



Natural Language Processing

Info 159/259

Lecture 15: Review (Oct 11, 2018)

David Bamman, UC Berkeley

Big ideas

- Classification
 - Naive Bayes, Logistic regression, feedforward neural networks, CNN.
- Where does NLP data come from?
 - Annotation process
 - Interannotator agreement
- Language modeling
 - Markov assumption, featurized, neural
- Probability/statistics in NLP
 - Chain rule of probability, independence, Bayes' rule

Big ideas

- Lexical semantics
 - Distributional hypothesis
 - Distributed representations
 - Subword embedding models
 - Contextualized word representations (ELMO)
- Evaluation metrics (accuracy, precision, recall, F score, perplexity, parseval)
- Sequence labeling
 - POS, NER
 - Methods: HMM, MEMM, CRF, RNN, BiRNN
- Trees
 - Phrase-structure parsing, CFG, PCFG
 - CKY for recognition, parsing

Big ideas

- What defines the models we've seen so far? What formally distinguishes an HMM from an MEMM? How do we train those models?
- For all of the problems we've seen (sentiment analysis, POS tagging, phrase structure parsing), how do we evaluate the performance of different models?
- If faced with a new NLP problem, how would you decide between the alternatives you know about? How would you adapt an MEMM, for example, to a new problem?

Midterm

- In class next Tuesday
- Mix of multiple choice, short answer, long answer
- Bring 1 cheat sheet (1 page, both sides)
- Covers all material from lectures and readings

Multiple choice

1. John builds a system to detect lies in speeches by politicians. To evaluate his system, he halts development and runs every speech given by a national political candidate through the system, during one week. The system hypothesizes that ten statements are lies, which John has an expert human fact-checker check. Nine of the ten are found to be lies, and one is found to be true. “My system achieves 90% accuracy!” exclaims John. Your response:
 - A. Yes! Great job, John!
 - B. No, John, your system achieves 90% F-measure.
 - C. No, John, your system achieves 90% recall.
 - D. No, John, your system achieves 90% precision.

Multiple choice

2. Pointwise mutual information is a measure of association between:
- A. Two random variables
 - B. A random variable and one of its values
 - C. A word and document label
 - D. Two values of two random variables

Short answer

What is regularization and why is it important?

Short answer

For sequence labeling problems like POS tagging and named entity recognition, what are two strengths of using a bidirectional LSTM over an HMM? What's one weakness?

Long answer

1. Here is a fragment of a PCFG:

Production	log P(Production)	Notes
S → NP VP	-0.70	
S → Noun VP	-2.30	
VP → Verb Noun	-0.70	
VP → Verb ProperNoun	-1.50	
VP → buffalo	-4.00	
NP → ProperNoun Noun	-0.70	
Verb → buffalo	-4.00	"to bully or harrass"
ProperNoun → buffalo	-2.30	A city in New York
Noun → buffalo	-3.00	bison

Table 1: Fragment of PCFG

- (a) Use that grammar and (optionally) the chart below to find the most probable parse of *buffalo buffalo buffalo*. Draw the parse tree and its probability.
- (b) If the sentence is ambiguous, list the interpretations, explaining their differences in plain English.

2. **@kimkierkegaardashian** is a Twitter account that mashes up the language of celebrity Kim Kardashian and philosopher Søren Kierkegaard.

The figure displays three tweets from the Twitter account @KimKierkegaard, which is a mashup of Kim Kardashian and Søren Kierkegaard. The tweets are as follows:

- KimKierkegaard** @KimKierkegaard · Sep 8
It is hard to try on clothes without some question arising as to my relationship to the eternal
2 replies, 444 retweets, 898 likes
- KimKierkegaard** @KimKierkegaard · Sep 8
The perfect white tee reminds you that being nothing in this world is the condition for being something in the next
2 replies, 171 retweets, 497 likes
- KimKierkegaard** @KimKierkegaard · Sep 8
People despair about being lonely and therefore get married. But is this love? I should say it is self-love. Happy anniversary babe
1 reply, 149 retweets, 591 likes

Figure 2: **@kimkierkegaardashian**

- (a) Assume independent language models have been trained on the tweets of Kim Kardashian (generating language model \mathcal{L}_{Kim}) and the writings of Søren Kierkegaard (generating language model $\mathcal{L}_{\text{Søren}}$). Using concepts from class, how could you use \mathcal{L}_{Kim} and $\mathcal{L}_{\text{Søren}}$ to create a new language model $\mathcal{L}_{\text{Kim+Søren}}$ to generate tweets like those above?

2. **@kimkierkegaardashian** is a Twitter account that mashes up the language of celebrity Kim Kardashian and philosopher Søren Kierkegaard.

The image shows three tweets from the @kimkierkegaardashian Twitter account, each featuring a small profile picture of Kim Kardashian and the account name followed by the date (Sep 8). The tweets are separated by horizontal lines.

- KimKierkegaardashian @KimKierkegaard · Sep 8**
It is hard to try on clothes without some question arising as to my relationship to the eternal
2 444 898 ✉
- KimKierkegaardashian @KimKierkegaard · Sep 8**
The perfect white tee reminds you that being nothing in this world is the condition for being something in the next
2 171 497 ✉
- KimKierkegaardashian @KimKierkegaard · Sep 8**
People despair about being lonely and therefore get married. But is this love? I should say it is self-love. Happy anniversary babe
1 149 591 ✉

Figure 2: **@kimkierkegaardashian**

(b) How would you control that model to sound more like Kierkegaard than Kardashian?

2. **@kimkierkegaardashian** is a Twitter account that mashes up the language of celebrity Kim Kardashian and philosopher Søren Kierkegaard.

KimKierkegaardashian @KimKierkegaard · Sep 8
It is hard to try on clothes without some question arising as to my relationship to the eternal

2 444 898

KimKierkegaardashian @KimKierkegaard · Sep 8
The perfect white tee reminds you that being nothing in this world is the condition for being something in the next

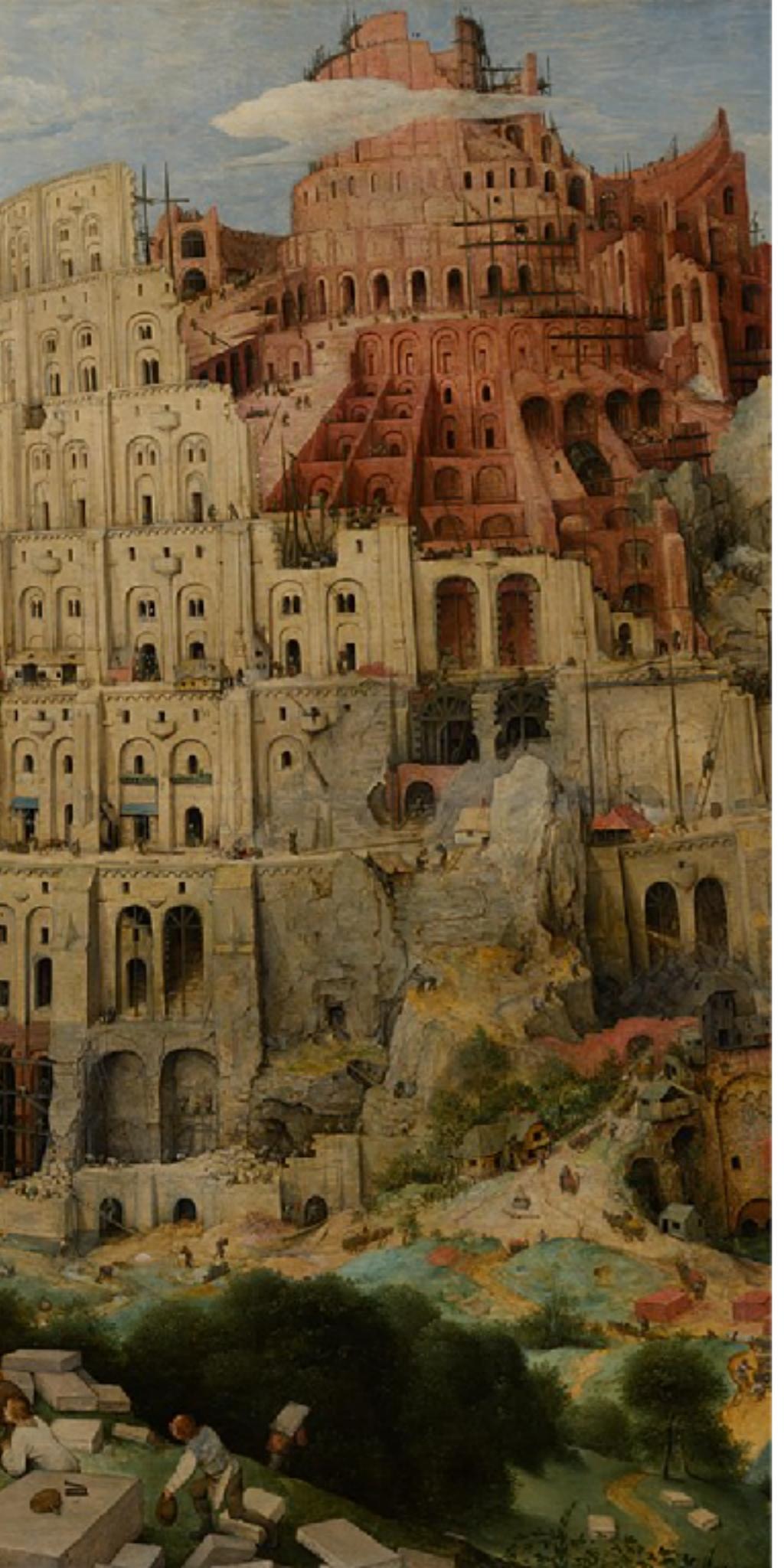
2 171 497

KimKierkegaardashian @KimKierkegaard · Sep 8
People despair about being lonely and therefore get married. But is this love? I should say it is self-love. Happy anniversary babe

1 149 591

Figure 2: **@kimkierkegaardashian**

(c.) Assume you have access to the full Twitter archive of **@kimkierkegaardashian**. How could you choose the best way to combine \mathcal{L}_{Kim} and $\mathcal{L}_{\text{Søren}}$? How would you operationalize “best”?



Classification

A mapping h from input data $\textcolor{magenta}{x}$ (drawn from instance space \mathcal{X}) to a label (or labels) $\textcolor{magenta}{y}$ from some enumerable output space \mathcal{Y}

$\textcolor{magenta}{\mathcal{X}}$ = set of all documents
 $\textcolor{magenta}{\mathcal{Y}}$ = {english, mandarin, greek, ...}

$\textcolor{magenta}{x}$ = a single document
 $\textcolor{magenta}{y}$ = ancient greek

Text categorization problems

task	x	y
language ID	text	{english, mandarin, greek, ...}
spam classification	email	{spam, not spam}
authorship attribution	text	{jk rowling, james joyce, ...}
genre classification	novel	{detective, romance, gothic, ...}
sentiment analysis	text	{positive, negative, neutral, mixed}

Bayes' Rule

Prior belief that $Y = y$
(before you see any data)

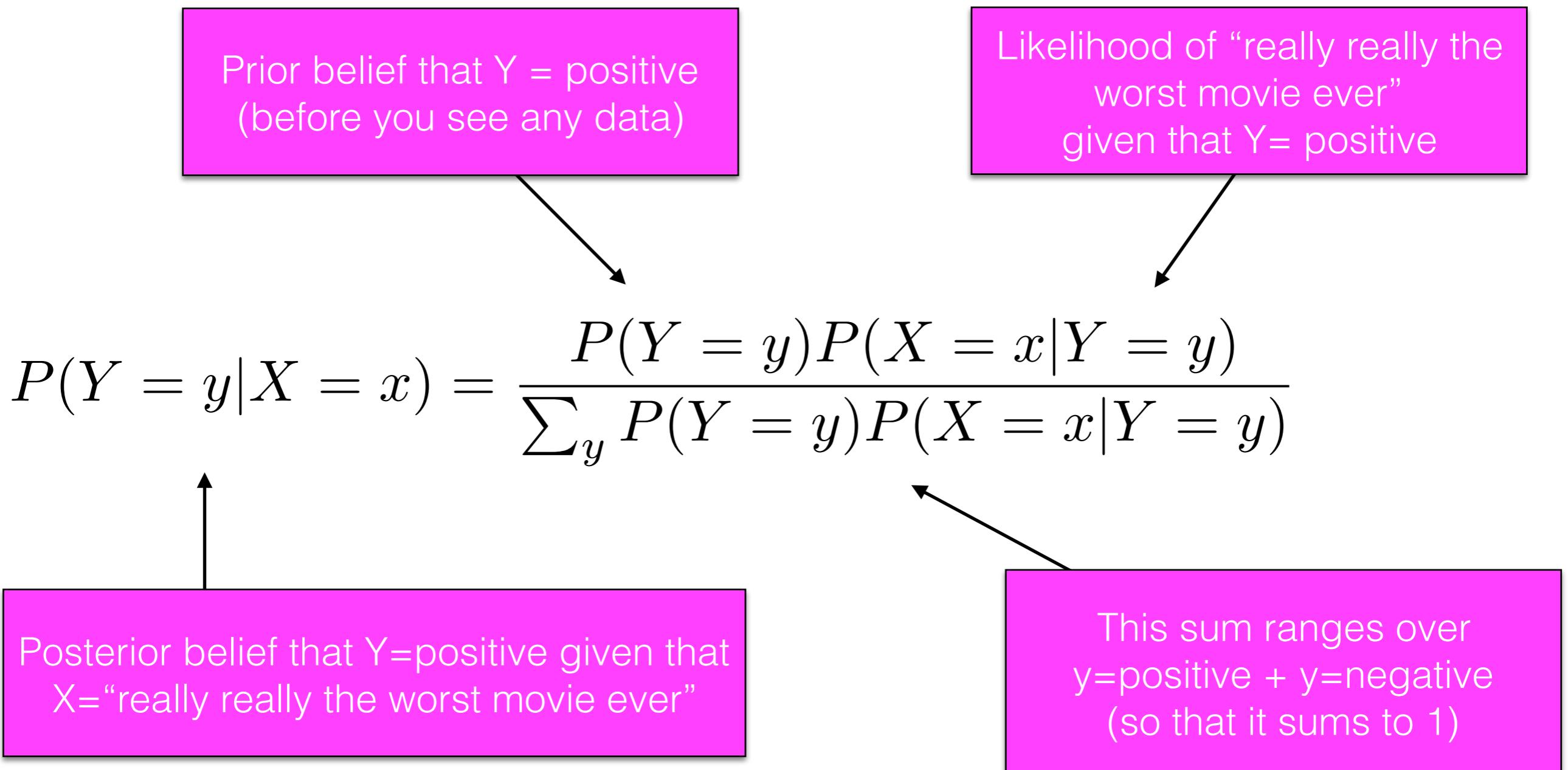
Likelihood of the data
given that $Y=y$

$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$



Posterior belief that $Y=y$ given that $X=x$

Bayes' Rule



Logistic regression

$$P(y = 1 \mid x, \beta) = \frac{1}{1 + \exp\left(-\sum_{i=1}^F x_i \beta_i\right)}$$

output space $\mathcal{Y} = \{0, 1\}$

x = feature vector

Feature	Value
the	0
and	0
bravest	0
love	0
loved	0
genius	0
not	0
fruit	1
BIAS	1

β = coefficients

Feature	β
the	0.01
and	0.03
bravest	1.4
love	3.1
loved	1.2
genius	0.5
not	-3.0
fruit	-0.8
BIAS	-0.1

Features

- As a discriminative classifier, logistic regression doesn't assume features are independent like Naive Bayes does.
- Its power partly comes in the ability to create richly expressive features with out the burden of independence.
- We can represent text through features that are not just the identities of individual words, but any feature that is scoped over **the entirety of the input**.

features

contains like

has word that shows up in positive sentiment dictionary

review begins with “I like”

at least 5 mentions of positive affectual verbs (like, love, etc.)

Stochastic g.d.

- Batch gradient descent reasons over every training data point for each update of β . This can be slow to converge.
- Stochastic gradient descent updates β after each data point.

