

Natural Language Understanding

CSE 415: Introduction to Artificial Intelligence
 University of Washington
 Winter, 2019

© S. Tanimoto and University of Washington, 2017
 [Some of these slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley.]

Problem: Ambiguities

How do we interpret these sentences? Syntactically, all correct, but there are two possible meanings.

- Headlines:

- Enraged Cow Injures Farmer With Ax
- Hospitals Are Sued By 7 Foot Doctors
- Ban on Nude Dancing on Governor's Desk
- Iraqi Head Seeks Arms
- Local HS Dropouts Cut in Half
- Juvenile Court to Try Shooting Defendant
- Stolen Painting Found by Tree
- Kids Make Nutritious Snacks



- Why are these funny?

Nat. Lang. Understanding/Processing

3

Parsing as Search



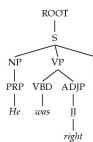
Nat. Lang. Understanding/Processing

4

Grammar: PCFGs

or the sum of the scores where score = $-\log(\text{prob})$

- Natural language grammars are very ambiguous!
- PCFGs are a formal probabilistic model of trees
 - Each "rule" has a conditional probability (like an HMM)
 - Tree's probability is the product of all rules used**
- Parsing: Given a sentence, find the best tree – search!



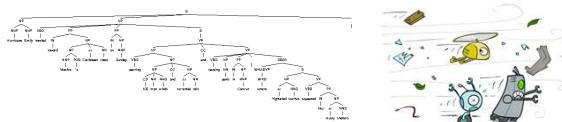
ROOT → S	375/420
S → NP VP.	320/392
NP → PRP	127/539
VP → VBD ADJP	32/401
.....	

We want most likely (highest prob lowest score) tree to use for our syntax tree.

Nat. Lang. Understanding/Processing

5

Syntactic Analysis



Hurricane Emily howled toward Mexico's Caribbean coast on Sunday packing 135 mph winds and torrential rain and causing panic in Cancun, where frightened tourists squeezed into musty shelters.

[Demo: Berkeley NLP Group Parser <http://tomato.banatao.berkeley.edu:8080/parser/parser.html>]

Nat. Lang. Understanding/Processing

6

Dialog Systems



Nat. Lang. Understanding/Processing

7

ELIZA



Nat. Lang. Understanding/Processing

8

- A “psychotherapist” agent (Weizenbaum, ~1964)
- Led to a long line of chatterbots
- How does it work:
 - Trivial NLP: string match and substitution
 - Trivial knowledge: tiny script / response database
 - Example: matching “I remember __” results in “Do you often think of __?”
- Can fool some people some of the time?

[Demo: <http://nlp-addiction.com/eliza>]

Watson



"a camel is a horse designed by a committee"

WikiAnswers

a camel is a horse designed by a committee

The Phrase Finder

A camel is a horse designed by committee

Discussion Forum

A camel is a horse designed by committee

Re: A camel is a horse designed by committee

Does anyone know the origin of this maxnet? I heard it way back at the United Nations, which is checklist of committees. It may have originated there, but I'd like an authoritative explanation. Thanks.

Re: A camel is a horse designed by committee

* Re: A camel is a horse designed by committee

If a camel is a horse designed by committee

If a camel is a horse designed by committee

Nat. Lang. Understanding/Processing

9

What's in Watson?



- A question-answering system (IBM, 2011)
- Designed for the game of Jeopardy
- How does it work:
 - Sophisticated NLP: deep analysis of questions, noisy matching of questions to potential answers
 - Lots of data: onboard storage contains a huge collection of documents (e.g. Wikipedia, etc.), exploits redundancy
 - Lots of computation: 90+ servers
 - Can beat all of the people all of the time?

Nat. Lang. Understanding/Processing

10

Machine Translation



Nat. Lang. Understanding/Processing

11

Machine Translation

automatic translation between languages.

"Il est impossible aux journalistes de rentrer dans les régions tibétaines"
Bruno Philip, correspondant du "Monde" à Pékin, a déclaré que les journalistes de l'AFP qui ont été évacués de la ville de Lhassa par le commandement du Qinghai "n'étaient pas dans l'illégalité".
Les faits Le dalaï-lama dénonce l'"ordre" imposé au Tibet depuis sa fuite en 1959 et la répression de la rébellion

"It is impossible for journalists to enter Tibetan areas"
Philip Bruno, correspondent for "Le Monde" in Beijing, said that journalists of the AFP who had been deported from the city of Lhasa by the command of the Qinghai "were not illegal".
Facts The Dalai Lama denounces the "order" imposed on Tibet since he fled Tibet in 1959
Video Anniversary of the Tibetan rebellion: China is great!

