

	i want to eat chinese food lunch spend							
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	1	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	19	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

$$V = 1446$$

$$P^*(i|i) = \frac{6}{2533 + 1446} = 0.0015$$

$$P^*(want|i) = \frac{828}{2533 + 1446} = 0.208$$

$$P^*(to|i) = \frac{1}{2533 + 1446} = 0.00025$$

$$P^*(eat|i) = \frac{10}{2533 + 1446} = 0.0025$$

$$P^*(chinese|i) = \frac{1}{2533 + 1446} = 0.00025 = P(food|i) = P(lunch|i)$$

$$P^*(spend|i) = \frac{3}{2533 + 1446} = 0.00075$$

$$C^*(i|i) = C^*(i|i) = \frac{6 \times 2533}{2533 + 1446} = 3.7995$$

$$C^*(want|i) = \frac{828 \times 2533}{2533 + 1446} = 526.864$$

$$C^*(to|i) = \frac{1 \times 2533}{2533 + 1446} = 0.6333 = C^*(chinese|i)$$

$$C^*(eat|i) = \frac{10 \times 2533}{2533 + 1446} = 6.3325$$

$$C^*(chinese|i) = C^*(spend|i) = \frac{3 \times 2533}{2533 + 1446} = 1.8998$$

Jolplace Smoothing:

- Sharp change in counts & probabilities occur
- In

Smoothing:

- 1) Unsmoothed
- 2) Laplace - Add -on
- 3) Add - k
- 4) Good Turing
- 5) Combining

estimation — Backoff

Good Turing Smoothing:

- N_c - Count of things we've seen c Times

I am I am I am I do not eat

I am 3

I am 2

$$N_1 = 3$$

I am 2

$$N_2 = 2$$

do 1

$$N_3 = 1$$

not 1

eat 1

N_c - frequency of frequency c

N - no. of bigrams seen - V^2 = 9

N_0 - no. of bigrams with count 0 = 0

N_c - no. of bigrams with count c

$$N_0 = V^2 - N$$

$$c^* = \frac{(c+1) \cdot N_{c+1}}{N_c}$$

c - original (real) word count

$(c+1) \cdot N_{c+1}$ - The probability mass for words with frequency $c+1$

c^* - new (adjusted) word count.

* MLE count for N_c is c.

Eg:

10 carp $N_1 = 3$ $N = 18$

3 perch $N_2 = 1$

2 whitefish $N_3 = 1$

1 trout $N_{10} = 1$

1 salmon $N_0 = 2$

1 eel

MLE count of hitherto-unseen species (catfish or bass) is 0

$$P_{GT}^*(\text{things with frequency 0 in training}) = \frac{N_1}{N} = \frac{3}{18}$$

$$P_{GT}^*(\text{bass}) = P_{GT}^*(\text{catfish}) = \frac{3}{18} * \frac{1}{N_b}$$

	unseen (bass or catfish)	trout
c	0	1
MLE P	$\hat{p} = \frac{0}{18} = 0$	$\frac{1}{18}$
c^*	$P_{GT}^*(\text{unseen}) = \frac{N_1}{N} = \frac{3}{18} = 0.17$	$c^*(\text{trout}) = 2 \times \frac{N_3}{N} = 2 \times \frac{1}{3} = 0.67$
$G_T P_{GT}^*$	$P_{GT}^*(\text{unseen}) = \frac{N_1}{N} = \frac{3}{18} = 0.17$	$P_{GT}^*(\text{trout}) = \frac{0.67}{18} = \frac{1}{27} = 0.037$

c^* computation isn't there for unseen i.e., for count 0

Revised count.

Consider vocabulary as $\{a, b, c\}$

Consider Corpus: ba b a a c b c a c a c

Find list of possible bigrams

Find out count of seen (observed) & unseen (unobserved) bigram
& Their frequencies using Good Turing.

	a	b	c	
a	1	1	3	(ba, ab, ba, aa, ac, cb, b, ca, ca, ac)
b	2	0	1	aa = 1 ba = 2 ca = ?
c	2	1	0	ab = 1 bb = 0 cb = 1 ac = 3 bc = 1 cc = 0
	3	5	$N = 12$	
	5	3		$N_0 = 2$ $N = 11$
	4	2		$N_1 = 4$
	2	1		$N_2 = 2$
	1	0		$N_3 = 1$

