



تمرین سوم درس پردازش زبان طبیعی

نام تهیه کننده: ملیکا نوبختیان

نسخه: ۱

۱- بخش تئوری

۱-۱- سوال اول

تعداد Verb = ۴ تعداد Modal = ۴ تعداد Noun = ۶

Emission probability

کلمات	Noun	Modal	Verb
Mark	1/6	0	2/4
Can	0	3/4	0
Watch	2/6	0	2/4
Will	1/6	1/4	0
Tom	2/6	0	0

Sentences:

<S> Mark can watch. <E>

N M V

<S> Will can mark watch. <E>

N M V N

<S> Can Tom watch? <E>

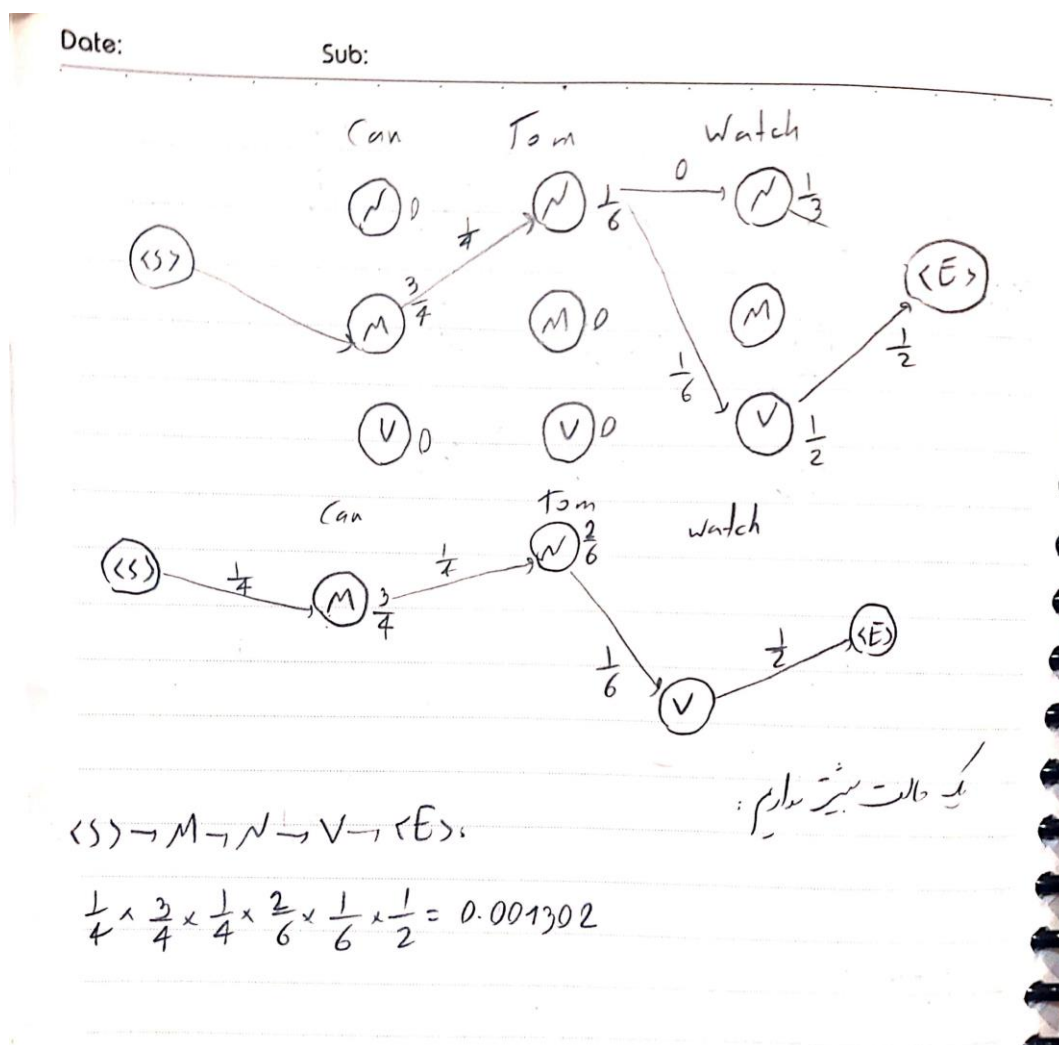
