





# Parts of Speech

- Perhaps starting with Aristotle in the West (384–322 BCE) the idea of having parts of speech
  - lexical categories, word classes, “tags”, POS
- Dionysius Thrax of Alexandria (c. 100 BCE): 8 parts of speech
  - Still with us! But his 8 aren’t exactly the ones we are taught today
    - *Thrax*: noun, verb, article, adverb, preposition, conjunction, participle, pronoun
    - *School grammar*: noun, verb, adjective, adverb, preposition, conjunction, pronoun, interjection

## Open class (lexical) words

### Nouns

#### Proper

*IBM*  
*Italy*

#### Common

*cat / cats*  
*snow*

### Verbs

#### Main

*see*  
*registered*

Adjectives *old older oldest*

Adverbs *slowly*

### Numbers

*122,312*  
*one*

*... more*

## Closed class (functional)

Determiners *the some*

Conjunctions *and or*

Pronouns *he its*

### Modals

*can*  
*had*

Prepositions *to with*

Particles *off up*

*... more*

Interjections *Ow Eh*



# Open vs. Closed classes

- Open vs. Closed classes
  - Closed:
    - determiners: *a, an, the*
    - pronouns: *she, he, I*
    - prepositions: *on, under, over, near, by, ...*
    - Why “closed”?
  - Open:
    - Nouns, Verbs, Adjectives, Adverbs.



# POS Tagging

- Words often have more than one POS: *back*
  - The back door = JJ
  - On my back = NN
  - Win the voters back = RB
  - Promised to back the bill = VB
- The POS tagging problem is to determine the POS tag for a particular instance of a word.



# POS Tagging

- Input: Plays well with others
- Ambiguity: NNS/VBZ UH/JJ/NN/RB IN NNS
- Output: Plays/VBZ well/RB with/IN others/NNS
- Uses:

- MT: reordering of adjectives and nouns (say from Spanish to English)
- Text-to-speech (how do we pronounce “lead”?)
- Can write regexps like (Det) Adj\* N+ over the output for phrases, etc.
- Input to a syntactic parser

Penn Treebank  
POS tags



# The Penn TreeBank Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two</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>	VCN	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, sing.	<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>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</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>			



## Penn Treebank tags

The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

**There/EX** are/VBP 70/CD children/NNS **there/RB**

Preliminary/JJ findings/NNS were/VBD **reported/VBN** in/IN today/NN 's/**POS** New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.





# POS tagging performance

- How many tags are correct? (Tag accuracy)
  - 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
  - Partly easy because
    - Many words are unambiguous
    - You get points for them (*the*, *a*, etc.) and for punctuation marks!



# Deciding on the correct part of speech can be difficult even for people

- Mrs/NNP Shaefer/NNP never/RB got/VBD around/RP to/TO joining/VBG
- All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ around/RB 250/CD



# How difficult is POS tagging?

- About 11% of the word types in the Brown corpus are ambiguous with regard to part of speech
- But they tend to be very common words. E.g., *that*
  - I know *that* he is honest = IN
  - Yes, *that* play was nice = DT
  - You can't go *that* far = RB
- 40% of the word tokens are ambiguous



# Sources of information

- What are the main sources of information for POS tagging?
  - Knowledge of neighboring words
    - Bill saw that man yesterday
    - NNP NN DT NN NN
    - VB VB(D) IN VB NN
  - Knowledge of word probabilities
    - *man* is rarely used as a verb....
- The latter proves the most useful, but the former also helps



# More and Better Features → Feature-based tagger

- Can do surprisingly well just looking at a word by itself:
  - Word                      the: the → DT
  - Lowercased word      Importantly: importantly → RB
  - Prefixes                unfathomable: un- → JJ
  - Suffixes                Importantly: -ly → RB
  - Capitalization        Meridian: CAP → NNP
  - Word shapes            35-year: d-x → JJ
- Then build a classifier to predict tag
  - Maxent  $P(t|w)$ :      93.7% overall / 82.6% unknown



# Overview: POS Tagging Accuracies

- Rough accuracies:

- Most freq tag:

~90% / ~50%

- Trigram HMM:

~95% / ~55%

- Maxent  $P(t|w)$ :

93.7% / 82.6%

- TnT (HMM++):

96.2% / 86.0%

- MEMM tagger:

96.9% / 86.9%

- Bidirectional dependencies:

97.2% / 90.0%

- Upper bound:

~98% (human agreement)

Most errors  
on unknown  
words



# POS tagging as a sequence classification task

- We are given a sentence (an “observation” or “sequence of observations”)
  - Secretariat is expected to race tomorrow
  - She promised to back the bill
- What is the best sequence of tags which corresponds to this sequence of observations?
- Probabilistic view:
  - Consider all possible sequences of tags
  - Out of this universe of sequences, choose the tag sequence which is most probable given the observation sequence of  $n$  words  $w_1...w_n$ .



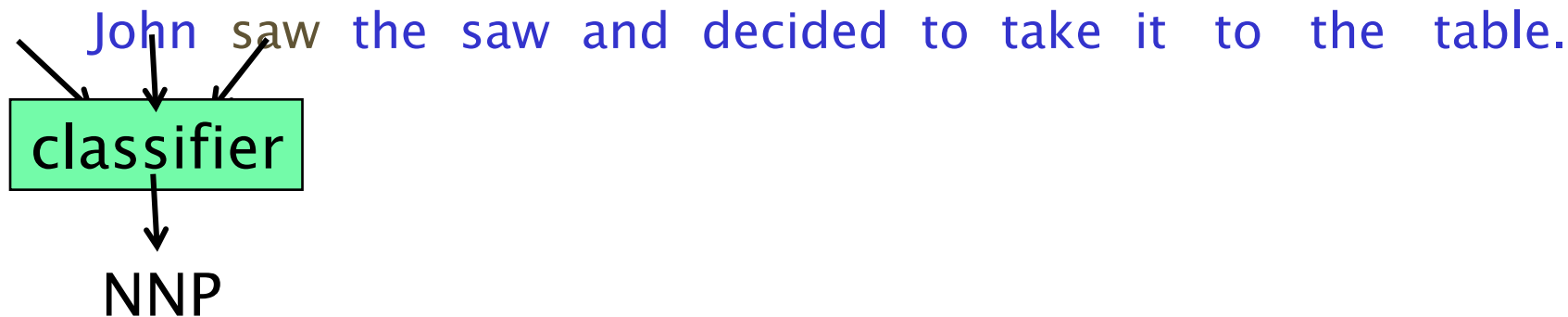
# How do we apply classification to sequences?





# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).





# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



VBD



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



DT



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

classifier

NN



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

classifier

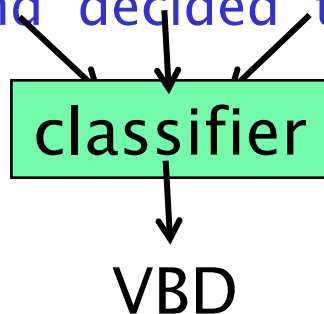
CC



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

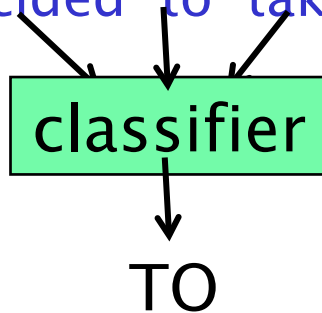




# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.

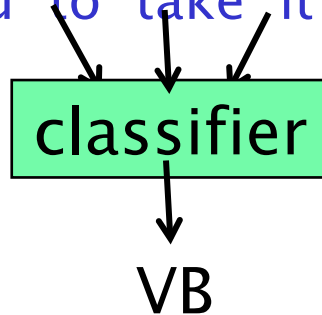




# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.







# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



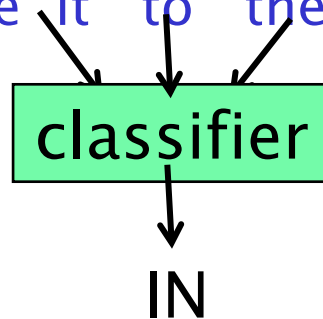
PRP



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.





# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



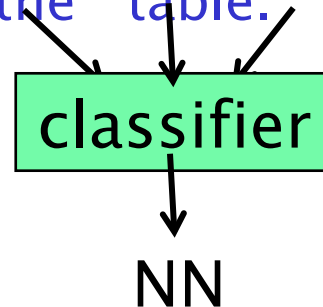
DT



# Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.





# Sequence Labeling as Classification

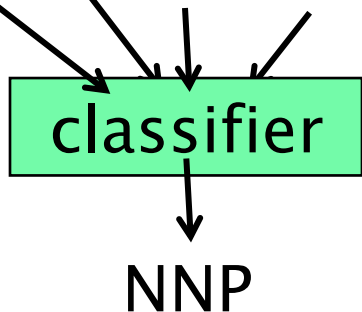
## Using Outputs as Inputs

- Better input features are usually the **categories** of the surrounding tokens, but these are not available yet.
- Can use category of either the preceding or succeeding tokens by going forward or back and using previous output.



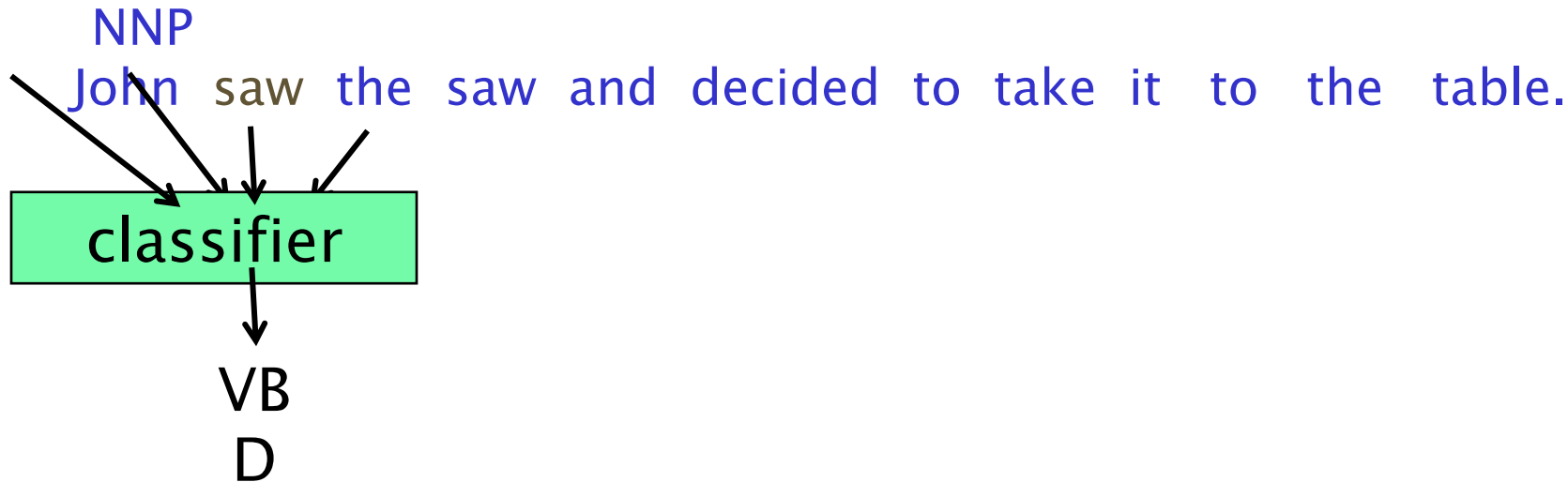
# Forward Classification

John saw the saw and decided to take it to the table.



