

Week 5 – Ngrams LM and Sequence Tagging

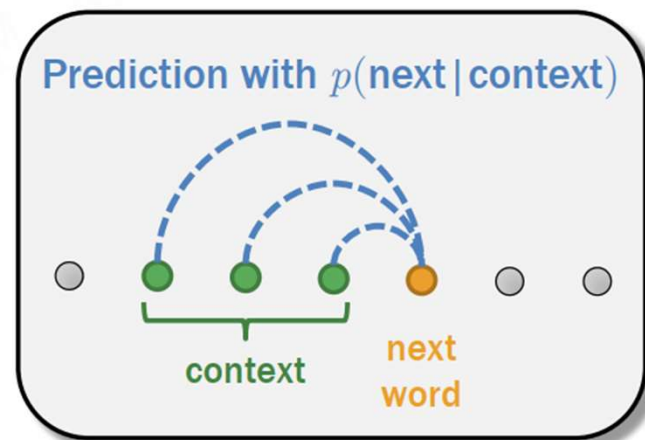
EGCO467 Natural Language and Speech Processing

n-gram models

- Language models whose context is truncated to at most $n-1$ previous words

$$P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k | \overbrace{w_{k-n+1} \dots w_{k-1}}^{n-1 \text{ prev words}} \overbrace{w_{k-1}}^{\text{next word}})$$

- Unigram** ($n=1$): $P(w_1 \dots w_N) = \prod_{k=1}^N p(w_k)$
- Bigram** ($n=2$): $P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k | w_{k-1})$
- Trigram** ($n=3$): $P(w_1 \dots w_N) = p(w_1)p(w_2|w_1) \prod_{k=3}^N p(w_k | w_{k-2}, w_{k-1})$



Formal Definition of LM

- m words: w_1, w_2, \dots, w_m
- LM is probability $P(w_1, w_2, \dots, w_m)$ for the whole text
- I like music <- high probability
- music like I <- low probability

Making a LM

- s1: a brown fox jumps over a lazy dog
- s2: the dog jumps over the fence

1. build vocab: {a, brown, fox, jumps, over, lazy, dog, the, over, fence}
2. $P(w)$ = how many time w appear / how many word total in corpus

$$N = 14$$

$$P(a) = 2/14, P(dog) = 2/14, P(brown) = 1/14,$$

Unigram LM

- Probability of any sequence $P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2) \dots P(w_n)$
- E.g. $P(a, brown, fox) = P(a)P(brown)P(fox) = (2/14) * (1/14) * (1/14)$

$$\text{Unigram } (n=1): P(w_1 \dots w_N) = \prod_{k=1}^N p(w_k)$$

$$\text{Bigram } (n=2): P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k | w_{k-1})$$

Bigram LM

- E.g. let's say that in corpus "a" occurs 4 times in: "**a** brown fox", "**a** brown cookie", "**a** big apple", "**a** big book"
- $P(\text{brown}|\text{a}) = 2/4$. Because we have two instances of "a brown" out of all 4 instances where "a ..." appears
- So $P(w_2|w_1) = \text{count of } w_1 \text{ followed by } w_2 / \text{count of } w_1 \text{ followed by } x$
= count of w_1 followed by w_2 / count of w_1
- Do this for all word pairs that occur in corpus.

Example

- s1: a brown fox jumps over a lazy dog
- s2: the dog jumps over the fence
- vocab: {a, brown, fox, jumps, over, lazy, dog, the, over, fence}
- $P(\text{brown}|\text{a}) = |\{ (\text{a brown}) \}| / |\{ (\mathbf{a} \text{ brown}), (\mathbf{a} \text{ lazy}) \}| = 0.5$
- $P(\text{jumps}|\text{fox}) = |\{ (\text{fox jumps}) \}| / |\{ (\mathbf{fox} \text{ jumps}) \}| = 1.0$
- $P(\text{over}|\text{the}) = \mathbf{the over}$ never occurs, so $P = 0$
- $P(\text{fox}|\text{brown}) = 1.0$

Example

- $P(a,brown,fox) = P(a)*P(brown|a)*P(fox|brown) = (2/14)*0.5*1.0$

Unigram ($n=1$): $P(w_1 \dots w_N) = \prod_{k=1}^N p(w_k)$

Bigram ($n=2$): $P(w_1 \dots w_N) = p(w_1) \prod_{k=2}^N p(w_k|w_{k-1})$

Trigram LM

- $P(w_3|w_2, w_1) = \text{count of sequence } (w_1, w_2, w_3) \text{ in corpus} / \text{count of sequence } (w_1, w_2) \text{ in corpus}$
- Do this for all possible values of w_1, w_2, w_3

n-grams

- $p(w_n | w_{n-1}, w_{n-2}, \dots, w_1)$
- The number of conditions grows with n
- Longer condition \rightarrow less chance to appear in corpus \rightarrow count = 0
- Maximum used in practice is only trigram
- For $n > 3$, we use NN to estimate the probability

LSTM Language Model

- At time k , predict k th word, give the previous $k-1$ words (context).
- Append the predicted word to the context.
- Go the next time step.