M N V

<S> Tom will mark watch. <E>

N M V N

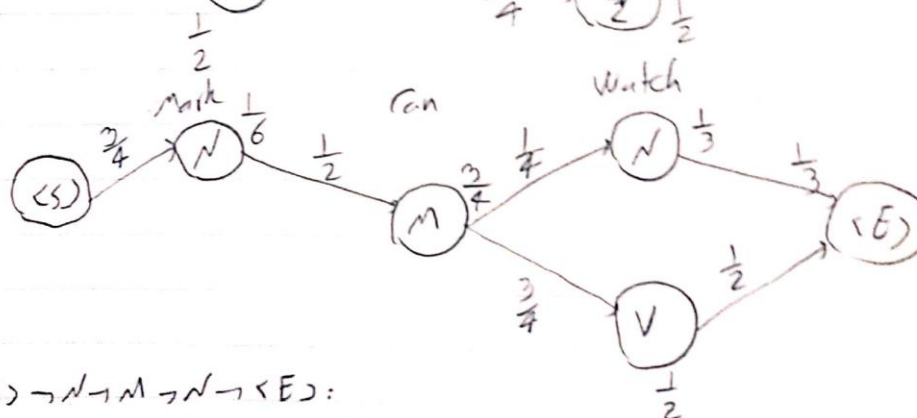
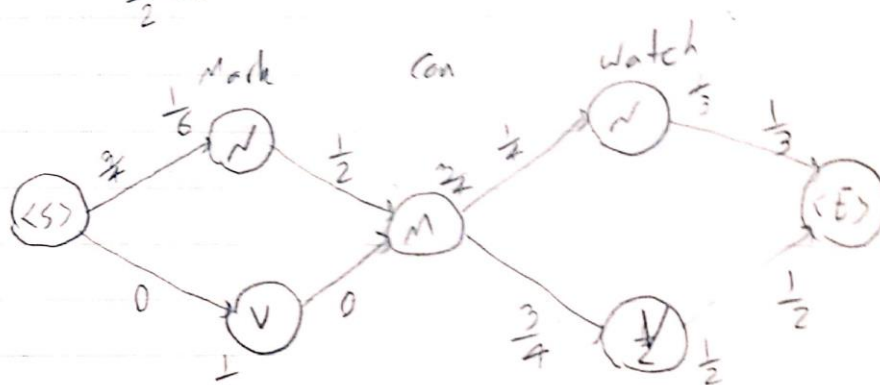
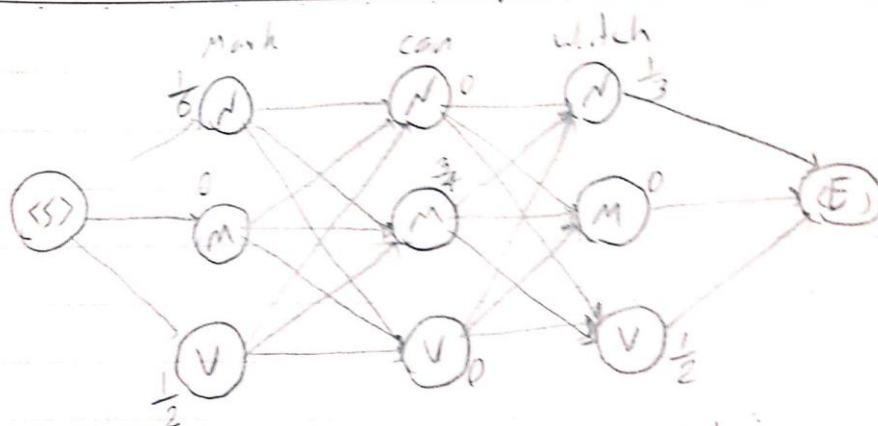
Transition Probability:

	Noun	Modal	Verb	<E>
<S>	3/4	1/4	0	0
Noun	0	3/6	1/6	2/6
Modal	1/4	0	3/4	0
Verb	2/4	0	0	2/4



Date:

Sub:



$\langle S \rangle \rightarrow N \rightarrow M \rightarrow N \rightarrow \langle E \rangle$:

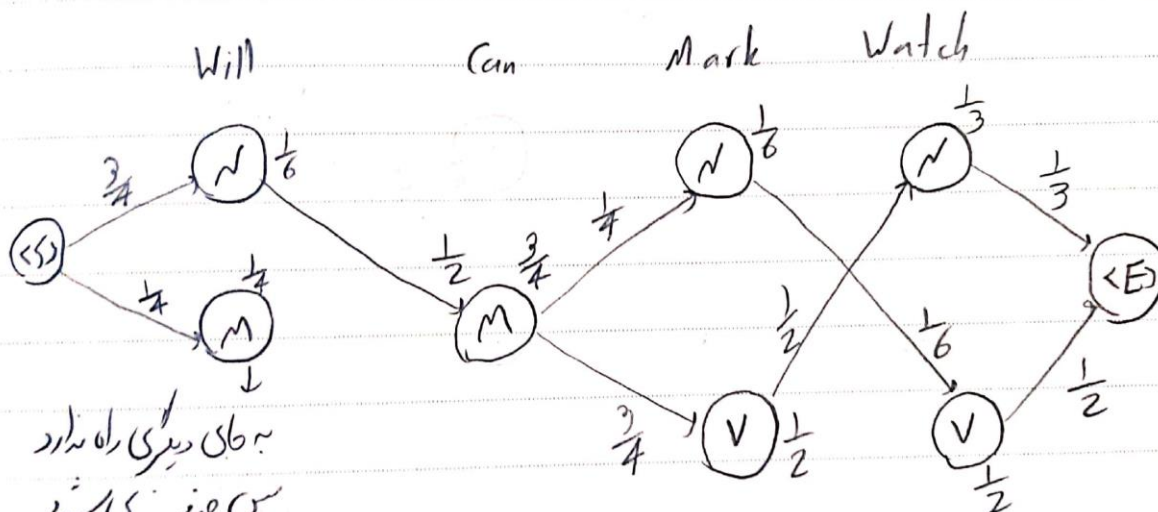
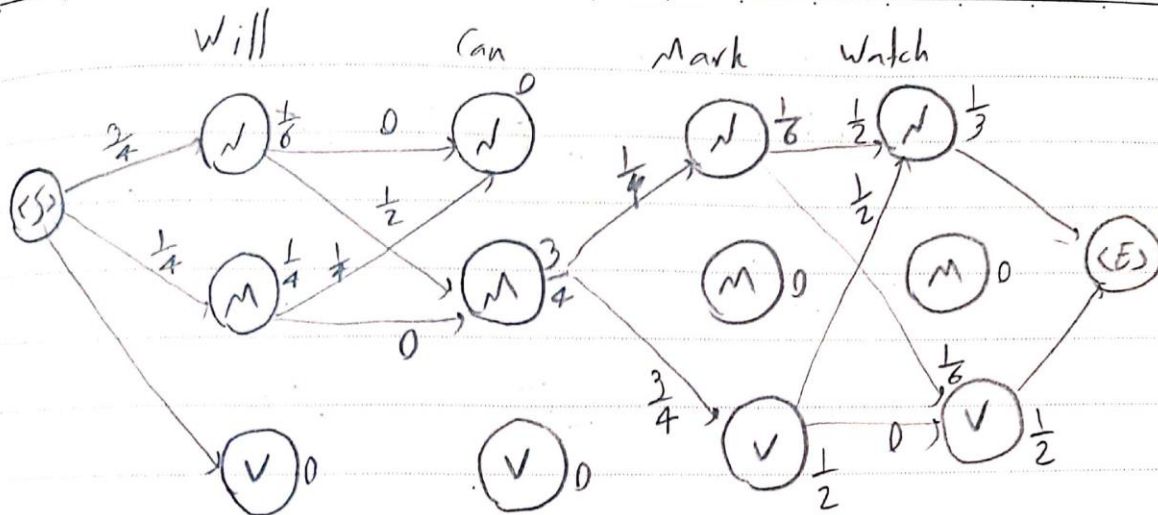
$$\frac{3}{4} \times \frac{1}{6} \times \frac{1}{2} \times \frac{3}{4} \times \frac{1}{4} \times \frac{1}{3} \times \frac{1}{3} = 0.001302$$

$\langle S \rangle \rightarrow N \rightarrow M \rightarrow V \rightarrow \langle E \rangle$:

$$\frac{3}{4} \times \frac{1}{6} \times \frac{1}{2} \times \frac{3}{4} \times \frac{3}{4} \times \frac{1}{2} \times \frac{1}{2} = 0.00878$$

Date:

Sub:



<S> → N → M → N → V → <E>:

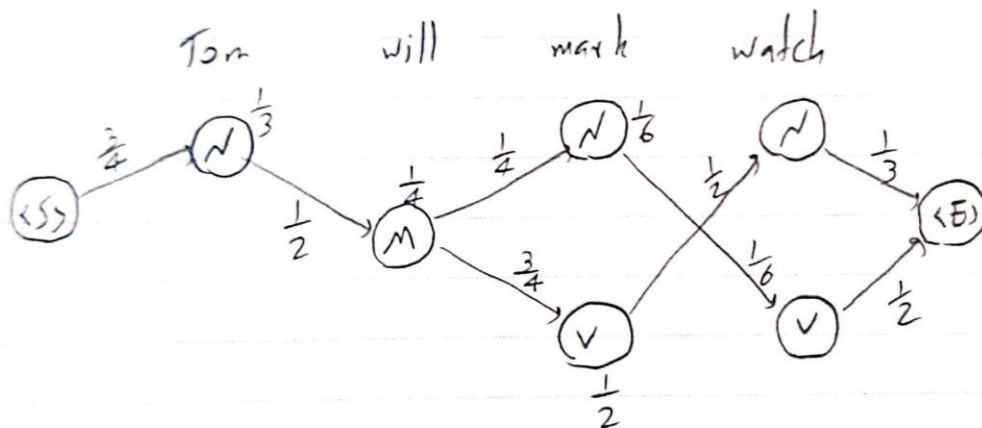
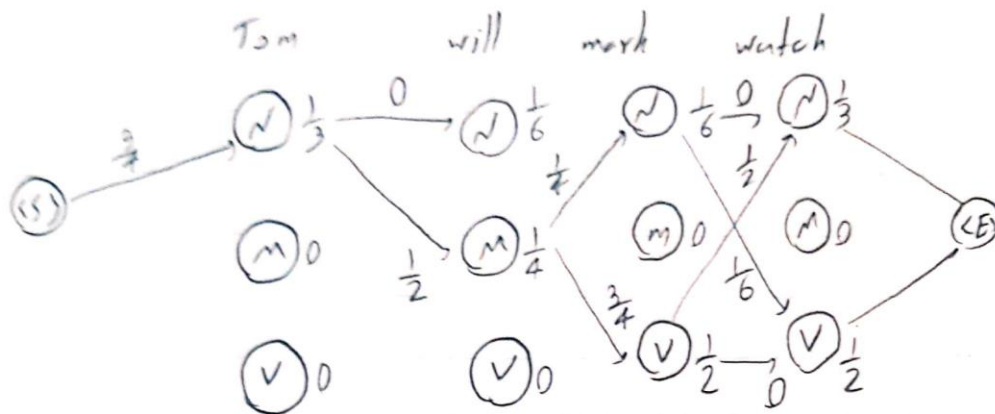
$$\frac{3}{4} \times \frac{1}{6} \times \frac{1}{2} \times \frac{3}{4} \times \frac{1}{4} \times \frac{1}{6} \times \frac{1}{6} \times \frac{1}{2} \times \frac{1}{2} = 0.0000813$$

<S> → N → M → V → N → <E>: ✓

$$\frac{3}{4} \times \frac{1}{6} \times \frac{1}{2} \times \frac{3}{4} \times \frac{3}{4} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{3} \times \frac{1}{3} = 0.000976$$

Date:

Sub:



$$\langle S \rangle \rightarrow N \rightarrow M \rightarrow N \rightarrow V \rightarrow \langle E \rangle:$$

$$\frac{3}{4} \times \frac{1}{3} \times \frac{1}{2} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{6} \times \frac{1}{6} \times \frac{1}{2} = 0.000108$$

$$\langle S \rangle \rightarrow N \rightarrow M \rightarrow V \rightarrow N \rightarrow \langle E \rangle:$$

$$\frac{3}{4} \times \frac{1}{3} \times \frac{1}{2} \times \frac{1}{4} \times \frac{3}{4} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{3} = 0.00097 \checkmark$$

۲-۱- سوال دوم

Date:

Sub:

(۲)