Algorithm 2 Logistic regression stochastic gradient descent

```
1: Data: training data  $x \in \mathbb{R}^F, y \in \{0, 1\}$ 
2:  $\beta = 0^F$ 
3: while not converged do
4:   for  $i = 1$  to N do
5:      $\beta_{t+1} = \beta_t + \alpha (y_i - \hat{p}(x_i)) x_i$ 
6:   end for
7: end while
```

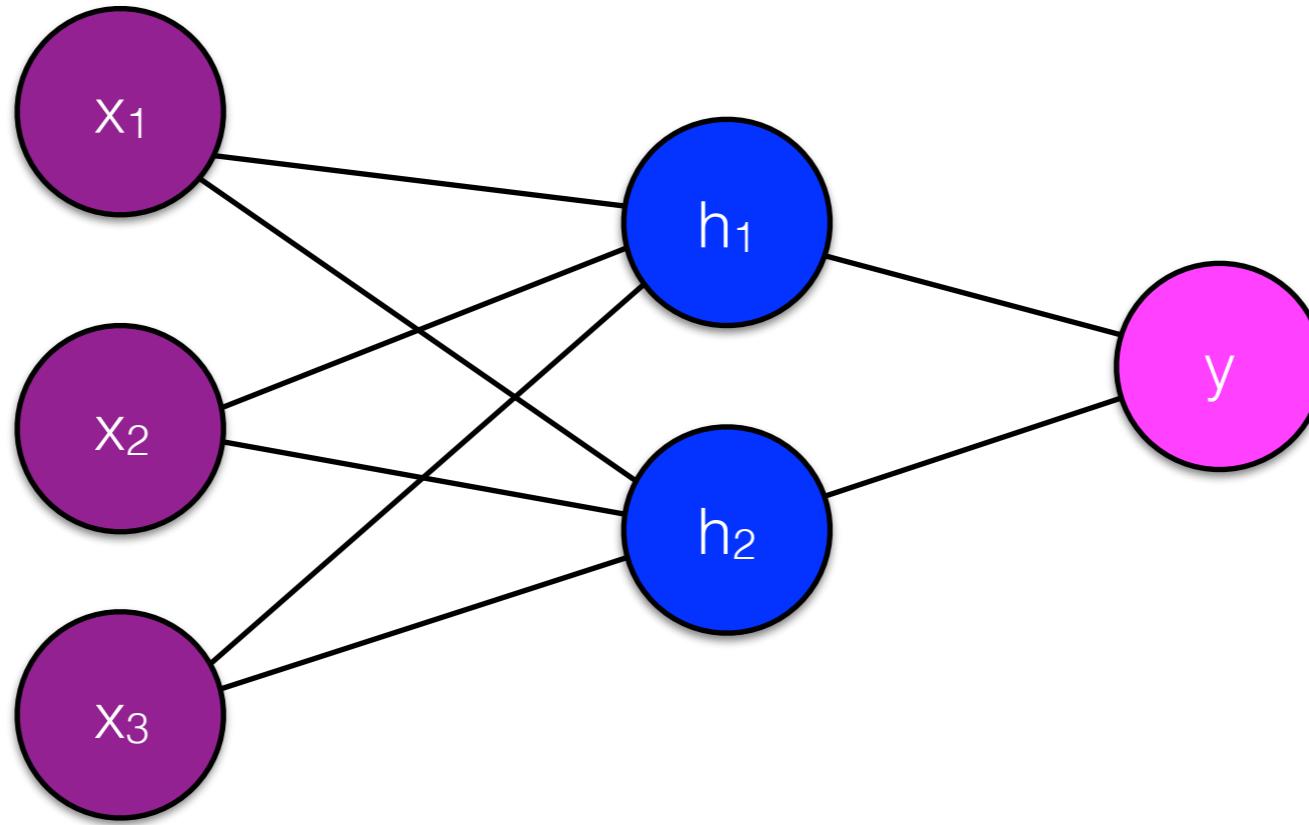
L2 regularization

$$\ell(\beta) = \underbrace{\sum_{i=1}^N \log P(y_i | x_i, \beta)}_{\text{we want this to be high}} - \underbrace{n \sum_{j=1}^F \beta_j^2}_{\text{but we want this to be small}}$$

- We can do this by changing the function we're trying to optimize by adding a penalty for having values of β that are high
- This is equivalent to saying that each β element is drawn from a Normal distribution centered on 0.
- n controls how much of a penalty to pay for coefficients that are far from 0 (optimize on development data)

W

V



$$\hat{y} = \sigma \left[V_1 \left(\sigma \left(\sum_i^F x_i W_{i,1} \right) \right) + V_2 \left(\sigma \left(\sum_i^F x_i W_{i,2} \right) \right) \right]$$

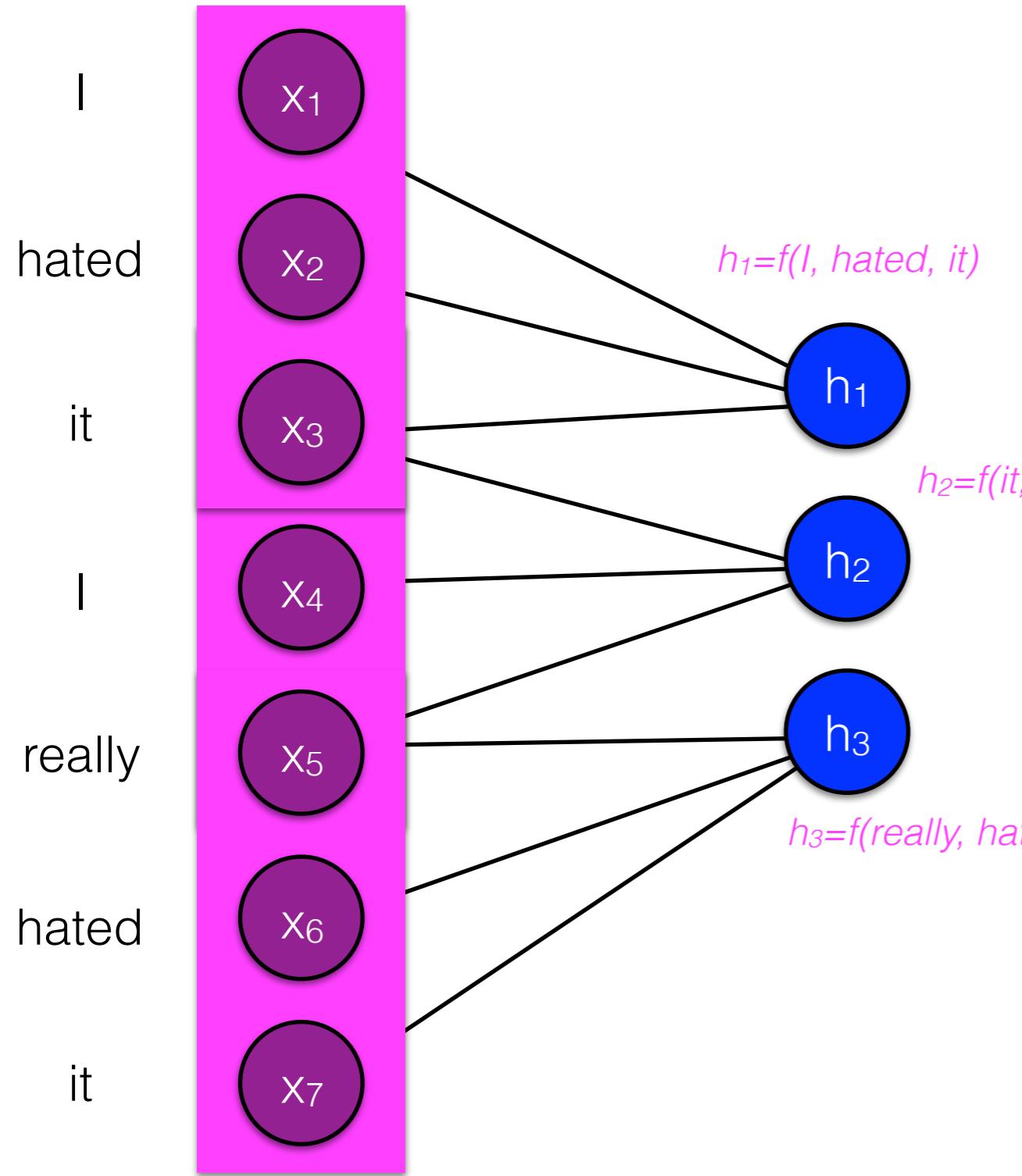
we can express y as a function only of the input x and the weights W and V

$$\hat{y} = \sigma \left[V_1 \underbrace{\left(\sigma \left(\sum_i^F x_i W_{i,1} \right) \right)}_{h_1} + V_2 \underbrace{\left(\sigma \left(\sum_i^F x_i W_{i,2} \right) \right)}_{h_2} \right]$$

This is hairy, but **differentiable**

Backpropagation: Given training samples of $\langle x, y \rangle$ pairs, we can use stochastic gradient descent to find the values of W and V that minimize the loss.

Convolutional networks



$$X \quad \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array}$$

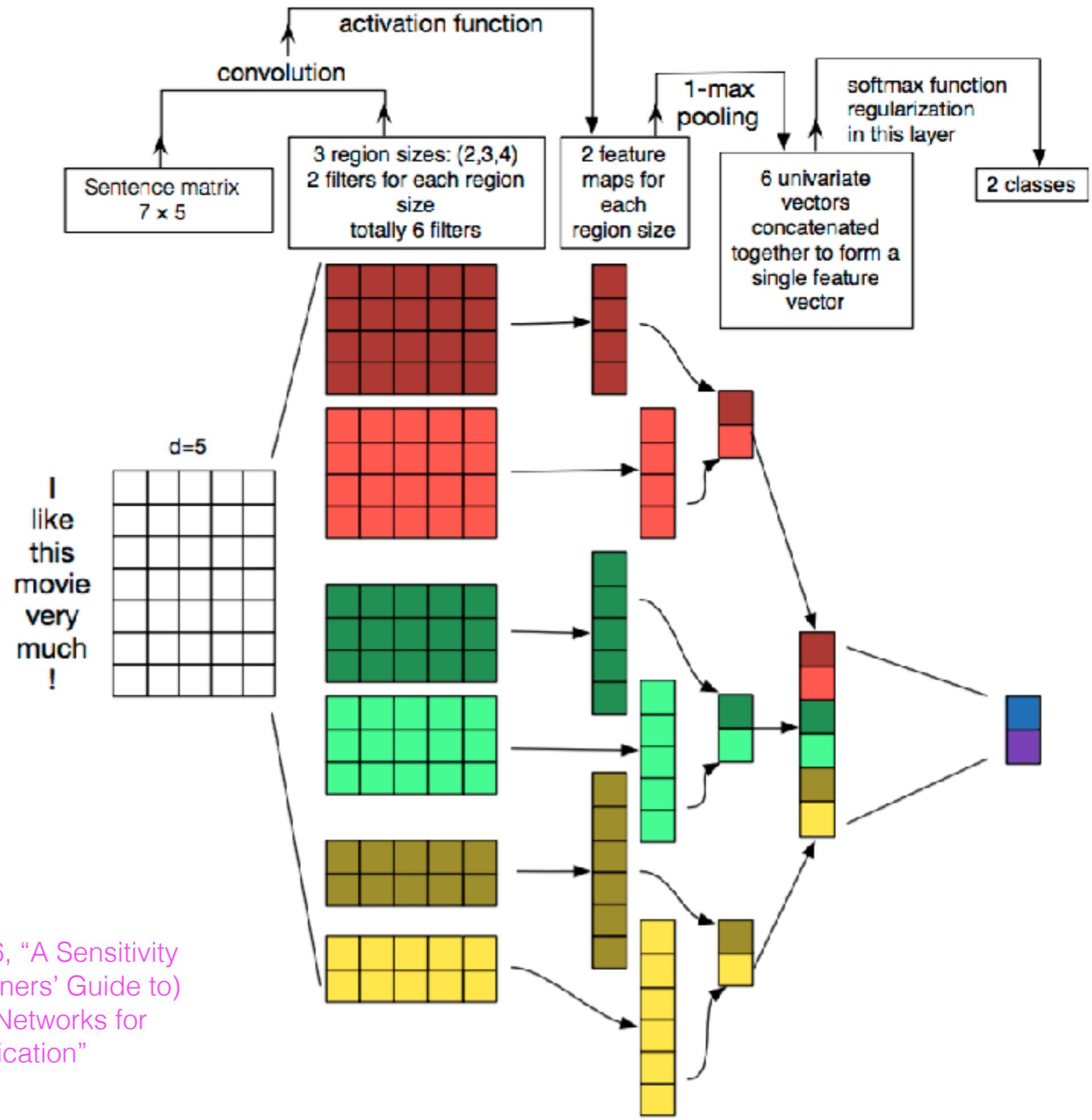
$$W \quad \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array}$$

$$h \quad \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array}$$

$$h_1 = \sigma(x_1 W_1 + x_2 W_2 + x_3 W_3)$$

$$h_2 = \sigma(x_3 W_1 + x_4 W_2 + x_5 W_3)$$

$$h_3 = \sigma(x_5 W_1 + x_6 W_2 + x_7 W_3)$$



Zhang and Wallace 2016, “A Sensitivity Analysis of (and Practitioners’ Guide to) Convolutional Neural Networks for Sentence Classification”

Language Model

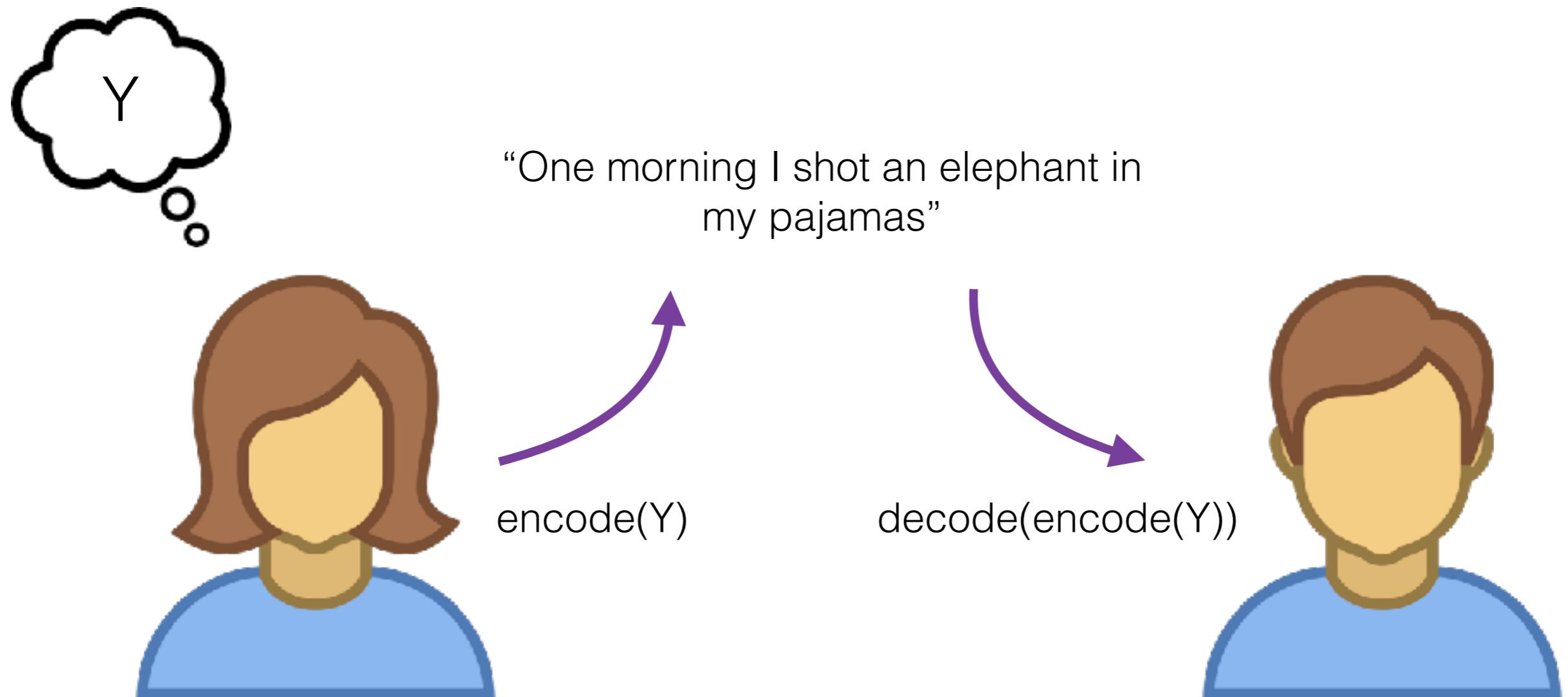
- Language models provide us with a way to quantify the likelihood of sequence — i.e., **plausible** sentences.

OCR

To see great Pompey passe the streets of Rome:
And when you saw his Chariot but appeare,
Haue you not made an Vniuersall shout,
That Tyber trembled vnderneath her bankes
To heare the replication of your sounds,
Made in her Concaue Shores?

- to fee great Pompey paffe the Areets of Rome:
- to see great Pompey passe the streets of Rome:

Information theoretic view



Shannon 1948

Noisy Channel

	X	Y
ASR	speech signal	transcription
MT	target text	source text
OCR	pixel densities	transcription

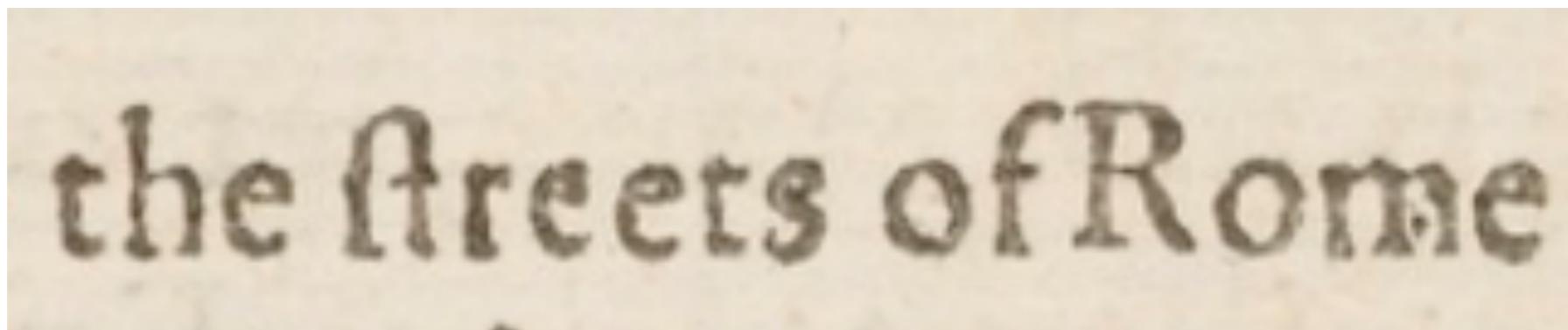
$$P(Y | X) \propto \underbrace{P(X | Y)}_{\text{channel model}} \underbrace{P(Y)}_{\text{source model}}$$

OCR

To see great Pompey passe the streets of Rome :
And when you saw his Chariot but appeare,
Haue you not made an Vniuersall shout,
That Tyber trembled vnderneath her bankes
To heare the replication of your sounds,
Made in her Concaue Shores ?

$$P(Y | X) \propto \underbrace{P(X | Y)}_{\text{channel model}} \underbrace{P(Y)}_{\text{source model}}$$

OCR



the streets of Rome

$$P(Y | X) \propto \underbrace{P(X | Y)}_{\text{channel model}} \underbrace{P(Y)}_{\text{source model}}$$

Language Model

- Language modeling is the task of estimating $P(w)$
- Why is this hard?