- Translate text from one language to another
- Recombines fragments of example translations
- Challenges:
 - What fragments? [learning to translate]
 - How to make efficient? [fast translation search]

Nat. Lang. Understanding/Processing

12

The Problem with Dictionary Lookups

There does not exist a one to one conversion of meaning for one word in one language to another. How do you handle the ambiguity?

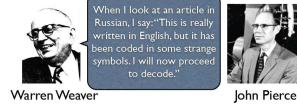
顶部	/top/roof/
顶端	/summit/peak/ top /apex/
顶头	/coming directly towards one/ top /end/
盖	/lid/ top /cover/canopy/build/Gai/
盖帽	/surpass/ top /
极	/extremely/pole/utmost/ top /collect/receive/
尖峰	/peak/ top /
面	/fade/side/surface/aspect/ top /face/flour/
擔心	/ top /topping/

Example from Douglas Hofstadter

Nat. Lang. Understanding/Processing

13

MT: 60 Years in 60 Seconds

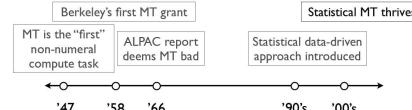


When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."

"Machine Translation" presumably means going by textable source text to useful target text... In this context, there has been no machine translation...

Warren Weaver

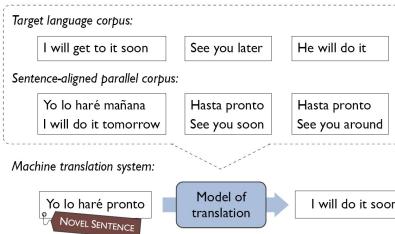
John Pierce



Nat. Lang. Understanding/Processing

14

Data-Driven Machine Translation



Nat. Lang. Understanding/Processing

15

Learning to Translate

Finding repeating words and try to find correspondance in their translation.

CLASSIC SOUPS

		Sm.	Lg.
大排湯	57. House Chicken Soup (Chicken, Celery)	1.50	2.75
雞蛋湯	Potato, Onion, Carrot)	1.85	3.25
雞肉湯	Chicken Rice Soup	1.85	3.25
香茅魚湯	Chicken & Rice (Miso)Soup	NA	2.75
冬瓜湯	Tomato Clear Egg Drop Soup	1.65	2.95
冬筍湯	Regular WontonSoup	1.10	2.10
酸辣湯	Hol & Sour Soup	1.10	2.10
蛋花湯	Egg Drop WontonSoup	1.10	2.10
文慶湯	Egg Drop WontonMix	1.10	2.10
酸辣湯	Tofu Vegetable Soup	NA	3.50
鴨血湯	Chicken Corn Cream Soup	NA	3.50
蟹味王	Crab Meat Corn Cream Soups	NA	3.50
海鮮湯	Seafood Soup	NA	3.50

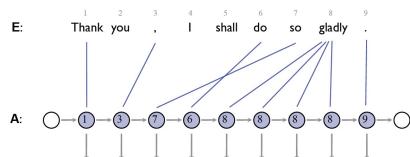
Example from Adam Lopez

Nat. Lang. Understanding/Processing

16

Hidden Markov Model An HMM Translation Model

True state = X in language to be translated into observations/evidence E = words in language to be translated.



Model Parameters

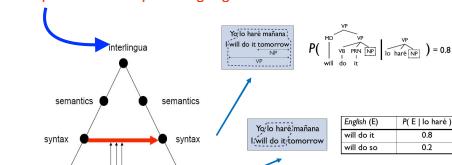
Emissions: $P(F_i = \text{Gracias} | E_{A_i} = \text{Thank})$ Transitions: $P(A_3 = 3 | A_1 = 1)$

Nat. Lang. Understanding/Processing

17

Levels of Transfer

independent concept of language

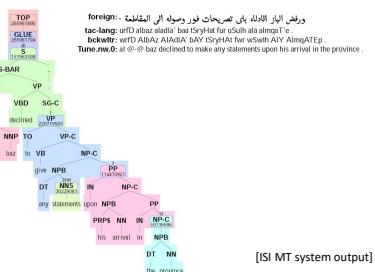


Try to translate at word level, phrase level, syntax level...

Nat. Lang. Understanding/Processing

18

Example: Syntactic MT Output



Nat. Lang. Understanding/Processing

[ISI MT system output]

Document Analysis with LSA: Outline

- Motivation
- Bag-of-words representation
- Stopword elimination, stemming, reference vocabulary
- Vector-space representation
- Document comparison with the cosine similarity measure
- Latent Semantic Analysis

Break down documents into its key words and phrases.

Nat. Lang. Understanding/Processing

22

Motivation

- Document analysis is a highly active area, very relevant to information science, the World Wide Web, and search engines.
- Algorithms for document analysis span a wide range of techniques, from string processing to large matrix computations.
- One application: automatic essay grading.