C	C_{GT}^*		P_{GT}^*		$P_{GT}^* \text{ (based on no. of words)}$
	P_{MLE}	C_{GT}^*	P_{GT}^*		
$C=0$	0.	0	$\frac{4}{11}$	$\frac{N_1}{N}$	$P_{GT}^*(bb) = P_{GT}^*(cc) = \frac{2}{11} \frac{N_1}{N}$
$C=1$	$\frac{1}{11}$	1	$2 \times \frac{2}{4} \times \frac{1}{11}$ $= \frac{1}{11}$		$P_{GT}^*(aa) = P_{GT}^*(ab) = P_{GT}^*(bc)$ $= P_{GT}^*(cb) = \frac{1}{11}$
$C=2$	$\frac{2}{11}$	$3 \times \frac{1}{2}$	$\frac{3}{2} \times \frac{1}{11} = \frac{3}{22}$		$P_{GT}^*(ba) = P_{GT}^*(ca)$ $= \frac{3}{22}$
	$\frac{4}{11}$	$\frac{4}{11}$	$\frac{6}{22}$	$= \frac{8+8+6}{22}$	

replace:

	a	b	c
a	1	1	3
b	2	0	1
c	2	1	0

unsmoothed.

$$c_i^* = C(w_{i+1}) * c(w_{i-1})$$

$$P_{\text{replace}} \# (w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i) + 1}{C(w_{i-1}) + V}$$

	a	b	c
a	2	2	4
b	3	1	2
c	3	2	1

$$V = 9$$

Add one smoothed.

$$P^*(a|a) = \frac{2}{5+9} = \frac{2}{14} = \frac{1}{7}$$

$$P^*(b|a) = \frac{3}{5+9} = \frac{3}{14} = \frac{1}{7}$$

$$P^*(c|a) = \frac{4}{5+9} = \frac{4}{14} = \frac{2}{7}$$

$$P^*(a|b) = \frac{3}{3+9} = \frac{3}{12} = \frac{1}{4}$$

$$P^*(b|b) = 0 = P^*(c|c)$$

$$P^*(c|b) = \frac{1}{3+9} = \frac{1}{12}$$

$$P^*(a|c) = \frac{2}{3+9} = \frac{2}{12}$$

$$P^*(b|c) = \frac{1}{3+9} = \frac{1}{12}$$

$$\mathbf{N} \quad \frac{(c+1)N}{N+V}$$

a	b	c
6	3	4

	a	b	c
a	2	2	4
b	3	1	3
c	3	2	1

	a	b	c
a	1.25	1.25	2.5
b	1.5	0.5	1
c	1.71	1.14	0.54

Add one $\frac{(I+1)N}{N+V}$
 reconstructed
 bigram count = $\frac{2 \times 5}{5+3}$
 c*

	a	b	c
a	0.25	0.25	0.5
b	0.5	0.0	0.33
c			

~~CPM~~ SLIA

$$P(w_{i+1} | w_i) = \frac{C(w_{i-1}, w_i) + 1}{C(w_{i-1}) + V}$$

$$P(a|a) = \frac{C(a, a) + 1}{C(a) + V}$$

$$= \frac{2}{5+3}$$

$$= 2/8$$

$$= 0.25$$

Interpolation:

- Simple linear interpolation:
- We combine different order N-grams by linearly interpolating all models.
- To combine estimate Trigram probability $P(w_n | w_{n-1}, w_{n-2})$
- :

$$P(w_n | w_{n-1}, w_{n-2}) = \lambda_1 P(w_n | w_{n-1}, w_{n-2}) \\ + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n)$$

backoff N-gram modelling: Katz backoff

$$P_{\text{Katz}}(z|x,y) = \begin{cases} P^*(z|x,y) & \text{if } c(x,y,z) \geq 0 \\ \alpha(x,y) P_{\text{NG}}(z|y) & \text{else if } c(y,z) \geq 0 \\ P^*(z) & \text{otherwise} \end{cases}$$

$$P_{\text{NG}}(z|y) = \begin{cases} P^*(z|y) & \text{if } c(y,z) > 0 \\ \alpha(y) P^*(z) & \text{otherwise} \end{cases}$$

