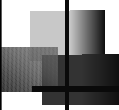


6.863J Natural Language Processing

Lecture 7: The Red Pill or the Blue Pill, Episode 2: part-of-speech tagging



Instructor: Robert C. Berwick
berwick@ai.mit.edu

The Menu Bar



• Administrivia:

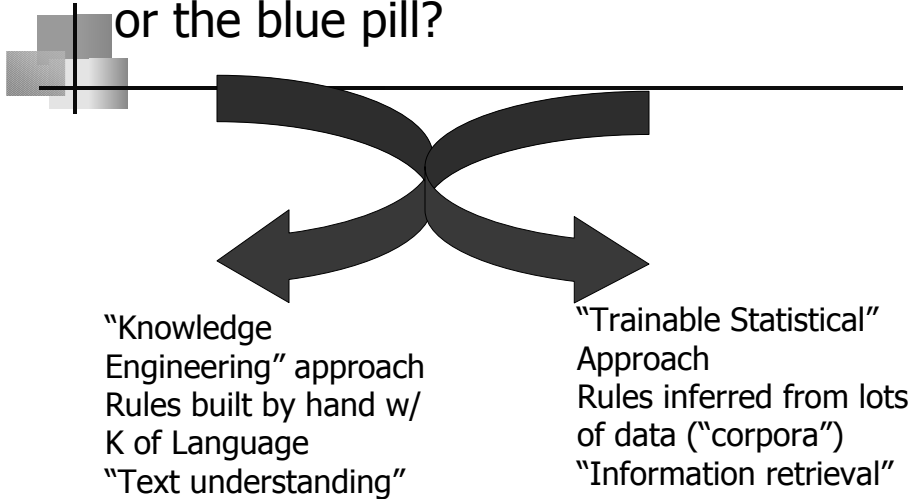
- Schedule alert: Lab1b due today
- Lab 2b released, this Weds (later today)

Agenda:

Red vs. Blue:

- Part of speech 'tagging' via statistical models
- Part of speech tagging via rules
- Ch. 6 & 8 in Jurafsky

The Great Divide in NLP: the red pill or the blue pill?



6.863J/9.611J SP04 Lecture 7

Two approaches

1. ~~Statistical model~~
 2. Deterministic baseline tagger composed with a cascade of fixup transducers
- These two approaches are the guts of Lab 2
(lots of others methods: decision trees, ...)

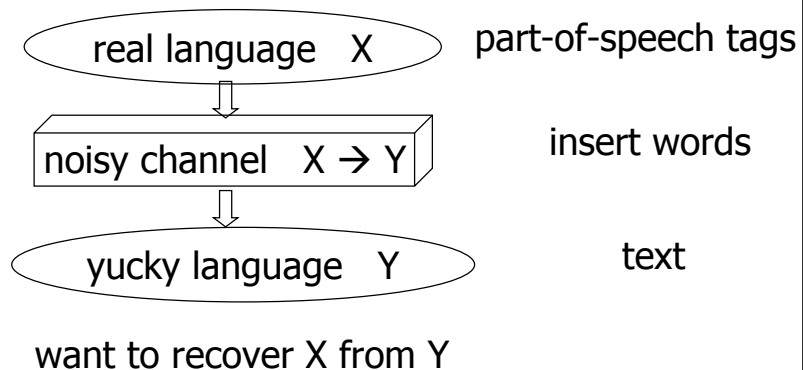
6.863J/9.611J SP04 Lecture 7

The problem

- In unseen data, we wish to find the part of speech tags
- The set of part of speech tags are decided by experts

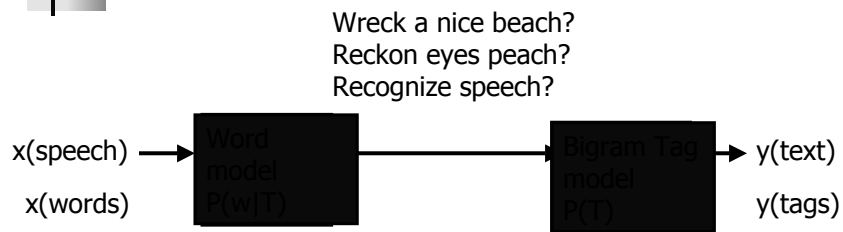
6.863J/9.611J SP04 Lecture 7

Noisy Channel Muddle (statistical)



6.863J/9.611J SP04 Lecture 7

A picture: the statistical, noisy channel view



6.863J/9.611J SP04 Lecture 7

Formulation, in general

6.863J/9.611J SP04 Lecture 7

General probabilistic decision problem

- E.g.: data = bunch of text
 - label = language
 - label = topic
 - label = author
- E.g.2: (sequential prediction)
 - label = translation or summary of entire text
 - label = part of speech of current word
 - label = identity of current word (ASR) or character (OCR)

6.863J/9.611J SP04 Lecture 7

Language models – statistical view

- Application to speech recognition (and parsing, generally)
 - x = Input (speech/words)
 - y = output (text/Tags)
 - We want to find $\max P(y|x)$ Problem: we don't know the tags – that is what we want to find!
 - Solution: We have an estimate of $P(y)$ [the language model] and $P(x|y)$ [the prob. of some sound/words given text/Tags = an acoustic model or Tag model]

6.863J/9.611J SP04 Lecture 7



Finding inner form given outside:

From Bayes' law, we have,
 $\max P(y|x) = \max P(x|y) \cdot P(y) =$
 $\max \text{Pr acoustic model} \times \text{lang model}$
(hold $P(x)$ fixed, i.e., $P(x|y) \cdot P(y) / P(x)$, but
max is same for both)

So, in tagging case, we have a word model
instead - so we find $\max P(\text{tags}|w)$ from:
 $\max P(\text{words}|\text{tags}) \cdot P(\text{tags})$

6.863J/9.611J SP04 Lecture 7



HMM for POS tagging

- In a Hidden Markov model, it is hypothesized that there is an underlying finite state machine (not directly observable, hence hidden) that changes state with each input element
- For us, the hidden states are the tags, and the input elements are the words