Nat. Lang. Understanding/Processing

23

Representations for Documents

- Text string
- Image (I.e., .jpg, .gif, and .png files)
- linguistically structured files: PostScript, Portable Doc. Format (PDF), XML.
- Vector: e.g., bag-of-words
- Hypertext, hypermedia



Nat. Lang. Understanding/Processing

24

Fundamental Problems

- Representation*
- Lexical Analysis (tokenizing)*
- Information Extraction*
- Comparison (similarity, distance)*
- Classification (e.g., for net-nanny service)*
- Indexing (to permit fast retrieval)
- Retrieval (querying and query processing)

*important for AI

Natural Language Understanding/Processing

25

Bag-of-Words Representation

A **multiset** is a collection like a set, but which allows **duplicates** (any number of copies) of elements.
 $\{a, b, c\}$ is a set. (It is also a multiset.)
 $\{a, a, b, c, c\}$ is not a set, but it is a multiset.
 $\{c, a, b, a, c, c\}$ is the same multiset. (Order does matter).
A multiset is also called a **bag**.



Natural Language Understanding/Processing

26

Bag-of-Words (continued)

Let document D =
"The big fox jumped over the big fence."
The bag representation is:
{ big, big, fence, fox, jumped, over, the, the }

For notational consistency, we use alphabetical order.
Also, we omit punctuation and normalize the case.

The ordering information in the document is lost. But
this is OK for some applications.

Natural Language Understanding/Processing

27

Eliminating Stopwords

In information retrieval and some other types of document analysis, we often begin by deleting words that don't carry much meaning or that are so common that they do little to distinguish one document from another. Such words are called **stopwords**.

Examples: (articles) a, an, the; (quantifiers) any, some, only, many, all, no; (pronouns) I, you, it, he, she, they, me, him, her, them, his, hers, their, which; (prepositions) above, at, behind, below, beside, for, in, into, of, on, onto, over, under; (verbs) am, are, be, can, could, do, does, did, going, go, have, has, had, do, did, can, can't, will, would, might, may, must; (conjunctions) and, but, if, nor, and, neither, nor, either, or; (other) yes, perhaps, first, last, there, where, when.



Natural Language Understanding/Processing

28

Stemming

In order to detect similarities among words, it often helps to perform stemming. We typically stem a word by removing its suffixes, leaving the basic word, or "uninflecting" the word

- apples → apple
- cacti → cactus
- swimming → swim
- swam → swim

regularizing our set of words!

Natural Language Understanding/Processing

29

Reference Vocabulary

A counterpart to stopwords is the **reference vocabulary**. These are the words that ARE allowed in document representations.

These are **all stemmed**, and are not stopwords. There might be several hundred or even thousands of terms in a reference vocabulary for real document processing.

Natural Language Understanding/Processing

30

Vector representation

As we use same reference vocabulary for all the documents we hope to compare, all our vector representations "line up" entry in one refers to same word as that same entry in vector for different document

Assume we have a reference vocabulary of words that might appear in our documents.

{apple, big, cat, dog, fence, fox, jumped, over, the, zoo}

We represent our bag

{ big, big, fence, fox, jumped, over, the, the }

by giving a vector (list) of occurrence count of each reference term in the document:

[0, 2, 0, 0, 1, 1, 1, 2, 0]

If there are n terms in the reference vocabulary, then each document is represented by a point in an n-dimensional space.

Natural Language Understanding/Processing

31

Indexing

Create links from terms to documents or document parts

- (a) concordance
- (b) table of contents
- (c) book index
- (d) index for a search engine
- (e) database index for a relation (table)

Natural Language Understanding/Processing

32

Concordance

A **concordance** for a document is a sort of dictionary that lists, for each word that occurs in the document the sentences or lines in which it occurs.

"document":
A concordance for a **document** is a sort of dictionary that lists, for each word that occurs in the **document** the

"occurs":
that lists, for each word that **occurs** in the document the sentences or lines in which it **occurs**.

Natural Language Understanding/Processing 33

Search Engine Index

Query terms are organized into a large table or tree that can be quickly searched.
(e.g., large hash-table in memory, or a B-Tree with its top levels in memory).

Associated with each term is a list of occurrences, typically consisting of Document IDs or URLs.

Natural Language Understanding/Processing 34

Document Comparison

i.e. using our vectors for each document, how much alike are they?
Geometrically, what's the angle between the two vectors?

Typical problems:

- Determine whether two documents are slightly different versions of the same document.
(applications: search engine hit filtering, plagiarism detection).
- Find the longest common subsequence for a pair of documents. (can be useful in genetic sequencing).
- Determine whether a new document should be placed into the same category as a model document. (essay grading, automatic response generation, etc.)