Different

Verbs : VB, VBP, VBZ, VBD, VBG, VBN

- base, present - non-3rd, plural - 3rd, past - ing, -en

Nouns : NNP, NNPS, NN, NNS

- proper/common, singular/plural (singular includes mass+generic)

Adjectives : JJ, JJR, JJS (base, comparative, superlative)

Adverbs : RB, RBR, RBS, RP (base, comparative, superlative, particle)

Pronouns : PRP, PRP\$ (personal, possessive)

Interrogatives : WP, WP\$, WDT, WRB (compare to : PRP, PRP\$, DT, RB)

Other Closed Class : CC, CD, DT, PDT, IN, MD

Punctuation : # \$, .

Examples:

The grand jury commented on a no. of other topics

- Eg: Determiner : the and a
- The adjective grand & other
- The common nouns jury, number & topics,
- The past tense verb commented

1) $P(t_n | t_{n-1}) \rightarrow$ tag transition probability
prior probability

2) $P(w_i | t_i) \rightarrow$ observation likelihood probability / estimated

HMM part of speech tagging

- We want, out of all sequences of n tags t_1, \dots, t_n , the single sequence such that $P(t_1, \dots, t_n | w_1, \dots, w_n)$ is highest

$$\hat{t}_1^n = \arg \max_{t_1^n} P(t_1^n | w_1^n)$$

- $\hat{\cdot}$ means 'our estimate is best one'

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

$$P(t_1^n | w_1^n) = \frac{P(w_1^n | t_1^n) \cdot P(t_1^n)}{P(w_1^n)}$$

=

HMM simplifying assumptions

- It is too hard to compute $\hat{t}_1^n = \arg \max_{t_1^n} P(w_1^n | t_1^n)$

- Assumption 1: The probability of a word is dependent only on part-of-speech, not on the surrounding words or other tags around it.

$$P(w_i^n | t_i^n) \approx \prod_{i=1}^n P(w_i | t_i)$$

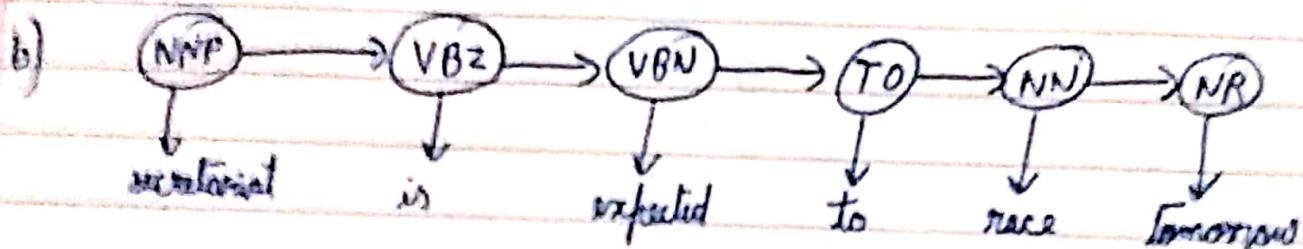
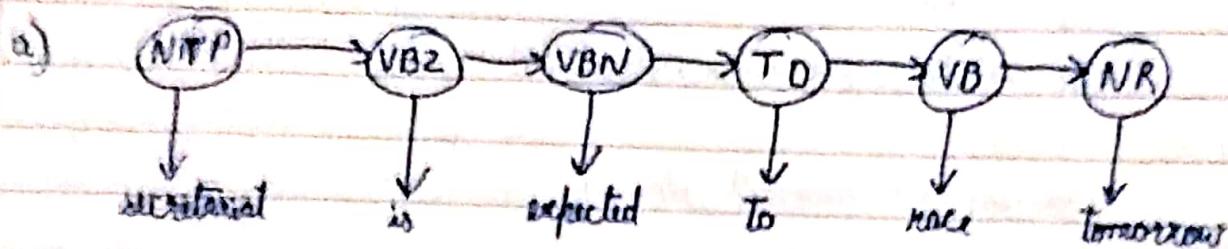
- Assumption 2: The probability of a tag appearing is dependent on the previous tag (bigram assumption)

$$P(t_i) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

word likelihoods, log transition probabilities

$$\hat{t}_i = \arg \max P(t_i | w_i) \approx \arg \max \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

