

1

#### Bookkeeping

### Final exam is Tuesday in class

- Final exam review Thursday
- Project Phase II due 12/7 at 11:59 PM
- Project final paper due 12/15 at 11:59 PM

#### Speech and Language Processing

 Getting computers to do reasonably intelligent things with human language is the domain of Computational Linguistics (or Natural Language Processing or Human Language Technology)









7

#### Applications

- Applications of NLP can be broken down into Small and Big
- Small applications include many things you never think about:
  - Hyphenation
  - Spelling correction
  - OCR
  - Grammar checkers
- Big applications include:
  - Machine translation
  - Question answering
  - Conversational speech

#### Applications

- There's another kind: Medium
  - Speech in closed domains
  - Question answering in closed domains
  - Question answering for factoids
  - Information extraction from news-like text
  - Generation and synthesis in closed/small domains.

9

#### NLP Research

- In between the linguistics and the big applications are a host of hard problems.
  - Robust Parsing
  - Word Sense Disambiguation
  - Semantic Analysis
  - Etc
- Not too surprisingly, solving these problems involves:
  - Choosing the right logical representations
  - Managing hard search problems
  - Dealing with uncertainty
  - Using machine learning to train systems to do what we need

#### Why Study NLP?

- A hallmark of human intelligence.
- To interact with computing devices using human (natural) languages
  - Building intelligent robots (AI)
  - Enabling voice-controlled operation
- To access (large amount of) information and knowledge stored in the form of human languages quickly
  - Text is the largest repository of human knowledge and is growing quickly.
    - Emails, news articles, web pages, IM, scientific articles, insurance claims, customer complaint letters, transcripts of phone calls, technical documents, government documents, patent portfolios, court decisions, contracts, ...

11

#### NL and NLP

- "Natural" languages = human languages
  - English, Russian, Wolof, ...
- Natural Language Processing: any form of dealing with NL computationally
- Many, many sub-areas; from an AI perspective, 2 are most crucial:
  - Natural Language Understanding: understanding the meaning (semantics) of spoken or written text
  - Natural Language Generation: Producing meaningful, relevant language

#### Fundamental NLP Tasks

- Semantics: map sentence to corresponding "meaning representation" (e.g., logic)
  - give(John, Book, Sally)
  - Quantification: Everybody loves somebody
- Word Sense Disambiguation
  - orange juice vs. orange coat
  - Where do word senses come from?
- Co-reference resolution:
  - The dog found a cookie; He ate it.
- Implicit "text" what is left unsaid?
  - Joe entered the restaurant, ordered, ate and left. The owner said "Come back soon!"
  - Joe entered the restaurant, ordered, ate and left. The owner said "Hey, I'll call the police!"



#### Applied NLP

- Machine translation
- Spelling/grammar correction
- Information Retrieval/extraction
- Data mining
- Document classification
- Question answering
- Conversational agents

#### You See It Daily

- Question answering: Siri, OK Google, Cortana, Alexa
- spelling/grammar correction
- Automated response systems
- To get input for
  - Information Retrieval
  - Data mining
  - Document classification
- Machine translation





#### Semantics: Meaning



18

# Semantics What kinds of things can we not do well with the tools we have already looked at? Retrieve information in response to unconstrained questions: e.g., travel planning Play the "chooser" side of 20 Questions Provide accurate information about a topic Support human-computer interfaces Read a newspaper article and answer questions about it These tasks require that we also consider semantics: the meaning of our tokens and their sequences

#### Semantics



• What are the semantic primitives?















#### Human Languages

- You know ~50,000 words of primary language, each with several meanings
- Six year old knows ~13000 words
- First 16 years we learn 1 word every 90 min of waking time
- Mental grammar generates sentences
  - virtually every sentence is novel!
- 3 year olds already have 90% of grammar
- ~6000 human languages none of them simple!

