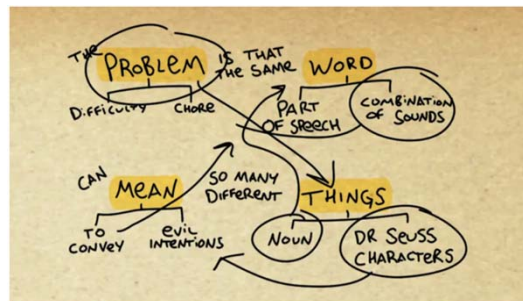


# 600.465 Connecting the dots - I (NLP in Practice)

Delip Rait Master subtitle style  
delip@jhu.edu



## What is "Text"?

*Shall I compare thee to a summer's day?  
Thou art more lovely and more temperate;  
Rough winds do shake the darling buds of May,  
And summer's lease hath all too short a date;  
Sometime too hot the eye of heaven shines,  
And often is his gold complexion dimm'd;  
And every fair from fair sometime declines,  
By chance or nature's changing course untrimm'd;  
But thy eternal summer shall not fade,  
Nor lose possession of that fair thou ow'st;  
Nor shall Death brag thou wander'st in his shade,  
When in eternal lines to time thou grow'st:  
So long as men can breathe or eyes can see,  
So long lives this, and this gives life to thee.*

## What is "Text"?

*Shall I compare thee to a summer's day?*

**A relation-preserving isomorphism**  
If one object consists of a set  $X$  with a binary relation  $R$  and the other object consists of a set  $Y$  with a binary relation  $S$  then an isomorphism from  $X$  to  $Y$  is a bijective function  $f: X \rightarrow Y$  such that  
 $S(f(u), f(v)) \iff R(u, v)$   
 $S$  is reflexive, intransitive, symmetric, antisymmetric, transitive, total, trichotomous, a partial order, total order, strict weak order, total preorder (weak order), an equivalence relation, or a relation with any other special properties, if and only if  $R$  is.  
For example,  $R$  is an ordering  $\leq$  and  $S$  an ordering  $\sqsubseteq$ , then an isomorphism from  $X$  to  $Y$  is a bijective function  $f: X \rightarrow Y$  such that  
 $f(u) \sqsubseteq f(v) \iff u \leq v$ .  
Such an isomorphism is called an *order isomorphism* or (less commonly) an *isotone isomorphism*.  
If  $X = Y$  we have a *relation-preserving automorphism*.

**An operation-preserving isomorphism**  
Suppose that on these sets  $X$  and  $Y$ , there are two binary operations  $\ast$  and  $\circ$  which happen to constitute the groups  $(X, \ast)$  and  $(Y, \circ)$ . Note that the operators operate on elements from the domain and range, respectively, of the "one-to-one" and "onto" function  $f$ . There is an isomorphism from  $X$  to  $Y$  if the bijective function  $f: X \rightarrow Y$  happens to produce results, that sets up a correspondence between the operator  $\ast$  and the operator  $\circ$ .  
 $f(u) \circ f(v) = f(u \ast v)$   
for all  $u, v \in X$ .

## What is "Text"?

## "Real" World

- Tons of data on the web
- A lot of it is text
- In many languages
- In many genres

Language by itself is complex.  
The Web further complicates language.

Click to edit Master subtitle style

11/30/11

### NLP for fun and profit

- Making NLP more accessible
  - Provide APIs for common NLP tasks

```
var text = document.get(...);
var entities = agent.markNE(text);
```

- **Big \$\$\$\$**
- Backend to intelligent processing of text

### Desideratum: Multilinguality

- Except for feature ex be language agnostic

Language	Millions of users
English	326.6
Chinese	244.9
Spanish	133.3
Japanese	95.1
Portuguese	83.5
German	75.7
Arabic	65.4
French	55.9
Russian	54.7
Korean	26.4
All the rest	358.8

Source: Internet World Stats - www.internetworldstats.com/stats7.htm  
 Estimated internet users are 1,305,514,816 on June 30, 2010  
 Copyright © 2000 - 2010, Minicraft Marketing Group  
 Level1

### In this lecture

- Understand how to solve and ace in NLP tasks
- Learn general methodology or approaches
- End-to-End development using an example task
- Overview of (un)common NLP tasks

### Case study: Named Entity Recognition

A hand-drawn diagram on a piece of paper. A red cube is drawn in the center. Three arrows point towards the cube from the words 'THING', 'PLACE', and 'PERSON' written around it.