$P(\text{"It was the best of times, it was the worst of times"})$

Markov assumption

bigram model
(first-order markov)

$$\prod_i^n P(w_i \mid w_{i-1}) \times P(\text{STOP} \mid w_n)$$

trigram model
(second-order markov)

$$\prod_i^n P(w_i \mid w_{i-2}, w_{i-1})$$

$$\times P(\text{STOP} \mid w_{n-1}, w_n)$$

Smoothing LM

- Additive smoothing; Laplace smoothing
- Interpolating LMs of different orders
- Kneser-Ney
- Stupid backoff

Featurized LMs

- We can use multi class logistic regression for language modeling by treating the vocabulary as the output space

$$\mathcal{Y} = \mathcal{V}$$

Featurized LMs

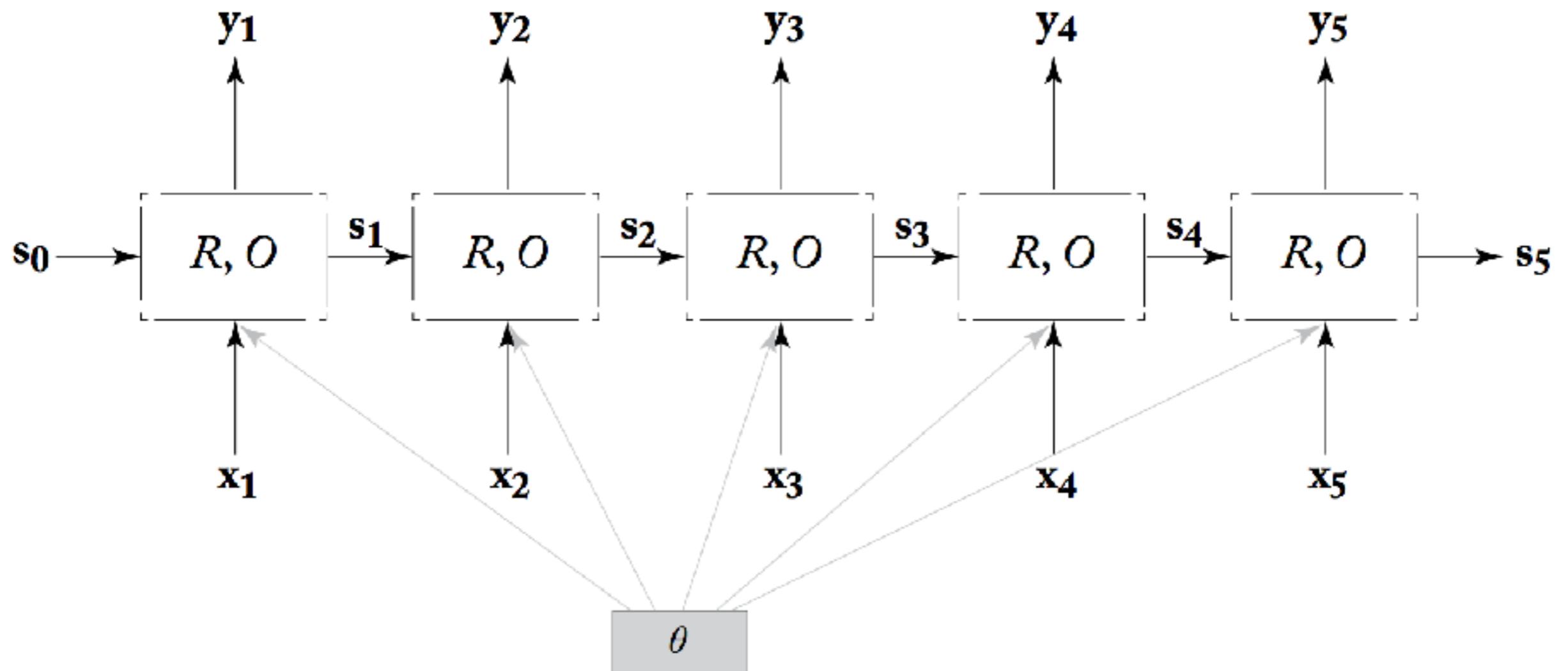
$$P(w_i = \text{dog} \mid w_{i-2} = \text{and}, w_{i-1} = \text{the})$$

	Feature	Value
second-order features	$w_{i-2}=\text{the} \wedge w_{i-1}=\text{the}$	0
	$w_{i-2}=\text{and} \wedge w_{i-1}=\text{the}$	1
	$w_{i-2}=\text{bravest} \wedge w_{i-1}=\text{the}$	0
	$w_{i-2}=\text{love} \wedge w_{i-1}=\text{the}$	0
first-order features	$w_{i-1}=\text{the}$	1
	$w_{i-1}=\text{and}$	0
	$w_{i-1}=\text{bravest}$	0
	$w_{i-1}=\text{love}$	0
BIAS		1

Richer representations

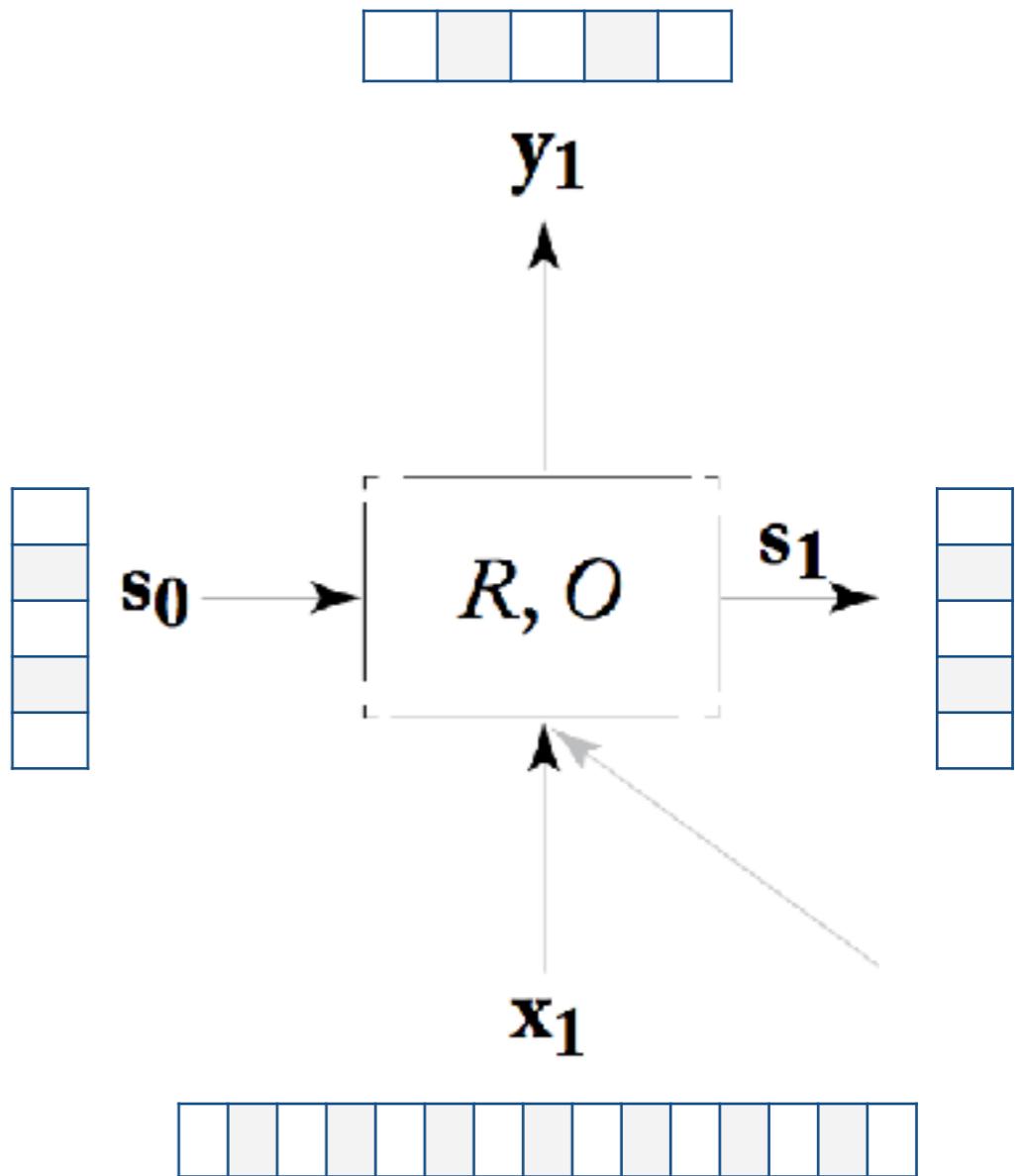
- Log-linear models give us the flexibility of encoding richer representations of the **context** we are conditioning on.
- We can reason about any observations from the entire history and not just the local context.

Recurrent neural network



Recurrent neural network

- Each time step has two inputs:
 - x_i (the observation at time step i); one-hot vector, feature vector or **distributed representation**.
 - s_{i-1} (the output of the previous state); base case: $s_0 = 0$ vector



Distributed representations

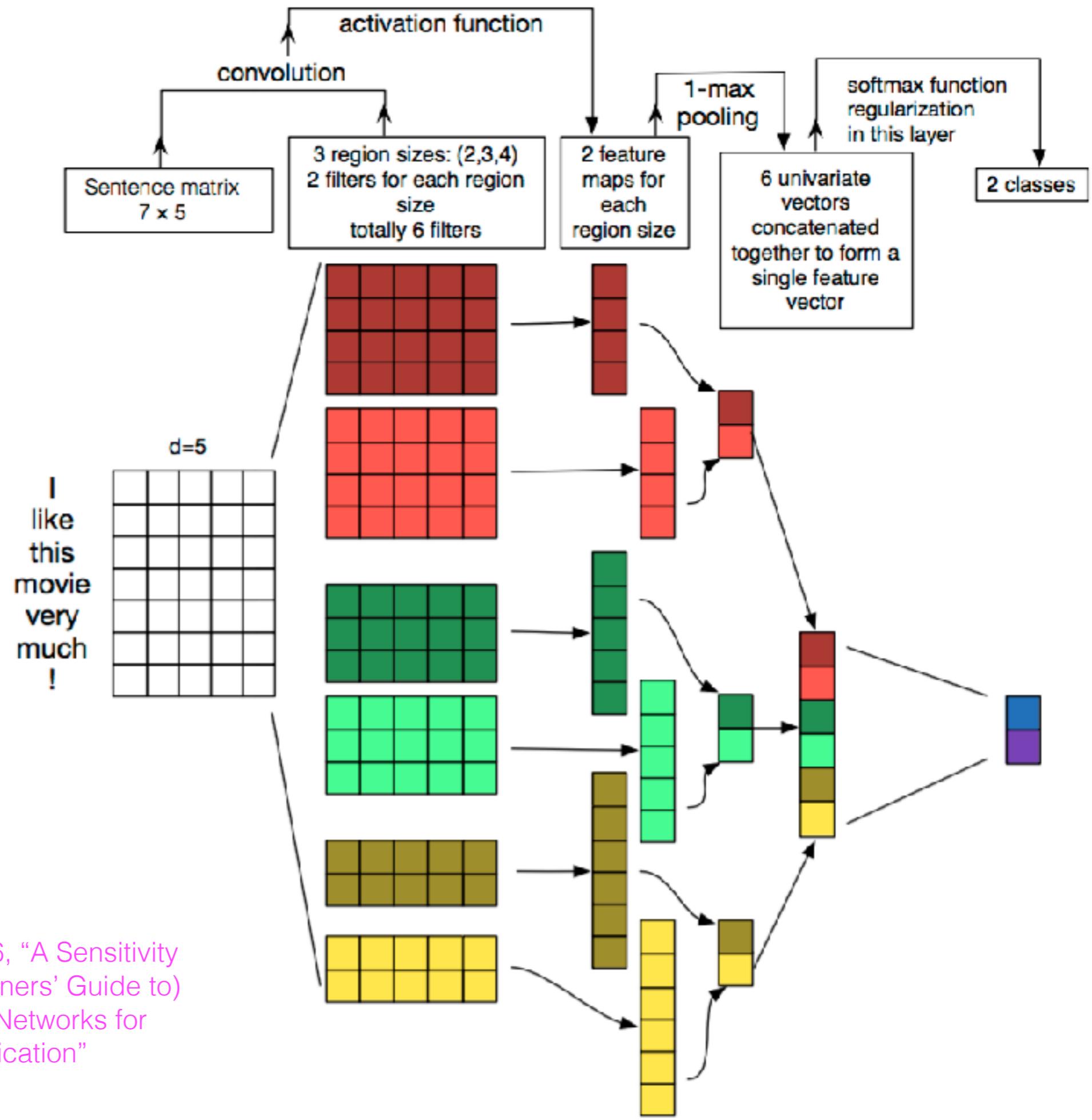
- Low-dimensional, dense word representations are extraordinarily powerful (and are arguably responsible for much of gains that neural network models have in NLP).
- Lets your representation of the input share statistical strength with words that behave similarly in terms of their distributional properties (often **synonyms** or words that belong to the same **class**).

Dense vectors from prediction

- Learning low-dimensional representations of words by framing a predicting task: using context to predict words in a surrounding window
- Transform this into a **supervised** prediction problem; similar to language modeling but we're ignoring order within the context window

Using dense vectors

- In neural models (CNNs, RNNs, LM), replace the V -dimensional sparse vector with the much smaller K -dimensional dense one.
- Can also take the derivative of the loss function with respect to those representations to optimize for a particular task.



Zhang and Wallace 2016, “A Sensitivity Analysis of (and Practitioners’ Guide to) Convolutional Neural Networks for Sentence Classification”

Subword models

- Rather than learning a single representation for each word type w , learn representations z for the set of ngrams \mathcal{G}_w that comprise it [Bojanowski et al. 2017]
- The word itself is included among the ngrams (no matter its length).
- A word representation is the sum of those ngrams

$$w = \sum_{g \in \mathcal{G}_w} z_g$$

FastText

$$e(\text{where}) =$$

$e(<\text{wh}>)$
+ $e(\text{whe})$
+ $e(\text{her})$
+ $e(\text{ere})$
+ $e(\text{re}>)$

3-grams

+ $e(<\text{whe}>)$
+ $e(\text{wher})$
+ $e(\text{here})$
+ $e(\text{ere}>)$

4-grams

+ $e(<\text{wher}>)$
+ $e(\text{where})$
+ $e(\text{here}>)$

5-grams

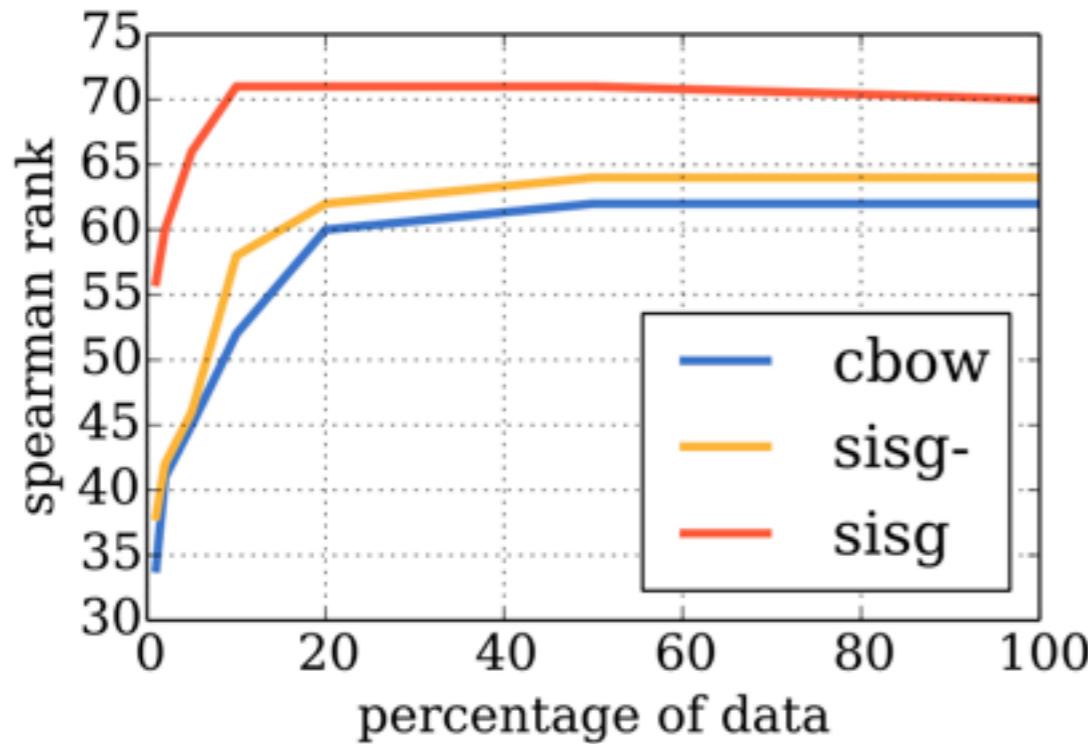
+ $e(<\text{where}>)$
+ $e(\text{where}>)$