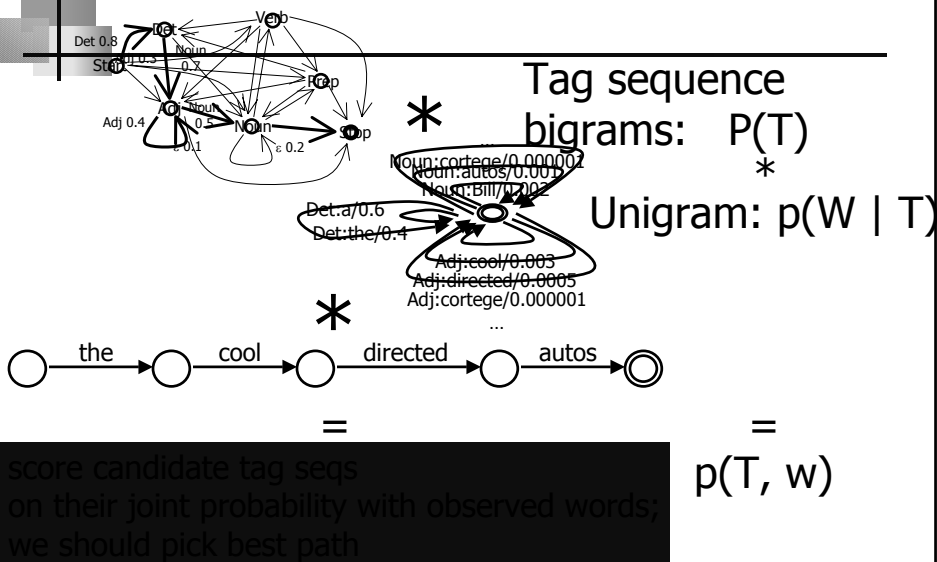
6.863J/9.611J SP04 Lecture 7

Hidden Markov Model tagging for POS

- Prob(Tag sequence) – based on n-grams: train on marked up, tagged text
- Prob(W|T) – unigram prob, based on tagged text
- Prob(T, w) computed from Viterbi trellis computation - max over all possible tag sequence paths, and 'emission' probabilities of word, tag combination
- Unseen tag sequence

6.863J/9.611J SP04 Lecture 7

Cartoon form Review

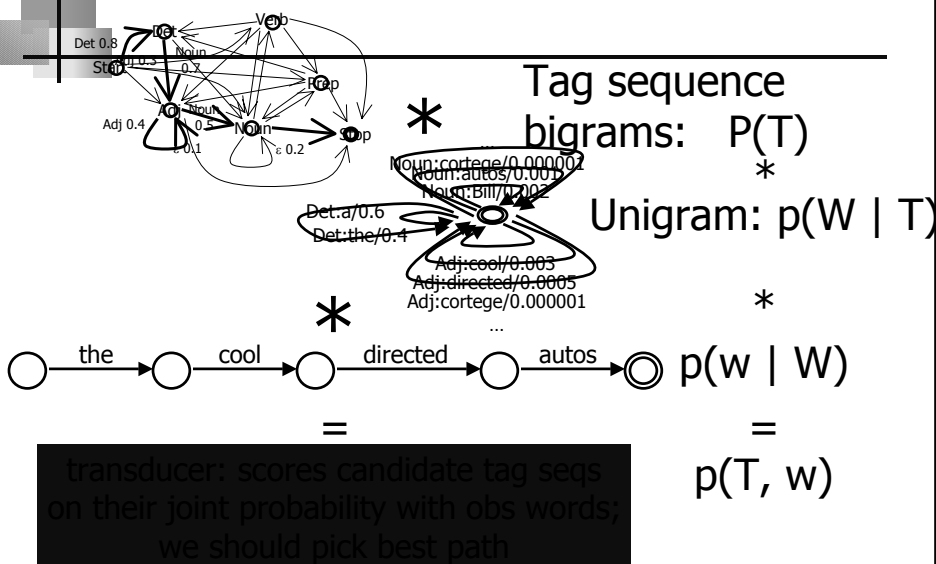


HMM construction

- Hidden state transition model governs observed word sequences
- Transitions probabilistic
- Estimate transition probabilities from an annotated corpus state 's' = tag state
 - $P(s_j | s_{j-1}, w_j)$
 - Based just on prior state and current word seen (hence Markovian assumption)
- At runtime, find maximum likelihood path through the network, using a max-flow algorithm (Viterbi)

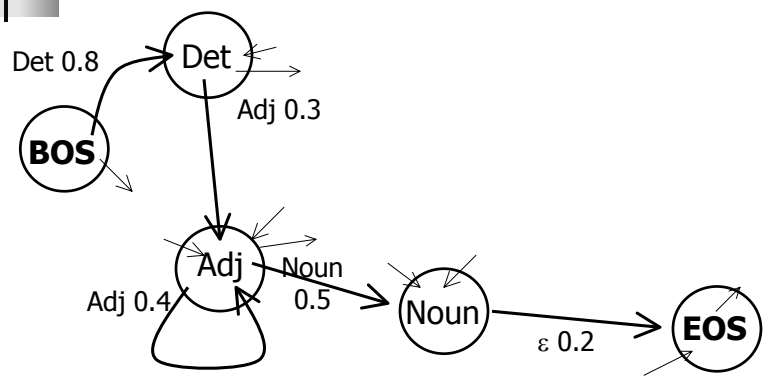
6.863J/9.611J SP04 Lecture 7

Cartoon form Review



P(T) - Tag bigram picture

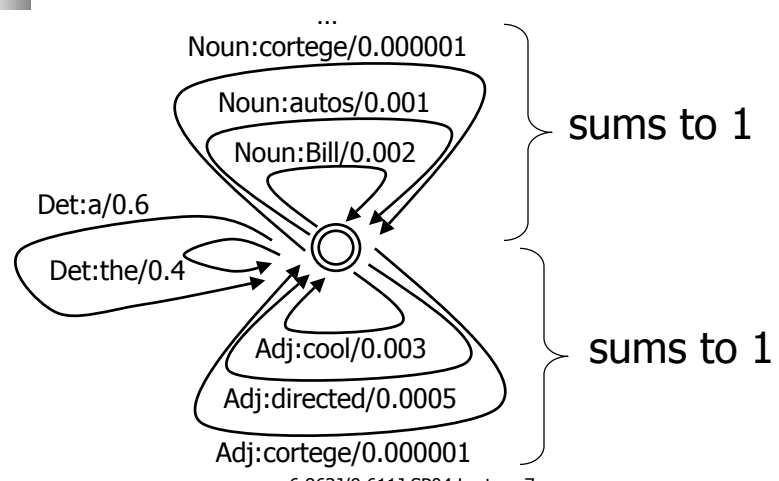
$p(\text{tag seq})$



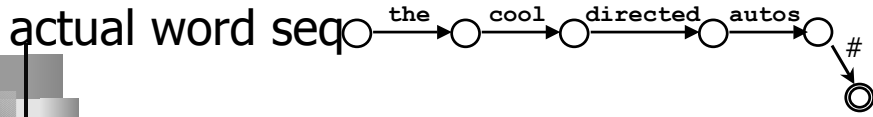
BOS Det Adj Adj Noun EOS = $0.8 * 0.3 * 0.4 * 0.5 * 0.2$