### Case study: Named Entity Recognition

- Demo: <http://viewer.opencalais.com>
- How do we build something like this?
- How do we find out well we are doing?
- How can we improve?


Click to edit Master subtitle style

11/30/11

### Case study: Named Entity Recognition


- Collect data to learn from
  - Sentences with words marked as PER, ORG, LOC, NONE
- How do we get this data?

### Pay the experts



Linguistic Data Consortium

The Linguistic Data Consortium supports language-related education, research and technology development by creating and sharing linguistic resources: data, tools and standards.




### Wisdom of the crowds

#### Make Money by working on HITs

HITs - Human Intelligence Tasks - are individual tasks that you work on. [Find HITs here.](#)

**As a Mechanical Turk Worker you:**

- Can work from home
- Choose your own work hours
- Get paid for doing good work




[Find HITs Now](#)

#### Get Results from Mechanical Turk Workers

Ask workers to complete HITs - Human Intelligence Tasks - and get results using Mechanical Turk. [Request Now](#)

**As a Mechanical Turk Requester you:**

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results.



[Get Started](#)

### Getting the data: Annotation

- Time consuming
- Costs \$\$\$
- Need for quality control
  - Inter-annotator agreement
  - Kappa score (Krippendorff, 1980)
- Smarter ways to annotate
  - Get fewer annotations: Active Learning
  - Rationales (Zaidan, Eisner & Piatko, 2007)

Click to edit Master subtitle style

11/30/11

### Our recipe ...

- 1. Formalize some insights
- 2. Study the formalism mathematically
- 3. Develop & implement algorithms
- 4. Test on real data

### NER: Designing features

Only
France
and
Great
Britain
backed
Fischler
's
proposal
.

Click to edit the outline text format  
 Need to segment sentences  
 Tokenize the sentences  
 Level  
 - Third Outline  
 Level  
 • Not as trivial as you think  
 • Fourth Outline  
 Level  
 • Original text itself might be in an ugly HTML  
 Level  
 • Cleaneval! Outline  
 Level

### Preprocessing

### NER: Designing features

Only	IS_CAPITALIZED
France	IS_CAPITALIZED
and	
Great	IS_CAPITALIZED
Britain	IS_CAPITALIZED
backed	
Fischler	IS_CAPITALIZED
's	
proposal	
.	

### NER: Designing features

Only	IS_CAPITALIZED	IS_SENT_START
France	IS_CAPITALIZED	
and		
Great	IS_CAPITALIZED	
Britain	IS_CAPITALIZED	
backed		
Fischler	IS_CAPITALIZED	
's		
proposal		
.		

### NER: Designing features

Only	IS_CAPITALIZED	IS_SENT_START
France	IS_CAPITALIZED	
and		
Great	IS_CAPITALIZED	
Britain	IS_CAPITALIZED	
backed		
Fischler	IS_CAPITALIZED	
's		
proposal		
.		

### NER: Designing features

Only	IS_CAPITALIZED	IS_SENT_START	
France	IS_CAPITALIZED		IN_LEXICON_LOC
and			
Great	IS_CAPITALIZED		
Britain	IS_CAPITALIZED		IN_LEXICON_LOC
backed			
Fischler	IS_CAPITALIZED		
's			
proposal			
.			

### NER: Designing features

These are extracted during preprocessing!

Only	POS=RB	IS_CAPITALIZED	...	PREV_WORD=_NONE_	...
France	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=only	...
and	POS=CC		...	PREV_WORD=france	...
Great	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=and	...
Britain	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=great	...
backed	POS=VBD		...	PREV_WORD=britain	...
Fischler	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=backed	...