Adapted from Martin Nowak 2000 – Evolutionary biology of language – Phil.Trans. Royal Society London

## Human spoken language Most complicated mechanical motion of the body Movements must be accurate to within half mm synchronized within hundredths of a second We can understand up to 50 phonemes/sec (normal speech 10-15ph/sec) but if sound is repeated 20 times /sec we hear continuous buzz! All aspects of language processing are involved and manage to keep apace

28

#### One Huge Challenge in NLP: Ambiguity

- Lexical:
  - Label (noun or verb)?
  - London (Jack or capital of the UK)?
- Syntactic (examples from newspaper headlines):
  - Prepositional Phrase Attachment: Ban on Nude Dancing on Governor's Desk
  - Word Sense Disambiguation: Iraqi Head Seeking Arms
  - Syntactic Ambiguity (what's modifying what): Juvenile Court to Try Shooting Defendant
- Semantic ambiguity:
  - "snake poison"
- Rampant metaphors:
  - "prices went up"

#### Some Ambiguous Headlines

- Juvenile Court to Try Shooting Defendant
- Teacher Strikes Idle Kids
- Kids Make Nutritious Snacks
- Bush Wins on Budget, but More Lies Ahead
- Hospitals are Sued by 7 Foot Doctors

Source: Marti Hearst, i256, at UC Berkeley

30

#### Dealing with Ambiguity

- Four possible approaches:
  - Formal approaches -- Tightly coupled interaction among processing levels; knowledge from other levels can help decide among choices at ambiguous levels.
  - 2. Pipeline processing that ignores ambiguity as it occurs and hopes that other levels can eliminate incorrect structures.
  - 3. Probabilistic approaches based on making the most likely choices
  - 4. Don't do anything, maybe it won't matter







#### What Are Words?

- Hard to get agreement
- (Human) Language-dependent
- White-space separation is a sometimes okay (for written English longform)
- Social media? Spoken vs. written? Other languages?

What Are Words?

pişirdiler

They cooked it.

What Are Words?

pişmişlermişlerdi

They had it cooked it.

What Are Words?

):





#### Why Language is Hard

- NLP is Al-complete
- Abstract concepts are difficult to represent
- LOTS of possible relationships among concepts
- Many ways to represent similar concepts
- Tens of hundreds or thousands of features/dimensions

#### Why Language is Easy

- Highly redundant
- Relatively crude methods provide fairly good results
- Lots of subject matter experts!

42

#### Some of the Tools

- A mixed bag, at various levels...
  - Tokenizers
  - Regular Expressions and Finite State Automata
  - Part of Speech taggers
  - Grammars
  - Parsers
  - N-Grams
  - Semantic Analysis

#### What will it take?

- Models of computation (state machines)
- Formal grammars
- Knowledge representation
- Search algorithms
- Dynamic programming
- Logic
- Machine learning
- Probability theory



#### Text annotation tasks

- Classify the entire document ("text categorization")
- Classify word tokens individually
- Classify word tokens in a sequence (i.e., order matters)
- Identify phrases ("chunking")
- Syntactic annotation (parsing)
- Semantic annotation
- Text generation

46

#### Parts of Speech Tagging

- Part-of-Speech (POS) taggers identify nouns, verbs, adjectives, noun phrases, etc.
- More recent work uses machine learning to create taggers from labeled examples



#### Named Entities (NE) Tagging

- Persons, places, companies
  - "Proper nouns"
  - One of most common information extraction tasks
  - Combination of rules and dictionary
- Example rules:
  - Capitalized word not at beginning of sentence
  - Two capitalized words in a row
  - One or more capitalized words followed by Inc
  - Dictionaries of common names, places, major corporations.
    - Sometimes called "gazetteer"

48

#### **Reference Resolution**

- Discourse Knowledge what have we just said?
  - <u>Paula</u> is here. <u>She</u> is ready.
- Domain/World Knowledge
  - U: I would like to register in a CMSC Course.
  - S: Which number?
  - U: 647.
  - S: Which section?
  - U: Which section is in the evening?
  - S: section 1.
  - U: Then that one.

#### Word Sense Resolution

- Many words have several meanings or senses
- We need to resolve which of the senses of an ambiguous word is invoked in a particular use of the word
- I made her duck. (meanings?)