Unigram replacement model

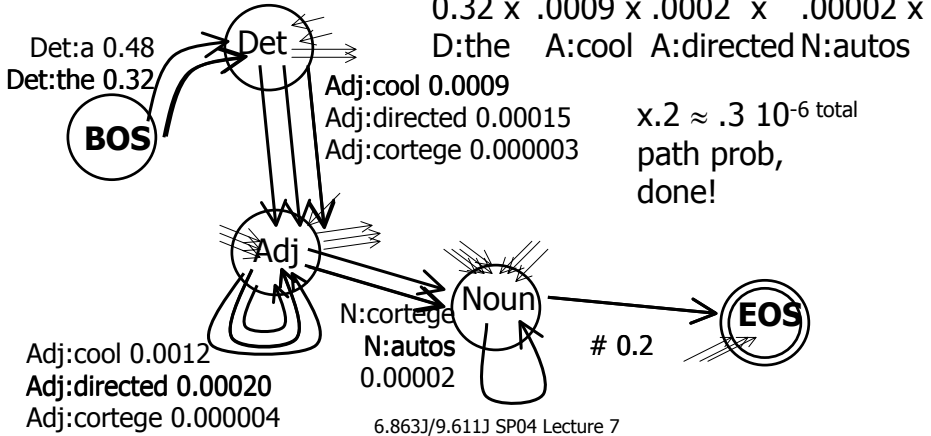
$P(\text{word} | \text{tag})$



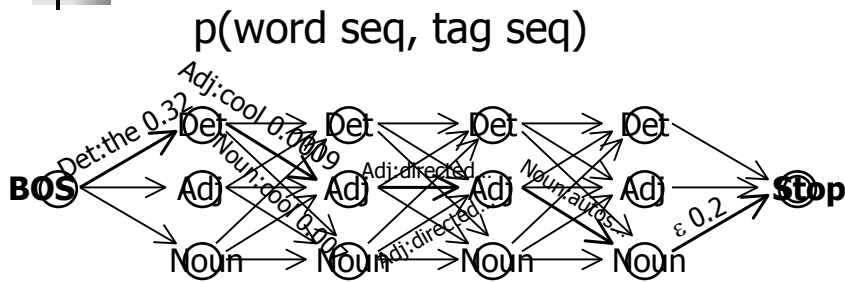
Compose with



$$p(\text{word seq, tag seq}) = p(\text{tag seq}) * p(\text{word seq} | \text{tag seq})$$



Unroll the fsa - All paths together form 'trellis'

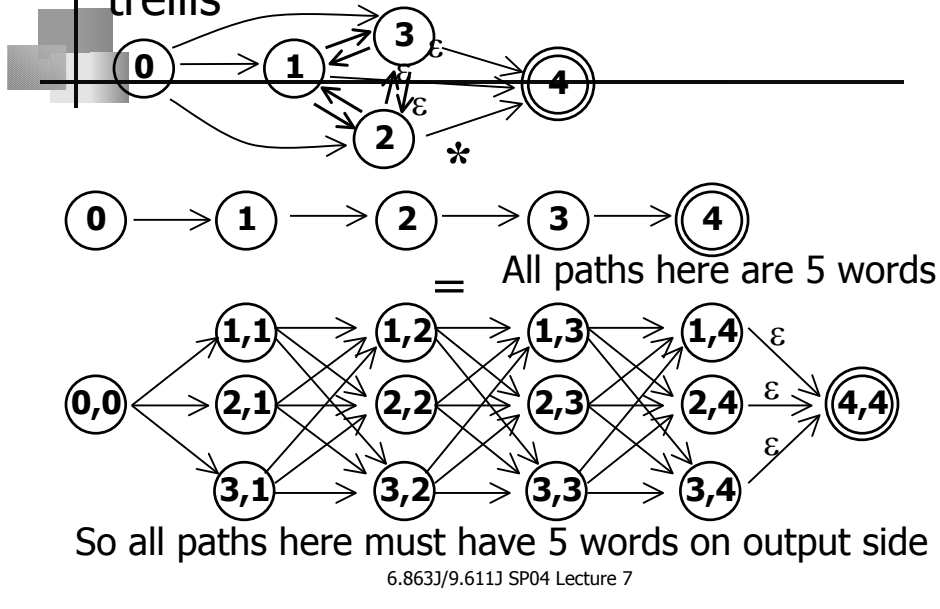


WHY?

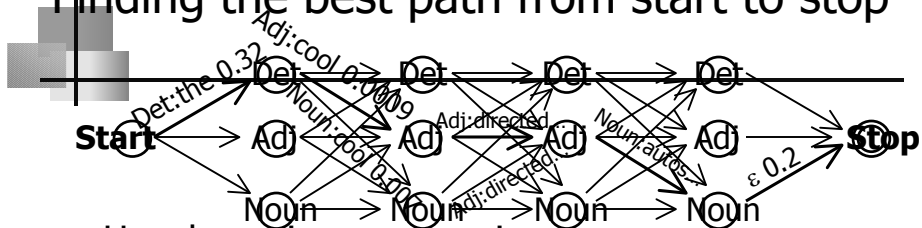
The best path:

BOS Det Adj Adj Noun **EOS** = $0.32 * 0.0009 \dots$
 the cool directed autos

Cross-product construction forms trellis



Finding the best path from start to stop



- Use dynamic programming
- What is best path from Start to each node?
 - Work from left to right
 - Each node stores its best path from Start (as probability plus one backpointer)
- Special acyclic case of Dijkstra's shortest-path algorithm
- Faster if some arcs/states are absent

Method to find max probability path: Viterbi algorithm

- For each path reaching state s at step (word) w , we compute a path probability. We call the max of these viterbi(s,w)
- [Base step] Compute $\text{viterbi}(0,0)=1$
- [Induction step] Compute $\text{viterbi}(s',w+1)$, assuming we know $\text{viterbi}(s,w)$ for all s

6.863J/9.611J SP04 Lecture 7

Viterbi recursion

$$\text{path-prob}(s'|s,w) = \text{viterbi}(s,w) * a[s,s']$$

probability of path to s' through s max path score * transition prob
for state s at word w $s \rightarrow s'$

$$\text{viterbi}(s',w+1) = \max_{s \text{ in STATES}} \text{path-prob}(s' | s,w)$$

↓
Bi-gram POS
probability

6.863J/9.611J SP04 Lecture 7

Method...