Natural Language Understanding/Processing 35

Cosine Similarity Function

Cosine as in the trig function; essentially we are using a dot product to determine the angle between the two documents' vectors

Document 1:
"All Blues. First the key to last night's notes."

Document 2:
"How to get your message across. Restate your key points first and last."

Reference vocabulary:
{ across, blue, first, key, last, message, night, note, point, restate, zebra }

Natural Language Understanding/Processing 36

Cosine Similarity (cont)

Document 1 reduced:
blue first key last night note

Document 2 reduced:
message across restate key point first last

Document 1 vector representation:
[0, 1, 1, 1, 0, 1, 1, 0, 0, 0]

Document 2 vector representation:
[1, 0, 1, 1, 1, 0, 0, 1, 1, 0]

Natural Language Understanding/Processing 37

Cosine Similarity (cont)

Dot product (same as "inner product")

$$[0, 1, 1, 1, 0, 1, 1, 0, 0, 0] \cdot [1, 0, 1, 1, 1, 0, 0, 1, 1, 0]$$

$$= 0 \cdot 1 + 1 \cdot 0 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 0 \cdot 1 + 1 \cdot 0 + 1 \cdot 0 + 0 \cdot 1 + 0 \cdot 1 + 0 \cdot 0 = 3$$

Normalized:

$$\cos \theta = (v_1 \cdot v_2) / (\| v_1 \| \| v_2 \|)$$

$\| v \| = \sqrt{v \cdot v}$ $\cos \theta = \frac{3}{\sqrt{6} \sqrt{7}} \approx 0.4629.$
 $0 \approx 62.4 \text{ deg.}$



Natural Language Understanding/Processing 38

Properties of the Cosine Similarity

$\cos \theta = 0$ means that the document vectors are orthogonal and the documents have no reference vocabulary occurrences in common.

$\cos \theta = 1$ means that the documents are either identical or the vectors point in the same direction in the n-dim space. That is, the documents share the same distribution of occurrences of the reference terms.

Natural Language Understanding/Processing

39

Latent Semantic Analysis

A problem with the cosine similarity function:
Unless both documents use the same term for something, the similarity is not recognized.

"Computer learning environments have a great future."

"Educational technology offers wonderful potential."

cosine similarity is 0.

Natural Language Understanding/Processing

40

LSA (continued)

With Latent Semantic Analysis, the vector for each document is first transformed into a vector in another space – a "semantic space" in which related terms get mapped to the same element or set of elements.

After that, the cosine similarity between the new vectors will be greater, if the documents share RELATED terms.

Natural Language Understanding/Processing

41

LSA (continued)

The semantic space for LSA is obtained from a set of documents given in advance.

The space is created using matrix factorization via the Singular Value Decomposition (SVD) method.

This is computationally costly, but modern computers are powerful enough to do it.

For more details, see Chapter 16 of *Introduction to Python for Artificial Intelligence*.

Natural Language Understanding/Processing

42

Singular Value Decomposition

Given term-document matrix A, having t rows and d columns, find TSD such that:

$A = TSD$
T is a t by t orthonormal matrix
D is a d by d orthonormal matrix
S is an m by m diagonal matrix, where m is the rank of A.

```
import LinearAlgebra as LA
(TSD) = LA.singular_value_decomposition(A)
```

Natural Language Understanding/Processing

43

Latent Semantic Model

Given TSD, form a reduced (and generalized) product $T_r S_r D_r$ by deleting the rows and columns of S that contain the $n - k$ smallest diagonal values. Then eliminate the last $n - k$ columns of T to get T_r and eliminate the last $n - k$ rows of D to get D_r .

$A_r = T_r S_r D_r$

To compare two documents in the latent semantic space, first map the documents into the space and then compute their cosine similarity.
 $doc_1' = D_r doc_1 ; doc_2' = D_r doc_2 ; cossim(doc_1', doc_2')$

Natural Language Understanding/Processing

44

Example

d_1 = "the brown weasel followed the fox and stole the eggs"
 d_2 = "behind the fence the thief fled with half a dozen"
 d_3 = "artificial limbs can offer full mobility"

Documents used to create a semantic space:
"the lazy brown fox jumped over the fence"
"the thief jumped the lazy fence and fled"
"artificial intelligence is full of surprises"

$\text{cossim}(d_1, d_2) = 0$ Without LSA, d_1 and d_2 seem dissimilar.
 $\text{cossim}(d_1', d_2') = 1$ With LSA, they are completely similar.
 $\text{cossim}(d_1, d_3) = \text{cossim}(d_1', d_3') = 0$ But LSA does not make d_3 any more similar to the others.