

Sequence Labeling

CS 690N, Spring 2017

Advanced Natural Language Processing

<http://people.cs.umass.edu/~brenocon/anlp2017/>

Brendan O'Connor

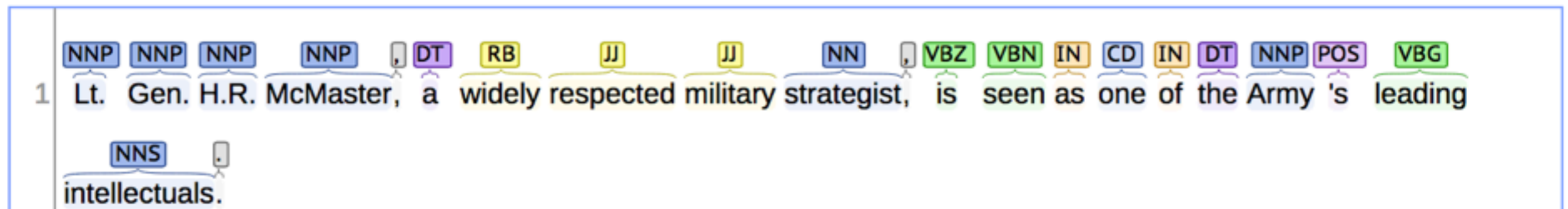
College of Information and Computer Sciences

University of Massachusetts Amherst

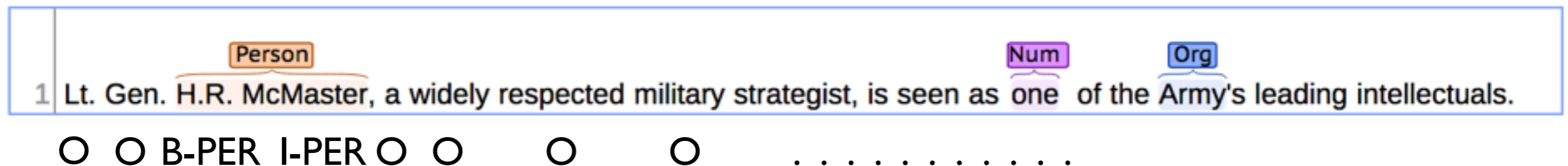
- Sequence labeling problems
 - Part of speech tags
- Models: HMMs, CRFs

- Sequence labeling: from $x_1..x_n$, predict tags $y_1..y_n$
- Named entity recognition:
an example of *span recognition*
- BIO tags allow treatment as a sequence labeling problem

Part-of-Speech:

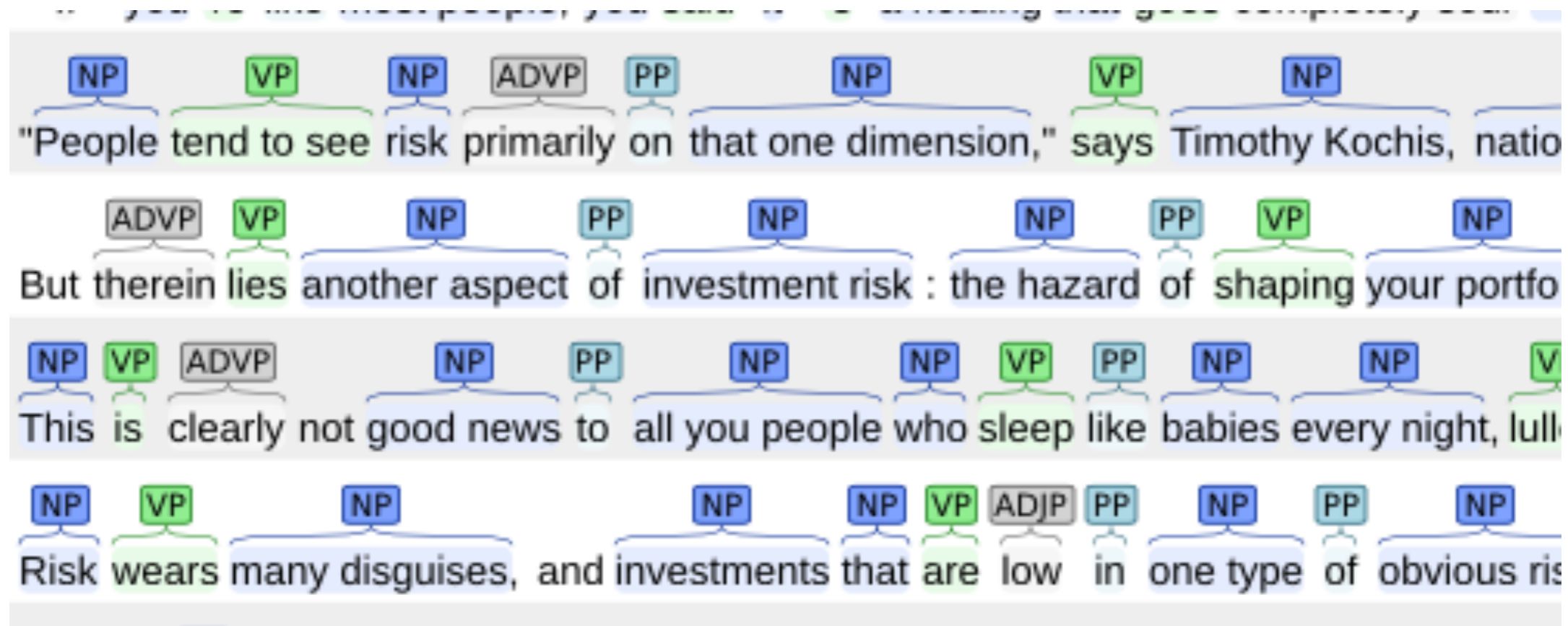


Named Entity Recognition:



More span labeling tasks

- Syntactic chunking



- Biological entities

Characterization of undifferentiated spc human cell type ES cells and differentiated EBs anat by antibodies. All monoclonals initially selected for their abilities to recognize recombinant proteins in direct ELISAs.

A subset were also tested by Western Blot analysis using recombinant proteins and cell lysate to confirm epitope.

The best clone was later screened for its applications for immunocytochemistry and flow cytometry us

spc Human anatomy peripheral blood component platelets were used for screening spc mouse anti-spc human gene CD9 antibody.