- This is *almost* correct...but again, we need to also factor in the *unigram* prob of a state s' 'emitting' a word w given an observed surface word w , or $b(o_w)$ at tag state s'
- So the correct formula for the path prob is:
$$\text{path-prob}(s'|s,w) = \text{viterbi}(s,w) * a[s,s'] * b_{s'}(o_w)$$

bigram unigram

6.863J/9.611J SP04 Lecture 7

Or as in your text...p. 179 (NB: t used instead of w for index)

function VITERBI(*observations* of len T , *state-graph*) **returns** *best-path*

$\text{num-states} \leftarrow \text{NUM-OF-STATES}(\text{state-graph})$

Create a path probability matrix $\text{viterbi}[\text{num-states}+2, T+2]$

$\text{viterbi}[0,0] \leftarrow 1.0$

for each time step t **from** 0 **to** T **do**

for each state s **from** 0 **to** num-states **do**

for each transition s' from s specified by *state-graph*

$\text{new-score} \leftarrow \text{viterbi}[s, t] * a[s,s'] * b_{s'}(o_t)$

if $((\text{viterbi}[s', t+1] = 0) \parallel (\text{new-score} > \text{viterbi}[s', t+1]))$

then

$\text{viterbi}[s', t+1] \leftarrow \text{new-score}$

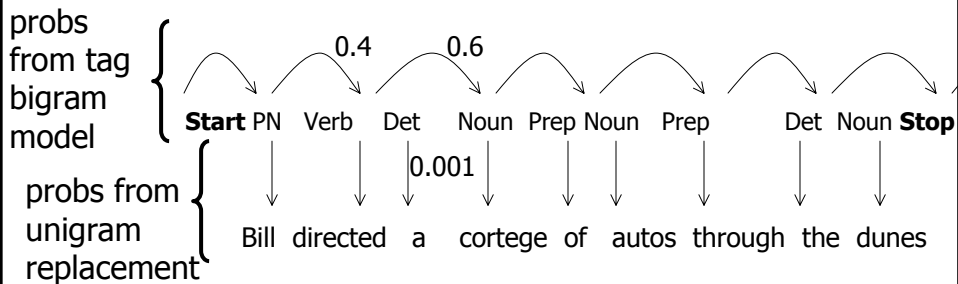
$\text{back-pointer}[s', t+1] \leftarrow s$

Backtrace from highest probability state in the final column of $\text{viterbi}[]$ and return path

6.863J/9.611J SP04 Lecture 7

Summary

- We are modeling $p(\text{word seq}, \text{tag seq})$
- The tags are hidden, but we see the words
- Is tag sequence X likely with these words?
- This is a "Hidden Markov Model":



- Find X that maximizes probability product

6.863J/9.611J SP04 Lecture 7

How much data is needed?

- System performance bears a roughly log-linear relationship to the training data quantity, at least up to about 1.2 million words
- Obtaining 1.2 million words of training data requires transcribing and annotating approximately 200 hours of broadcast news programming, or if annotating text, this would amount to approximately 1,777 average-length Wall Street Journal articles

6.863J/9.611J SP04 Lecture 7



Other tasks

6.863J/9.611J SP04 Lecture 7



If you think POS tagging is not relevant, then...

- Microsoft announced plans to include “Smart Tags” in its browser and other products. This is a feature that automatically inserts hyperlinks from concepts in text to related web pages chosen by Microsoft.
- The best way to make automatic hyperlinking unbiased is to base it on an unbiased source of web pages, such as Google.

6.863J/9.611J SP04 Lecture 7

How to do this?

- The main technical problem is to find pieces of text that are good concept anchors... like names!
- So: Given a text, find the starting and ending points of all the names. Depending on our specific goals, we can include the names of people, places, organizations, artifacts (such as product names), etc.
- . Once we have the anchor text, we can send it to a search engine, retrieve a relevant URL (or set of URLs, once browsers can handle multi-way hyperlinks), and insert them into the original text on the fly.

6.863J/9.611J SP04 Lecture 7

Example – “Message understanding” (MUC)

ST1-MUC3-0011

SANTIAGO, 18 MAY 90 (RADIO COOPERATIVA NETWORK) -- [REPORT] [JUAN ARAYA]
[TEXT]

EDMUNDO VARGAS CARRENO, CHILEAN FOREIGN MINISTRY UNDER SECRETARY, HAS STATED THAT THE BRYANT TREATY WITH THE UNITED STATES WILL BE APPLIED IN THE LETELIER CASE ONLY TO COMPENSATE THE RELATIVES OF THE FORMER CHILEAN FOREIGN MINISTER MURDERED IN WASHINGTON AND THE RELATIVES OF HIS U.S. SECRETARY, RONNIE MOFFIT. THE CHILEAN FOREIGN UNDER SECRETARY MADE THIS STATEMENT IN REPLY TO U.S. NEWSPAPER REPORTS STATING THAT THE TREATY WOULD BE PARTIALLY RESPECTED.