#### Word Sense Resolution

- Many words have several meanings or senses
- We need to resolve which of the senses of an ambiguous word is invoked in a particular use of the word
- I made her duck. (meanings?)
  - 1. I cooked waterfowl for her benefit (to eat)
  - 2. I cooked waterfowl belonging to her
  - 3. I created the (plaster?) duck she owns
  - 4. I caused her to quickly lower her head or body
  - 5. I waved my magic wand and turned her into undifferentiated waterfowl
- Again, discourse and world knowledge

Duck example Jurafsky & Martin "Speech and Language Processing"

#### Issues with Statistical Parsing

- Statistical parsers still make plenty of errors
- Tree banks are language specific
- Tree banks are genre specific
  - Train on WSJ → fail on the Web
  - standard distributional assumption
- Unsupervised, un-lexicalized parsers exist
  - But performance is substantially weaker

52

#### **Big Applications**

- POS tagging, parsing and word sense disambiguation are all mediumsized enabling applications.
  - They don't actually do anything that anyone actually cares about.
  - MT and QA are things people seem to care about.

#### How Difficult is Morphology?

- Examples from Woods et. al. 2000
  - delegate (de + leg + ate) take the legs from
  - caress (car + ess) female car
  - cashier (cashy + er) more wealthy
  - lacerate (lace + rate) speed of tatting
  - ratify (rat + ify) infest with rodents
  - infantry (infant + ry) childish behavior





#### But "Reading" the Web is Tough

- Traditional IE is narrow
- IE has been applied to small, homogenous corpora
- No parser achieves high accuracy
- No named-entity taggers
- No supervised learning
- How about semi-supervised learning?





#### **Machine Translation** The automatic translation of texts between languages is one of the • oldest non-numerical applications in Computer Science. In the past 20 years or so, MT has gone from a niche academic • curiosity to a robust commercial industry. 巨大な銃規制集 会が米国を席巻 Huge gun-学生が主催する「私たちの生活 のための行進」イベントでは、 全国的に数十万人の抗議者が集 control rallies まります。 sweep US ◎4時間 米国とカナダ Student-led March For Our Lives events nationwide draw ARCH hundreds of thousands of protesters. ORDUR () 4h US & Canada IVES



#### Sentiment Analysis

• The field of sentiment analysis deals with categorization (or classification) of opinions expressed in textual documents

The TV is wonderful. Great size, great picture, easy interface. It makes a cute little song when you boot it up and when you shut it off. I just want to point out that the 43" does not in fact play videos from the USB. This is really annoying because that was one of the major perks I wanted from a new TV. Looking at the product description now, I realize that the feature list applies to the X758 series as a whole, and that each model's capabilities are listed below. Kind of a dumb oversight on my part, but it's equally stupid to put a description that does not apply on the listing for a very specific model.

Green color represents positive tone, red color represents negative tone, and product features and model names are highlighted in blue and brown, respectively.















#### Known LLM issues

- Bad reproducibility
- Copyright issues
- Can't explain what it's doing
- Can't remember things long term
- Confident bullshitter

Slide: Dr. Lara Martir

69

#### Conclusions

- NLP is harder than it might seem naively
- Many subtasks
- Statistical NLP is the dominant paradigm
  - supervised learning
  - corpus-based statistics (language models)
  - Some important limitations in practice!
- NL "understanding" has received very little attention

#### Our NLP/NLU Class: Course Goals

- Be introduced to some of the core problems and solutions of NLP (big picture)
- Learn different ways that success and progress can be measured in NLP
- Relate to statistics, machine learning, and linguistics
- Implement NLP programs

#### Our NLP/NLU Class: Course Topics

- Probability, classification, and the efficacy of simple counting methods
- Language modeling (n-gram models, smoothing heuristics, maxent/log-linear models, and distributed/vector-valued representations)
- Sequences of latent variables (e.g., hidden Markov models, some basic machine translation alignment)
- Trees and graphs, as applied to syntax and semantics
- Some discourse-related applications (coreference resolution, textual entailment)
- Special and current topics (e.g., fairness and ethics in NLP)
- Modern, neural approaches to NLP, such as recurrent neural networks and transformers (e.g., BERT or GPT-2)