF)

c. Bootstrapping - Yarowsky Algorithm

- Given a small seed set $\Lambda^{(0)}$ of labelled instances of each sense, and a much larger unlabelled corpus $V^{(0)}$
- The algorithm first trains an initial decision-list classifier on the seed set $\Lambda^{(0)}$
- Use this classifier to label the unlabeled corpus $V^{(0)}$
- Select the examples in $V^{(0)}$ that it is most confident about, removes them, and adds them to the training set - called it $\Lambda^{(1)}$
- Iterate by applying the classifier to the unlabeled set $V^{(1)}$ extracting a new training set $\Lambda^{(2)}$ and so on.

- With each iteration, the training corpus grows and the untagged corpus shrinks.
- The process is repeated until there is some sufficiently low error-rate on the training set, or until no further examples from the untagged corpus are above threshold.
- The key to bootstrapping - an accurate initial set of seeds. One way to generate the initial seeds is to hand-label a small set of examples or to use a heuristic to select accurate seeds.
- Yarowsky used the One Sense per Collocation heuristic : which relies on, certain words or phrases strongly associated with the target sense, do not occur with the other sense.