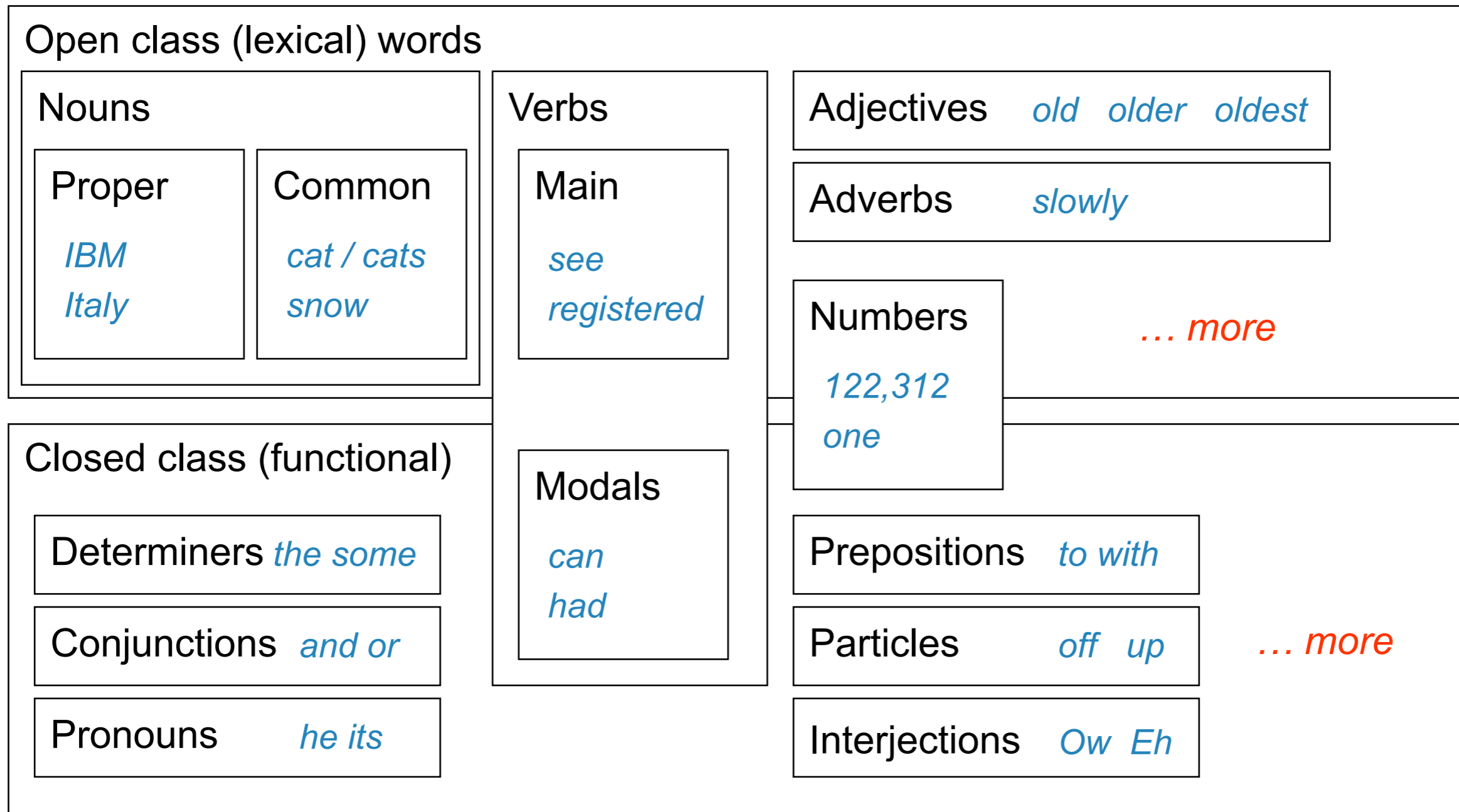
c line MCF-7 cells were used for screening spc mouse anti-spc human gene E-Cadherin and gene PODXL (gene or protein podocalyxin-like) an

c line MG-63 cells were used for screening spc mouse anti-spc human gene GATA1 (gene or protein GATA binding protein 1) antibody.

What's a part-of-speech (POS)?

- Syntax = how words compose to form larger meaning-bearing units
- POS = syntactic categories for words
 - You could substitute words within a class and have a syntactically valid sentence.
 - Give information how words can combine.
 - I saw the dog
 - I saw the cat
 - I saw the {table, sky, dream, school, anger, ...}
- (Phrasal/constituent categories generalize this idea. POS tags are constrained to single words.)

Open vs closed classes



ikr smh he asked fir yo last

name so he can add u on

fb lololol

ikr	smh	he	asked	fir	yo	last
!	G	O	V	P	D	A
name	so	he	can	add	u	on
N	P	O	V	V	O	P
fb	lololol					
^	!					

Why do we want POS?

- Useful for many syntactic and other NLP tasks.
 - Phrase identification (“chunking”)
 - Named entity recognition
 - Full parsing
 - Sentiment
- Especially when there’s a low amount of training data

POS patterns: sentiment

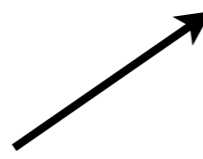
- Turney (2002): identify bigram phrases, from unlabeled corpus, useful for sentiment analysis.

Table 1. Patterns of tags for extracting two-word phrases from reviews.

	First Word	Second Word	Third Word (Not Extracted)
1.	JJ	NN or NNS	anything
2.	RB, RBR, or RBS	JJ	not NN nor NNS
3.	JJ	JJ	not NN nor NNS
4.	NN or NNS	JJ	not NN nor NNS
5.	RB, RBR, or RBS	VB, VBD, VBN, or VBG	anything



(plus sentiment PMI stuff)



POS patterns: sentiment

- Turney (2002): identify bigram phrases, from unlabeled corpus, useful for sentiment analysis.

Table 1. Patterns of tags for extracting two-word phrases from reviews.