### NER: Designing features

Only	POS=RB	IS_CAPITALIZED	...	PREV_WORD=_NONE_
France	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=only
and	POS=CC		...	PREV_WORD=france
Great	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=and
Britain	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=great
backed	POS=VBD		...	PREV_WORD=britain
Fischler	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=backed

### NER: Designing features

Only	POS=RB	IS_CAPITALIZED	...	PREV_WORD=_NONE_	...
France	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=only	...
and	POS=CC		...	PREV_WORD=france	...
Great	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=and	...
Britain	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=great	...
backed	POS=VBD		...	PREV_WORD=britain	...
Fischler	POS=NNP	IS_CAPITALIZED	...	PREV_WORD=backed	...
's	POS=XX		...	PREV_WORD=fischler	...

### NER: Designing features

- Can you think of other features?

WORD	HAS_DIGITS
PREV_WORD	IS_HYPHENATED
NEXT_WORD	IS_ALLCAPS
PREV_BIGRAM	FREQ_WORD
NEXT_BIGRAM	RARE_WORD
POS	USEFUL_UNIGRAM_PER
PREV_POS	USEFUL_BIGRAM_PER
NEXT_POS	USEFUL_UNIGRAM_LOC
PREV_POS_BIGRAM	USEFUL_BIGRAM_LOC
NEXT_POS_BIGRAM	USEFUL_UNIGRAM_ORG
IN_LEXICON_PER	USEFUL_BIGRAM_ORG
IN_LEXICON_LOC	USEFUL_SUFFIX_PER
IN_LEXICON_ORG	USEFUL_SUFFIX_LOC
IS_CAPITALIZED	USEFUL_SUFFIX_ORG

Click to edit Master subtitle style

11/30/11

Click to edit Master subtitle style

11/30/11

## NER: How else can we improve?

- Unlabeled data!

With more unlabeled data

instance 1: ... headquartered in (Washington State)<sup>L</sup> ...  
 instance 2: ... (Mr. Washington)<sup>P</sup>, the vice president of ...  
 instance 3: ... headquartered in (Kazakhstan) ...  
 instance 4: ... flew to (Kazakhstan) ...  
 instance 5: ... (Mr. Smith), a partner at Steptoe & Johnson ...  
 test: ... (Robert Jordan), a partner at ...  
 test: ... flew to (China) ...

*example from Jerry Zhu*

## NER : Challenges

- Domain transfer  
 WSJ → NYT  
 WSJ → Blogs ??  
 WSJ → Twitter ???!
- Tough nut: Organizations
- Non textual data?

Entity Extraction is a Boring Solved Problem – or is it?  
 (Vilain, Su and Lubar, 2007)

## NER: Related application

- Extracting real estate information from Craigslist Ads

Our oversized one, two and three bedroom apartment homes with floor plans featuring 1 and 2 baths offer space unlike any competition. Relax and enjoy the views from your own private balcony or patio, or feel free to entertain, with plenty of space in your large living room, dining area and eat-in kitchen. The lovely pool and sun deck make summer fun a splash. Our location makes commuting a breeze – Near MTA bus lines, the Metro station, major shopping areas, and for the little ones, an elementary school is right next door.

## NER: Related Application

- BioNLP: Annotation of chemical entities

Type	Description	Example
CM	chemical compound	citric acid
RN	chemical reaction	1,3-dimethylation
CJ	chemical adjective	pyrazolic
ASE	enzyme	methylase
CPR	chemical prefix	1,3-

*Corbet, Batchelor & Teufel, 2007*

## Shared Tasks: NLP in practice

- Shared Task
  - Everybody works on a (mostly) common dataset
  - Evaluation measures are defined
  - Participants get ranked on the evaluation measures
  - Advance the state of the art
  - Set benchmarks
- Tasks involve common hard problems or new interesting problems