Disambiguating "race"



$$\begin{aligned}
 a) \quad & P(VB | TO) * P(NR | VB) * P(race | VB) \\
 & \approx 0.83 * 0.0027 * 0.00012 \\
 & \approx \underline{\underline{2.69 \times 10^{-7}}}
 \end{aligned}$$

$$\begin{aligned}
 b) \quad & P(NN | TO) * P(NR | NN) * P(race | NN) \\
 & \approx 0.00067 * 0.0012 * 0.00057 \\
 & \approx \underline{\underline{3.21 \times 10^{-10}}}
 \end{aligned}$$

Start probability, $P(\pi) = 1$

Draw part of

$S \rightarrow NP \ VP$

$NP \rightarrow \text{pronoun} |$

Proper-noun |

Det.-nominal

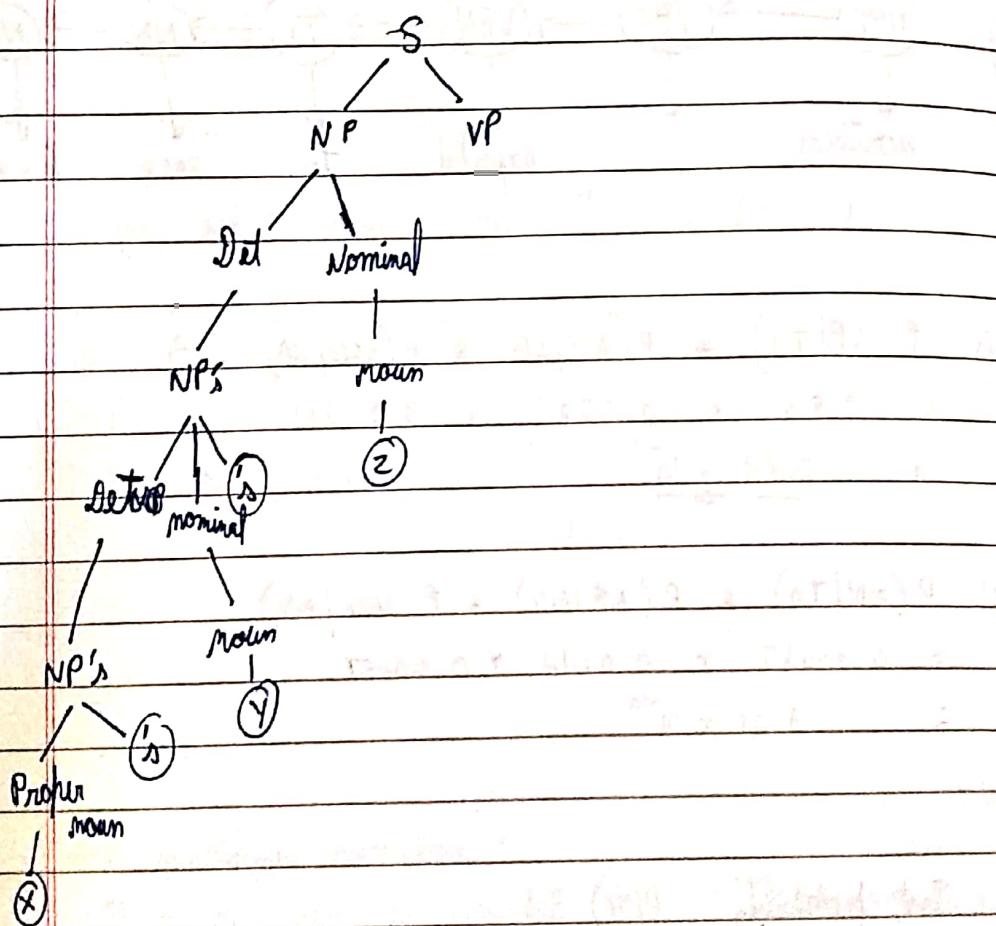
$\text{Det} \rightarrow NP's$

$\text{Det} \rightarrow \text{the} | \text{a} | \text{an} | \text{this} | \text{these} | \text{that}$

$\text{Nominal} \rightarrow \text{Nominal Noun} |$

Noun

→ ~~Draw~~



$S \rightarrow NP \ VP$

$NP \rightarrow \text{Pronoun} | \text{Proper noun} | \text{Det. Nominal}$

$\text{Nominal} \rightarrow \text{Nominal noun} | \text{Noun} | \text{Nominal PP (pp)}$

