# Lecture 1:
# Course Introduction

CMPSCI 585, Fall 2015
Introduction to Natural Language Processing
http://people.cs.umass.edu/~brenocon/inlp2015/

Brendan O'Connor

# What

- Learn fundamental principles and methods in natural language processing
  - Hands-on implementation experience
  - Appreciation of basic linguistic issues
  - Know about useful open-source NLP tools
  - Know when NLP works and when it doesn't
  - Be able to read current research papers in NLP
- "AI systems"

# How

- Lots of math!
  - This course: more than traditional CS, but less than traditional physics, engineering
- Lots of data!
- Lots of code that *implements* math & algorithms
  - Skill: translating from math to code
  - Skill: debugging math/linguistic/algorithm code
- A little bit of linguistics (goes a long way)

3

# Requirements

- (20%) Participation and short exercises
- (30%) Problem sets
  - Written: math and concepts
  - Programming: Python programs
- (20%) Midterm (in-class, Oct 20)
- (30%) Final projects (groups of 1-2)
  - Choose a topic, or select a suggested topic
  - Project Proposal
  - Progress Report
  - In-class presentations (last week)
  - Final Report

4

# Logistics

- Main course website for assignments and links to everything. This is ground truth if there are contradictions http://people.cs.umass.edu/~brenocon/inlp2015/
- Piazza for announcements, discussions and restricted access files
- Moodle for homework submissions (sometimes) and to see grades

- Waitlist situation

- To check:
  - SPIRE-registered students should have Piazza invites. Check @umass.edu email if you don't!
  - Email me if you can't access Piazza.
- Homework #0 due Thursday!

5

# Readings

- Readings will be provided as PDFs on Piazza
    - Draft chapters of Jurafsky and Martin, *Speech and Language Processing*
    - We will use selections from both
        - 2nd edition (published)
        - 3rd edition (unpublished)
    - Other readings on occasion as well

6

# Related courses at UMass

- Computational Linguistics:  Ling 409, 492B (Bhatt, Dillon)

  https://sites.google.com/site/umasslx409/home
  http://www.umass.edu/linguist/courses/detail.php?cid=571

- Speech:  Ling 592B (Yu)

  http://courses.umass.edu/linguist592b-kmyu/category/info.html

- Information Retrieval:  CS 446, 646 (Allan, Croft)

  http://ciir.cs.umass.edu/cmpsci446/
  http://ciir.cs.umass.edu/~allan/cs646/

# NLP is interdisciplinary

Algorithms

Linguistics

Statistics +
Machine Learning

Cognitive Science

Artificial Intelligence

8

# "Can Machines Think?"

- British mathematician and founding figure in computer science
- <u>Alan Turing (1950)</u>
- How do we know when we have AI?
- "Imitation Game"

# NLP imagined

# NLP today

- Speech interfaces
- Machine translation
- Sentiment analysis
- Search engines
- ...

- [This course: document text analysis]

# NLP today: Speech interfaces

# "Rao's coffee in Amherst, Massachusetts"

# "Rao's coffee in Amherst, Massachusetts"

# NLP today: Question answering



IBM Watson

*Wanted for general evilness, last seen at the Tower of
Barad-Dur. It's a giant eye, folks, kinda hard to miss*

# NLP today: Question answering

IBM Watson

25 engineers, 4 years, 200 subsystems,
2,880 CPU cores, 15 TB storage

# NLP today: Question answering

From <u>IBM Journal of Research and Development, 2012</u>

**Table 1**    DeepQA technology performance on public benchmark sets. (ACE: automatic content extraction; RTE: recognizing textual entailment.)

| NLP task | Evaluation set | Project start | State of art | Watson |
|---|---|---|---|---|
| Parsing | Wikipedia** accuracy | 84.4 | 81.1 Charniak parser [19] | 88.7 |
| Entity disambiguation | Wikipedia disambiguation $F_1$ | 72.5 | 81.9 Hoffart et al. [42] | 92.5 |
| Relation detection | ACE 2004 $F_1$ | 45.8 | 72.1 Zhang et al. [43] | 73.2 |
| Textual entailment | RTE-6 2010 $F_1$ | 34.6 | 48.0 PKUTM [44] | 48.8 |