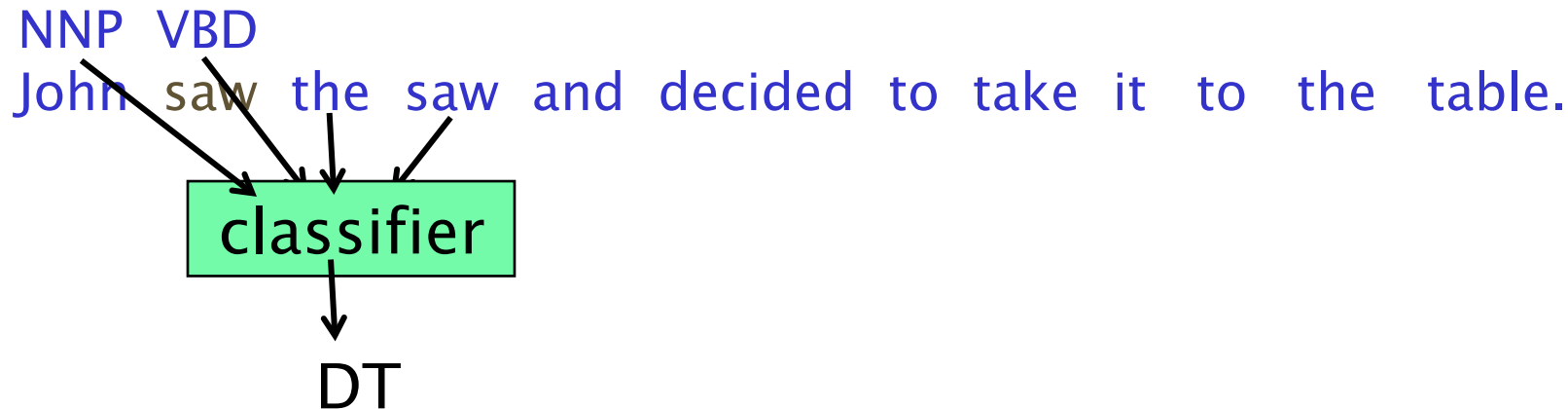


# Forward Classification





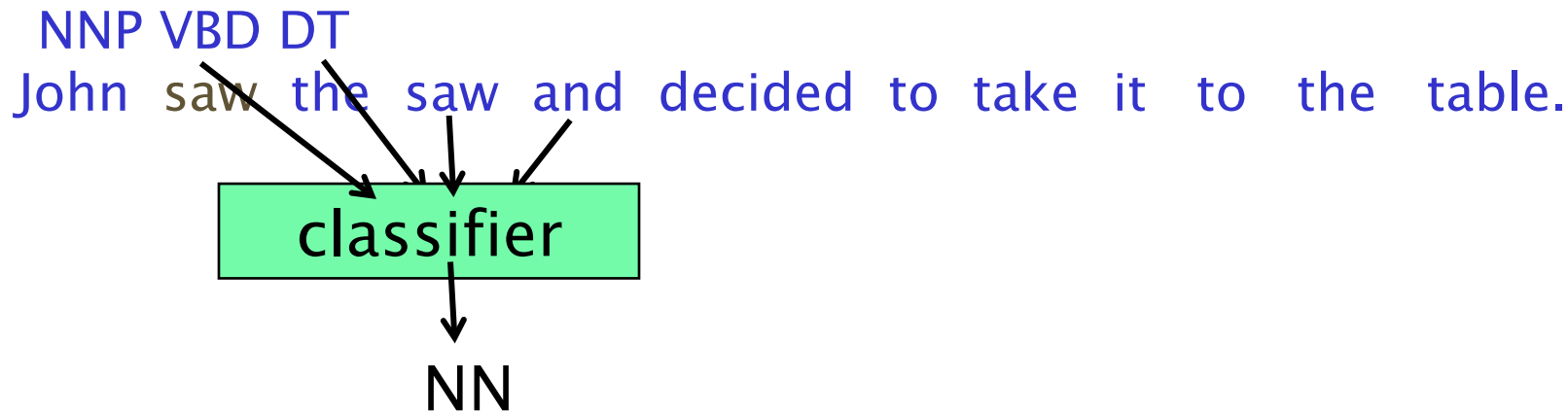
# Forward Classification





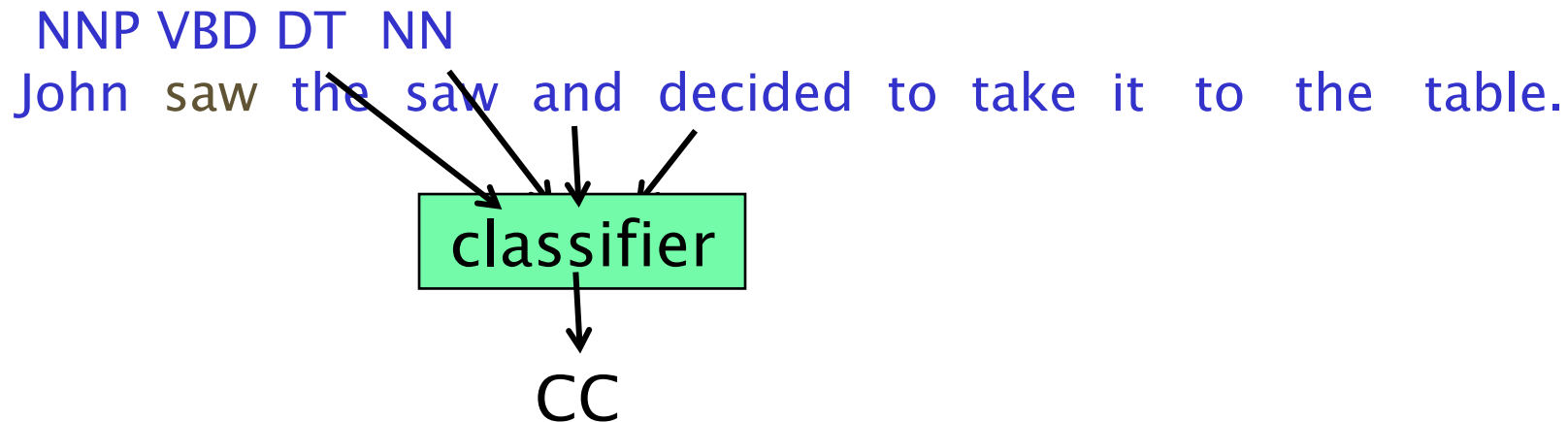


# Forward Classification



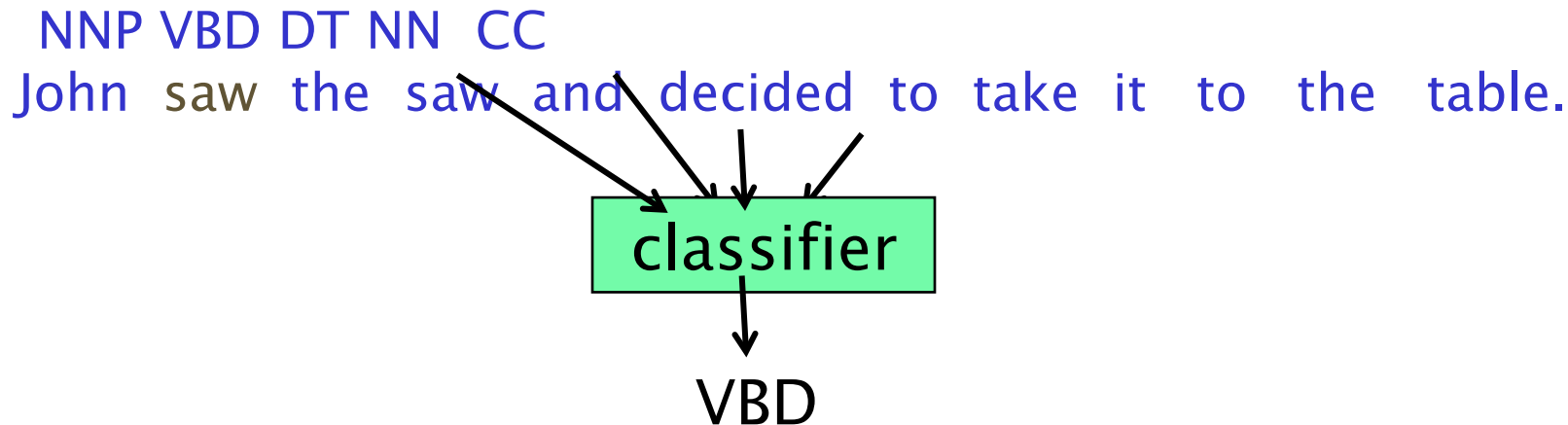