FOLLOWING ARE VARGAS CARRENO'S STATEMENTS AT A NEWS CONFERENCE HE HELD IN BUENOS AIRES BEFORE CONCLUDING HIS OFFICIAL VISIT TO ARGENTINA:

6.863J/9.611J SP04 Lecture 7

0. MESSAGE ID	TST: MUC3 001
1. MESSAGE TEMPLATE	1
2. INCIDENT: DATE	18 MAY 90
3. INCIDENT: LOCATION	UNITED STATES: WASHINGTON D.C. (CITY)
4. INCIDENT: TYPE	ATTACK
5. INCIDENT: STAGE OF EXECUTION	ACCOMPLISHED
6. INCIDENT: INSTRUMENT ID	-
7. INCIDENT: INSTRUMENT TYPE	-
8. PERP: INCIDENT CATEGORY	STATE-SPONSORED VIOLENCE
9. PERP: INDIVIDUAL ID	-
10. PERP: ORGANIZATION ID	"CHILEAN GOVERNMENT"
11. PERP: ORGANIZATION CONFIDENCE REPORTED AS FACT:	"CHILEAN GOVERNMENT"
12. PHYS TGT: ID	-
13. PHYS TGT: TYPE	-
14. PHYS TGT: NUMBER	-
15. PHYS TGT: FOREIGN NATION	-
16. PHYS TGT: EFFECT OF INCIDENT	-
17. PHYS TGT: TOTAL NUMBER	-
18. HUM TGT: NAME	"ORLANDO LETELIER" "RONNIE MOFFIT"
19. HUM TGT: DESCRIPTION	"FORMER CHILEAN FOREIGN MINISTER": "ORLANDO LETELIER" "U.S. SECRETARY" / "ASSISTANT" / "SECRETARY": "RONNIE MOFFIT"
20. HUM TGT: TYPE	GOVERNMENT OFFICIAL: "ORLANDO LETELIER"
21.	6.863J/9.611J SP04 Lecture 7

Recognizing domain patterns

- Match domain patterns against complex phrase heads
 - <company> <form><joint venture> with <company>
 - <company> capitalized at <currency>

"Bridgestone Sports Co. said Friday it has set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be shipped to Japan.

"The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990."

Relationship:	TIE-UP
Entities:	"Bridgestone Sports Co." "a local concern" "a Japanese trading house"
JV Company:	--
Capitalization:	--

Relationship:	TIE-UP
Entities:	--
JV Company:	"Bridgestone Sports Taiwan Co."
Capitalization:	20000000 TWD

What about part of speech tagging here?

- Advantages
 - Ambiguity can be potentially reduced (but we shall see in our laboratory if this is true)
 - Avoid errors due to incorrect categorization of rare senses e.g., “has been” as noun
- Disadvantages
 - Errors taggers make often those you’d most want to eliminate
 - High performance requires training on similar genre
 - Training takes time

6.863J/9.611J SP04 Lecture 7

Proper names...

- Proper names are particularly important for extraction systems
- Because typically one wants to extract events, properties, and relations about some particular object, and that object is usually identified by its name

6.863J/9.611J SP04 Lecture 7



...A challenge...

- Problems though...
 - proper names are huge classes and it is difficult, if not impossible to enumerate their members
 - Hundreds of thousands of names of locations around the world
 - Many of these names are in languages other than the one in which the extraction system is designed

6.863J/9.611J SP04 Lecture 7

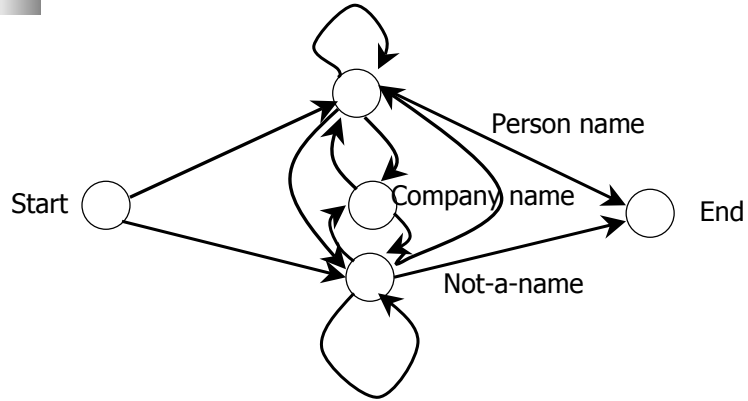


How are names extracted?

- (Hidden) Markov Model
- Hypothesized that there is an underlying finite state machine (not directly observable, hence hidden) that changes state with each input element
- probability of a recognized constituent is conditioned not only on the words seen, but the state that the machine is in at that moment
- "John" followed by "Smith" is likely to be a person, while "John" followed by "Deere" is likely to be a company (a manufacturer of heavy farming and construction equipment).

6.863J/9.611J SP04 Lecture 7

HMM statistical name tagger



6.863J/9.611J SP04 Lecture 7

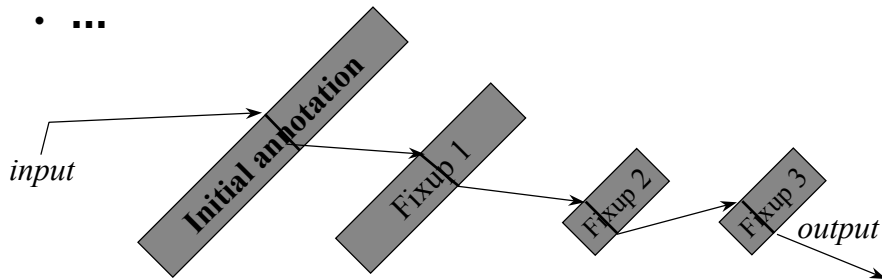
Method 2: Rule system (but learned)

- Error-based tagging

6.863J/9.611J SP04 Lecture 7

Another FST Paradigm: Successive Fixups

- Like successive markups but *alter*
- Morphology
- Phonology
- Part-of-speech tagging
- ...