6-grams

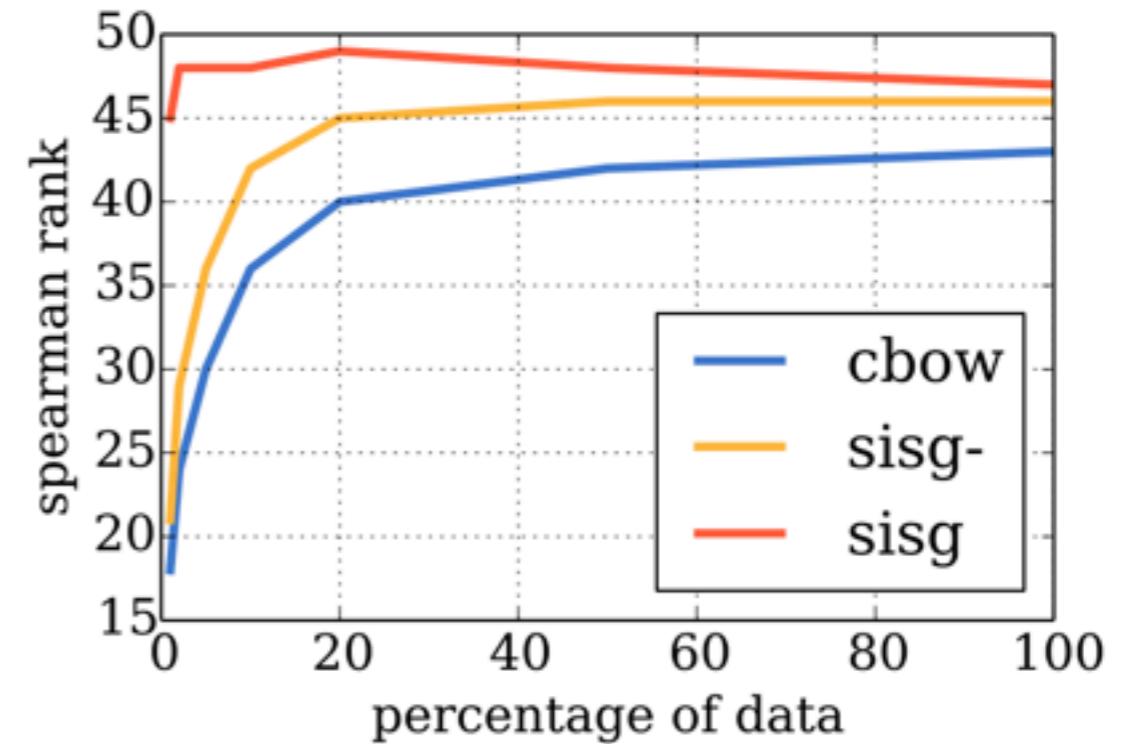
+ $e(<\text{where}>)$

word

$e(*)$ = embedding for *



(a) DE-GUR350



(b) EN-RW

1% =
~20M tokens

100% =
~1B tokens

- Subword models need less data to get comparable performance.

ELMo

- Peters et al. (2018), “Deep Contextualized Word Representations” (NAACL)
- Big idea: transform the representation of a word (e.g., from a static word embedding) to be sensitive to its local context in a sentence and optimized for a specific NLP task.
- Output = word representations that can be plugged into just about any architecture a word embedding can be used.

ELMo

TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lcc et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model baselines across six benchmark NLP tasks. The performance metric varies across tasks – accuracy for SNLI and SST-5; F₁ for SQuAD, SRL and NER; average F₁ for Coref. Due to the small test sizes for NER and SST-5, we report the mean and standard deviation across five runs with different random seeds. The “increase” column lists both the absolute and relative improvements over our baseline.

Parts of speech

- Parts of speech are categories of words defined **distributionally** by the morphological and syntactic contexts a word appears in.

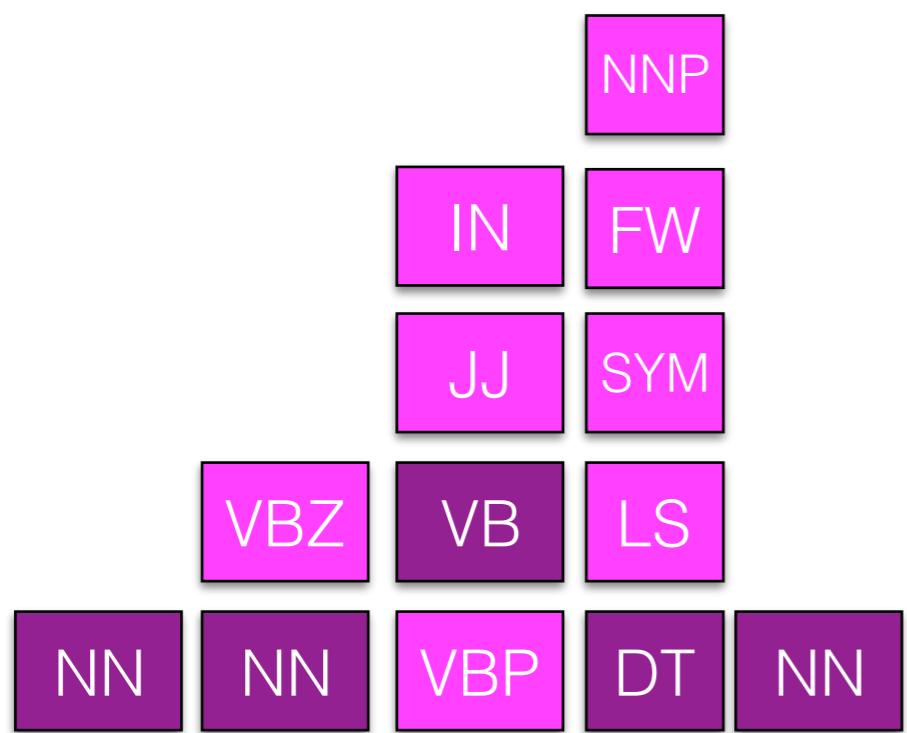
	-s	-ed	-ing
walk	walks	walked	walking
slice	slices	sliced	slicing
believe	believes	believed	believing
of	*ofs	*ofed	*ofing
red	*reds	*redded	*reding

Kim saw the	elephant	before we did
	dog	
	idea	
	*of	
	*goes	

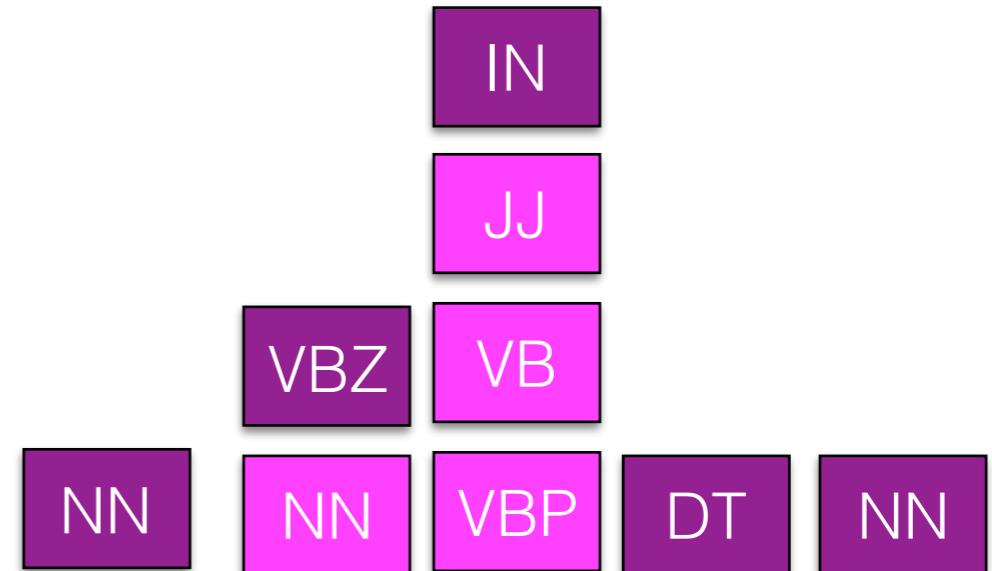
Open class	Nouns	fax, affluenza, subtweet, bitcoin, cronut, emoji, listicle, mocktail, selfie, skort
	Verbs	text, chillax, manspreading, photobomb, unfollow, google
	Adjectives	crunk, amazeballs, post-truth, woke
	Adverbs	hella, wicked
Closed class	Determiner	
	Pronouns	
	Prepositions	English has a new preposition, because internet [Garber 2013; Pullum 2014]
	Conjunctions	

POS tagging

Labeling the tag that's correct
for the context.



Fruit flies like a banana



Time flies like an arrow

(Just tags in evidence within the Penn Treebank — more are possible!)

Why is part of speech tagging useful?

Sequence labeling

$$x = \{x_1, \dots, x_n\}$$

$$y = \{y_1, \dots, y_n\}$$

- For a set of inputs x with n sequential time steps, one corresponding label y_i for each x_i
- Model correlations in the labels y .

HMM

$$P(x_1,\dots,x_n,y_1,\dots,y_n) \approx \prod_{i=1}^{n+1} P(y_i \mid y_{i-1}) \prod_{i=1}^n P(x_i \mid y_i)$$

Hidden Markov Model

Prior probability of label sequence

$$P(y) = P(y_1, \dots, y_n)$$

$$P(y_1, \dots, y_n) \approx \prod_{i=1}^{n+1} P(y_i \mid y_{i-1})$$

- We'll make a first-order Markov assumption and calculate the joint probability as the product the individual factors conditioned **only on the previous tag**.

Hidden Markov Model

$$P(x \mid y) = P(x_1, \dots, x_n \mid y_1, \dots, y_n)$$

$$P(x_1, \dots, x_n \mid y_1, \dots, y_n) \approx \prod_{i=1}^N P(x_i \mid y_i)$$

- Here again we'll make a strong assumption: the probability of the word we see at a given time step is only dependent on its label

Parameter estimation

$$P(y_t \mid y_{t-1})$$

$$\frac{c(y_1, y_2)}{c(y_1)}$$

MLE for both is just counting
(as in Naive Bayes)

$$P(x_t \mid y_t)$$

$$\frac{c(x, y)}{c(y)}$$

Decoding

- Greedy: proceed left to right, committing to the best tag for each time step (given the sequence seen so far)

Fruit flies like a banana

NN VB IN DT NN

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

create a path probability matrix $viterbi[N+2,T]$

for each state s **from** 1 **to** N **do** ; initialization step

$viterbi[s,1] \leftarrow a_{0,s} * b_s(o_1)$

$backpointer[s,1] \leftarrow 0$

for each time step t **from** 2 **to** T **do** ; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s,t] \leftarrow \max_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)$

$backpointer[s,t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s',t-1] * a_{s',s}$

$viterbi[q_F,T] \leftarrow \max_{s=1}^N viterbi[s,T] * a_{s,q_F}$; termination step

$backpointer[q_F,T] \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s,T] * a_{s,q_F}$; termination step

return the backtrace path by following backpointers to states back in time from $backpointer[q_F,T]$

Figure 10.8 Viterbi algorithm for finding optimal sequence of tags. Given an observation sequence and an HMM $\lambda = (A, B)$, the algorithm returns the state path through the HMM that assigns maximum likelihood to the observation sequence. Note that states 0 and q_F are non-emitting.

MEMM

General maxent form

$$\arg \max_y P(y \mid x, \beta)$$

Maxent with first-order Markov
assumption: Maximum Entropy
Markov Model

$$\arg \max_y \prod_{i=1}^n P(y_i \mid y_{i-1}, x)$$

Features

$$f(t_i, t_{i-1}; x_1, \dots, x_n)$$

Features are scoped over
the previous predicted
tag and the entire
observed input

feature	example
$x_i = \text{man}$	1
$t_{i-1} = \text{JJ}$	1
$i=n$ (last word of sentence)	1
$x_i \text{ ends in } -ly$	0

Viterbi decoding

Viterbi for HMM: max joint probability

$$P(y)P(x \mid y) = P(x, y)$$

$$v_t(y) = \max_{u \in \mathcal{Y}} [v_{t-1}(u) \times P(y_t = y \mid y_{t-1} = u)P(x_t \mid y_t = y)]$$

Viterbi for MEMM: max conditional probability

$$P(y \mid x)$$

$$v_t(y) = \max_{u \in \mathcal{Y}} [v_{t-1}(u) \times P(y_t = y \mid y_{t-1} = u, x, \beta)]$$

MEMM Training

$$\prod_{i=1}^n P(y_i \mid y_{i-1}, x, \beta)$$

Locally normalized — at each time step,
each conditional distribution sums to 1

Label bias

NN TO VB

will to fight

$$\prod_{i=1}^n P(y_i \mid y_{i-1}, x, \beta)$$

Because of this local normalization, $P(\text{TO} \mid \text{context})$ will always be 1 if $x = \text{"to"}$

Label bias

NN TO VB

will to fight

That means our prediction for *to* can't help us disambiguate will. We lose the information that MB + TO sequences rarely happen.

Conditional random fields

- We can solve this problem using global normalization (over the entire sequences) rather than locally normalized factors.

MEMM

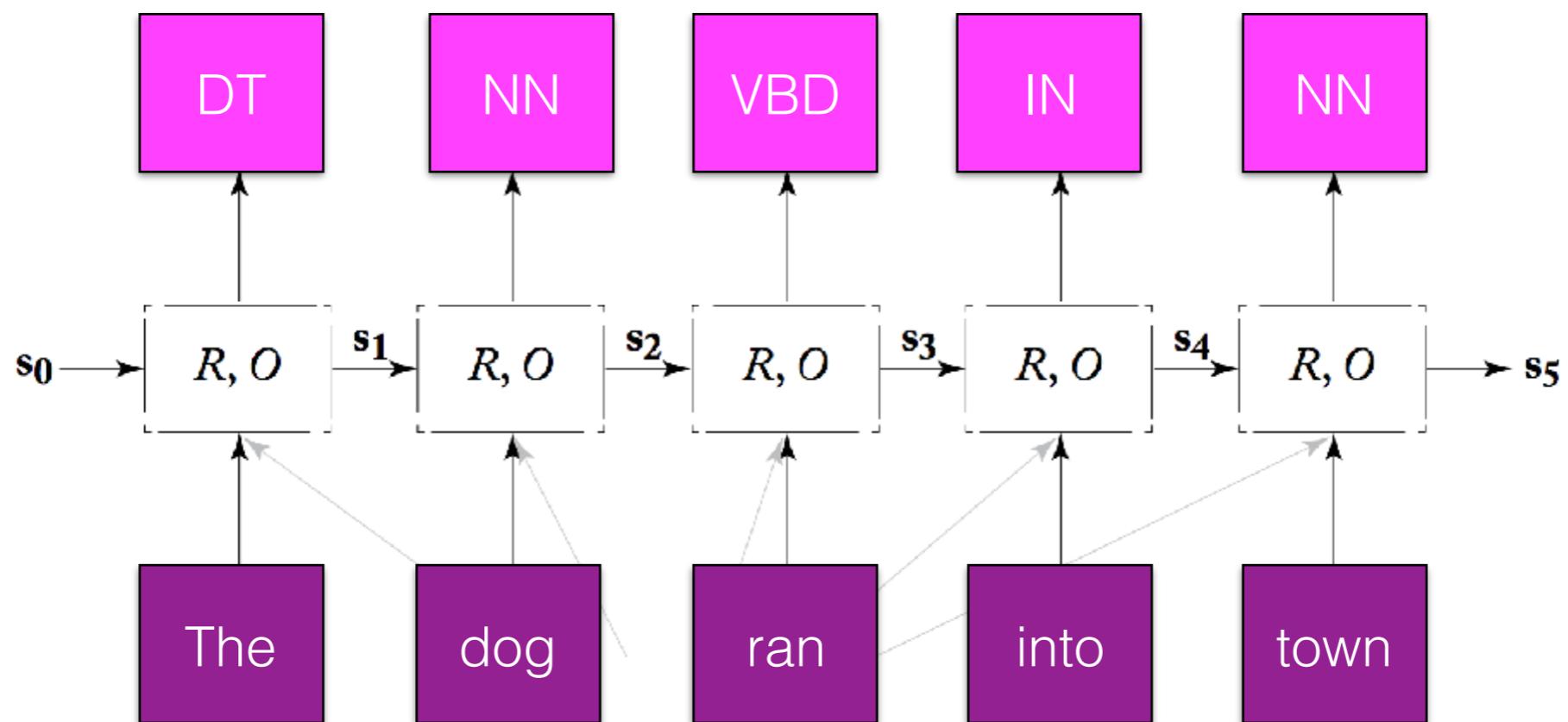
$$P(y \mid x, \beta) = \prod_{i=1}^n P(y_i \mid y_{i-1}, x, \beta)$$

CRF

$$P(y \mid x, \beta) = \frac{\exp(\Phi(x, y)^\top \beta)}{\sum_{y' \in \mathcal{Y}} \exp(\Phi(x, y')^\top \beta)}$$

Recurrent neural network

- For POS tagging, predict the **tag** from **y** conditioned on the context



Bidirectional RNN

- A powerful alternative is make predictions conditioning both on the **past** and the **future**.
- Two RNNs
 - One running left-to-right
 - One right-to-left
- Each produces an output vector at each time step, which we concatenate

Evaluation

- A critical part of development new algorithms and methods and demonstrating that they work

Experiment design

	training	development	testing
size	80%	10%	10%
purpose	training models	model selection	evaluation; never look at it until the very end

Accuracy

$$\frac{1}{N} \sum_{i=1}^N I[\hat{y}_i = y_i]$$

$I[x] \begin{cases} 1 & \text{if } x \text{ is true} \\ 0 & \text{otherwise} \end{cases}$

True (y)

Predicted (\hat{y})

	NN	VBZ	JJ
NN	100	2	15
VBZ	0	104	30
JJ	30	40	70

Precision

Precision(NN) =

$$\frac{\sum_{i=1}^N I(y_i = \hat{y}_i = \text{NN})}{\sum_{i=1}^N I(\hat{y}_i = \text{NN})}$$

Precision: proportion
of predicted class
that are actually that
class.