# Forward Classification



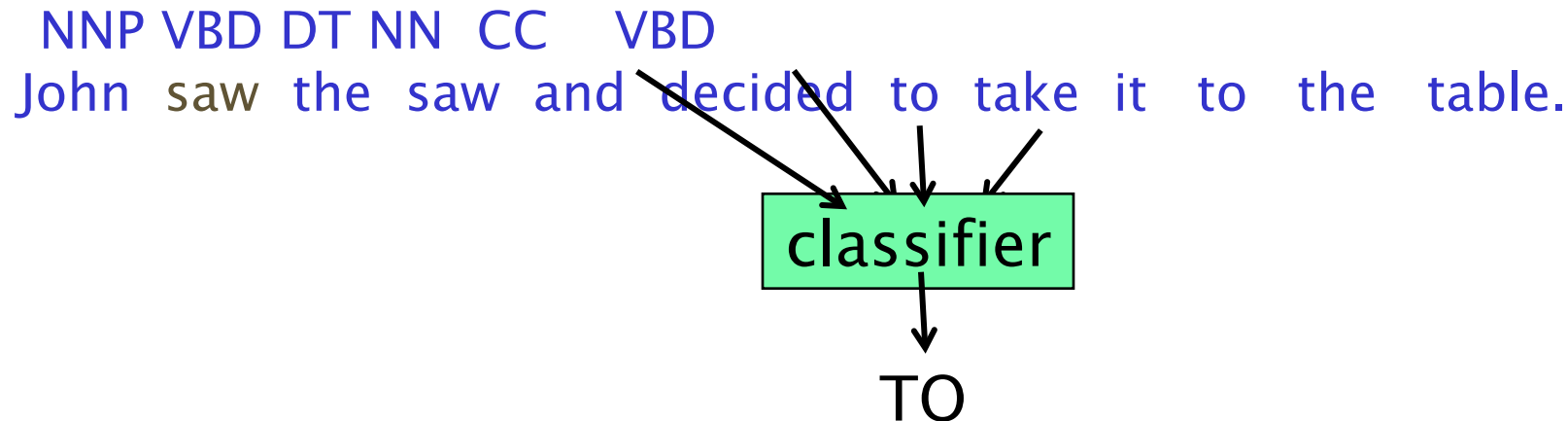


# Forward Classification





# Forward Classification





# Forward Classification

NNP VBD DT NN CC VBD TO  
John saw the saw and decided to take it to the table.

