

Part-of-Speech Tagging

A Canonical Finite-State Task

6.00.465 - Intro to NLP - J. Eisner

1

The Tagging Task

Input: the lead paint is unsafe

Output: the/Det lead/N paint/N is/V unsafe/Adj

- Uses:
 - text-to-speech (how do we pronounce "lead"?)
 - can write regexps like (Det) Adj* N+ over the output
 - preprocessing to speed up parser (but a little dangerous)
 - if you know the tag, you can back off to it in other tasks

6.00.465 - Intro to NLP - J. Eisner

2

Why Do We Care?

Input: the lead paint is unsafe

Output: the/Det lead/N paint/N is/V unsafe/Adj

- The first statistical NLP task
- Been done to death by different methods
- Easy to evaluate (how many tags are correct?)
- Canonical finite-state task
 - Can be done well with methods that look at local context
 - Though should "really" do it by parsing!

6.00.465 - Intro to NLP - J. Eisner

3

Degree of Supervision

- Supervised:** Training corpus is tagged by humans
- Unsupervised:** Training corpus isn't tagged
- Partly supervised:** Training corpus isn't tagged, but you have a dictionary giving possible tags for each word
- We'll start with the supervised case and move to decreasing levels of supervision.

6.00.465 - Intro to NLP - J. Eisner

4

Current Performance

Input: the lead paint is unsafe

Output: the/Det lead/N paint/N is/V unsafe/Adj

- How many tags are correct?
 - About 97% currently
 - But baseline is already 90%
 - Baseline is performance of stupidest possible method
 - Tag every word with its most frequent tag
 - Tag unknown words as nouns

6.00.465 - Intro to NLP - J. Eisner

5

What Should We Look At?

correct tags

PN	Verb	Det	Noun	Prep	Noun	Prep	Det	Noun
Bill	directed	a	cortege	of	autos	through	the	dunes
PN	Adj	Det	Noun	Prep	Noun	Prep	Det	Noun
Verb	Verb	Noun	Verb	Adj				
				Prep				
				...?				

some possible tags for each word (maybe more)

Each unknown tag is **constrained** by its word and by the tags to its immediate left and right. But those tags are unknown too ...

6.00.465 - Intro to NLP - J. Eisner

6

What Should We Look At?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
 Bill directed a cortege of autos through the dunes

PN Adj Det Noun Prep Noun Prep Det Noun
 Bill directed a cortege of autos through the dunes

Verb Verb Noun Verb
 directed a cortege of autos through the dunes

Adj
 some possible tags for
 Prep each word (maybe more)
 ...?

Each unknown tag is **constrained** by its word and by the tags to its immediate left and right. But those tags are unknown too ...

6.00.465 - Intro to NLP - J. Eisner

7

What Should We Look At?

correct tags

PN Verb Det Noun Prep Noun Prep Det Noun
 Bill directed a cortege of autos through the dunes

PN Adj Det Noun Prep Noun Prep Det Noun
 Bill directed a cortege of autos through the dunes

Verb Verb Noun Verb
 directed a cortege of autos through the dunes

Adj
 some possible tags for
 Prep each word (maybe more)
 ...?

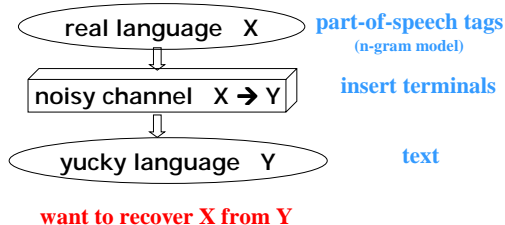
Each unknown tag is **constrained** by its word and by the tags to its immediate left and right. But those tags are unknown too ...

6.00.465 - Intro to NLP - J. Eisner

8

Three Finite-State Approaches

- Noisy Channel Model (statistical)



6.00.465 - Intro to NLP - J. Eisner

9

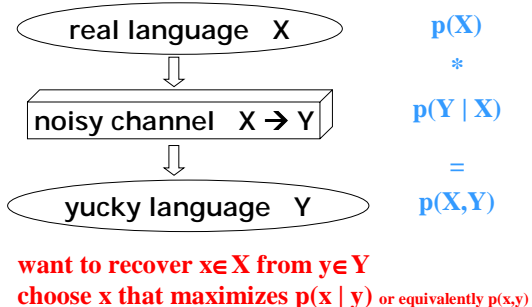
Three Finite-State Approaches

- Noisy Channel Model (statistical)
- Deterministic baseline tagger composed with a cascade of fixup transducers
- Nondeterministic tagger composed with a cascade of finite-state automata that act as filters