IBM Watson

Imperfect NLP is still useful

# Ambiguity: why NLP is hard

# Ambiguity: why NLP is hard

- Juvenile Court to Try Shooting Defendant

# Ambiguity: why NLP is hard

- Juvenile Court to Try Shooting Defendant
- Hospitals Are Sued by 7 Foot Doctors

# Ambiguity: why NLP is hard

- Juvenile Court to Try Shooting Defendant
- Hospitals Are Sued by 7 Foot Doctors
- Alice saw Bob with a telescope.

# Ambiguity: why NLP is hard

- Juvenile Court to Try Shooting Defendant
- Hospitals Are Sued by 7 Foot Doctors
- Alice saw Bob with a telescope.
- Our company is training workers.

# Ambiguity: why NLP is hard

- Juvenile Court to Try Shooting Defendant
- Hospitals Are Sued by 7 Foot Doctors
- Alice saw Bob with a telescope.
- Our company is training workers.
- They found that in order to attract settlers -- and make a profit from their holdings -- they had to offer people farms, not just tenancy on manorial estates.

# Levels of linguistic structure

Characters

| A | l | i | c | e |   | t | a | l | k | e | d |   | t | o |   | B | o | b | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Levels of linguistic structure

Morphology

Characters

talk -ed [VerbPast]

Alice talked to Bob.

18

# Levels of linguistic structure

| Words |
| Morphology |
| Characters |

Alice | talked | to | Bob | .

talk | -ed | [VerbPast]

| A | l | i | c | e | | t | a | l | k | e | d | | t | o | | B | o | b | . |

# Levels of linguistic structure

| Syntax: Part of Speech |
|---|

| Words |
|---|

| Morphology |
|---|

| Characters |
|---|

Noun    VerbPast    Prep    Noun    Punct

`Alice` `talked` `to` `Bob` `.`

`talk` `-ed`  [VerbPast]

`Alice talked to Bob.`

18

# Levels of linguistic structure

Syntax: Constituents

Syntax: Part of Speech

Words

Morphology

Characters

S
VP
NP PP .
Noun VerbPast Prep Noun Punct

Alice | talked | to | Bob | .

talk | -ed | [VerbPast]

Alice talked to Bob.

18

# Levels of linguistic structure

Discourse

Semantics

Syntax: Constituents

Syntax: Part of Speech

Words

Morphology

Characters

CommunicationEvent(e)    SpeakerContext(s)
Agent(e, Alice)          TemporalBefore(e, s)
Recipient(e, Bob)

S
├─ NP
│  └─ Noun
└─ VP
   ├─ VerbPast
   └─ PP
      ├─ Prep
      └─ Noun
   .
   └─ Punct

Alice    talked    to    Bob    .

talk -ed    [VerbPast]

Alice talked to Bob.

18

# NLP today: Machine translation

# NLP today: Machine translation

# NLP today: Trend analysis

Data: news articles

Dependency parsing to identify events



British officials in Tehran and London have been meeting discretely with their Iranian counterparts

GBR                                                                                          IRN

Machine learning from text:

## (1) **Event class dictionaries**

"diplomacy"
> arrive in, visit, meet with, travel to, leave, hold with, meet, meet in, fly to, be in, arrive for talk with, say in, arrive with, head to, hold in, due in, leave for, make to, arrive to,

"verbal conflict"
> accuse, blame, say, break with, sever with, blame on, warn, call, attack, rule with, charge, say←ccomp come from, say ←ccomp, suspect, slam, accuse government ←poss,

"material conflict"
> kill in, have troops in, die in, be in, wound in, have soldier in, hold in, kill in attack in, remain in, detain in, have in, capture in, stay in, about ←pobj troops in, kill, have troops