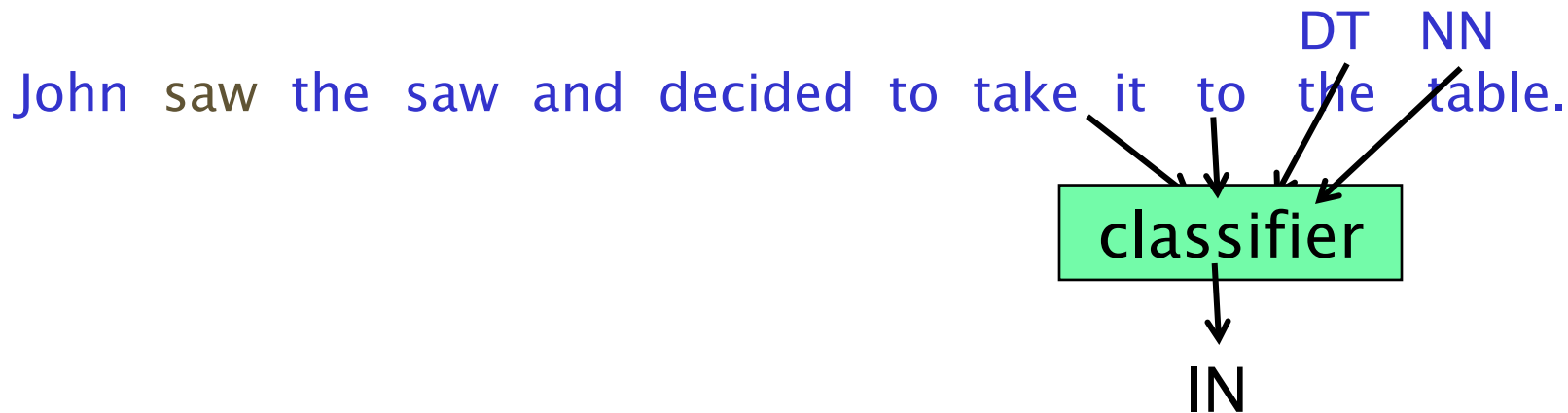
classifier

VB



# Backward Classification

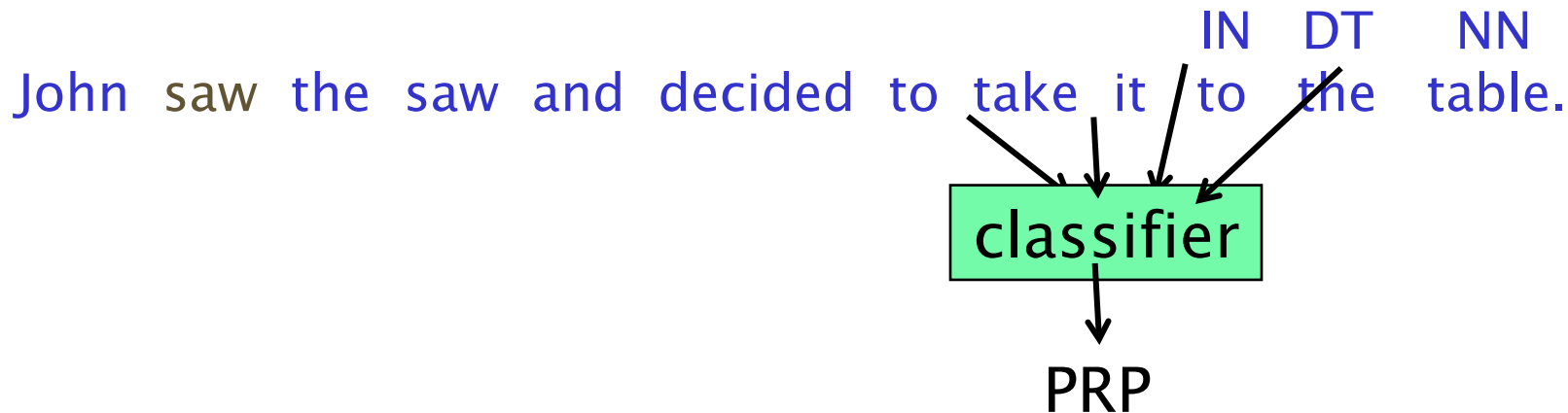
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

- Disambiguating “to” in this case would be even easier backward.



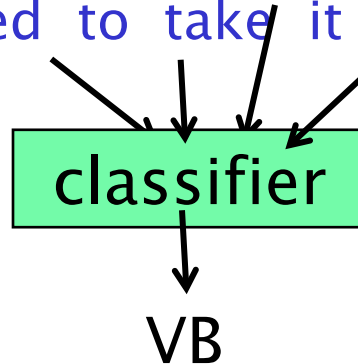


# Backward Classification

- Disambiguating “to” in this case would be even easier backward.

John saw the saw and decided to take it to the table.

PRP IN DT NN

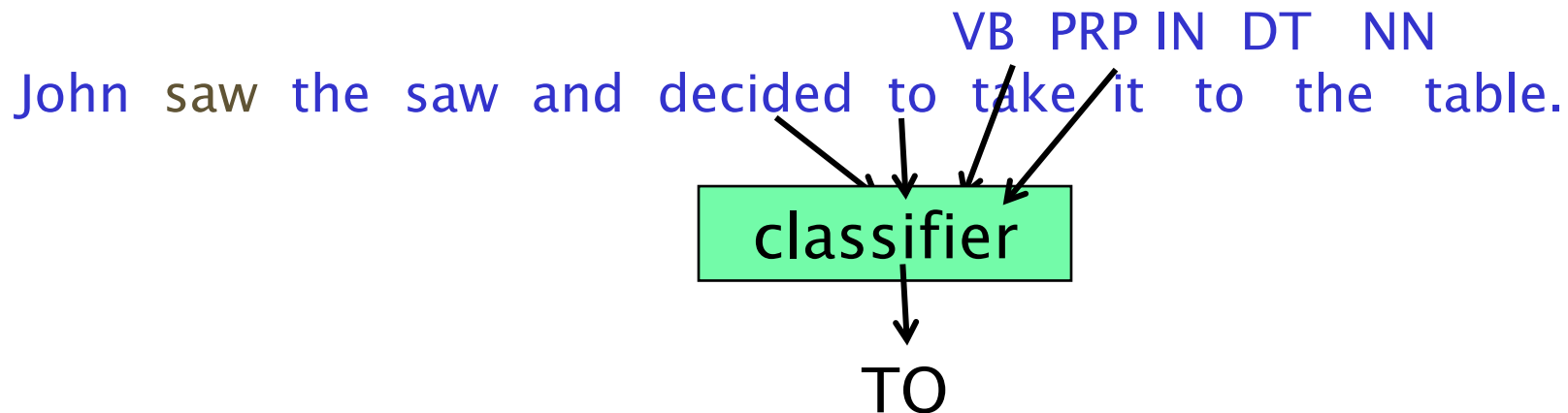






# Backward Classification

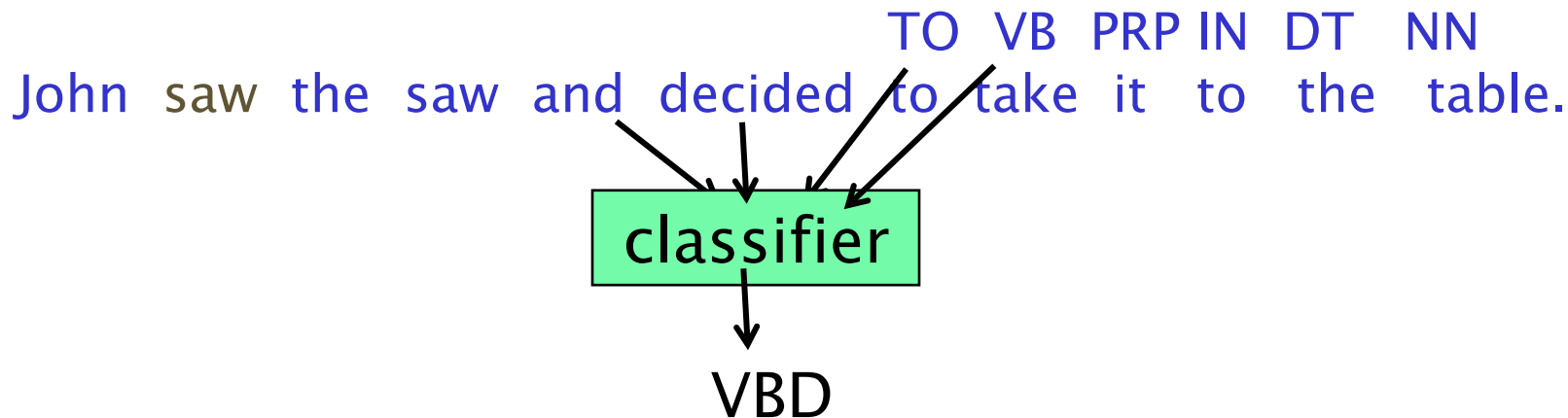
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