ابتدا باید گزاره‌ها را برای ساختن جدول به فرمائی قابل تبدیل کنیم:

$$S' \rightarrow S \quad 1$$

$$S \rightarrow NP \quad VP$$

$$S \rightarrow SC \quad S \quad 0.2 \quad SC \rightarrow S \quad conj \quad 1$$

$$NP \rightarrow Noun \mid Det \quad Noun \mid NP \quad PP \mid \underbrace{NPC \quad NP}_{0.2}$$

$$NPC \rightarrow NP \quad conj \quad 1$$

$$VP \rightarrow Verb \mid Verb \quad NP \mid \underbrace{Verb \quad 2NP}_{0.1} \mid VP \quad PP$$

$$2NP \rightarrow NP \quad NP \quad 1$$

$$PP \rightarrow P \quad NP$$

$$Noun \rightarrow John \mid Jack \mid pie \mid cream \mid cake$$

$$Verb \rightarrow eat \mid eats \mid drinks \quad P \rightarrow with \mid by$$

$$Det \rightarrow a \mid the \quad conj \rightarrow and \mid or$$

$$S' \rightarrow NP \quad VP \mid SC \quad S \quad \text{حذف} \quad S' \rightarrow S$$

1x0.8 1x0.2

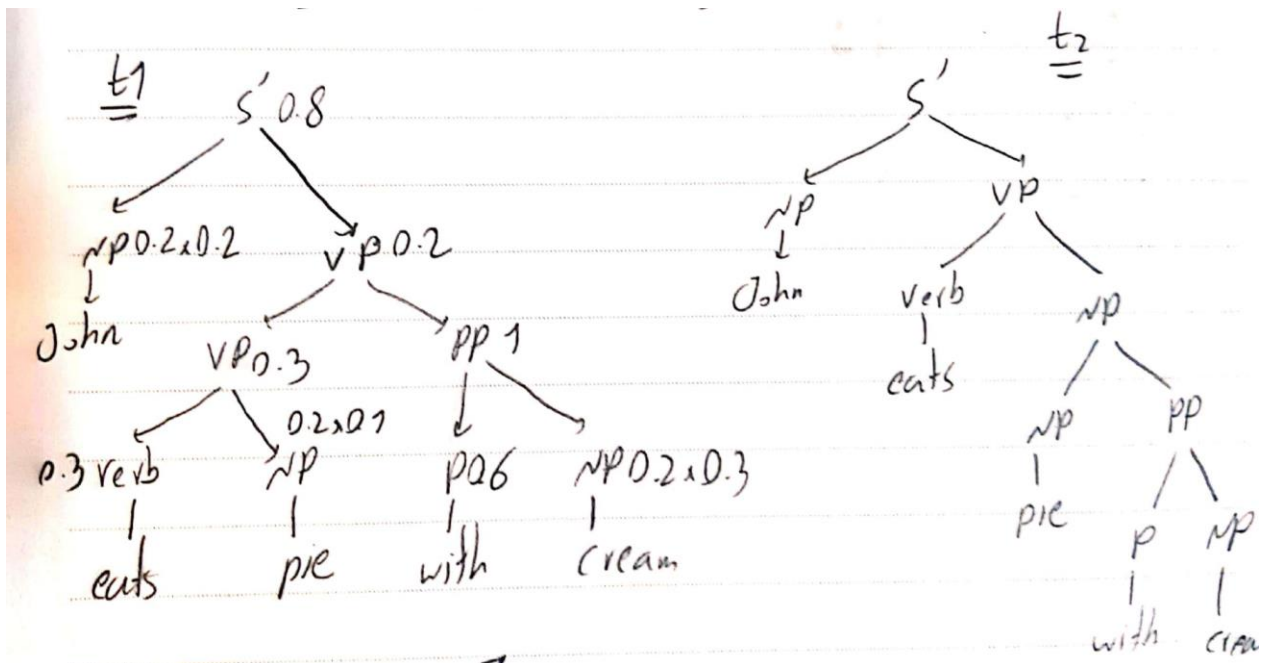
Elipon

Date:

Sub:

$NP \rightarrow Noun$ عزت
 $NP \rightarrow Det Noun \mid NP PP \mid NP C NP$
 John | Jack | pie | cream | cake
 $P(NP \rightarrow Noun) = P(Noun \rightarrow John)$
 $0.2 \times 0.2 \quad 0.2 \times 0.3 \quad 0.2 \times 0.1 \quad 0.2 \times 0.3 \quad 0.2 \times 0.1$
 $VP \rightarrow Verb$ ص
 $VP \rightarrow Verb NP \mid Verb 2NP \mid VP PP$
 eat | eats | drinks
 $0.4 \times 0.2 \quad 0.4 \times 0.3 \quad 0.4 \times 0.5$

