

Natural Language Processing CS224N/Ling280



Christopher Manning
Spring 2008
Lecture 1



Course logistics in brief

- Instructor: Christopher Manning
- TAs: Paul Baumstarck, Pichuan Chang
- Time: MW 11:00–12:15. (Section: ?? F 11:00–12:15)
- Handouts:
 - Course syllabus, lecture 1, **assignment 1**
- Programming language: Java 1.5+
- Other information: see the webpage.
 - <http://cs224n.stanford.edu/>



This class

- Assumes you come with some skills...
 - Some basic linear algebra, probability, and statistics; decent programming skills
 - But not everyone has the same skills
 - Assumes some ability to learn missing knowledge
- Teaches key theory and methods for statistical NLP: MT, parsing, semantics, etc.
 - Learn techniques which can be used in practical, robust systems that can (partly) understand human language
- But it's something like an "AI Systems" class:
 - A lot of it is hands on problem-based learning
 - Often practical issues are as important as theoretical niceties
 - We often combine a bunch of ideas



Natural language: the earliest UI

Dave Bowman: Open the pod bay doors, HAL.
HAL: I'm sorry Dave. I'm afraid I can't do that.



(cf. also false Maria in Metropolis – 1926)



Goals of the field of NLP

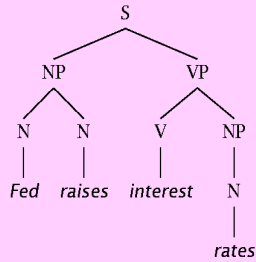
- Computers would be a lot more useful if they could handle our email, do our library research, chat to us ...
- But they are fazed by natural human languages.
 - Or at least their programmers are ... most people just avoid the problem and get into XML, or menus and drop boxes, or ...
- But someone has to work on the hard problems!
 - How can we tell computers about language?
 - Or help them learn it as kids do?
- In this course we seek to identify many of the open research problems in natural language



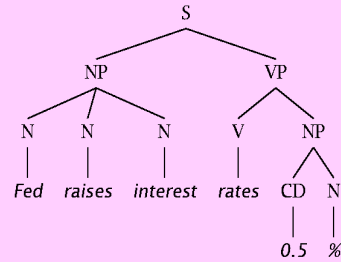
What/where is NLP?

- Goals can be very far reaching ...
 - True text understanding
 - Reasoning about texts
 - Real-time participation in spoken dialogs
- Or very down-to-earth ...
 - Finding the price of products on the web
 - Context sensitive spell-checking
 - Analyzing reading level or authorship statistically
 - Extracting facts or relations from documents
- These days, the latter predominate (as NLP becomes increasingly practical, it is increasingly engineering-oriented – also related to changes in approach in AI/NLP)

The bad effects of V/N ambiguities (2)



The bad effects of V/N ambiguities (3)



Why NLP is difficult: Newspaper headlines

- Ban on Nude Dancing on Governor's Desk
- Iraqi Head Seeks Arms
- Juvenile Court to Try Shooting Defendant
- Teacher Strikes Idle Kids
- Stolen Painting Found by Tree
- Local High School Dropouts Cut in Half
- Red Tape Holds Up New Bridges
- Clinton Wins on Budget, but More Lies Ahead
- Hospitals Are Sued by 7 Foot Doctors
- Kids Make Nutritious Snacks
- Minister Accused Of Having 8 Wives In Jail



LSAT / (former) GRE Analytic Section Questions

- Six sculptures – C, D, E, F, G, H – are to be exhibited in rooms 1, 2, and 3 of an art gallery.
 - Sculptures C and E may not be exhibited in the same room.
 - Sculptures D and G must be exhibited in the same room.
 - If sculptures E and F are exhibited in the same room, no other sculpture may be exhibited in that room.
 - At least one sculpture must be exhibited in each room, and no more than three sculptures may be exhibited in any room.
- If sculpture D is exhibited in room 3 and sculptures E and F are exhibited in room 1, which of the following may be true?
 - Sculpture C is exhibited in room 1
 - Sculpture H is exhibited in room 1
 - Sculpture G is exhibited in room 2
 - Sculptures C and H are exhibited in the same room
 - Sculptures G and F are exhibited in the same room



Reference Resolution

U: Where is **A Bug's Life** playing in **Mountain View**?
 S: A Bug's Life is playing at the **Century 16 theater**.
 U: When is **it** playing **there**?
 S: It's playing at 2pm, 5pm, and 8pm.
 U: I'd like 1 **adult** and 2 **children** for **the first show**.
 How much would **that** cost?

- Knowledge sources:
 - Domain knowledge
 - Discourse knowledge
 - World knowledge



Why is natural language computing hard?

- Natural language is:
 - highly ambiguous at all levels
 - complex and subtle use of context to convey meaning
 - fuzzy, probabilistic
 - involves reasoning about the world
 - a key part of people interacting with other people (a social system):
 - persuading, insulting and amusing them
- But NLP can also be surprisingly easy sometimes:
 - rough text features can often do half the job



Making progress on this problem...