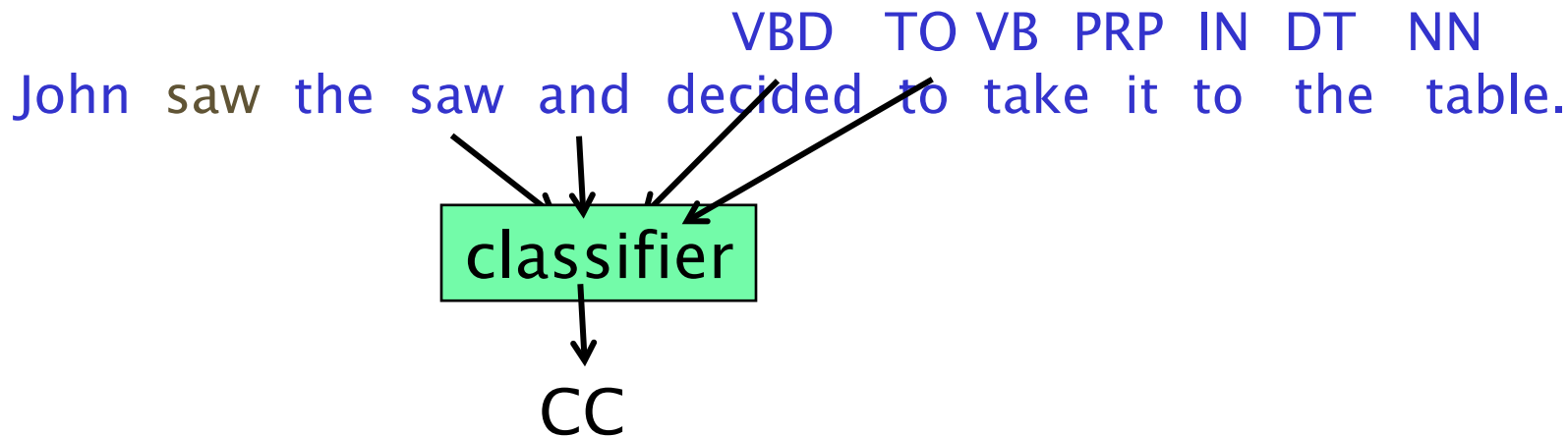
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

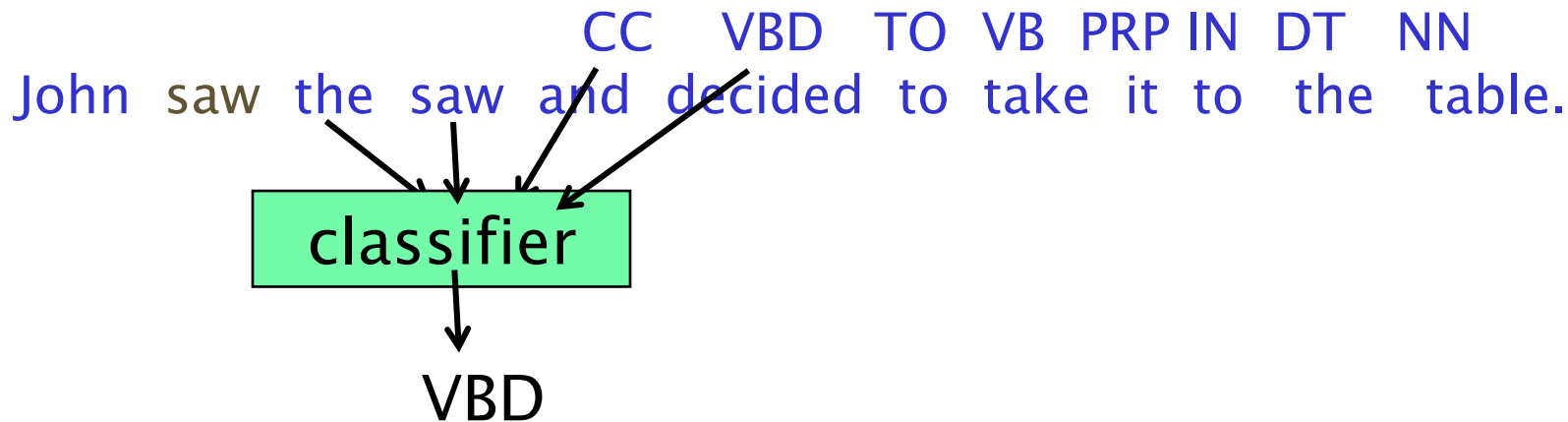
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

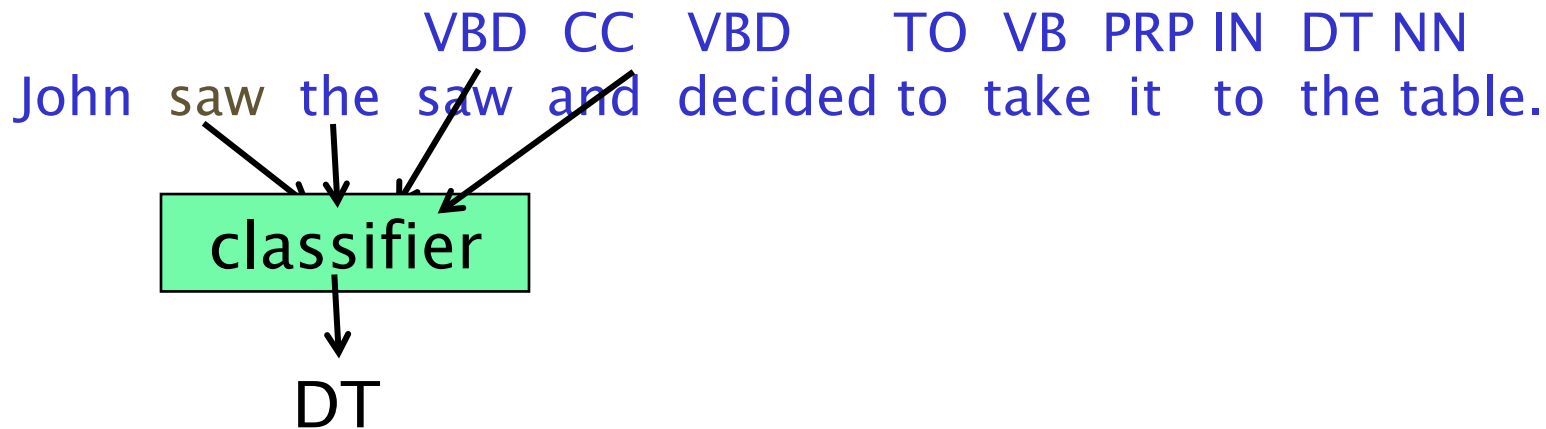
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