6.00.465 - Intro to NLP - J. Eisner

10

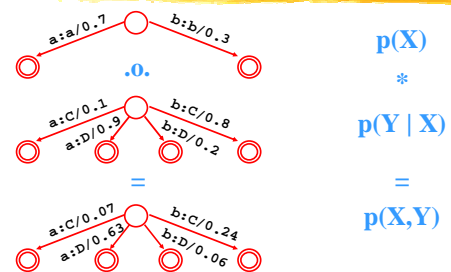
Review: Noisy Channel



6.00.465 - Intro to NLP - J. Eisner

11

Review: Noisy Channel

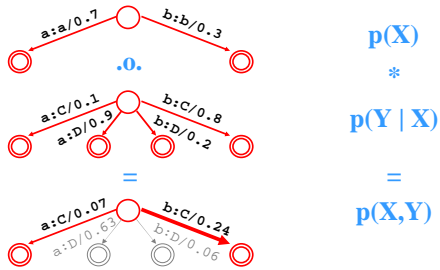


Note $p(x, y)$ sums to 1.
 Suppose $y = "C"$; what is best $"x"$?

6.00.465 - Intro to NLP - J. Eisner

12

Review: Noisy Channel

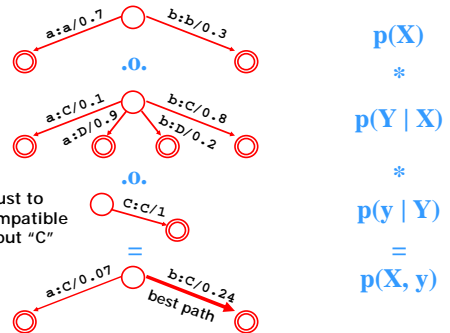


Suppose $y = "C"$; what is best $"x"$?

600.465 - Intro to NLP - J. Eisner

13

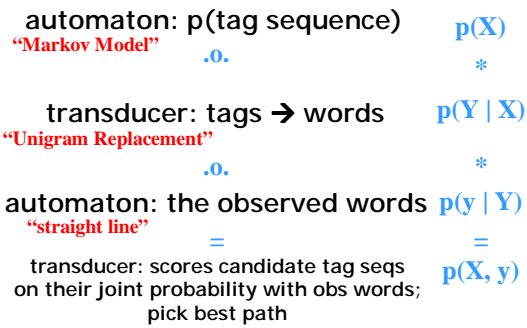
Review: Noisy Channel



600.465 - Intro to NLP - J. Eisner

14

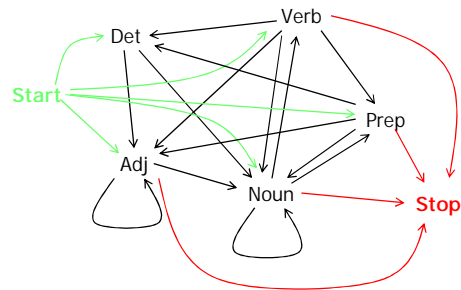
Noisy Channel for Tagging



600.465 - Intro to NLP - J. Eisner

15

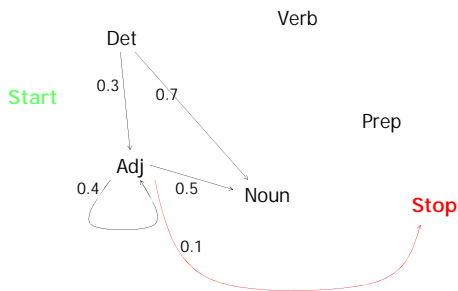
Markov Model (bigrams)



600.465 - Intro to NLP - J. Eisner

16

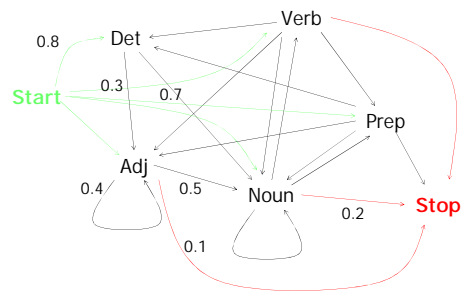
Markov Model



600.465 - Intro to NLP - J. Eisner

17

Markov Model

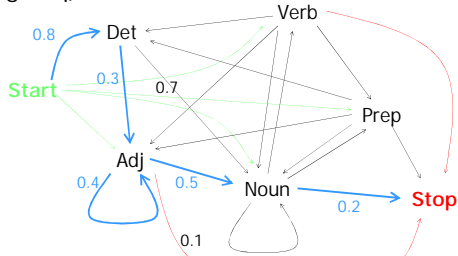


600.465 - Intro to NLP - J. Eisner

18

Markov Model

p(tag seq)