$\text{VP} \rightarrow \text{Verb} | \text{Verb NP} | \text{Verb NP PP} | \text{Verb PP}$

$\text{PP} \rightarrow \text{Preposition NP}$

PP postmodifiers
- any

Draw Top-down parse Tree for following

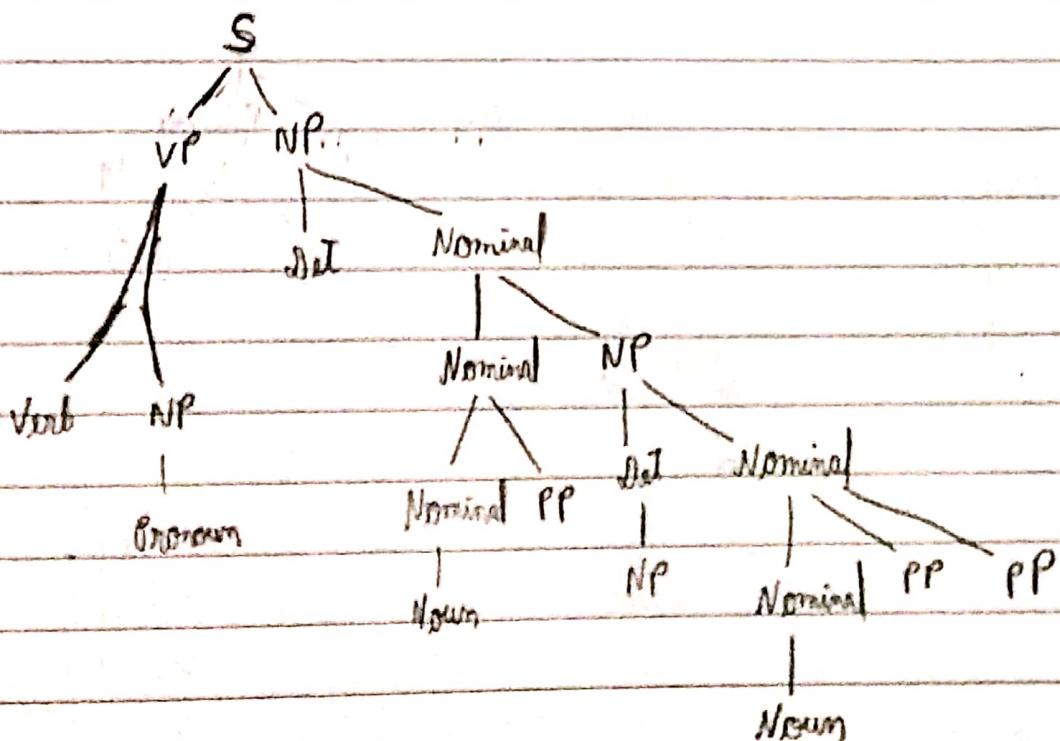
1) [All] [The] [morning flights] [from P1] [To P2] [having [before [to]]]
(adit Det Nominal PP PP GerundV GerundVP N)

Nominal → Nominal Gerund VP

Gerund VP → Gerund V NP | Gerund V PP | Gerund V / Gerund V NPP

Gerund V → being / preferring / arriving / having / ...

[Show] [me] [The] [meal] [on] [flight Boeing -780] [from P1] [To P2]
Verb Interrog Det noun NP PP PP



$S \rightarrow NP \ VP$

$NP \rightarrow \text{Pre-Head Modifiers Nominal Post-head Modifiers}$

$\text{Pre-Head-Modifiers} \rightarrow \text{Predet Det Post-Determiners}$

$\text{Post-Determiners} \rightarrow (\text{Card}) (\text{Ord}) (\text{Quant}) (\text{AP})$

$NP \rightarrow \text{Pronoun} / \text{Proper-Name} / \text{Det Nominal}$

$\text{Det} \rightarrow \text{NPs'}$

$\text{Nominal} \rightarrow \text{Nominal Noun} / \text{Noun} / \text{Nominal NP}$

$\text{Pre-head Modifiers} \rightarrow (\text{Card}) (\text{Ord}) (\text{Quant}) (\text{AP})$

