

تمرین سری سوم
NLP

یاسمین مدنی

۹۷۵۳۲۲۶۵

فهرست

سوال ۱.....	۲
مرحله اول.....	۲
مرحله دوم.....	۳
مرحله سوم.....	۴
سوال ۲.....	۶
سوال ۳.....	۸
سوال عملی.....	۹

سوال ۱

مرحله اول

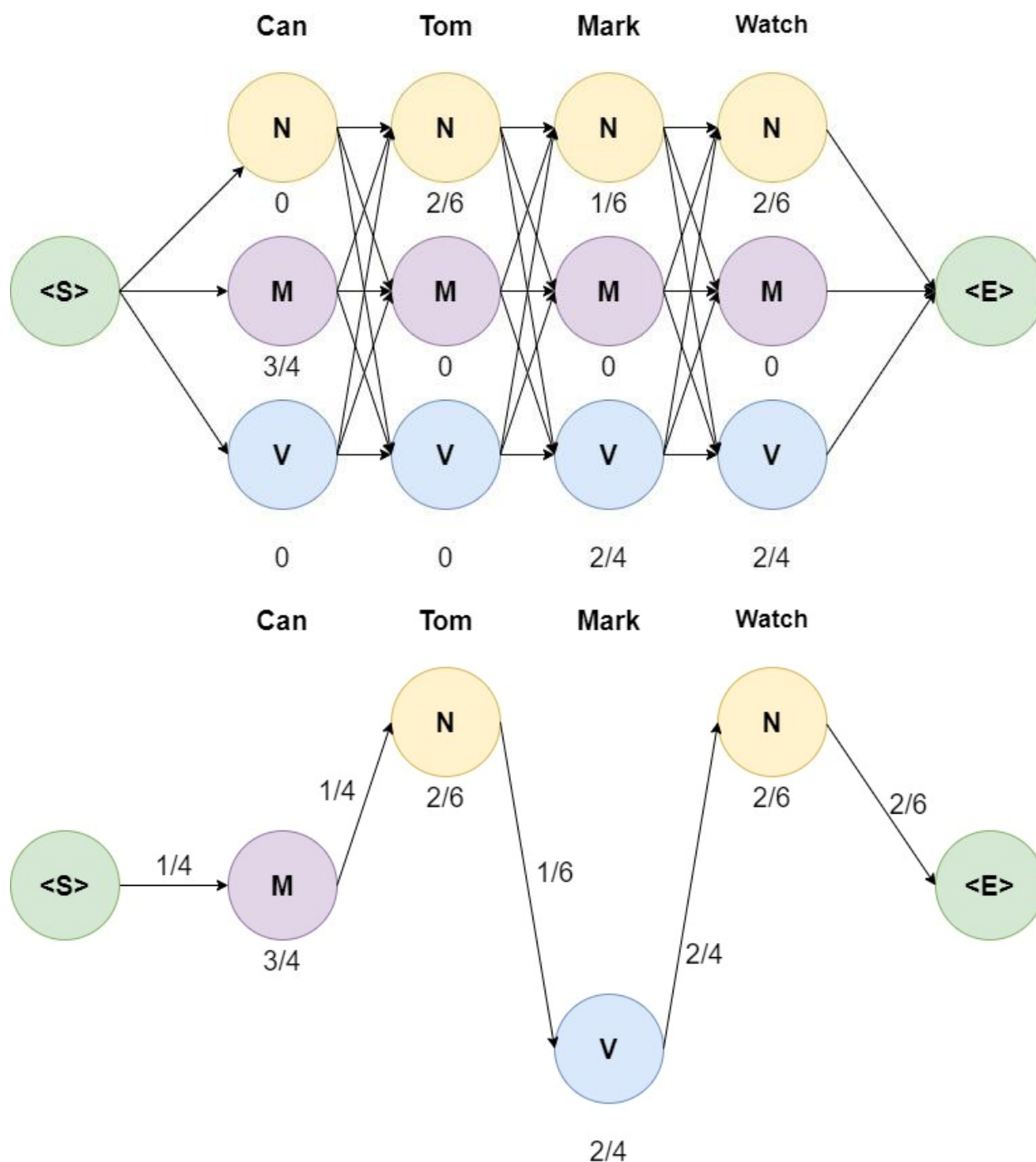
word	Noun	Modal	Verb
Mark	1/6	0	2/4
Tom	2/6	0	0
Can	0	3/4	0
Watch	2/6	0	2/4
Will	1/6	1/4	0