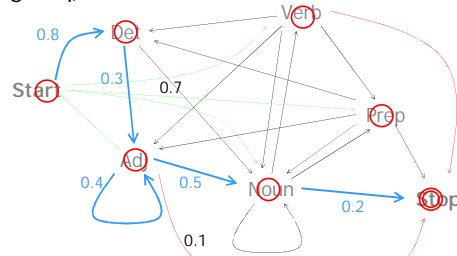
Start Det Adj Adj Noun Stop = $0.8 * 0.3 * 0.4 * 0.5 * 0.2$

600.465 - Intro to NLP - J. Eisner

19

Markov Model as an FSA

p(tag seq)



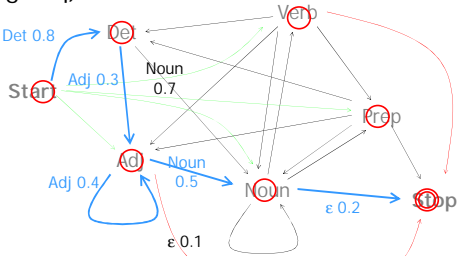
Start Det Adj Adj Noun Stop = $0.8 * 0.3 * 0.4 * 0.5 * 0.2$

600.465 - Intro to NLP - J. Eisner

20

Markov Model as an FSA

p(tag seq)



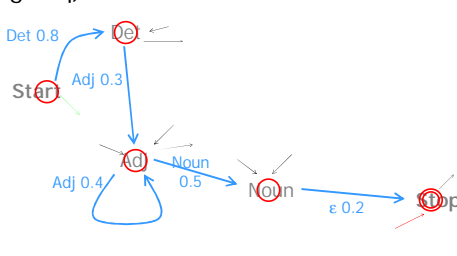
Start Det Adj Adj Noun Stop = $0.8 * 0.3 * 0.4 * 0.5 * 0.2$

600.465 - Intro to NLP - J. Eisner

21

Markov Model (tag bigrams)

p(tag seq)



Start Det Adj Adj Noun Stop = $0.8 * 0.3 * 0.4 * 0.5 * 0.2$

600.465 - Intro to NLP - J. Eisner

22

Noisy Channel for Tagging

automaton: p(tag sequence)

"Markov Model"

.0.

p(X)

*

transducer: tags → words

"Unigram Replacement"

.0.

p(Y | X)

*

automaton: the observed words

"straight line"

.0.

p(y | Y)

*

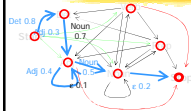
transducer: scores candidate tag seqs
on their joint probability with obs words;
pick best path

p(X, y)

600.465 - Intro to NLP - J. Eisner

23

Noisy Channel for Tagging



Noun: cortège/0.000007
Noun: autos/0.000007
Noun: cortège/0.000007
Det: a/0.5
Det: the/0.4
Adj: cool/0.003
Adj: directed/0.00005
Adj: cortège/0.000001

the cool directed autos

p(y | Y)

transducer: scores candidate tag seqs
on their joint probability with obs words;
we should pick best path

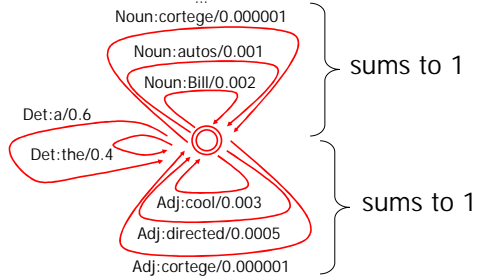
p(X, y)