$\text{Nominal} \rightarrow \text{Nominal PP (PP) (PP)}$

$\text{Nominal} \rightarrow \text{Nominal Gerund VP}$

$\text{Nominal} \rightarrow \text{Nominal RelClause}$

$\text{RelClause} \rightarrow (\text{who} / \text{that}) VP$

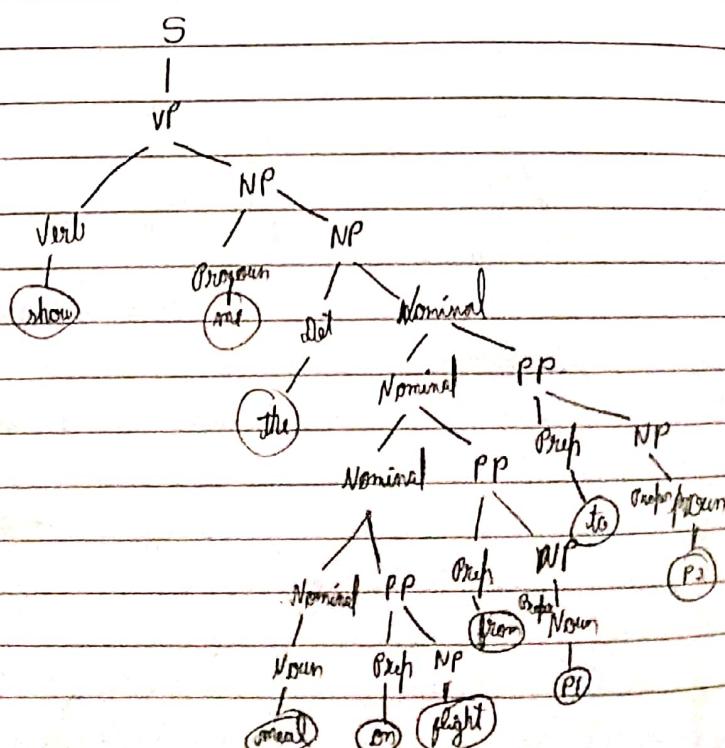
$\text{Gerund VP} \rightarrow \text{Gerund V} / \text{Gerund V NP} / \text{Gerund V PP} / \text{Gerund V NP PP}$

$VP \rightarrow \text{Verb} / \text{Verb NP} / \text{Verb PP} / \text{Verb NP PP}$

$PP \rightarrow \text{Preposition NP}$

Note: Post-head-Modifiers \rightarrow PR-Modifier / Non-finite-Modifier / Rel-Clause

PP-Modifier \rightarrow



CC	coordinating conjunction	and, but, or
CD	cardinal no.	one, Two, Three
DT	determiner	a, an, the
EX	existential 'there'	There
FW	Foreign word	meat, coffee
IN	Preposition / sub-conj	of, in, by
JJ	adjective	yellow
JJR	adj. comparative	bigger
JJS	adj. superlative	wildest
LS	list item marker	1, 2, one
MD	Modal	can, should
NN	Noun, singular or mass	llama
NNS	Noun plural	llamas
NNP	Proper noun, singular	IBM
NNPS	Proper noun, plural	Caroline's
PDT	pre-determiner	all, both
POS	possessive ending	's
PP	personal pronoun	I, you, he
PPS	possessive personal pronoun	your, one's
RB	adverb	quickly, never
RBR	adverb, comparative	faster
RBS	adverb, superlative	fastest
RP	particle	up, off
SYM	symbol	+, -, x
TO	"to"	To
UH	interjection	oh, ph
VB	Verb, base form	eat
VBD	Verb, past tense	ate
VBG	Verb, gerund	eating
VBN	Verb, past participle	eaten
VBP	Verb, non-3sg pres	eat
VBZ	Verb, 3sg pres	eats
WDT	wh-determiner	which, that
WP	wh-pronoun	what, who

WP\$	possessive wh-	whose
WRB	wh-adverb	how, where

Modal Verbs

Auxiliaries

do

have

progressive be

passive be

modal < perfect < progressive < passive

modal perfect

could have been a contender

modal passive

will be married

perfect progressive

have been fearing

modal perfect passive

might have been prevented

Co-ordination