مرحله دوم

- <S> Mark(N) can(M) watch(V) <E>
- <S> Will(N) can(M) mark(V) watch(N)<E>
- <S> Can(M) Tom(N) watch(V) <E>
- <S> Tom(N) will(M) mark(V) watch(N)<E>

	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

مرحله سوم



$$P = 1/4 * 3/4 * 1/4 * 2/6 * 1/6 * 2/4 * 2/4 * 2/6 * 2/6 = 0.000072337$$

چون احتمال بزرگتر از ۰ است pos tag مناسب به صورت جدول زیر است.

Can	Modal
Tom	Noun
Mark	Verb
Watch	Noun

سوال ۲

نخست گرامر را به فرم نرمال چامسکی تبدیل می کنیم. از این رو بعضی از قواعد به صورت زیر باید تغییر کنند.

$S \rightarrow S \text{ S-Conj} \quad 0.2 \quad S\text{-Conj} \rightarrow \text{Conj } S \quad 1.0$

$NP \rightarrow \text{John } 0.04 \mid \text{Jack } 0.06 \mid \text{pie } 0.02 \mid \text{cream } 0.06 \mid \text{cake } 0.02 \quad 0.2$

$NP \rightarrow NP \text{ Conj-NP} \quad 0.2 \quad \text{Conj-NP} \rightarrow \text{Conj } NP \quad 1.0$

$VP \rightarrow \text{eat } 0.08 \mid \text{eats } 0.12 \mid \text{drinks } 0.2 \quad 0.4$

$VP \rightarrow \text{Verb } NP\text{-NP} \quad 0.1 \quad NP\text{-NP} \rightarrow NP \text{ NP} \quad 1.0$

John	eats	pie	with	cream	
Noun (0.2) / NP(0.04)	S (0.00384)	P (0.6)	-	S (0.00000041472)	John
	Verb (0.3) / VP (0.12)	VP (0.0018)	-	VP (0.00001296)	eats
		Noun (0.1) / NP (0.02)	-	NP (0.000144)	pie
			P (0.6)	PP (0.036)	with
				Noun (0.3) / NP (0.06)	cream

حالت های ممکن:

- $S \rightarrow NP VP \rightarrow \text{Noun VP} \rightarrow \text{John VP} \rightarrow \text{John Verb NP} \rightarrow \text{John eats NP} \rightarrow \text{John eats NP PP} \rightarrow \text{John eats Noun PP} \rightarrow \text{John eats pie PP} \rightarrow \text{John eats pie P NP} \rightarrow \text{John eats pie with NP} \rightarrow \text{John eats pie with Noun} \rightarrow \text{John eats pie with cream}$
- $S \rightarrow NP VP \rightarrow \text{Noun VP} \rightarrow \text{John VP} \rightarrow \text{John VP PP} \rightarrow \text{John Verb NP PP} \rightarrow \text{John eats NP PP} \rightarrow \text{John eats Noun PP} \rightarrow \text{John eats pie PP} \rightarrow \text{John eats pie P NP} \rightarrow \text{John eats pie with NP} \rightarrow \text{John eats pie with Noun} \rightarrow \text{John eats pie with cream}$

سوال ۳

یکی از مشکلات این روش این است که احتمالات را برای همه ی اجزا یکسان در نظر می گیرد و تفاوتی ندارد که کلمات در کجا ظاهر می شوند. کلمات خاص در زیرمجموعه های مختلف منجر به احتمالات متفاوت می شوند.