- The task is difficult! What tools do we need?
 - Knowledge about language
 - Knowledge about the world
 - A way to combine knowledge sources
- The answer that's been getting traction:
 - **probabilistic models** built from language data
 - P("maison" → "house") **high**
 - P("L'avocat général" → "the general avocado") **low**
- Some computer scientists think this is a new "A.I." idea
 - But really it's an old idea that was stolen from the electrical engineers....



Where do we head?

Look at subproblems, approaches, and applications at different levels

- Statistical machine translation
- Statistical NLP: classification and sequence models (part-of-speech tagging, named entity recognition, information extraction)
- Syntactic (probabilistic) parsing
- Building semantic representations from text. QA.
- (Unfortunately left out: natural language generation, phonology/morphology, speech dialogue systems, more on natural language understanding, There are other classes for some!)

Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯裔商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛保持高度戒备。



The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport .

The classic acid test for natural language processing.

Requires capabilities in both interpretation and generation.

About \$10 billion spent annually on human translation.

Mainly slides from Kevin Knight (at ISI)

Translation (human and machine)

对外经济贸易合作部今天提供的数据表明，今年至十一月中国实际利用外资四百六十九点五九亿美元，其中包括外商直接投资四百零七亿美元。

Ref According to the data provided today by the Ministry of Foreign Trade and Economic Cooperation, as of November this year, China has actually utilized 46.959 billion US dollars of foreign capital, including 40.007 billion US dollars of direct investment from foreign businessmen.

IBM4: the Ministry of Foreign Trade and Economic Cooperation, including foreign direct investment 40.007 billion US dollars today provide data include that year to November china actually using foreign 46.959 billion US dollars and Yamada/Knight: today's available data of the Ministry of Foreign Trade and Economic Cooperation shows that china's actual utilization of November this year will include 40.007 billion US dollars for the foreign direct investment among 46.959 billion US dollars in foreign capital

Machine Translation History

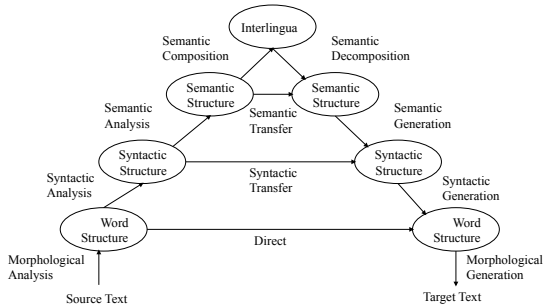
- 1950s: Intensive research activity in MT
- 1960s: Direct word-for-word replacement
- 1966 (ALPAC): NRC Report on MT
 - Conclusion: MT no longer worthy of serious scientific investigation.
- 1966-1975: 'Recovery period'
- 1975-1985: Resurgence (Europe, Japan)
 - Domain specific rule-based systems
- 1985-1995: Gradual Resurgence (US)
- 1995-2005: Statistical MT surges ahead

<http://ourworld.compuserve.com/homepages/WJHutchins/MTS-93.htm>

What happened between ALPAC and Now?

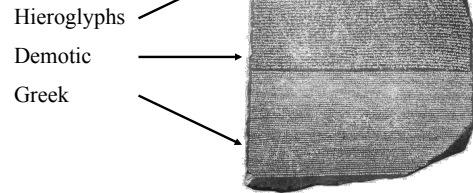
- Need for MT and other NLP applications confirmed
- Change in expectations
- Computers have become faster, more powerful
- WWW
- Political state of the world
- Maturation of Linguistics
- Availability of data
- Development of statistical and hybrid statistical/symbolic approaches

Three MT Approaches: Direct, Transfer, Interlingual (Vauquois triangle)



Statistical Solution

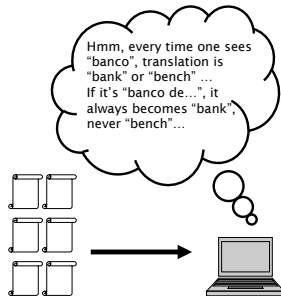
- Parallel Texts
 - Rosetta Stone



Statistical Solution

- Parallel Texts

- Instruction Manuals
- Hong Kong Legislation
- Macao Legislation
- Canadian Parliament Hansards
- United Nations Reports
- Official Journal of the European Communities



Warren Weaver

- “Also knowing nothing official about, but having guessed and inferred considerable about, the powerful new mechanized methods in cryptography—methods which I believe succeed even when one does not know what language has been coded—one naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. 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.’ ”
- Warren Weaver (1955:18, quoting a letter he wrote in 1947)

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok hihok yorok klok kantok ok-yurp

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok hihok yorok klok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anok plok sprok .	8a. lalok brok anok plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anok drok brok jok .	10a. lalok mok nok yorok ghrok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

process of elimination

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errok** **hihok** **yorok** **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

cognate?