John	eats	pie	with	cream	
NP, Noun	S', S	S', S	-	S', S	john
	VP, Verb	VP	-	VP	eats
		NP, Noun	-	NP	pie
			P	PP	with
				NP, Noun	cream



$$P(t_1) = 4,1472 \times 10^{-7}$$

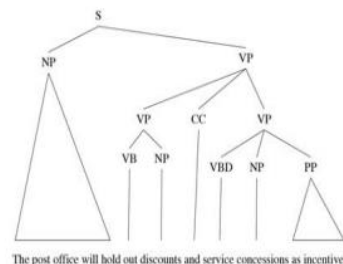
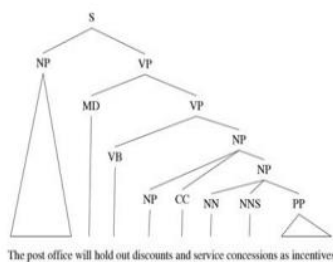
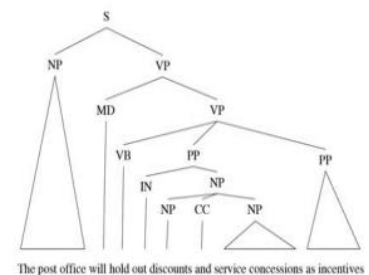
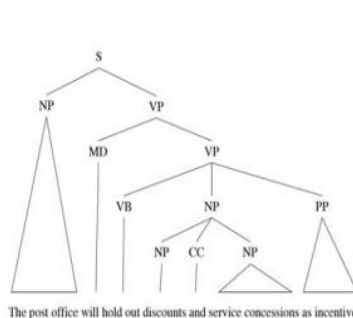
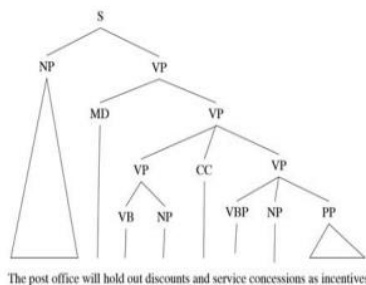
$$P(t_2) = 4,1472 \times 10^{-7}$$

امکان برابر دارند پس مرد در حالتی ته‌ایند (منا) شده

۳-۱- سوال سوم

اولین مشکل PCFG ها این است اهمیتی به lexical rules یا قواعد لغوی قائل نمی‌شوند به این دلیل ممکن است صرفاً به دلیل اینکه یک rule احتمال بالاتری دارد آن ساختار انتخاب خواهد شد بدون توجه به اینکه نگاهی به قواعد لغوی داشته باشد در حالی که این موارد مهم هستند و ساختار اصلی را مشخص می‌کنند. هم چنین ممکن است قواعدی که در حالت‌های مختلف parse کردن یک جمله استفاده می‌شوند متفاوت باشند ولی در نهایت احتمال هر دو جمله یکسان شوند چون احتمال ها یکی بوده‌اند. تصویر زیر حالت‌های مختلف parse کردن یک جمله را به ما نشان می‌دهد:

The post office will hold out discounts and service concessions as incentives.



با در نظر گرفتن lexicalized گرامرها این مشکل تا حدودی حل می‌شود چون در این گرامرها هر rule بر اساس کلمه‌ای که در حال حاضر داریم انتخاب می‌شود به این صورت که شاید احتمال یک rule با توجه به آن کلمه بیشتر یا کمتر باشد که در این صورت اطلاعات لغوی نیز به کار برده می‌شوند و پیش‌بینی دقیق‌تری خواهیم داشت. تصویر زیر احتمال یک rule با توجه به کلمات مختلف را نشان می‌دهد:

- $p(VP \rightarrow V \text{ NP NP}) = 0.00151$
- $p(VP \rightarrow V \text{ NP NP} \mid \text{said}) = 0.00001$
- $p(VP \rightarrow V \text{ NP NP} \mid \text{gave}) = 0.01980$ ”

مشکل دیگر PCFG ها bias داشتن به سمت درخت‌های کوتاه‌تر است و احتمال بیشتری را برای آن‌ها قائل می‌شود در حالی که ممکن است چنین فرضی درست نباشد. Lexicalized گرامرها تغییر چندانی در این وضع ایجاد نمی‌کنند ولی شاید اطلاعات لغوی وضعیت را کمی بهتر کند.