- John saw the man with the hat.

VP => VBD NP, NP => NP PP

- John saw the moon with the telescope.

VP => VBD NP PP

همانطور که در مثال بالا می بینیم مشخص نیست که کدام یک از گزینه ها فوق مورد نظر است. برای حل این مشکل میتوانیم از Lexical Information و از قواعد هد استفاده کنیم . اگر مشخصه متمایز در جای دیگری رخ دهد، می تواند به روش های دیگری تقسیم شود. همچنین سیستم های برچسب گذاری مختلف می توانند به مشکل کمک کنند.

سوال عملی

برای این پیاده سازی ابتدا نیاز داریم تا یک مجموعه برچسب گذاری شده را دانلود و ایمپورت کنیم.

```
1 nltk.download('treebank')
2 nltk.download('punkt')
3 nltk.download('universal_tagset')

[nltk_data] Downloading package treebank to /root/nltk_data...
[nltk_data] Package treebank is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
True

▼ Tagged sentences

[33] 1 tagged_sentences = nltk.corpus.treebank.tagged_sents()
```

برای هر اصطلاح، بسته به جمله ای که عبارت از آن استخراج شده است، فرهنگ لغت ویژگی ها را ایجاد می کنیم. این ویژگی ها می تواند شامل اطلاعاتی در مورد کلمات قبلی و بعدی و همچنین پیشوندها و پسوندها باشد.

```
1 def features(sentence, index):
2     return {
3         'word': sentence[index],
4         'is_first': index == 0,
5         'is_last': index == len(sentence) - 1,
6         'is_capitalized': sentence[index][0].upper() == sentence[index][0],
7         'is_all_caps': sentence[index].upper() == sentence[index],
8         'is_all_lower': sentence[index].lower() == sentence[index],
9         'prefix-1': sentence[index][0],
10        'prefix-2': sentence[index][1:2],
11        'prefix-3': sentence[index][2:3],
12        'suffix-1': sentence[index][-1],
13        'suffix-2': sentence[index][-2:],
14        'suffix-3': sentence[index][-3:],
15        'prev_word': '' if index == 0 else sentence[index - 1],
16        'next_word': '' if index == len(sentence) - 1 else sentence[index + 1],
17        'has_hyphen': '-' in sentence[index],
18        'is_numeric': sentence[index].isdigit(),
19        'capitals_inside': sentence[index][1:].lower() != sentence[index][1:]
20    }
```

اکنون، مانند یادگیری ماشین، باید مجموعه داده را برای آموزش و آزمایش تقسیم کنیم.

اکنون، شبکه عصبی ما بردارها را به عنوان ورودی می گیرد، بنابراین باید ویژگی های فرهنگ لغت خود را به بردار تبدیل کنیم.

```
1 part=int(0.8 * len(tagged_sentences))
2 x_train, y_train = transform_to_dataset( tagged_sentences[:part])
3 x_test, y_test = transform_to_dataset( tagged_sentences[part:])

[48] 1 classifier=Pipeline([
2         ('vectorizer', DictVectorizer(sparse=False)),
3         ('classifier', DecisionTreeClassifier(criterion='entropy'))
4     ])
5 classifier.fit(x_train[:10000], y_train[:10000])
```

با استفاده از برچسب تریگرام یا استفاده از مجموعه داده متفاوت و بزرگتر می توانیم دقت را افزایش دهیم .

▾ Sample for output of your PoS tagger

```
[52] 1 print(list(pos_tag(word_tokenize('This is my friend, John.'))))

[('This', 'DT'), ('is', 'VBZ'), ('my', 'NN'), ('friend', 'NN'), (',', ','), ('John', 'NNP'), ('.', '.')]

1 print(list(pos_tag(word_tokenize("let's go shopping"))))

[('let', 'VBD'), ("s", 'POS'), ('go', 'NN'), ('shopping', 'VBG')]
```