True (y)

Predicted (\hat{y})

	NN	VBZ	JJ
NN	100	2	15
VBZ	0	104	30
JJ	30	40	70

Recall

Recall(NN) =

$$\frac{\sum_{i=1}^N I(y_i = \hat{y}_i = \text{NN})}{\sum_{i=1}^N I(y_i = \text{NN})}$$

Recall: proportion of true class that are predicted to be that class.

True (y)

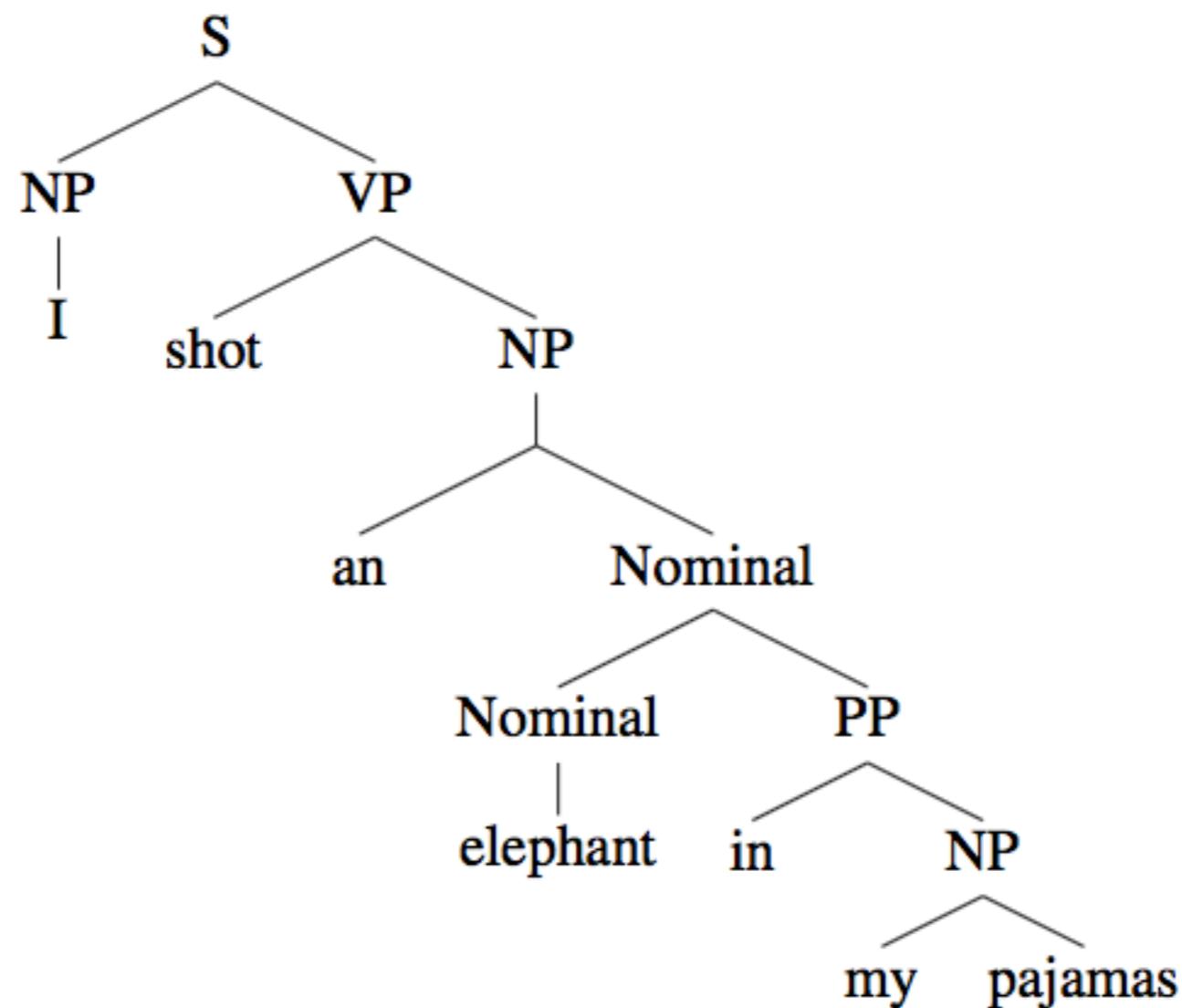
Predicted (\hat{y})

	NN	VBZ	JJ
NN	100	2	15
VBZ	0	104	30
JJ	30	40	70

F score

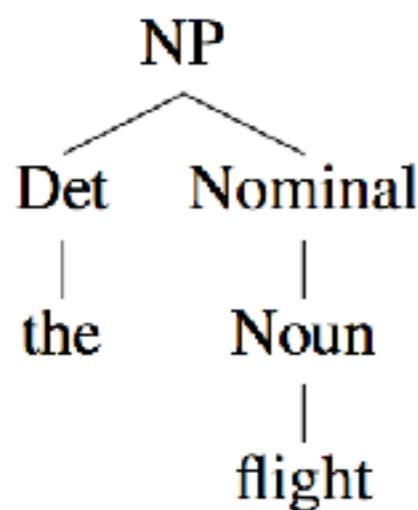
$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

Why is syntax important?

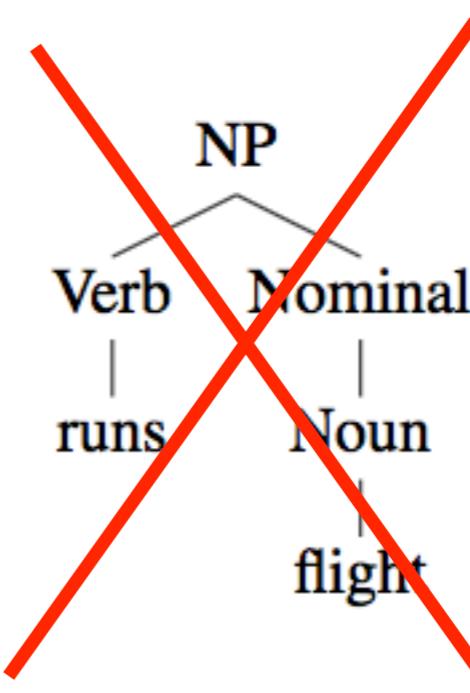


Context-free grammar

- A CFG gives a formal way to define what meaningful constituents are and exactly how a constituent is formed out of other constituents (or words). It defines **valid structure** in a language.

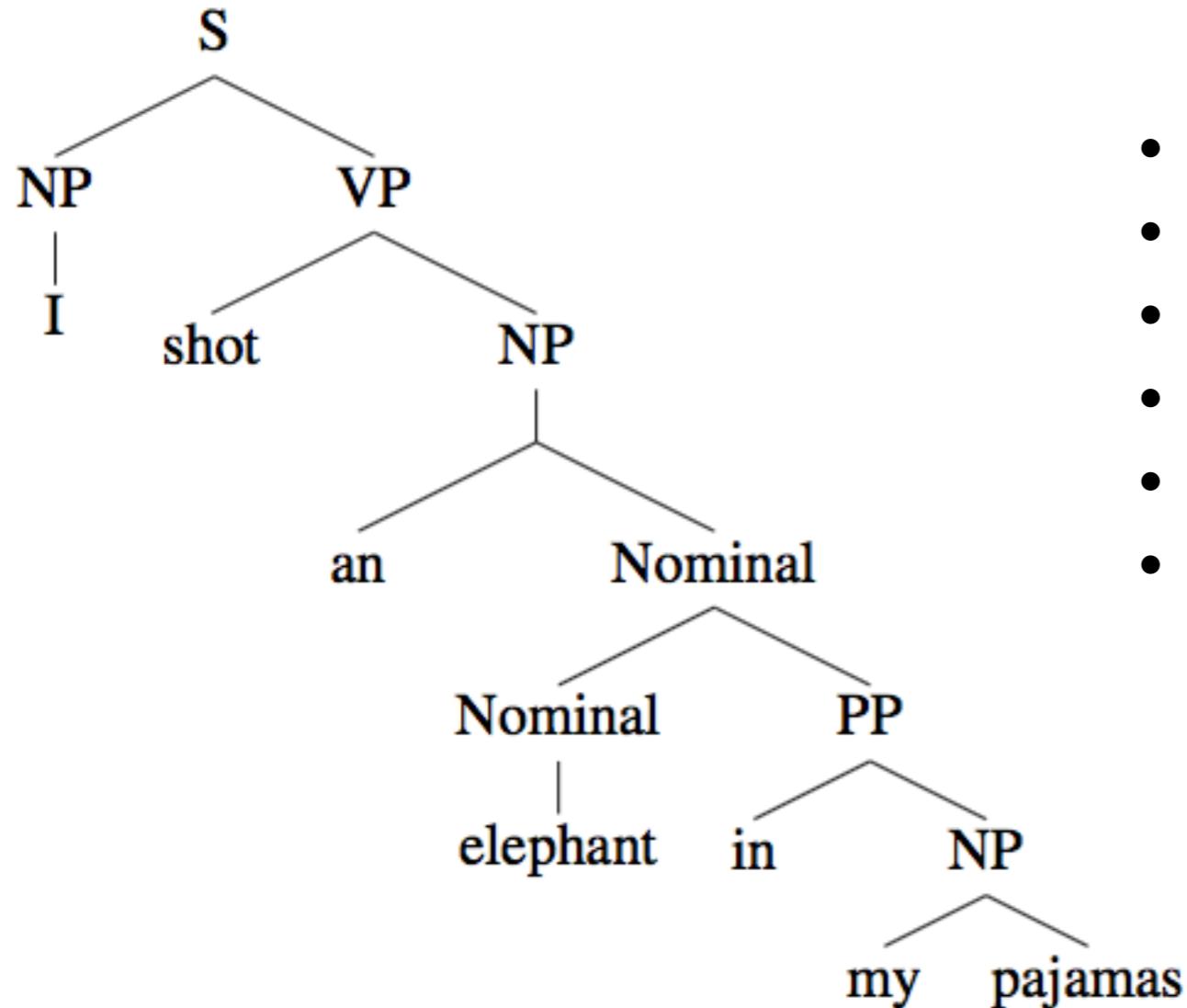


$NP \rightarrow Det\ Nominal$



$NP \rightarrow Verb\ Nominal$

Constituents



Every internal node is a phrase

- my pajamas
- in my pajamas
- elephant in my pajamas
- an elephant in my pajamas
- shot an elephant in my pajamas
- I shot an elephant in my pajamas

Each phrase could be replaced by another of the same type of constituent

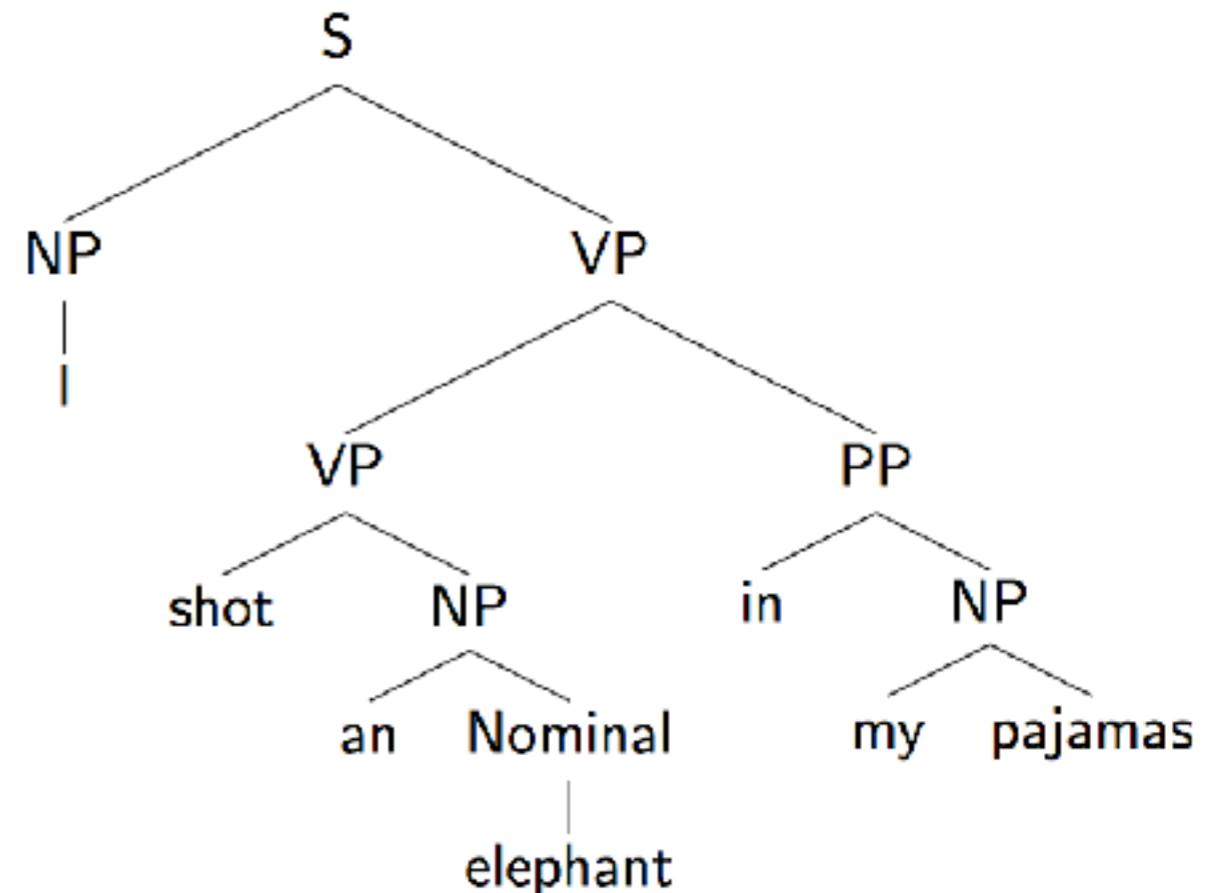
Evaluation

Parseval (1991):

Represent each tree as a collection of tuples:

$\langle l_1, i_1, j_1 \rangle, \dots, \langle l_n, i_n, j_n \rangle$

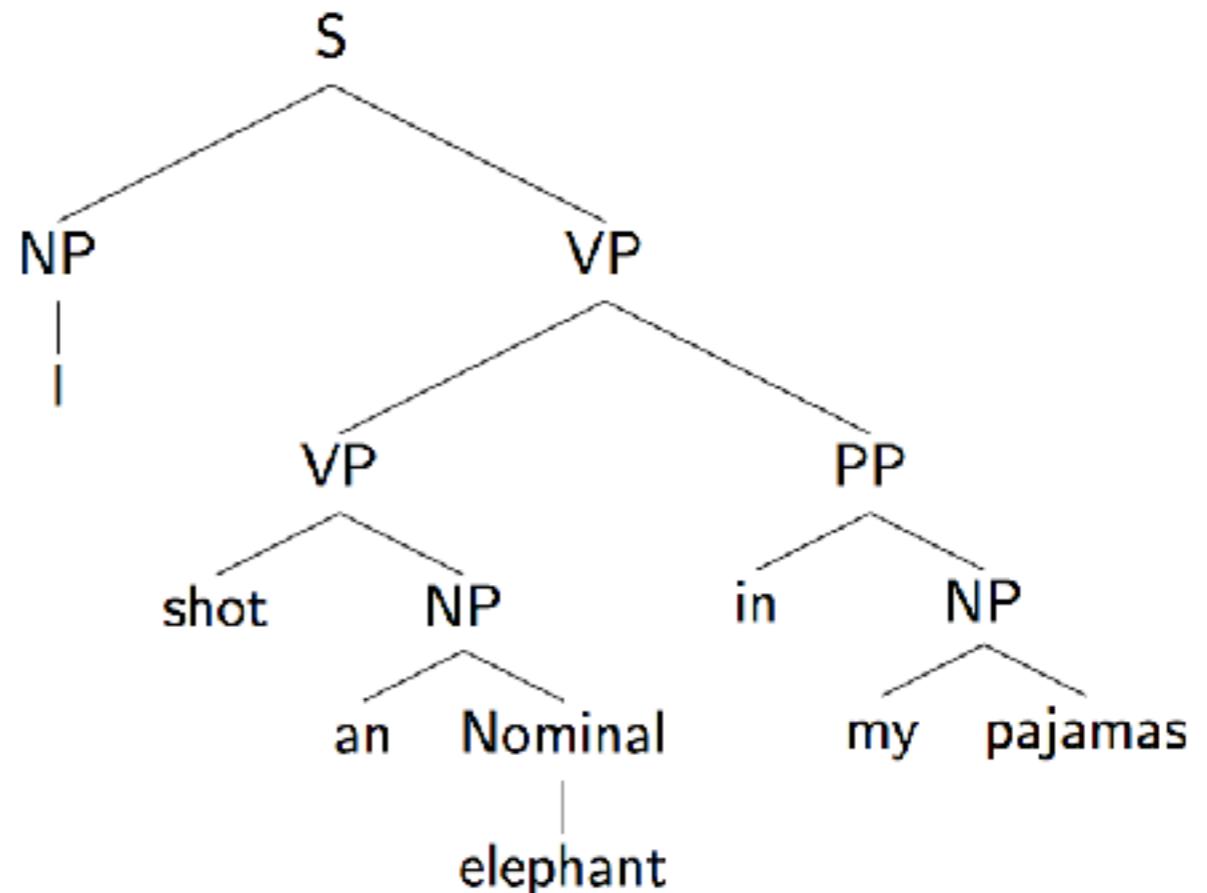
- l_k = label for kth phrase
- i_k = index for first word in zth phrase
- j_k = index for last word in kth phrase