600.465 - Intro to NLP - J. Eisner

24

Unigram Replacement Model

$p(\text{word seq} \mid \text{tag seq})$

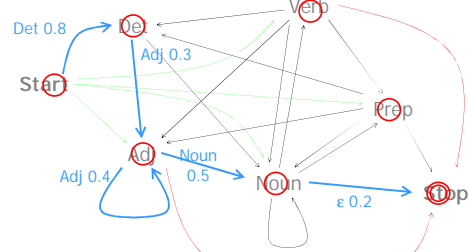


6.00.465 - Intro to NLP - J. Eisner

25

Compose

$p(\text{tag seq})$

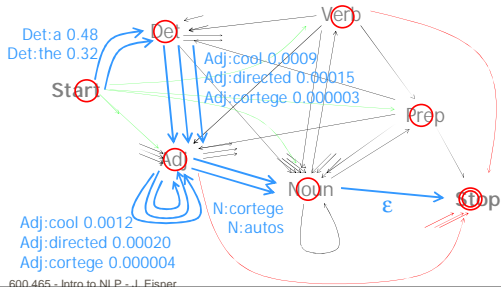


6.00.465 - Intro to NLP - J. Eisner

26

Compose

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

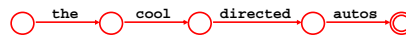


6.00.465 - Intro to NLP - J. Eisner

27

Observed Words as Straight-Line FSA

word seq

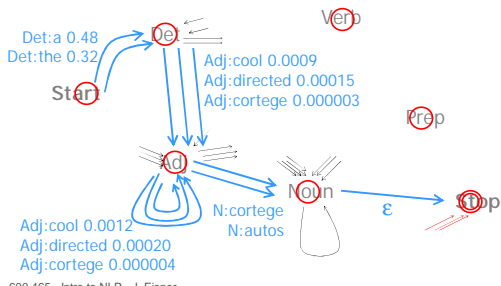


6.00.465 - Intro to NLP - J. Eisner

28

Compose with the cool directed autos

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

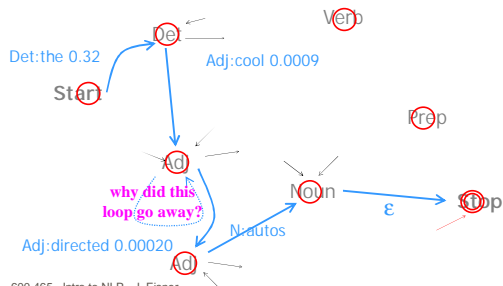


6.00.465 - Intro to NLP - J. Eisner

29

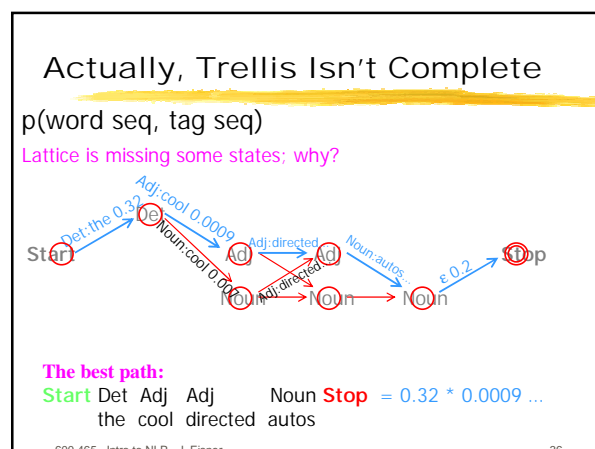
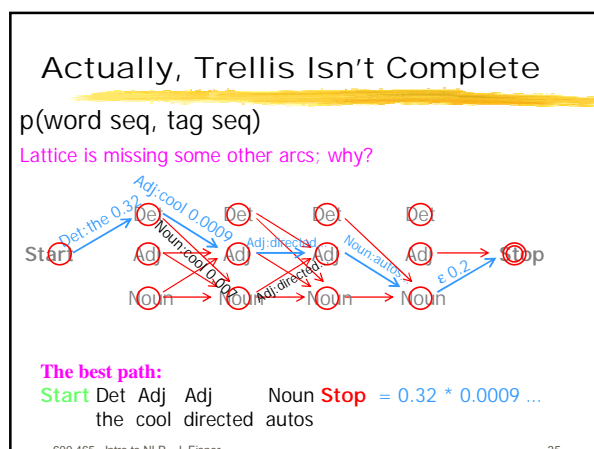
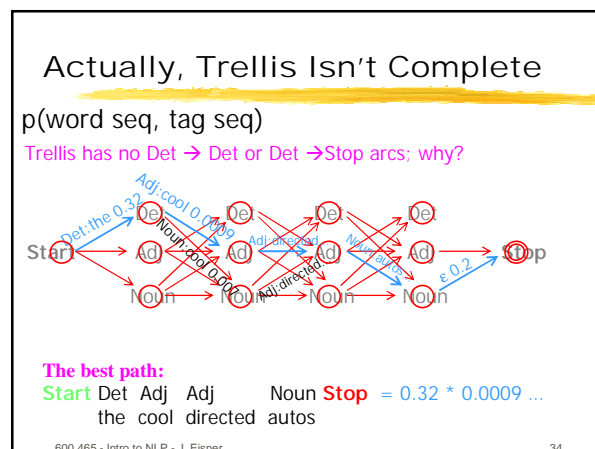
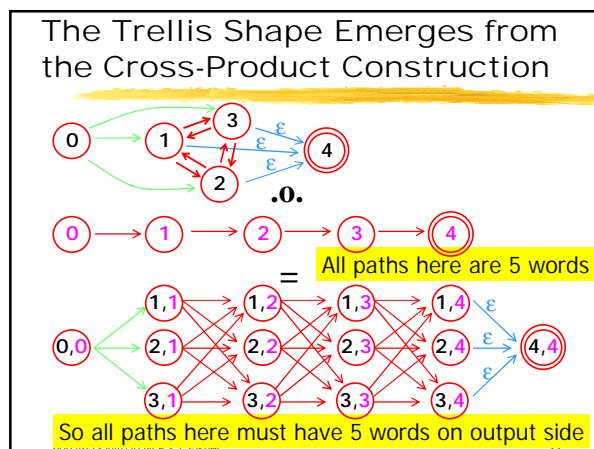