	First Word	Second Word	Third Word (Not Extracted)
1.	JJ	NN or NNS	anything
2.	RB, RBR, or RBS	JJ	not NN nor NNS
3.	JJ	JJ	not NN nor NNS
4.	NN or NNS	JJ	not NN nor NNS
5.	RB, RBR, or RBS	VB, VBD, VBN, or VBG	anything

Table 2. An example of the processing of a review that the author has classified as *recommended*.⁶

Extracted Phrase	Part-of-Speech Tags	Semantic Orientation
online experience	JJ NN	2.253
low fees	JJ NNS	0.333
local branch	JJ NN	0.421
small part	JJ NN	0.053
online service	JJ NN	2.780
printable version	JJ NN	-0.705
direct deposit	JJ NN	1.288
well other	RB JJ	0.237
inconveniently	RB VBN	-1.541
located		
other bank	JJ NN	-0.850
true service	JJ NN	-0.732

(plus sentiment PMI stuff)

POS patterns: simple noun phrases

- Quick and dirty noun phrase identification

<http://brenocon.com/JustesonKatz1995.pdf>

<http://brenocon.com/handler2016phrases.pdf>

Grammatical structure: Candidate strings are those multi-word noun phrases that are specified by the regular expression $((A | N)^+ | ((A | N)^*(NP)^?)(A | N)^*)N$,

Tag Pattern	Example
A N	<i>linear function</i>
N N	<i>regression coefficients</i>
A A N	<i>Gaussian random variable</i>
A N N	<i>cumulative distribution function</i>
N A N	<i>mean squared error</i>
N N N	<i>class probability function</i>
N P N	<i>degrees of freedom</i>

Table 5.2 Part of speech tag patterns for collocation filtering. These patterns were used by Justeson and Katz to identify likely collocations among frequently occurring word sequences.

POS Tagging: lexical ambiguity

Can we just use a tag dictionary
(one tag per word type)?

Types:		WSJ	Brown
Unambiguous	(1 tag)	44,432 (86%)	45,799 (85%)
Ambiguous	(2+ tags)	7,025 (14%)	8,050 (15%)

Most words types
are unambiguous ...

- Ambiguous wordtypes tend to be very common ones.
- I know **that** he is honest = IN (relativizer)
- Yes, **that** play was nice = DT (determiner)
- You can't go **that** far = RB (adverb)

POS Tagging: lexical ambiguity

Can we just use a tag dictionary
(one tag per word type)?

Types:

		WSJ	Brown
Unambiguous	(1 tag)	44,432 (86%)	45,799 (85%)
Ambiguous	(2+ tags)	7,025 (14%)	8,050 (15%)

Most words types
are unambiguous ...

Tokens:

Unambiguous	(1 tag)	577,421 (45%)	384,349 (33%)
Ambiguous	(2+ tags)	711,780 (55%)	786,646 (67%)

But not so for
tokens!

- Ambiguous wordtypes tend to be very common ones.
- I know **that** he is honest = IN (relativizer)
- Yes, **that** play was nice = DT (determiner)
- You can't go **that** far = RB (adverb)

Need careful guidelines (and do annotators always follow them?)

PTB POS guidelines, Santorini (1990)

4 Confusing parts of speech

This section discusses parts of speech that are easily confused and gives guidelines on how to tag such cases.

CD or JJ

Number-number combinations should be tagged as adjectives (JJ) if they have the same distribution as adjectives.

EXAMPLES: a 50–3/JJ victory (cf. a handy/JJ victory)

Hyphenated fractions *one-half*, *three-fourths*, *seven-eighths*, *one-and-a-half*, *seven-and-three-eighths* should be tagged as adjectives (JJ) when they are prenominal modifiers, but as adverbs (RB) if they could be replaced by *double* or *twice*.

EXAMPLES: one-half/JJ cup; cf. a full/JJ cup
one-half/RB the amount; cf. twice/RB the amount; double/RB the amount

Some other lexical ambiguities

- Prepositions (P) versus verb particles (T)
 - turn into/P a monster
 - take out/T the trash
 - check it out/T, what's going on/T, shout out/T

Test:

turn slowly into a monster

*take slowly out the trash

Careful annotator guidelines are necessary to define what to do in many cases.