Evaluation

I₁ shot₂ an₃ elephant₄ in₅ my₆ pajamas₇

- <S, 1, 7>
- <NP, 1, 1>
- <VP, 2, 7>
- <VP, 2, 4>
- <NP, 3, 4>
- <Nominal, 4, 4>
- <PP, 5, 7>
- <NP, 6, 7>



Evaluation

I₁ shot₂ an₃ elephant₄ in₅ my₆ pajamas₇

- <S, 1, 7>
- <NP, 1, 1>
- <VP, 2, 7>
- <VP, 2, 4>
- <NP, 3, 4>
- <Nominal, 4, 4>
- <PP, 5, 7>
- <NP, 6, 7>
- <S, 1, 7>
- <NP, 1, 1>
- <VP, 2, 7>
- <NP, 3, 7>
- <Nominal, 4, 7>
- <Nominal, 4, 4>
- <PP, 5, 7>
- <NP, 6, 7>

Evaluation

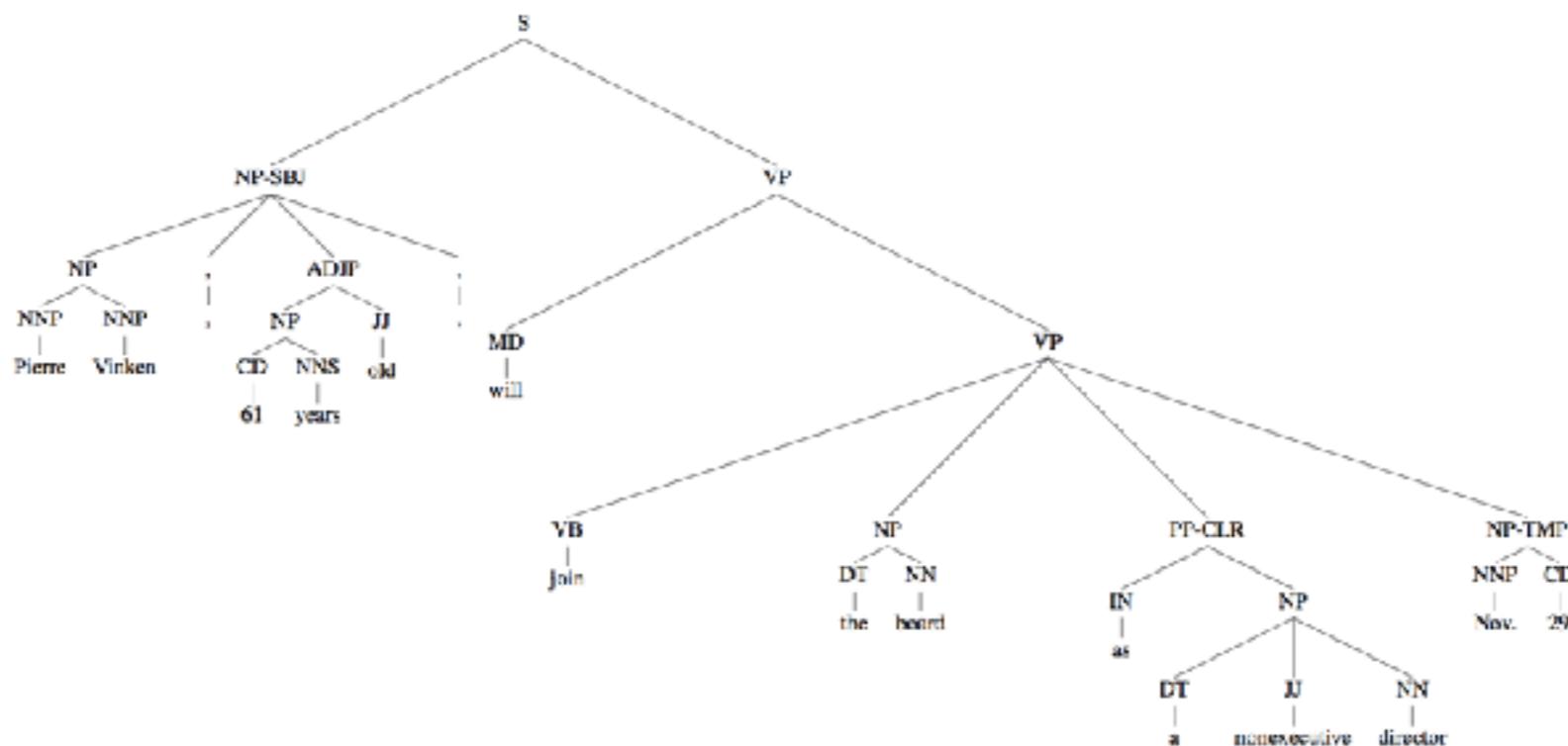
Calculate precision, recall, F1 from these collections of tuples

- Precision: number of tuples in predicted tree also in gold standard tree, divided by number of tuples in predicted tree
- Recall: number of tuples in predicted tree also in gold standard tree, divided by number of tuples in gold standard tree

Treebanks

- Rather than create the rules by hand, we can annotate sentences with their syntactic structure and then extract the rules from the annotations
- Treebanks: collections of sentences annotated with **syntactic structure**

Penn Treebank



NP	→	NNP NNP
NP-SBJ	→	NP , ADJP ,
S	→	NP-SBJ VP
VP	→	VB NP PP-CLR NP-TMP

Example rules extracted from this single annotation

PCFG

- Probabilistic context-free grammar: each production is also associated with a probability.
- This lets us calculate the probability of a parse for a given sentence; for a given parse tree T for sentence S comprised of n rules from R (each $A \rightarrow \beta$):

$$P(T, S) = \prod_i^n P(\beta | A)$$

Estimating PCFGs

$$\sum_{\beta} P(\beta \mid A) = \frac{C(A \rightarrow \beta)}{\sum_{\gamma} C(A \rightarrow \gamma)}$$

(equivalently)

$$\sum_{\beta} P(\beta \mid A) = \frac{C(A \rightarrow \beta)}{C(A)}$$

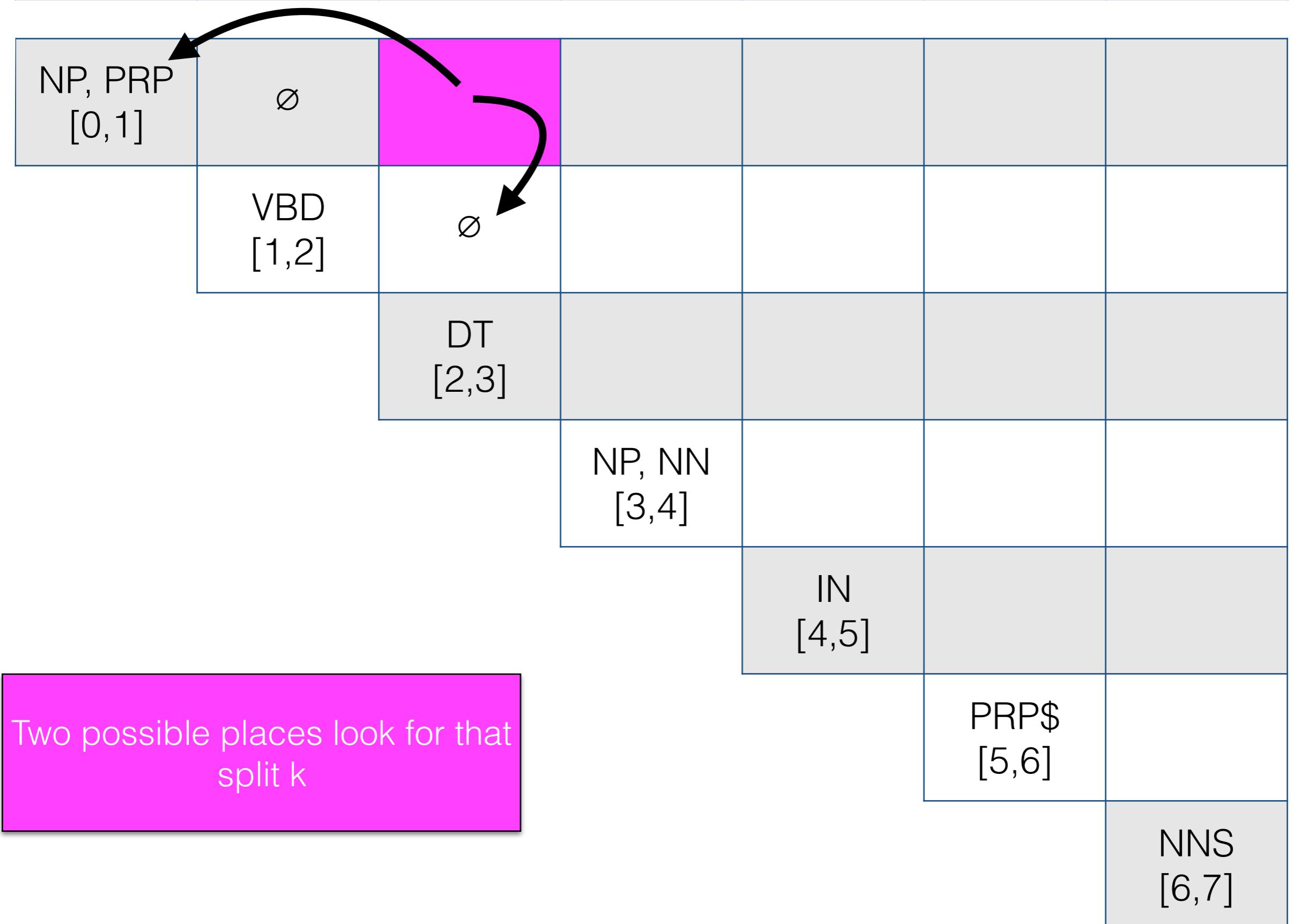
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
Does any rule generate PRP VBD?					PRP\$ [5,6]	
						NNS [6,7]

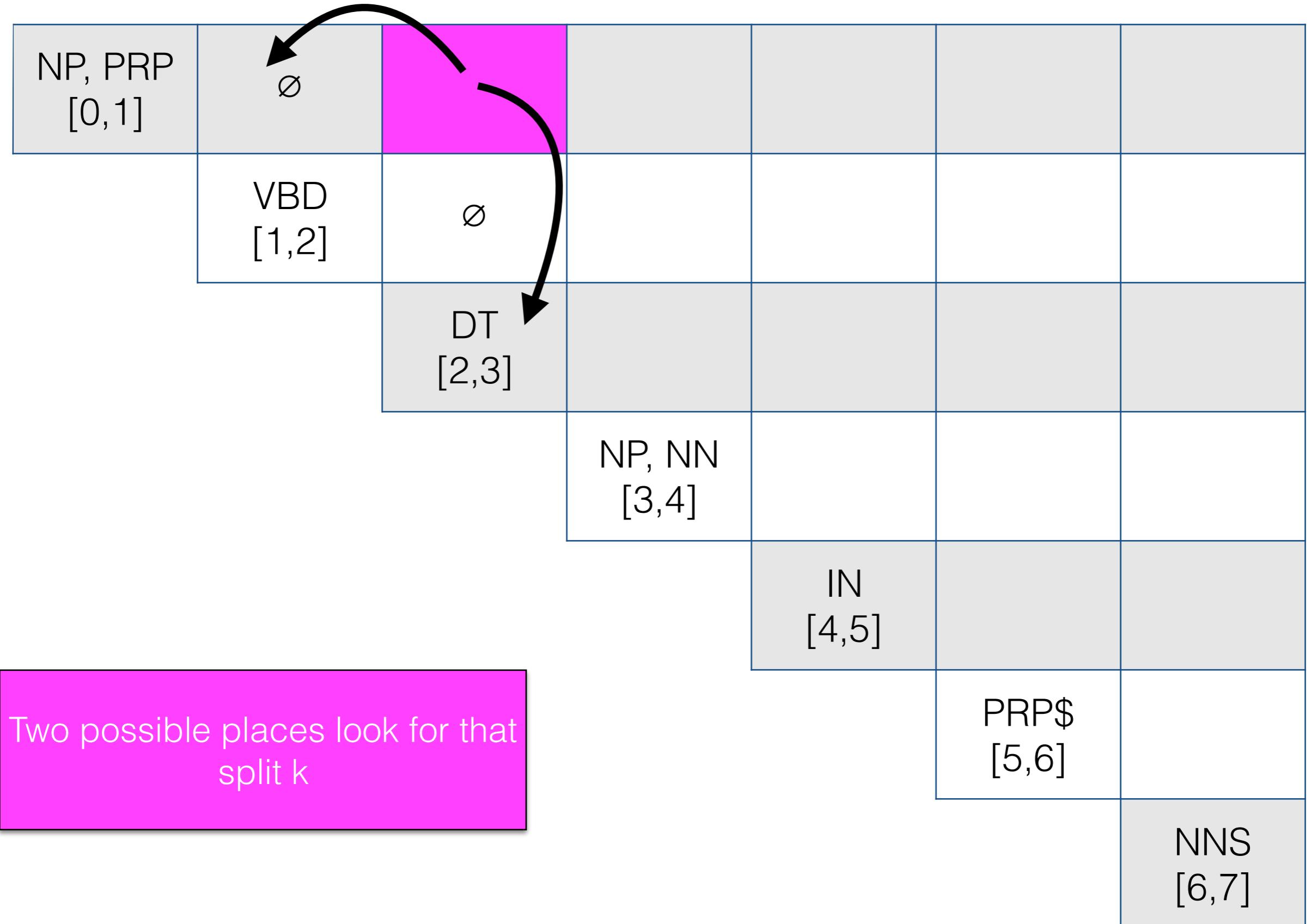
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset					
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
Does any rule generate VBD DT?					PRP\$ [5,6]	
						NNS [6,7]

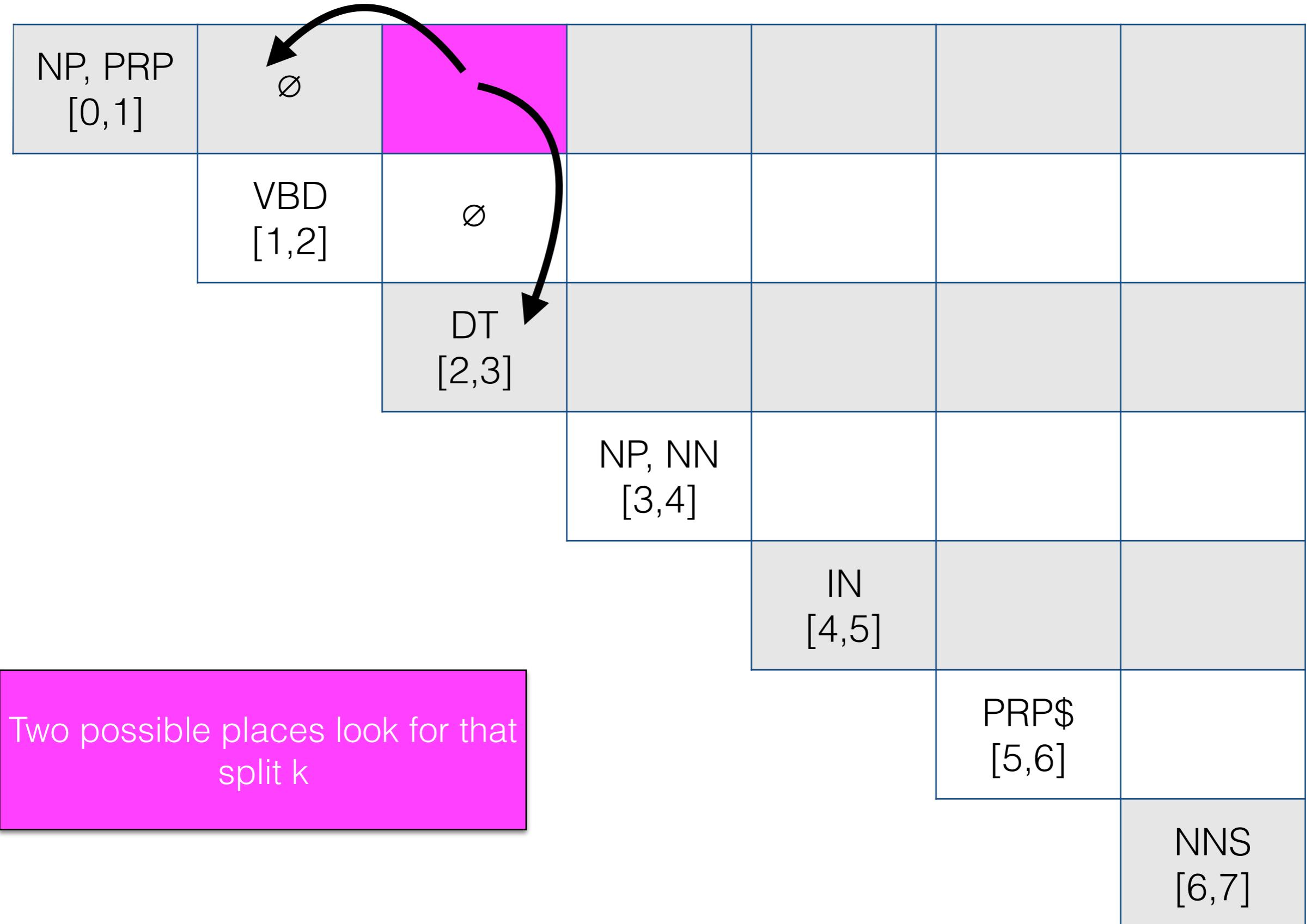
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