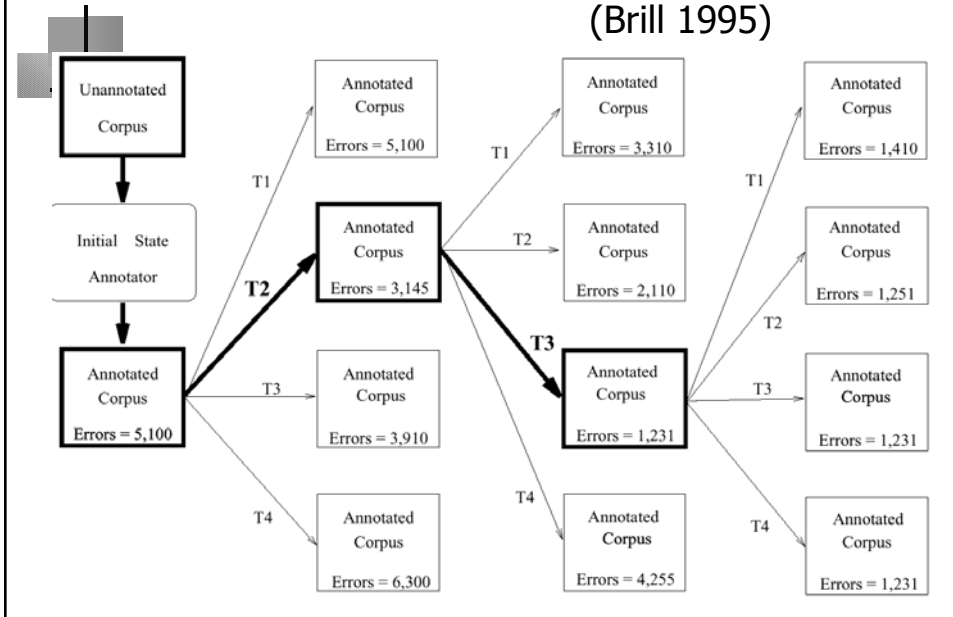
6.863J/9.611J SP04 Lecture 7

Brown/Upenn corpus tags

J. text,
p. 297
Fig 8.6
1M words
60K tag
counts

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(' or ")</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or ")</i>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, { , <</i>
PP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>([, } , ></i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... --)</i>
RP	Particle	<i>up, off</i>			

Transformation-Based Tagging (Brill 1995)



Transformations Learned

#	From	To	Condition
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

BaselineTag*
 NN @ → VB // TO _
 VBP @ → VB // ... _
 etc.

Compose this cascade of FSTs.

Get a big FST that does the initial tagging and the sequence of fixups "all at once."

Initial Tagging of OOV Words

#	Change Tag		Condition
	From	To	
1	NN	NNS	Has suffix -s
2	NN	CD	Has character .
3	NN	JJ	Has character -
4	NN	VBN	Has suffix -ed
5	NN	VBG	Has suffix -ing
6	??	RB	Has suffix -ly
7	??	JJ	Adding suffix -ly results in a word.
8	NN	CD	The word \$ can appear to the left.
9	NN	JJ	Has suffix -al
10	NN	VB	The word would can appear to the left.
11	NN	CD	Has character 0
12	NN	JJ	The word be can appear to the left.
13	NNS	JJ	Has suffix -us
14	NNS	VBZ	The word it can appear to the left.
15	NN	JJ	Has suffix -ble
16	NN	JJ	Has suffix -ic
17	NN	CD	Has character 1
18	NNS	NN	Has suffix -ss
19	??	JJ	Deleting the prefix un- results in a word
20	NN	JJ	Has suffix -ive

Laboratory 2

- Goals:
 1. Use both HMM and Brill taggers
 2. Find errors that both make
 3. Compare performance – use of kappa & 'confusion matrix'
 4. All the slings & arrows of corpora – use Wall Street Journal excerpts

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Instant Message 6.863J Syllabus Google BE768 Biologica Members WebMail Connections BizJournal SmartUpdate Mktplace Home

Bookmarks Location: http://www.ling.gu.se/~lager/Home/brilltagger_ui.html

Brill Tagger

Powered by μ -TBL Technology

Swedish
 English
 Russian

Text:

Secretariat is expected to race tomorrow

Trace

© Torbjörn Lager 1999, Russian tagger by Natalia Zinovjeva

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Instant Message 6.863J Syllabus Google BE768 Biologica Members WebMail Connections BizJournal SmartUpdate Mktplace Home

Bookmarks Location: <http://www.ling.gu.se/~lager/tagger.cgi?language=English&input=Secretariat+is+expected+to+race+tomorrow&trace=on>

Tokenization

Secretariat is expected to race tomorrow

Lexical lookup

Secretariat/NNP is/VBZ expected/VBN to/TO race/NN tomorrow/NN

Guessing

Contextual-rule application

Intermediate analysis:

Secretariat/NNP is/VBZ expected/VBN to/TO race/NN tomorrow/NN

Applied rule:

tag:NN>VB <- tag:TO&[-1].

Analysis

Secretariat/NNP is/VBZ expected/VBN to/TO race/VB tomorrow/NN

Transformation based tagging

- ~~Combines symbolic and stochastic approaches: uses machine learning to refine its tags, via several passes~~
- Analogy: painting a picture, use finer and finer brushes - start with broad brush that covers a lot of the canvas, but colors areas that will have to be repainted. Next layer colors less, but also makes fewer mistakes, and so on.
- Similarly: tag using broadest (most general) rule; then an narrower rule, that changes a smaller number of tags, and so on. (We haven't said how the rules are learned)
- First we will see how the TBL rules are *applied*

6.863J/9.611J SP04 Lecture 7

Contextual Rules

- *Change tag a to tag b when:*
- 1. The preceding (following) word is tagged z.
- 2. The word two before (after) is tagged z.
- 3. One of the two preceding (following) words is tagged z.
- 4. One of the three preceding (following) words is tagged z.
- 5. The preceding word is tagged z and the following word is tagged w.
- 6. The preceding (following) word is tagged z and the word two before (after) is tagged w.
- 7. The preceding (following) word x.

...

6.863J/9.611J SP04 Lecture 7

Lexical Rules

Change the tag of an unknown word (from X) to Y if:

- 1. Deleting the prefix (suffix) x , $|x| \leq 4$, results in a word (x is any string of length 1 to 4).
- 2. The first (last) (1,2,3,4) characters of the word are x .
- 3. Adding the character string x as a prefix (suffix) results in a word ($|x| \leq 4$).
- 4. Word w ever appears immediately to the left (right) of the word.
- 5. Character z appears in the word.

6.863J/9.611J SP04 Lecture 7

Example Lexical Rules

NN s fhassuf 1 NNS

change the tag of an unknown word from NN to NNS if it has suffix -s

webpages/NN → *webpages/NNS*

6.863J/9.611J SP04 Lecture 7

Example 2

NN - fchar JJ

change the tag of an unknown word from NN to JJ if it has character '-'

man-made, rule-based, three-year-old, etc.