Centauri/Arcturan [Knight, 1997]

Your assignment, put these words in order: { **jjat**, **arrat**, **mat**, **bat**, **olloat**, **at-yurp** }

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

zero fertility

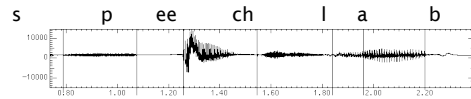
It's Really Spanish/English

Clients do not sell pharmaceuticals in Europe => Clientes no venden medicinas en Europa

1a. Garcia and associates .	7a. the clients and the associates are enemies .
1b. Garcia y asociados .	7b. los clients y los asociados son enemigos .
2a. Carlos Garcia has three associates .	8a. the company has three groups .
2b. Carlos Garcia tiene tres asociados .	8b. la empresa tiene tres grupos .
3a. his associates are not strong .	9a. its groups are in Europe .
3b. sus asociados no son fuertes .	9b. sus grupos estan en Europa .
4a. Garcia has a company also .	10a. the modern groups sell strong pharmaceuticals .
4b. Garcia tambien tiene una empresa .	10b. los grupos modernos venden medicinas fuertes .
5a. its clients are angry .	11a. the groups do not sell zanzanine .
5b. sus clientes estan enfadados .	11b. los grupos no venden zanzanina .
6a. the associates are also angry .	12a. the small groups are not modern .
6b. los asociados tambien estan enfadados .	12b. los grupos pequenos no son modernos .

Speech Recognition: Acoustic Waves

- Human speech generates a wave
 - like a loudspeaker moving
- A wave for the words "speech lab" looks like:



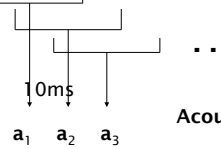
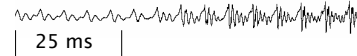
"l" to "a" transition:

Graphs from Simon Arnfield's web tutorial on speech, Sheffield: <http://www.psyc.leeds.ac.uk/research/cogn/speech/tutorial/>

43

Acoustic Sampling

- 10 ms frame (ms = millisecond = 1/1000 second)
- ~25 ms window around frame [wide band] to allow/smooth signal processing – it let's you see formants

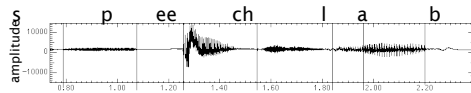


Result:
Acoustic Feature Vectors
(after transformation,
numbers in roughly R^{14})

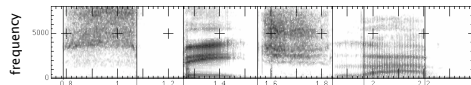
44

Spectral Analysis

- Frequency gives pitch; amplitude gives volume
 - sampling at ~8 kHz phone, ~16 kHz mic (kHz=1000 cycles/sec)



- Fourier transform of wave displayed as a spectrogram
 - darkness indicates energy at each frequency
 - hundreds to thousands of frequency samples



45

The Speech Recognition Problem

- The **Recognition Problem: Noisy channel model**
 - We started out with English words, they were encoded as an audio signal, and we now wish to decode.
 - Find most likely sequence w of "words" given the sequence of acoustic observation vectors a

- Use Bayes' rule to create a **generative model** and then decode
- $\text{ArgMax}_w P(w|a) = \text{ArgMax}_w P(a|w) P(w) / P(a)$
 $= \text{ArgMax}_w P(a|w) P(w)$

- Acoustic Model:** $P(a|w)$
- Language Model:** $P(w)$

A probabilistic theory of a language

46

Probabilistic Language Models

- Assign probability $P(w)$ to word sequence $w = w_1, w_2, \dots, w_k$
- Can't directly compute probability of long sequence – one needs to decompose it
- Chain rule provides a **history-based model**:
$$P(w_1, w_2, \dots, w_k) = P(w_1) P(w_2|w_1) P(w_3|w_1, w_2) \dots P(w_k|w_1, \dots, w_{k-1})$$
- Cluster** histories to reduce number of parameters
- E.g., just based on the last word (1st order Markov model):
$$P(w_1, w_2, \dots, w_k) = P(w_1|<s>) P(w_2|w_1) P(w_3|w_2) \dots P(w_k|w_{k-1})$$
- How do we estimate these probabilities?
 - We count word sequences in corpora
 - We "smooth" probabilities so as to allow unseen sequences

47

N-gram Language Modeling

- n -gram assumption clusters based on last $n-1$ words
 - $P(w_j|w_1, \dots, w_{j-1}) \approx P(w_j|w_{j-n+1}, \dots, w_{j-2}, w_{j-1})$
 - unigrams $\sim P(w_j)$
 - bigrams $\sim P(w_j|w_{j-1})$
 - trigrams $\sim P(w_j|w_{j-2}, w_{j-1})$
- Trigrams often interpolated with bigram and unigram:

$$\hat{P}(w_3 | w_1, w_2) = \lambda_3 \frac{F(w_3 | w_1, w_2)}{\sum_k F(w_k | w_1, w_2)} + \lambda_2 \frac{F(w_3 | w_2)}{\sum_k F(w_k | w_2)} + \lambda_1 \frac{F(w_3)}{\sum_k F(w_k)}$$

- the λ_i typically estimated by maximum likelihood estimation on held out data ($F(\cdot, \cdot)$ are relative frequencies)
- many other interpolations exist (another standard is a non-linear **backoff**)

48