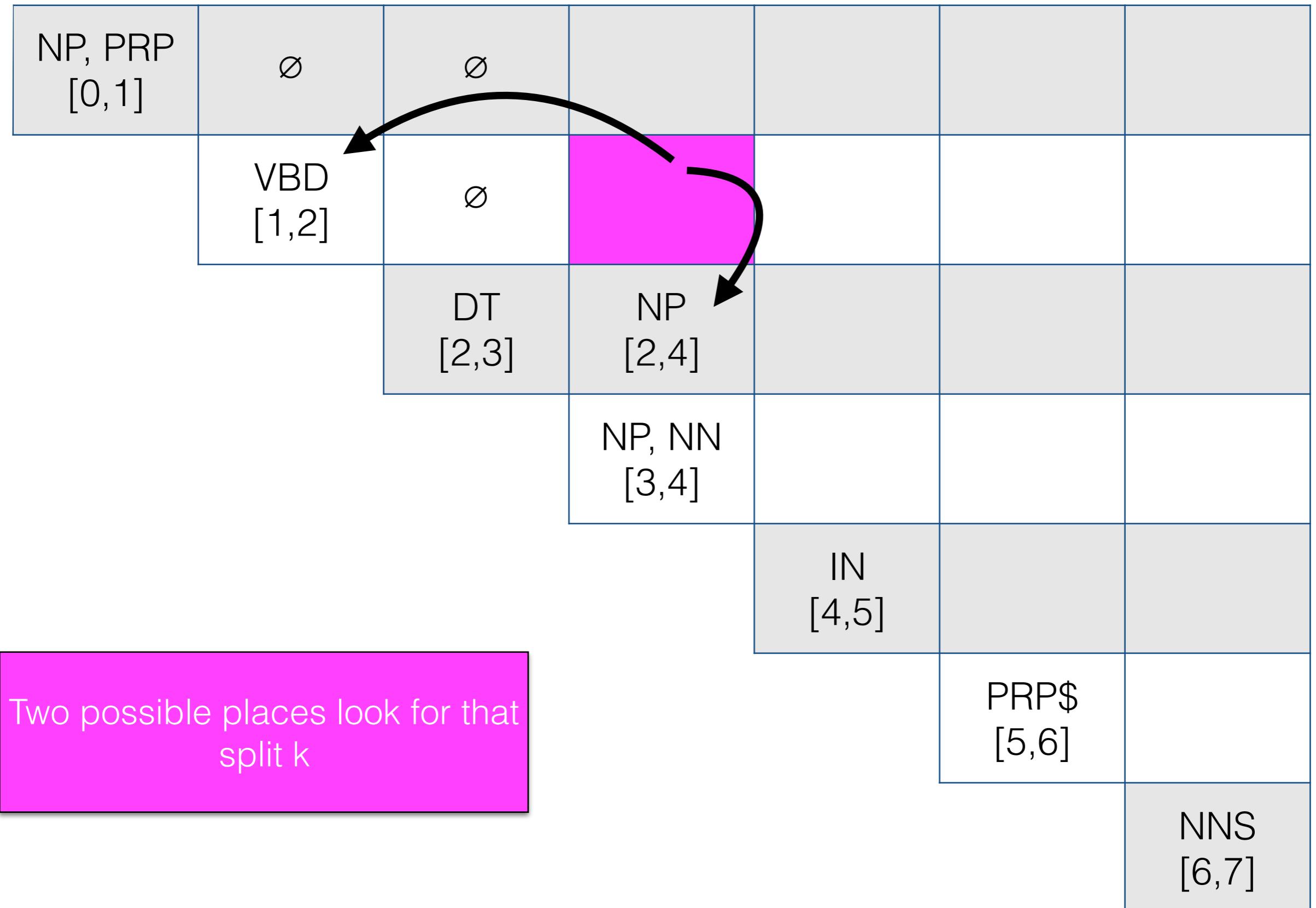
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



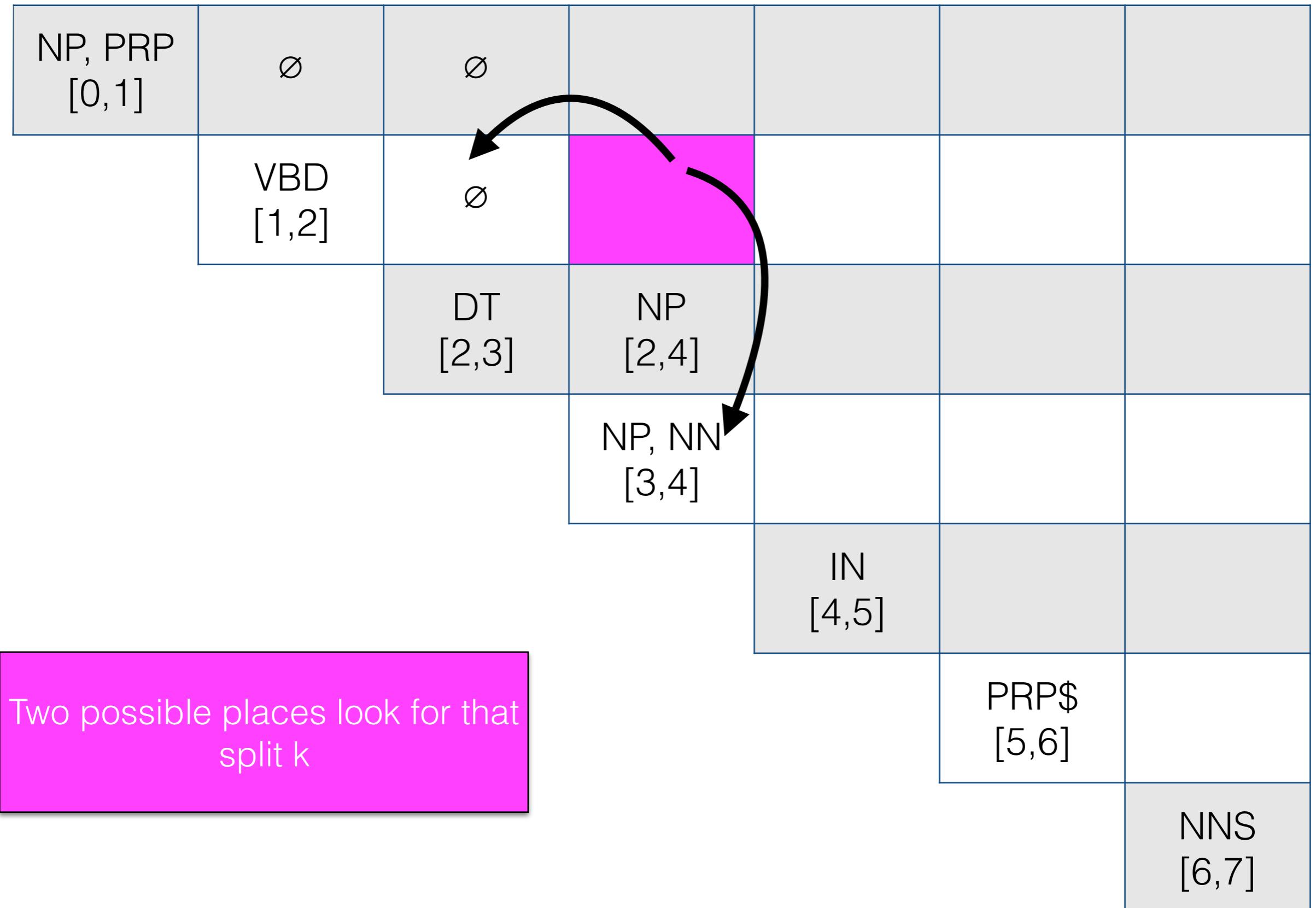
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	∅	∅				
	VBD [1,2]	∅				
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
Does any rule generate DT NN?					PRP\$ [5,6]	
						NNS [6,7]

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	∅	∅				
VBD [1,2]	∅		VP [1,4]			
	DT [2,3]		NP [2,4]			
		NP, NN [3,4]				
			IN [4,5]			
Three possible places look for that split k				PRP\$ [5,6]		
				NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
NP, PRP [0,1]	∅	∅	elephant			
VBD [1,2]	∅		VP [1,4]			
	DT [2,3]	NP [2,4]				
		NP, NN [3,4]				
			IN [4,5]			
Three possible places look for that split k			PRP\$ [5,6]			
			NNS [6,7]			

I	shot	an	elephant	in	my	pajamas
NP, PRP [0,1]	∅	∅	██████████			
VBD [1,2]	∅		VP [1,4]			
	DT [2,3]		NP [2,4]			
		NP, NN [3,4]				
Three possible places look for that split k			IN [4,5]		PRP\$ [5,6]	
					NNS [6,7]	

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset				
VBD [1,2]	\emptyset		VP [1,4]			
	DT [2,3]		NP [2,4]			
			NP, NN [3,4]			
				IN [4,5]		
Three possible places look for that split k				PRP\$ [5,6]		
				NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]			
	VBD [1,2]	Ø	VP [1,4]			
		DT [2,3]	NP [2,4]			
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	
	DT [2,3]	NP [2,4]		Ø	Ø	
		NP, NN [3,4]		Ø	Ø	
			IN [4,5]		Ø	
*elephant in	*in my				PRP\$ [5,6]	
*an elephant in	*elephant in my					
*shot an elephant in	*an elephant in my					
*I shot an elephant in	*shot an elephant in my					
	*I shot an elephant in my				NNS [6,7]	

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	
	DT [2,3]	NP [2,4]		Ø	Ø	
		NP, NN [3,4]		Ø	Ø	
			IN [4,5]		Ø	
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	
	DT [2,3]	NP [2,4]		Ø	Ø	
		NP, NN [3,4]		Ø	Ø	
			IN [4,5]	Ø	Ø	PP [4,7]
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

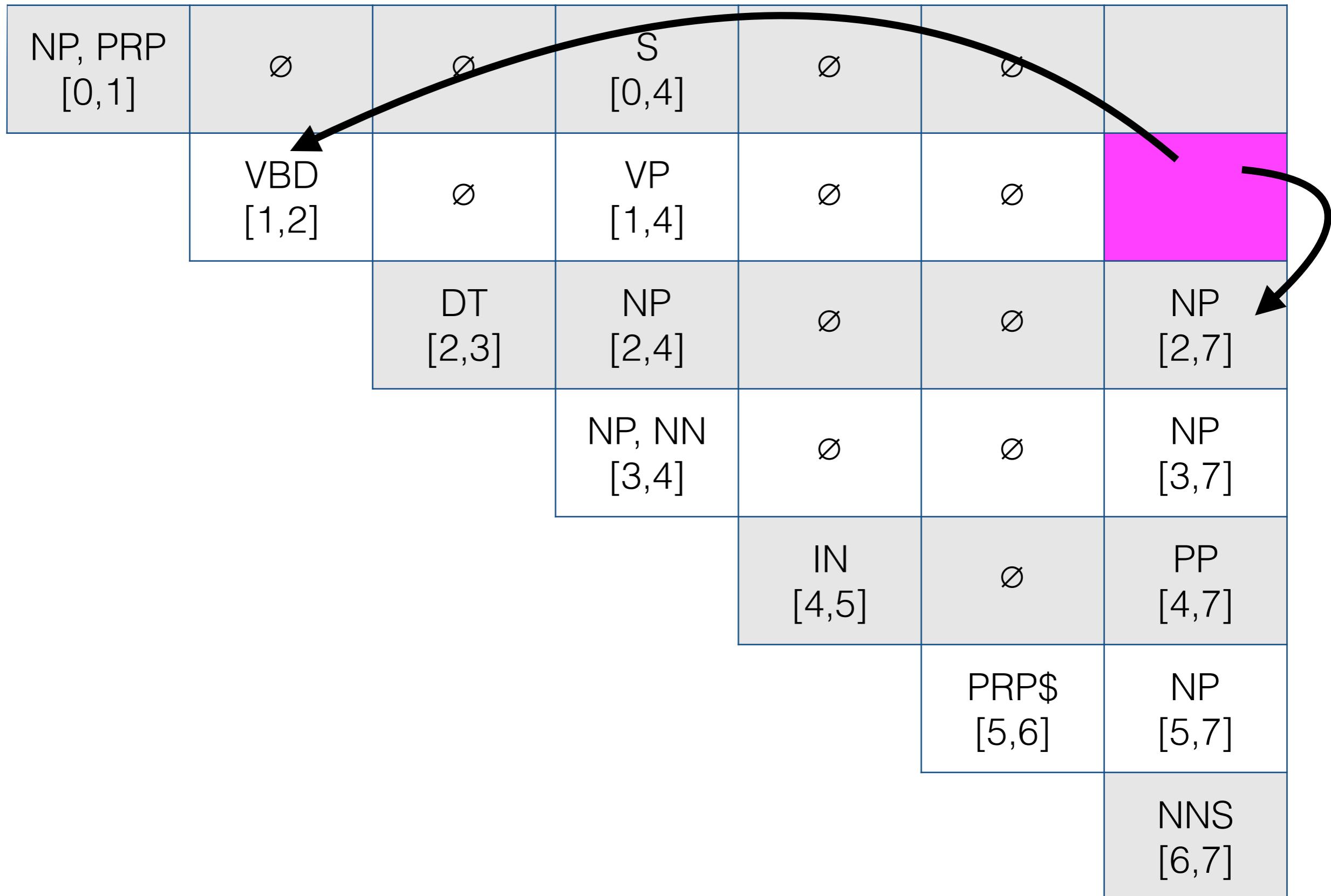
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	
	DT [2,3]	NP [2,4]		Ø	Ø	
		NP, NN [3,4]		Ø	Ø	NP [3,7]
			IN [4,5]		Ø	PP [4,7]
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

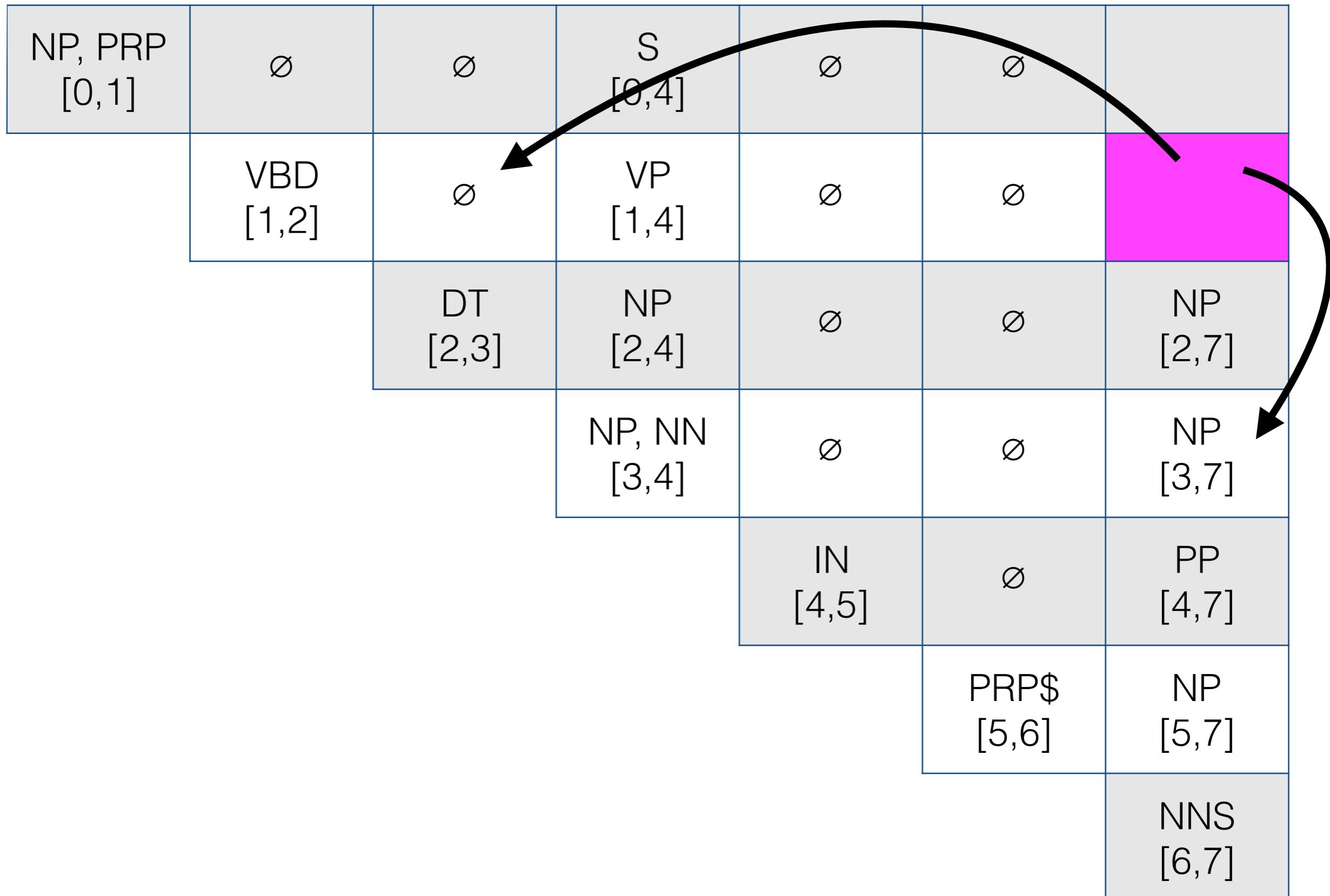
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	
	DT [2,3]	NP [2,4]		Ø	Ø	NP [2,7]
		NP, NN [3,4]		Ø	Ø	NP [3,7]
			IN [4,5]		Ø	PP [4,7]
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