6.863J/9.611J SP04 Lecture 7

Applying the rules

1. First label every word with its most-likely tag (as we saw, this gets 90% right...!) for example, in Brown corpus, *race* is most likely to be a Noun:

$$P(\text{NN}|\text{race}) = 0.98$$

$$P(\text{VB}|\text{race}) = 0.02$$

2. ...expected/VBZ to/T TO race/VB tomorrow/NN
...the/DT race/NN for/IN outer/JJ space/NN

3. Use transformational (learned) rules to change tags:

*Change **NN** to **VB** when the previous tag is **TO***

6.863J/9.611J SP04 Lecture 7

Initial Tagging of OOV Words

#	Change Tag		Condition
	From	To	
1	NN	NNS	Has suffix -s
2	NN	CD	Has character .
3	NN	JJ	Has character -
4	NN	VBN	Has suffix -ed
5	NN	VBG	Has suffix -ing
6	??	RB	Has suffix -ly
7	??	JJ	Adding suffix -ly results in a word.
8	NN	CD	The word \$ can appear to the left.
9	NN	JJ	Has suffix -al
10	NN	VB	The word would can appear to the left.
11	NN	CD	Has character 0
12	NN	JJ	The word be can appear to the left.
13	NNS	JJ	Has suffix -us
14	NNS	VBZ	The word it can appear to the left.
15	NN	JJ	Has suffix -ble
16	NN	JJ	Has suffix -ic
17	NN	CD	Has character 1
18	NNS	NN	Has suffix -ss
19	??	JJ	Deleting the prefix un- results in a word
20	NN	JJ	Has suffix -ive

OK, the proof is in the (supervised) learning pudding - How?

- 3 stages
 1. Start by labeling every word with most-likely tag
 2. Then examine every possible transformation, and selects one that results in most improved tagging
 3. Finally, re-tags data according to this rule
 4. Repeat 1-3 until some stopping criterion (no new improvement, or small improvement)
- Output is ordered list of transformations that constitute a tagging procedure

How this works

- Set of possible 'transforms' is infinite, e.g., "transform NN to VB if the previous word was *MicrosoftWindoze* & word *braindead* occurs between 17 and 158 words before *that*"
- To limit: start with small set of abstracted transforms, or *templates*

6.863J/9.611J SP04 Lecture 7

Templates used: Change *a* to *b* when...

The preceding (following) word is tagged **z**.
The word two before (after) is tagged **z**.
One of the two preceding (following) words is tagged **z**.
One of the three preceding (following) words is tagged **z**.
The preceding word is tagged **z** and the following word is tagged **w**.
The preceding (following) word is tagged **z** and the word two before (after) is tagged **w**.

Variables **a**, **b**, **z**, **w**, range over parts of speech

6.863J/9.611J SP04 Lecture 7

Examples of Contextual Rules

- NN VB PREVTAG TO
- change tag NN to tag VB when the preceding word is tagged TO
to/TO run/NN
would be changed to
to/TO run/VB
- VBP VB PREV1OR2OR3TAG MD
- Change tag VBP(verb, non-3rd person singular present) to VB(verb, base form) when one of the three preceding words is tagged MD (modal verb)

6.863J/9.611J SP04 Lecture 7

Method

1. Call `Get-best-transform` with list of potential templates; this calls
2. `Get-best-instance` which instantiates each template over all its variables (given specific values for where we are)
3. Try it out, see what score is (improvement over known tagged system -- supervised learning); pick best one locally

6.863J/9.611J SP04 Lecture 7

```

function TBL(corpus) returns transforms-queue
  INITIALIZE-WITH-MOST-LIKELY-TAGS(corpus)
  until end condition is met do
    templates ← GENERATE-POTENTIAL-RELEVANT-TEMPLATES
    best-transform ← GET-BEST-TRANSFORM(corpus, templates)
    APPLY-TRANSFORM(best-transform, corpus)
    ENQUEUE(best-transform-rule, transforms-queue)
  end
  return(transforms-queue)

function GET-BEST-TRANSFORM(corpus, templates) returns transform
  for each template in templates
    (instance, score) ← GET-BEST-INSTANCE(corpus, template)
    if (score > best-transform.score) then best-transform ← (instance, score)
  return(best-transform)

```

6.863J/9.611J SP04 Lecture 7

```

function GET-BEST-INSTANCE(corpus, template) returns transform
  for from-tag ← from tag-1 to tag-n do
    for to-tag ← from tag-1 to tag-n do
      for pos ← from 1 to corpus-size do
        if (correct-tag(pos) == to-tag && current-tag(pos) == from-tag)
          num-good-transforms(current-tag(pos-1))++
        elseif (correct-tag(pos) == from-tag && current-tag(pos) == from-tag)
          num-bad-transforms(current-tag(pos-1))++
      end
      best-Z ← ARGMAXt(num-good-transforms(t) - num-bad-transforms(t))
      if (num-good-transforms(best-Z) - num-bad-transforms(best-Z)
          > best-instance.Z) then
        best-instance ← "Change tag from from-tag to to-tag
                       if previous tag is best-Z"
  return(best-instance)

```

```

procedure APPLY-TRANSFORM(transform, corpus)
  for pos ← from 1 to corpus-size do
    if (current-tag(pos) == best-rule-from)
      && (current-tag(pos-1) == best-rule-prev)
      current-tag(pos) = best-rule-to

```

nonlexicalized rules learned by TBL tagger

#	Change tags		Condition	Example
	From	To		
1	NN	VB	Previous tag is TO	to/TO race/NN → VB
2	VBP	VB	One of the previous 3 tags is MD	might/MD vanish/VBP → VB
3	NN	VB	One of the previous 2 tags is MD	might/MD not reply/NN → VB
4	VB	NN	One of the previous 2 tags is DT	
5	VBD	VBN	One of the previous 3 tags is VBZ	

6.863J/9.611J SP04 Lecture 7

figure from Brill's thesis

Transformations Learned

#	Change Tag		Condition
	From	To	
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

BaselineTag*

NN @ → VB // TO _
VBP @ → VB // ... _
 etc.

Compose this cascade of FSTs.

Get a big FST that does the initial tagging and the sequence of fixups "all at once."