## (2) **Political dynamics**



Israeli–Palestinian Diplomacy

# NLP today: Story generation

## Earnings for OmniVision Technologies Expected to Fall

By Narrative Science

+ Comment Now   + Follow Comments
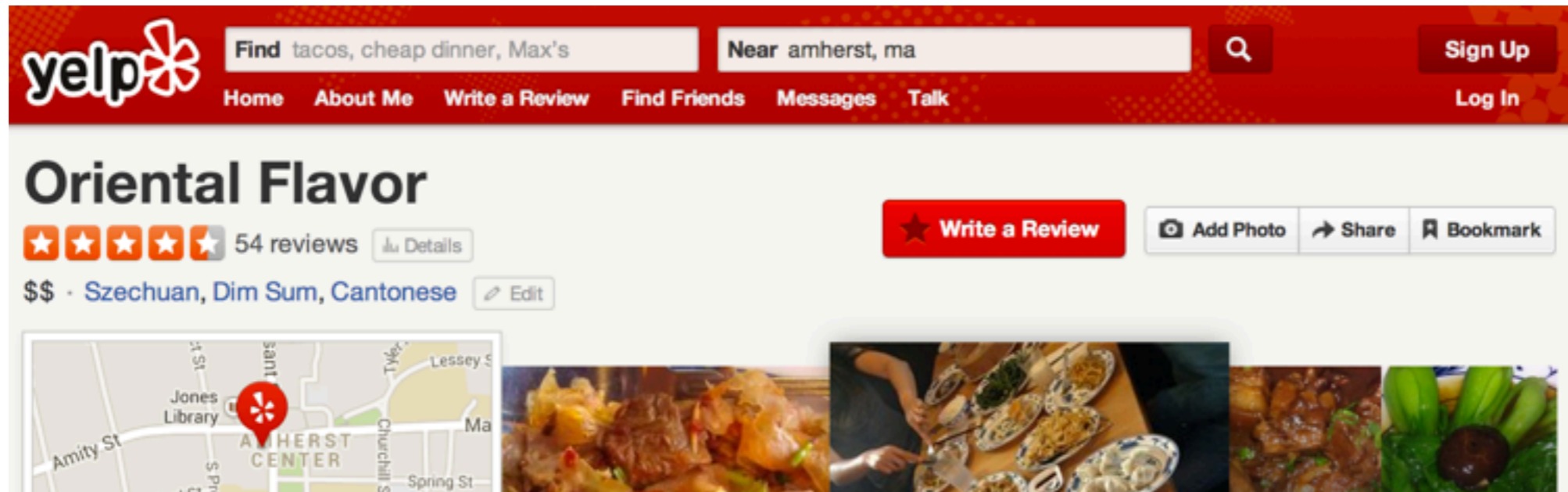
**Narrative Science**

FOLLOW

full bio →

Opinions expressed by Forbes Contributors are their own.

Wall Street is expecting lower profit for **OmniVision Technologies** when the company reports its first quarter results on Thursday, August 28, 2014. Analysts are expecting earnings per share of 39 cents after the company booked a profit of 42 cents a share a year earlier.

The consensus estimate has risen from 16 cents over the past three months. Analysts are expecting earnings of 99 cents per share for the fiscal year. Revenue is projected to eclipse the year-earlier total of $373.7 million by 2%, finishing at $381.5 million for the quarter. For the year, revenue is projected to come in at $1.39 billion.

http://www.forbes.com/sites/narrativescience/

# NLP today: Search/summarization

# NLP today: Search/summarization

# NLP today: Search/summarization



TOP SECRET//COMINT//REL TO USA, AUS, CAN, GBR, NZL

## Entity Extraction

- Have technology (thanks to R6) – for English, Arabic and Chinese
- Allow queries like:
- Show me all the word documents with references to IAEO
- Show me all documents that reference Osama Bin Laden
- Will allow a 'show me more like this' capability

TOP SECRET//COMINT//REL TO USA, AUS, CAN, GBR, NZL

NSA slides from Snowden leaks

http://www.theguardian.com/world/interactive/2013/jul/31/nsa-xkeyscore-program-full-presentation

- HW0
- See you on Thursday

26