۲- سوال عملی

برای اینکه بتوانیم تصمیم‌گیری درستی در درخت تصمیم برای تشخیص tag ها داشته باشیم، لازم است ابتدا چند feature مناسب برای هر کلمه در جمله به دست آوریم تا بر اساس آن‌ها تصمیم بگیریم که چه tag را به آن کلمه نسبت دهیم. تابع features این کار را انجام می‌دهد. ویژگی‌های مختلفی در این تابع استخراج می‌شوند که از جمله آن‌ها می‌توان به پیشوند و پسوندهای کلمه و اینکه کلمه آخر یا اول جمله است و یا اینکه حروف بزرگ دارد یا نه اشاره کرد:

```
def features(sentence, index):
    return {
        'word': sentence[index],
        'is_first': index == 0,
        'is_last': index == len(sentence) - 1,
        'is_before_period': True if index != len(sentence) - 1 and sentence[index + 1] == '.' else False,
        'is_capitalized': sentence[index][0].upper() == sentence[index][0],
        'is_all_caps': sentence[index].upper() == sentence[index],
        'is_all_lower': sentence[index].lower() == sentence[index],
        'prefix-1': sentence[index][0],
        'prefix-2': sentence[index][:2],
        'suffix-1': sentence[index][-1],
        'suffix-2': sentence[index][-2:],
        'prev_word': '' if index == 0 else sentence[index - 1],
        'next_word': '' if index == len(sentence) - 1 else sentence[index + 1],
        'last_word_in_sentence': sentence[-1],
        'is_numeric': sentence[index].isdigit(),
        'capitals_inside': sentence[index][1:].lower() != sentence[index][1:]
    }
```

برای اینکه بتوانیم از جملات treebank برای آموزش استفاده کنیم لازم است tag های آن را حذف کنیم که از طریق تابع زیر این کار را انجام می‌دهیم:

```
def untag(tagged_sentence):
    return [w for w, t in tagged_sentence]
```

حالا باید داده‌ها را به دو بخش train و test تقسیم کنیم که نسبت آن‌ها ۸۰ به ۲۰ است:

```
cutoff = int(.80 * len(tagged_sentences))
training_sentences = tagged_sentences[:cutoff]
test_sentences = tagged_sentences[cutoff:]
```

حالا لازم است که برای آموزش یک دیتاست بسازیم. برای این کار برای جمله‌های untag شده، feature های موردنیاز را استخراج می‌کنیم و tag کلمه در جمله نیز label آن خواهد بود:

```
def transform_to_dataset(tagged_sentences):
    X, y = [], []

    for tagged in tagged_sentences:
        for index in range(len(tagged)):
            X.append(features(untag(tagged), index))
            y.append(tagged[index][1])

    return X, y

X, y = transform_to_dataset(training_sentences)
```

حالا برای آموزش مدل یک پایپ لاین تعریف می‌کنیم. در این پایپ لاین در ابتدا یک vectorizer خواهیم داشت که ویژگی‌ها را به بردارهای مناسب نگاشت می‌کند تا بتوانیم با استفاده از آن‌ها مدل را آموزش دهیم. سپس این داده‌ها به classifier داده خواهند شد تا tag ها را پیش‌بینی کند. نتیجه روی داده‌های تست دقت ۸۹ درصدی را مطابق شکل زیر نشان می‌دهد:

```
clf = Pipeline([
    ('vectorizer', DictVectorizer(sparse=False)),
    ('classifier', DecisionTreeClassifier(criterion='entropy'))
])

clf.fit(X[:10000], y[:10000])

print ('Training completed')

X_test, y_test = transform_to_dataset(test_sentences)

print ("Accuracy:", clf.score(X_test, y_test))
```

```
Training completed
Accuracy: 0.8947552273067518
```

حالا با تابع pos_tag می‌توانیم هر جمله دلخواهی را tag کنیم:

```
def pos_tag(sentence):  
    tags = clf.predict([features(sentence, index) for index in range(len(sentence))])  
    return list(zip(sentence, tags))
```

```
pos_tag(word_tokenize('This is a test sentence!'))
```

```
[('This', 'DT'),  
 ('is', 'VBZ'),  
 ('a', 'DT'),  
 ('test', 'NN'),  
 ('sentence', 'NN'),  
 ('!', 'CD')]
```