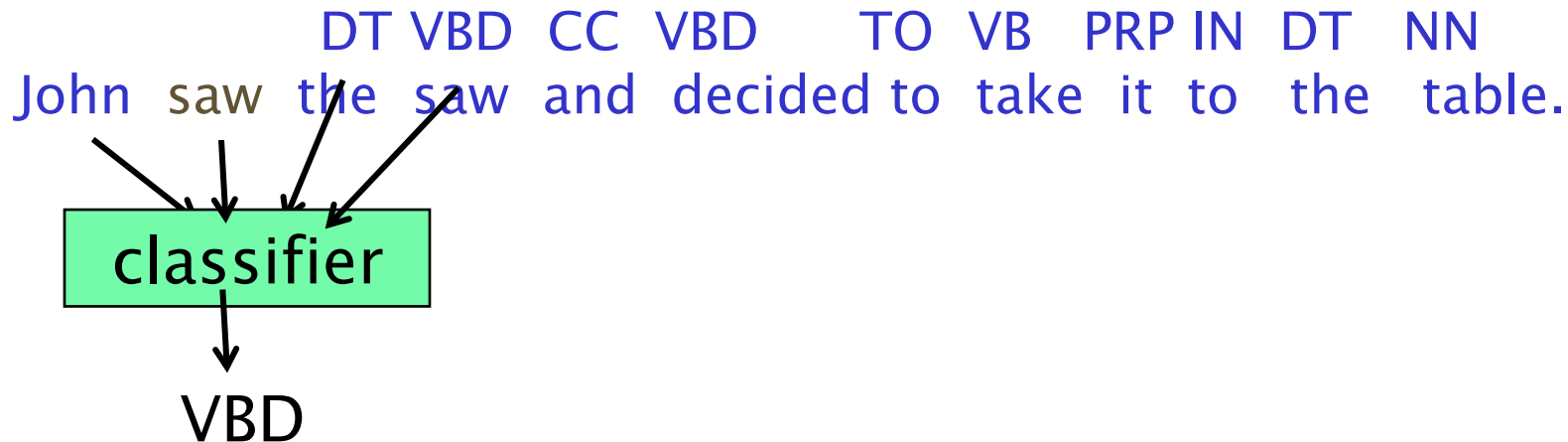
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

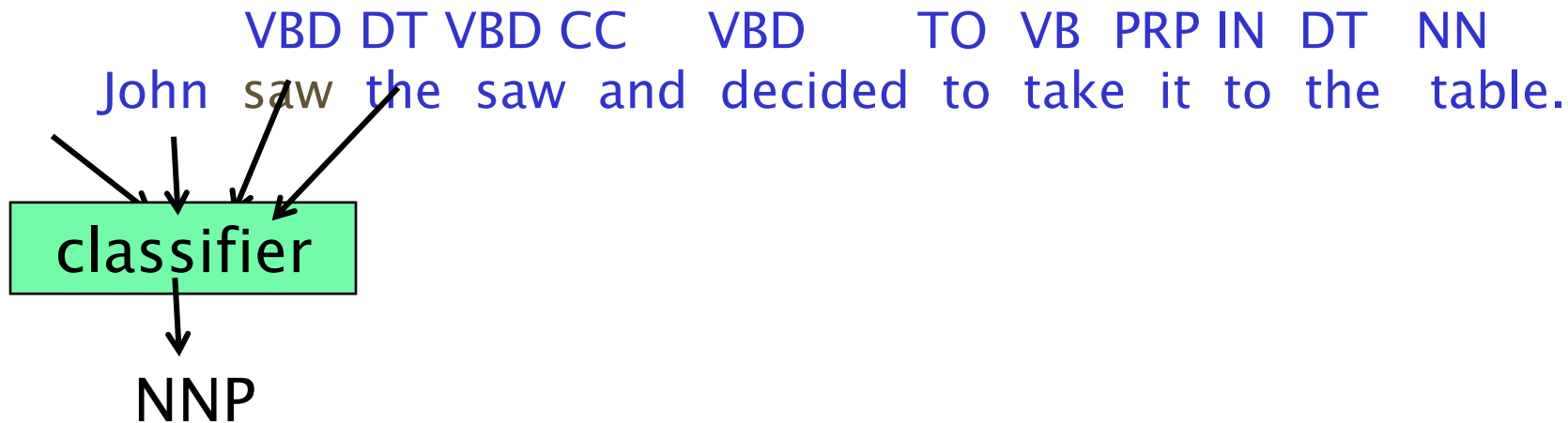
- Disambiguating “to” in this case would be even easier backward.





# Backward Classification

- Disambiguating “to” in this case would be even easier backward.