- http://repository.upenn.edu/cgi/viewcontent.cgi?article=1603&context=cis_reports
- http://www.ark.cs.cmu.edu/TweetNLP/annot_guidelines.pdf

Some other lexical ambiguities

- Prepositions (P) versus verb particles (T)
 - turn into/P a monster
 - take out/T the trash
 - check it out/T, what's going on/T, shout out/T
- this,that -- pronouns versus determiners
 - i just orgasmed over this/O
 - this/D wind is serious

Test:

turn slowly into a monster

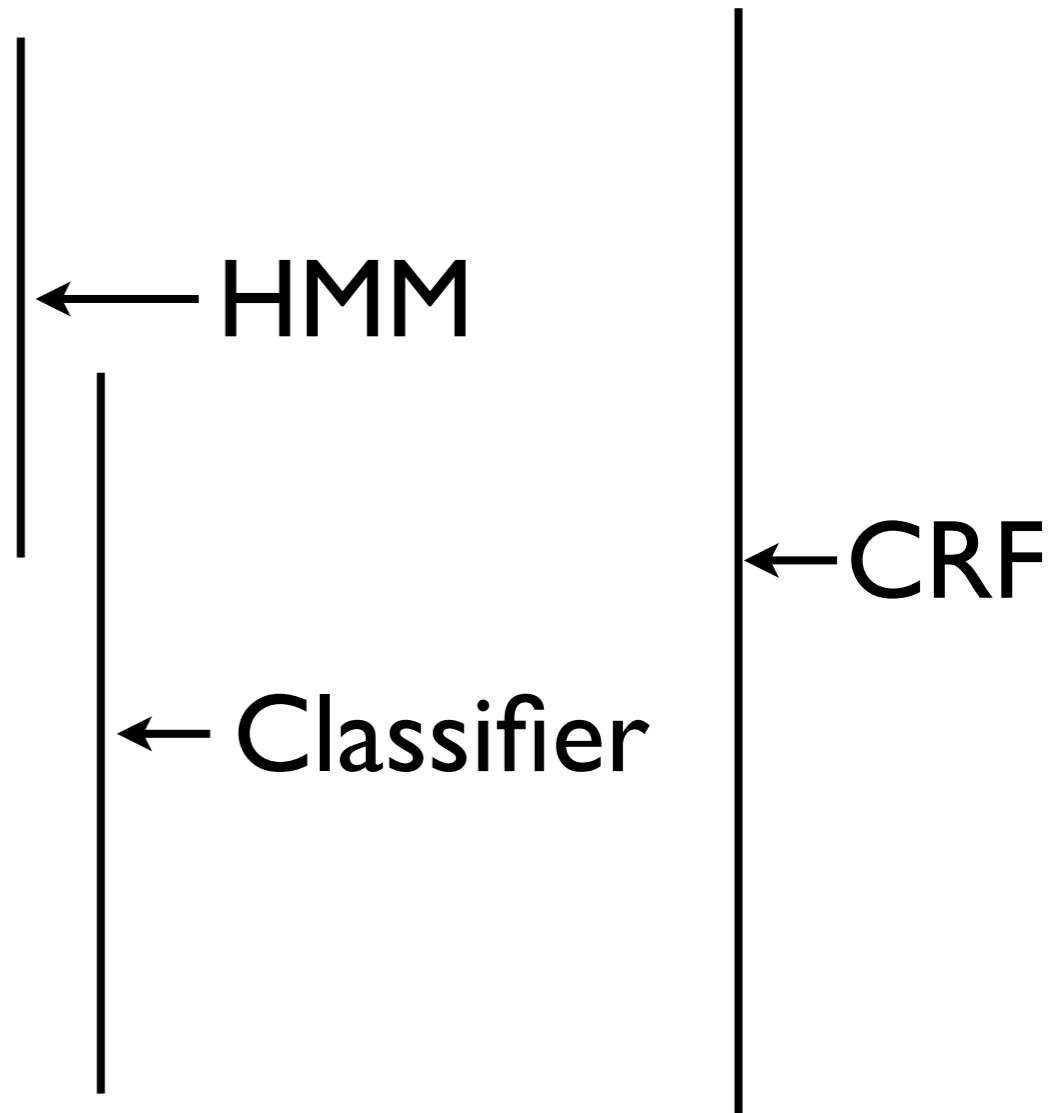
*take slowly out the trash

Careful annotator guidelines are necessary to define what to do in many cases.