7

Error analysis: what's hard for taggers

- Common errors (> 4%)
 - NN vs .NNP (proper vs. other nouns) vs. JJ (adjective): hard to distinguish prenominally; important to distinguish esp. for information extraction
 - RP vs. RB vs IN: all can appear in sequences immed. after verb
 - VBD vs. VBN vs. JJ: distinguish past tense, past participles (*raced vs. was raced vs. the out raced horse*)

6.863J/9.611J SP04 Lecture 7

What's hard

- Unknown words
 - Order 0 idea: equally likely over all parts of speech
 - Better idea: same distribution as 'Things seen once' estimator of 'things never seen' - theory for this done by Turing (again!)
 - *Hapax legomenon*
 - Assume distribution of unknown words is like this
 - But most powerful methods make use of how word is spelled
- See file in the course tagging dir on this

6.863J/9.611J SP04 Lecture 7

Or unknown language

- Все schastlivye sen'i pokhozhi brug na druga, kazhdaja neschastlivaja sem'ja neschastliva po-svoemu

6.863J/9.611J SP04 Lecture 7

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Instant Message 6.863J Syllabus Google BE768 Biologica Members WebMail Connections BizJournal SmartUpdate Mktplace Home

Bookmarks Location: http://www.ling.gu.se/~lager/Home/brilltagger_ui.html

Brill Tagger

Powered by μ-TBL Technology

Swedish English Russian

Text:

Trace

© Torbjörn Lager 1999, Russian tagger by Natalia Zinovjeva

Most powerful unknown word detectors

- 3 inflectional endings (*-ed, -s, -ing*); 32 derivational endings (*-ion, etc.*); capitalization; hyphenation
- More generally: should use morphological analysis! (and some kind of machine learning approach)
- How hard is this? We don't know - we actually don't know how children do this, either (they make mistakes)

6.863J/9.611J SP04 Lecture 7

Laboratory 2

- Goals:
 1. Use both HMM and Brill taggers
 2. Find errors that both make, relative to genre
 3. Compare performance – use of kappa & 'confusion matrix'
 4. All the slings & arrows of corpora – use Wall Street Journal excerpts, as well as 'switchboard' corpus

6.863J/9.611J SP04 Lecture 7

Brown/Upenn corpus tags

J. text,
p. 297
Fig 8.6
1M words
60K tag
counts

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(' or "')</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or "')</i>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, (, { , <</i>
PP\$	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>([, } , ></i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... --)</i>
RP	Particle	<i>up, off</i>			

Evaluation of systems

- The principal measures for information extraction tasks are recall and precision.
- Recall is the number of answers the system got right divided by the number of possible right answers
 - It measures how complete or comprehensive the system is in its extraction of relevant information
- Precision is the number of answers the system got right divided by the number of answers the system gave
 - It measures the system's correctness or accuracy
 - Example: there are 100 possible answers and the system gives 80 answers and gets 60 of them right, its recall is 60% and its precision is 75%.

A better measure - Kappa

- Takes baseline & complexity of task into account – if 99% of tags are Nouns, getting 99% correct no great shakes
- Suppose no “Gold Standard” to compare against?
- $P(A)$ = proportion of times hypothesis *agrees* with standard (% correct)
- $P(E)$ = proportion of times hypothesis and standard would be *expected* to agree by chance (computed from some other knowledge, or actual data)

6.863J/9.611J SP04 Lecture 7

Kappa [p. 315 J&M text]

- Note K ranges between 0 (no agreement, except by chance; to complete agreement, 1)
- Can be used even if no ‘Gold standard’ that everyone agrees on
- $K > 0.8$ is good

6.863J/9.611J SP04 Lecture 7



Kappa

- A = actual agreement; E = expected agreement
- consistency is more important than “truth”
- methods for raising consistency
 - style guides (often have useful insights into data)
 - group by task, not chronologically, etc.
 - annotator acclimatization

6.863J/9.611J SP04 Lecture 7



Coda on kids

C: “Mommy, nobody don’t like me”

A: No, say, “nobody likes me”

C: Nobody don’t likes me

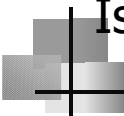
A: Say, “nobody likes me”

C: Nobody don’t likes me

[7 repetitions]

C: Oh! Nobody don’t like me!

6.863J/9.611J SP04 Lecture 7



Is that all there is?

6.863J/9.611J SP04 Lecture 7



What have we done so far?

- Only information we represent: is whether an item precedes (or follows) another
- Inventory of vocabulary items (classes)
- = Finite state machines

- Is there anything else in language???

6.863J/9.611J SP04 Lecture 7

Motivation

- What, How, and Why
- What: word *chunks* behave as units, like words or endings (morphemes), like *ing*
- How: we have to recover these from input
- Why: chunks used to discover *meaning*
- Parsing: mapping from *strings* to *structured representation*

6.863J/9.611J SP04 Lecture 7

Programming languages

```
printf ("/charset [%s"  
       (re_opcode_t) *(p - 1) == charset_not ? "^" : "");  
assert (p + *p < pend);  
for (c = 0; c < 256; c++)  
  if (c / 8 < *p && (p[1 + (c/8)] & (1 << (c % 8)))) {  
    /* Are we starting a range? */  
    if (last + 1 == c && ! inrange) {  
      putchar ('-');  
      inrange = 1;  
    }  
    /* Have we broken a range? */  
    else if (last + 1 != c && inrange) {  
      putchar (last);  
      inrange = 0;  
    }  
    if (! inrange)  
      putchar (c);  
    last = c;  
  }  
}
```

- Easy to parse.
- Designed that way!

6.863J/9.611J SP04 Lecture 7

Natural languages

```
printf "/charset %s", re opcode t *p - 1 == charset not ? "^" : "";  
assert p + *p < pend; for c = 0; c < 256; c++ if c / 8 < *p && p1 + c / 8  
& 1 << c % 8 Are we starting a range? if last + 1 == c && ! inrange  
putchar '-'; inrange = 1; Have we broken a range? else if last + 1 != c  
&& inrange putchar last; inrange = 0; if ! inrange putchar c; last =  
c;
```

- No {} () [] to indicate scope & precedence
- Lots of overloading (arity varies)
- Grammar isn't known in advance!
- What is the best formalism?

6.863J/9.611J SP04 Lecture 7

What can't linear relations represent?

- wine dark sea → (wine (dark sea)) *or*
((wine dark) sea) ?
- deep blue sky
- Can fsa's *represent* this?
- Not really: algebraically, *defined* as being associative (doesn't matter about concatenation order)

6.863J/9.611J SP04 Lecture 7