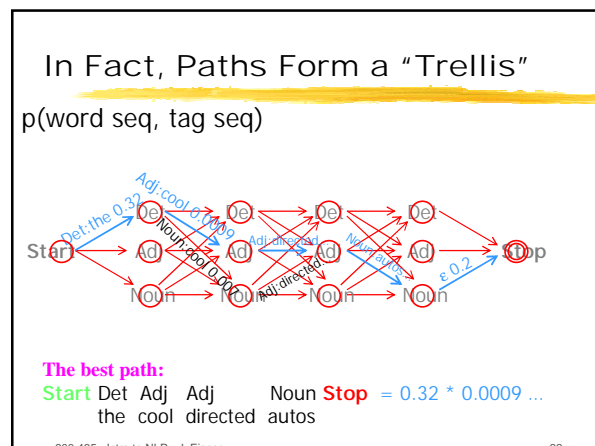
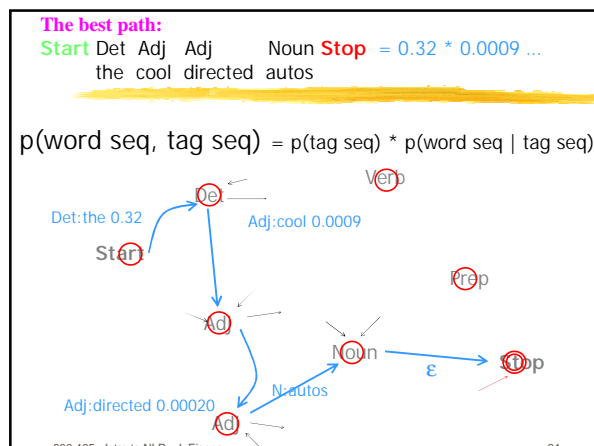
Compose with the cool directed autos

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



6.00.465 - Intro to NLP - J. Eisner

30



Find best path from Start to Stop



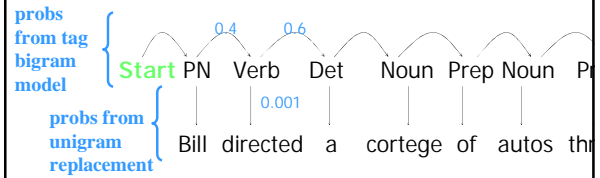
- Use dynamic programming – like prob. parsing:
 - What is best path from Start to *each* node?
 - Work from left to right
 - Each node stores its best path from Start (as probability plus one backpointer)
- Special acyclic case of Dijkstra's shortest-path alg.
- Faster if some arcs/states are absent

600.465 - Intro to NLP - J. Eisner

37

In Summary

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



- Find X that maximizes probability **product**

600.465 - Intro to NLP - J. Eisner

38

Another Viewpoint

- We are modeling $p(\text{word seq, tag seq})$
- Why not use chain rule + some kind of backoff?
- Actually, we are!