- http://repository.upenn.edu/cgi/viewcontent.cgi?article=1603&context=cis_reports
- http://www.ark.cs.cmu.edu/TweetNLP/annot_guidelines.pdf

How to build a POS tagger?

- Sources of information:
 - POS tags of surrounding words: syntactic context
 - The word itself
 - Features!
 - Word-internal information
 - External lexicons
 - Features from surrounding words



Sequence labeling

- Seq. labeling as classification
 - Each position m gets an independent classification

$$p(y_m \mid w_1 \dots w_n)$$

$$\arg \max_y \theta^T \mathbf{f}((\mathbf{w}, m), y)$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 1), \mathbf{N}) = \langle \text{they}, \mathbf{N} \rangle$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 2), \mathbf{V}) = \langle \text{can}, \mathbf{V} \rangle$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 3), \mathbf{V}) = \langle \text{fish}, \mathbf{V} \rangle.$$

Sequence labeling

- Seq. labeling as classification
 - Each position m gets an independent classification

$$p(y_m \mid w_1..w_n)$$

$$\arg \max_y \theta^T \mathbf{f}((\mathbf{w}, m), y)$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 1), \mathbf{N}) = \langle \text{they}, \mathbf{N} \rangle$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 2), \mathbf{V}) = \langle \text{can}, \mathbf{V} \rangle$$

$$\mathbf{f}((\mathbf{w} = \text{they can fish}, m = 3), \mathbf{V}) = \langle \text{fish}, \mathbf{V} \rangle.$$

- But syntactic (tag) context is sometimes necessary!
 - *The old man the boat* [garden path sentence]

- Seq. labeling as ***structured prediction***

$$\hat{\mathbf{y}}_{1:M} = \operatorname{argmax}_{\mathbf{y}_{1:M} \in \mathcal{Y}(\mathbf{w}_{1:M})} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}),$$

- **Hidden Markov model**

- Fully generative, simple sequence model
- Supports many operations
 - $P(\mathbf{w})$: Likelihood (generative model)
 - Forward algorithm
 - $P(\mathbf{y} \mid \mathbf{w})$: Predicted sequence (“decoding”)
 - Viterbi algorithm
 - $P(y_m \mid \mathbf{w})$: Predicted tag marginals
 - Forward-Backward algorithm
- The HMM is a type of log-linear model

Forward algorithm

- on board
- stopped here 2/21

Viterbi algorithm

- If the feature function decomposes into local features, dynamic programming gives global solution

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) \quad \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m).$$

- Decompose:

$$\max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) = \max_{\mathbf{y}_{1:M}} \sum_{m=1}^M \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

- Define Viterbi variables:

$$v_m(k) \triangleq \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n)$$

Viterbi algorithm

- If the feature function decomposes into local features, dynamic programming gives global solution

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) \quad \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m).$$

- Decompose:

$$\begin{aligned} \max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^M \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \max_{\mathbf{y}_{1:M}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_M, y_{M-1}, M) + \sum_{m=1}^{M-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \end{aligned}$$

- Define Viterbi variables:

$$v_m(k) \triangleq \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n)$$

Viterbi algorithm

- If the feature function decomposes into local features, dynamic programming gives global solution

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) \quad \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m).$$

- Decompose:

$$\begin{aligned} \max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^M \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \max_{\mathbf{y}_{1:M}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_M, y_{M-1}, M) + \sum_{m=1}^{M-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \max_{y_M} \max_{y_{M-1}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_M, y_{M-1}, M) + \max_{\mathbf{y}_{1:M-2}} \sum_{m=1}^{M-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m). \end{aligned}$$

- Define Viterbi variables:

$$v_m(k) \triangleq \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n)$$