# The Maximum Entropy Markov Model (MEMM)

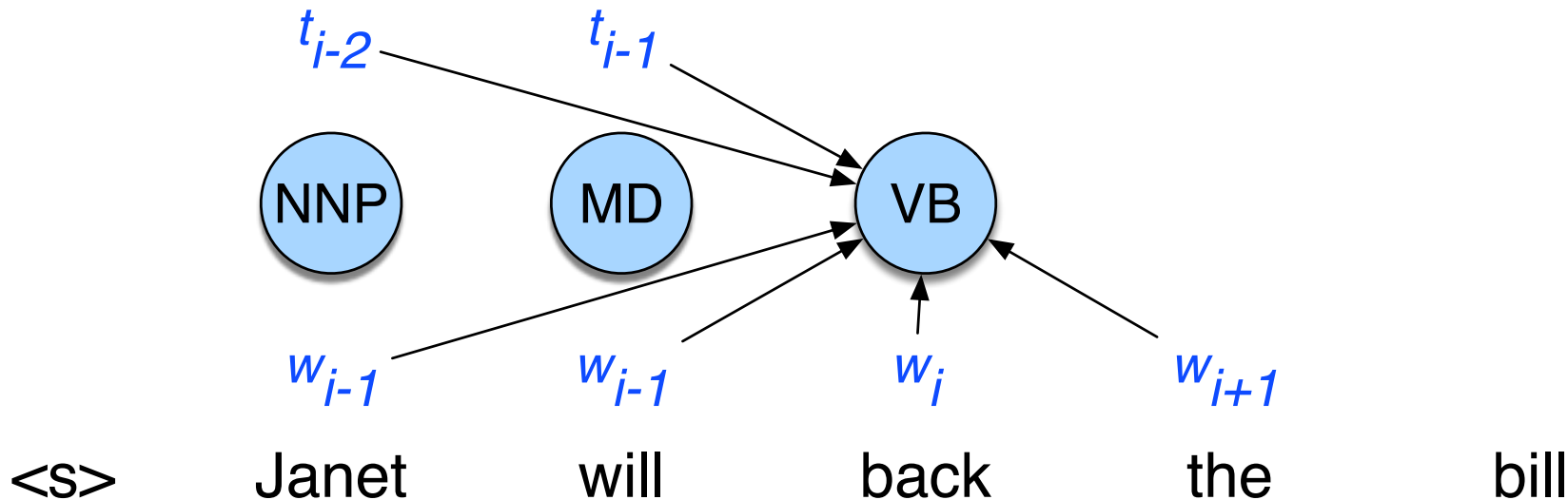
- A sequence version of the logistic regression (also called maximum entropy) classifier.
- Find the best series of tags:

$$\begin{aligned}\hat{T} &= \operatorname{argmax}_T P(T|W) \\ &= \operatorname{argmax}_T \prod_i P(t_i|w_i, t_{i-1})\end{aligned}$$



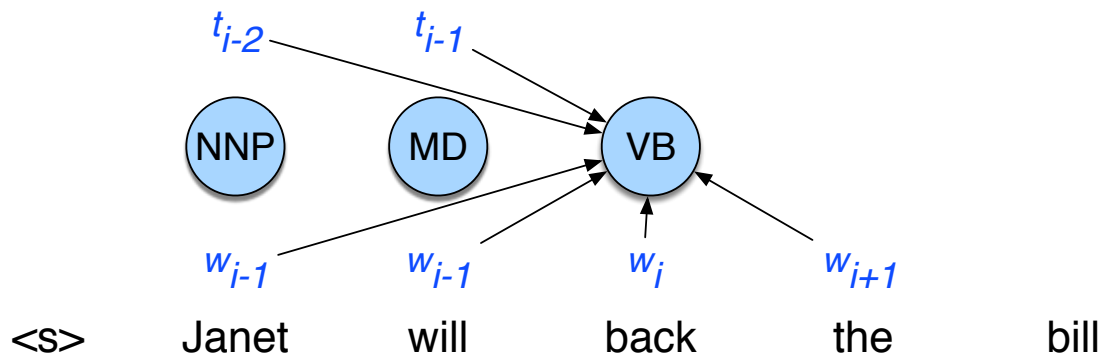


# The Maximum Entropy Markov Model (MEMM)





# Features for the classifier at each tag



$t_i = \text{VB}$  and  $w_{i-2} = \text{Janet}$

$t_i = \text{VB}$  and  $w_{i-1} = \text{will}$

$t_i = \text{VB}$  and  $w_i = \text{back}$

$t_i = \text{VB}$  and  $w_{i+1} = \text{the}$

$t_i = \text{VB}$  and  $w_{i+2} = \text{bill}$

$t_i = \text{VB}$  and  $t_{i-1} = \text{MD}$

$t_i = \text{VB}$  and  $t_{i-1} = \text{MD}$  and  $t_{i-2} = \text{NNP}$

$t_i = \text{VB}$  and  $w_i = \text{back}$  and  $w_{i+1} = \text{the}$



## More features

$w_i$  contains a particular prefix (from all prefixes of length  $\leq 4$ )

$w_i$  contains a particular suffix (from all suffixes of length  $\leq 4$ )

$w_i$  contains a number

$w_i$  contains an upper-case letter

$w_i$  contains a hyphen

$w_i$  is all upper case

$w_i$ 's word shape

$w_i$ 's short word shape

$w_i$  is upper case and has a digit and a dash (like *CFC-12*)

$w_i$  is upper case and followed within 3 words by Co., Inc., etc.



# MEMM computes the best tag sequence

$$\begin{aligned}\hat{T} &= \operatorname{argmax}_T P(T|W) \\ &= \operatorname{argmax}_T \prod_i P(t_i | w_{i-l}^{i+l}, t_{i-k}^{i-1}) \\ &= \operatorname{argmax}_T \prod_i \frac{\exp \left( \sum_i w_i f_i(t_i, w_{i-l}^{i+l}, t_{i-k}^{i-1}) \right)}{\sum_{t' \in \text{tagset}} \exp \left( \sum_i w_i f_i(t', w_{i-l}^{i+l}, t_{i-k}^{i-1}) \right)}\end{aligned}$$



# MEMM Decoding

- Simplest algorithm:

**function** GREEDY MEMM DECODING(words  $W$ , model  $P$ ) **returns** tag sequence  $T$

**for**  $i = 1$  **to**  $length(W)$

$$\hat{t}_i = \operatorname{argmax}_{t' \in T} P(t' \mid w_{i-l}^{i+l}, t_{i-k}^{i-1})$$

- What we use in practice: The **Viterbi** algorithm
- A version of the same dynamic programming algorithm we used to compute minimum edit distance.



# The Stanford Tagger

- Is a bidirectional version of the MEMM called a cyclic dependency network
- Stanford tagger:
  - <http://nlp.stanford.edu/software/tagger.shtml>