$$p(\text{Start PN Verb Det ... Bill directed a ...})$$

$$= p(\text{Start}) * p(\text{PN} | \text{Start}) * p(\text{Verb} | \text{Start PN}) * p(\text{Det} | \text{Start PN Verb}) * \dots$$

$$* p(\text{Bill} | \text{Start PN Verb ...}) * p(\text{directed} | \text{Bill, Start PN Verb Det ...})$$

$$* p(\text{a} | \text{Bill directed, Start PN Verb Det ...}) * \dots$$

600.465 - Intro to NLP - J. Eisner

39

Another Viewpoint

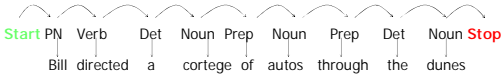
- We are modeling $p(\text{word seq, tag seq})$
- Why not use chain rule + some kind of backoff?
- Actually, we are!

$$p(\text{Start PN Verb Det ... Bill directed a ...})$$

$$= p(\text{Start}) * p(\text{PN} | \text{Start}) * p(\text{Verb} | \text{Start PN}) * p(\text{Det} | \text{Start PN Verb}) * \dots$$

$$* p(\text{Bill} | \text{Start PN Verb ...}) * p(\text{directed} | \text{Bill, Start PN Verb Det ...})$$

$$* p(\text{a} | \text{Bill directed, Start PN Verb Det ...}) * \dots$$



600.465 - Intro to NLP - J. Eisner

40

Three Finite-State Approaches

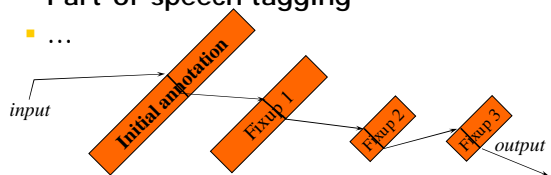
- Noisy Channel Model (statistical)
- Deterministic baseline tagger composed with a cascade of fixup transducers
- Nondeterministic tagger composed with a cascade of finite-state automata that act as filters

600.465 - Intro to NLP - J. Eisner

41

Another FST Paradigm: Successive Fixups

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



600.465 - Intro to NLP - J. Eisner

42

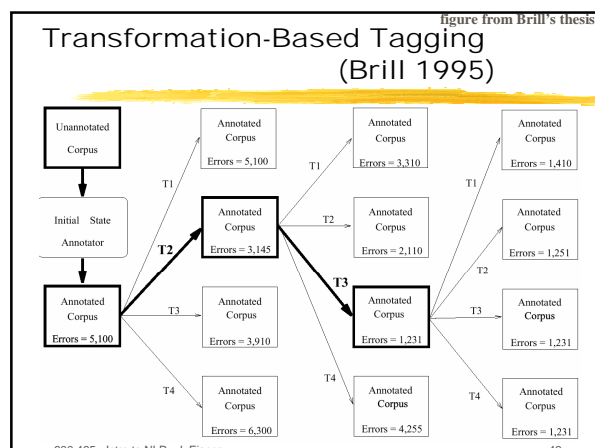


figure from Brill's thesis

Transformations Learned

#	Change	Tag	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>VBDZ</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	VBD	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	WDI	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	WDI	Next tag is <i>VBDZ</i>
17	IN	DT	Next tag is <i>NNP</i>
18	JJ	NNP	Next tag is <i>VBD</i>
19	IN	WDI	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.
 Gets a big FST that does the initial tagging and the sequence of fixups "all at once."

600.465 - Intro to NLP - J. Eisner 44

figure from Brill's thesis

Initial Tagging of OOV Words

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

600.465 - Intro to NLP - J. Eisner 45

- ## Three Finite-State Approaches
- Noisy Channel Model (statistical)
 - Deterministic baseline tagger composed with a cascade of fixup transducers
 - Nondeterministic tagger composed with a cascade of finite-state automata that act as filters
- 600.465 - Intro to NLP - J. Eisner 46

- ## Variations
- Multiple tags per word
 - Transformations to knock some of them out
 - How to encode multiple tags and knockouts?
 - Use the above for partly supervised learning
 - Supervised: You have a tagged training corpus
 - Unsupervised: You have an untagged training corpus
 - Here: You have an untagged training corpus and a dictionary giving possible tags for each word
- 600.465 - Intro to NLP - J. Eisner 47