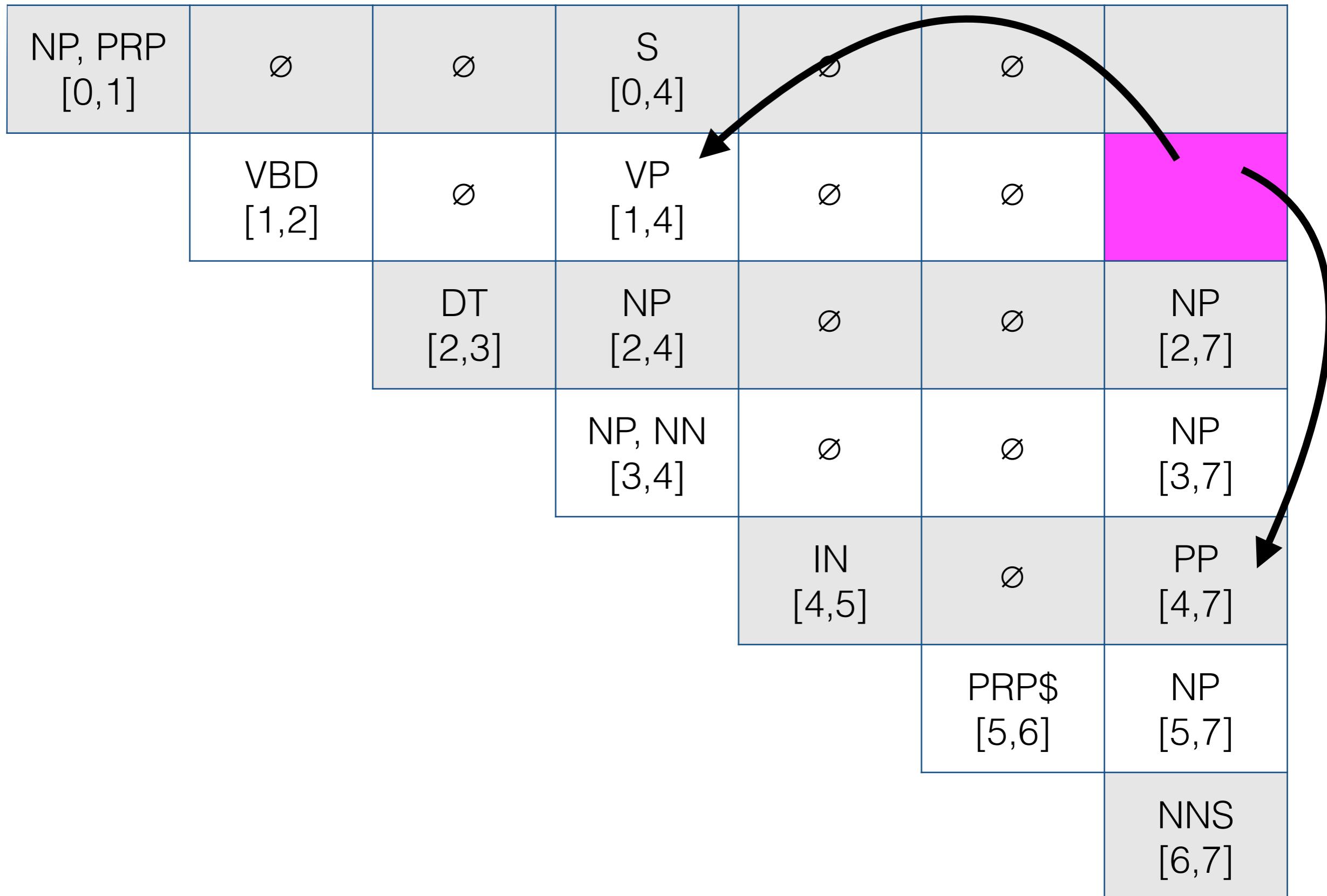
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



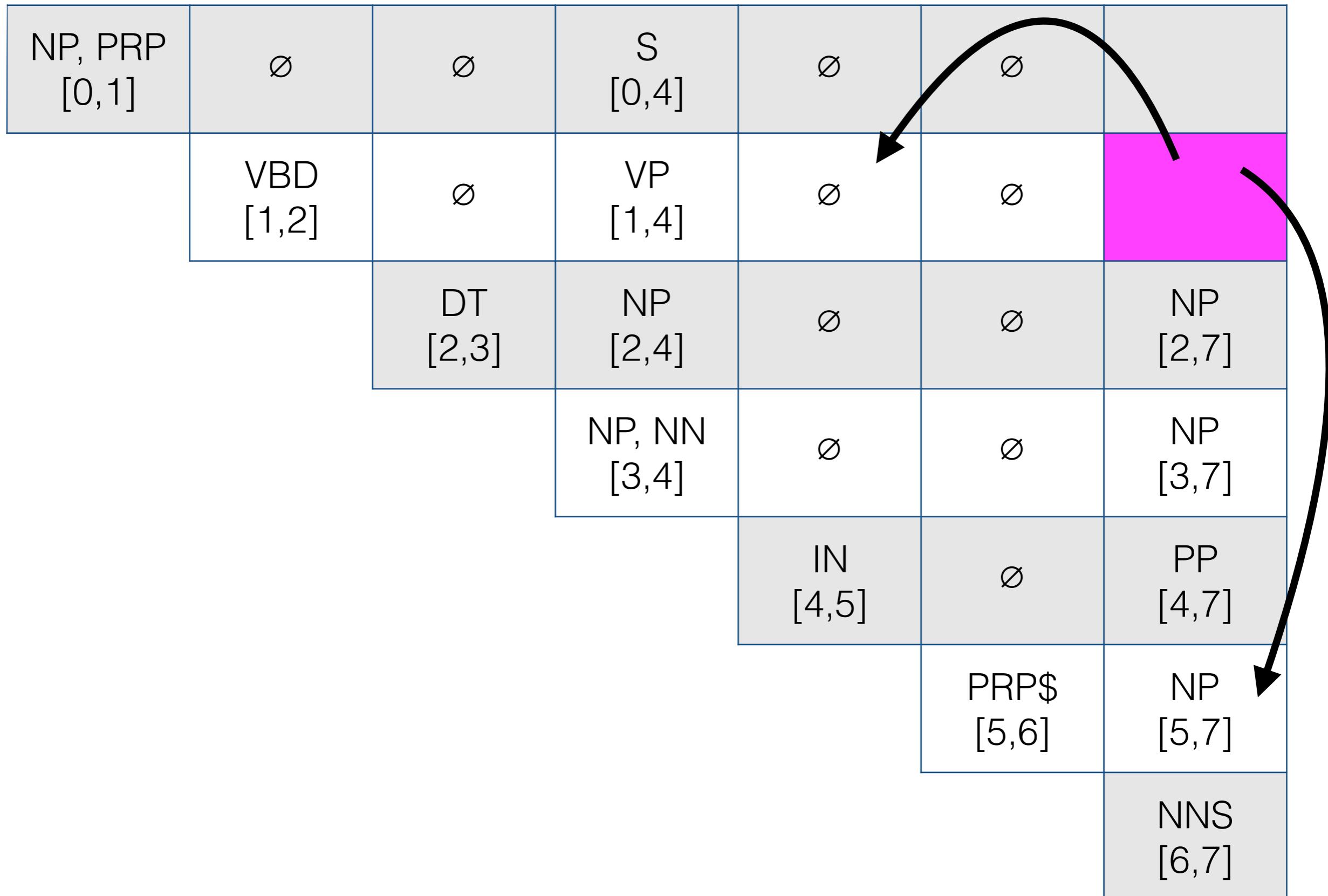
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



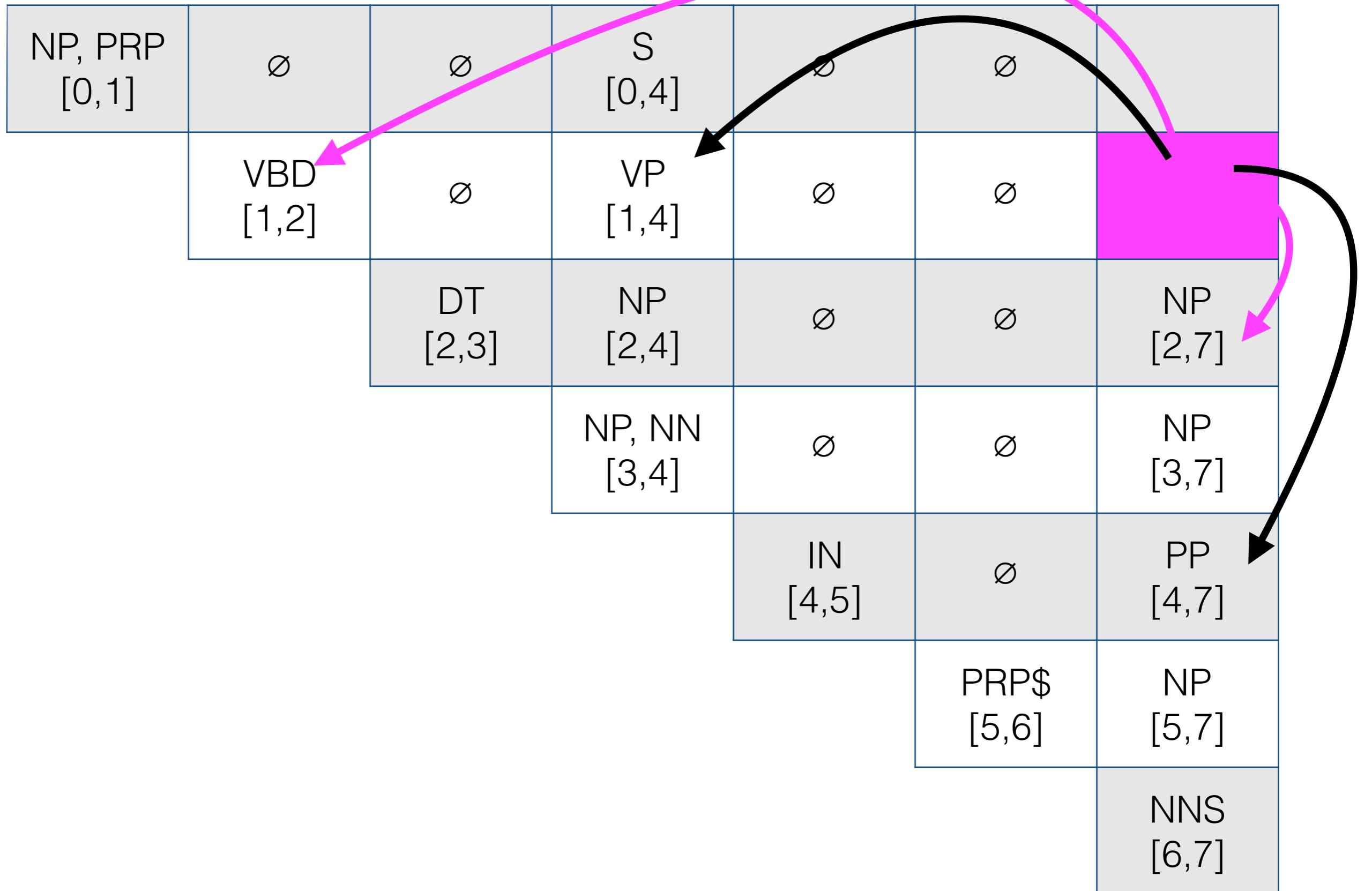
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	∅	∅	S [0,4]	∅	∅	
VBD [1,2]	∅	∅	VP [1,4]	∅	∅	
	DT [2,3]	NP [2,4]	∅	∅	∅	NP [2,7]
		NP, NN [3,4]	∅	∅	∅	NP [3,7]
			IN [4,5]	∅	∅	PP [4,7]
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
VBD [1,2]	Ø	Ø	VP [1,4]	Ø	Ø	VP ₁ , VP ₂ [1,7]
	DT [2,3]	NP [2,4]		Ø	Ø	NP [2,7]
		NP, NN [3,4]		Ø	Ø	NP [3,7]
			IN [4,5]		Ø	PP [4,7]
				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

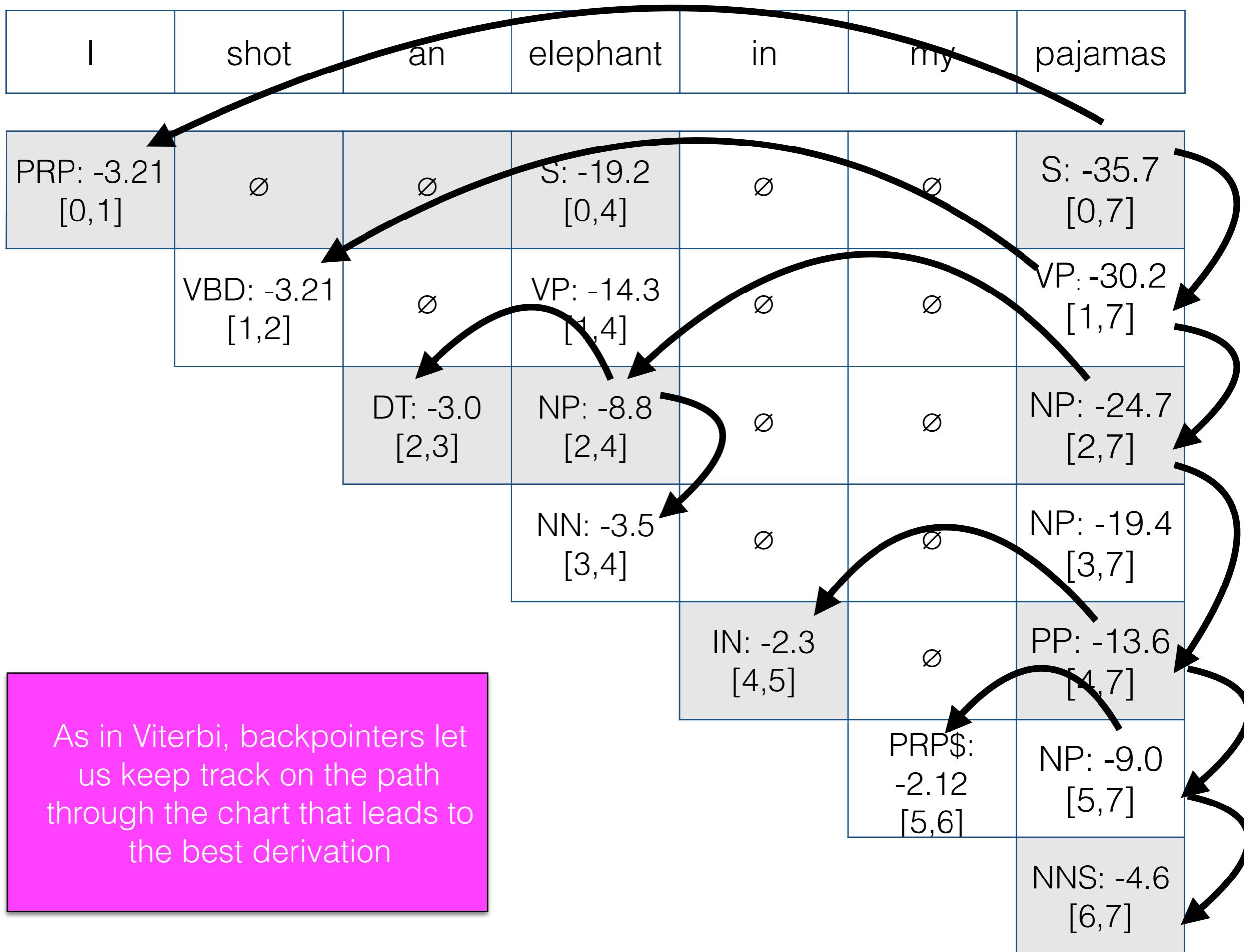
NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	\emptyset
VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	\emptyset	VP ₁ , VP ₂ [1,7]
	DT [2,3]	NP [2,4]	\emptyset	\emptyset	\emptyset	NP [2,7]
		NP, NN [3,4]	\emptyset	\emptyset	\emptyset	NP [3,7]
Possibilities:		IN [4,5]	\emptyset	\emptyset	\emptyset	PP [4,7]
			PRP\$ [5,6]	\emptyset	\emptyset	NP [5,7]
				NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	S ₁ , S ₂ [0,7]
VBD [1,2]	\emptyset	\emptyset	VP [1,4]	\emptyset	\emptyset	VP ₁ , VP ₂ [1,7]
	DT [2,3]	NP [2,4]		\emptyset	\emptyset	NP [2,7]
		NP, NN [3,4]		\emptyset	\emptyset	NP [3,7]
			IN [4,5]		\emptyset	PP [4,7]
Success! We've recognized a total of two valid parses				PRP\$ [5,6]	NP [5,7]	
					NNS [6,7]	

PCFGs

- A PCFG gives us a mechanism for assigning scores (here, probabilities) to different parses for the same sentence.
- But we often care about is finding **the single best parse** with the highest probability.
- We calculate the max probability parse using CKY by storing the probability of each phrase within each cell as we build it up.



Midterm

- In class next Tuesday
- Mix of multiple choice, short answer, long answer
- Bring 1 cheat sheet (1 page, both sides)
- Covers all material from lectures and readings