- $k=3$: *The quick* **brown**
- $k=4$: *The quick brown* **fox**
- $k=5$: *The quick brown fox* **jumps**
-

perplexity

- perplexity = *inverse normalized probability of test set*
- lower is better

$$PP(W) = \frac{1}{P(w_1, w_2, \dots, w_N)^{1/N}}$$

- W = test set
- w_1 to w_N = all words in the test set
- N = number of words in the test set

Probability of test set

Test Set

"Yesterday I went to the cinema"

"Hello, how are you?"

"The dog was wagging its tail"

High probability
Low perplexity

not surprised by test set

Fake/incorrect sentences

"Can you does it?"

"For wall a driving"

"She said me this"

Low probability
High perplexity

surprised by test set

Perplexity

- inverse normalized probability of test set

$$\begin{aligned} PP(W) &= \frac{1}{P(w_1, w_2, \dots, w_N)^{\frac{1}{N}}} \\ &= \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}} \end{aligned}$$

Example - trigram

- 05 - trigram NLTK.ipynb

Exercise

- Given the same corpus from the word2vec (03 - word2vec basic.ipynb)
- Count how many bigrams there are the corpus
- Calculate the probability of "of a"
- Modify the .ipynb in the previous slide to do this exercise.
- (ignore punctuation marks and lower/upper case)

We are about to study the idea of a computational process.

Computational processes are abstract beings that inhabit computers.

As they evolve, processes manipulate other abstract things called data.

The evolution of a process is directed by a pattern of rules

called a program. People create programs to direct processes. In effect,

we conjure the spirits of the computer with our spells.

Dealing with OOV words

- If corpus has: I am a police
- Test has: I am a farmer
- $p(\text{farmer} \mid \text{a}) = 0$
- $p(\text{I am a farmer}) = 0$
- perplexity = inf
- Use <unk> for unknown words (don't have in training data)

Sequence tagging problem

- Label each token, not the entire sequence
- Label whole sequence is classification
- Label each token is sequence tagging



Two main types of sequence tagging

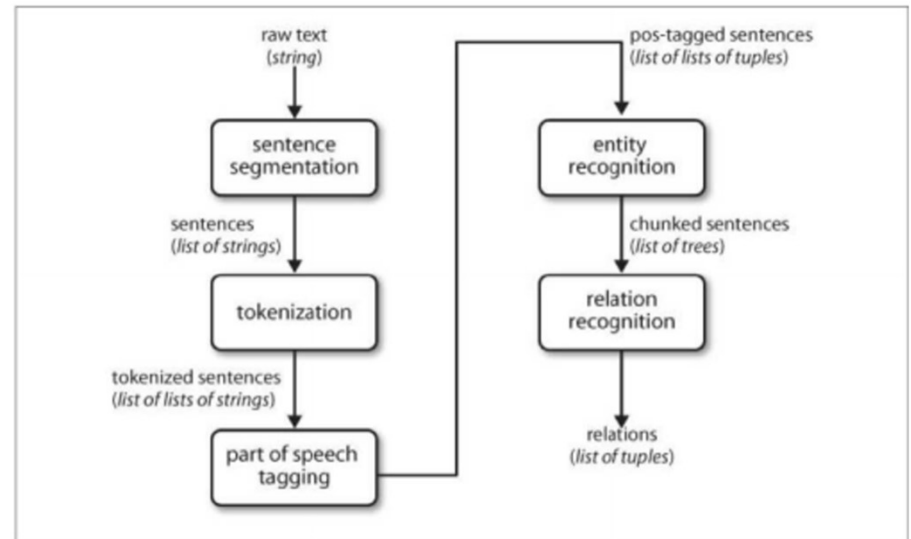


Elon Musk **PERSON** apparently wasn't aware that his company SpaceX had a **Facebook** **ORG** page. The SpaceX and **Tesla** **PRODUCT** CEO has responded to a comment on **Twitter** **ORG** calling for him to take down the SpaceX, Tesla and **Elon Musk** **ORG** official pages in support of the #deletefacebook movement by **first** **ORIGINAL** acknowledging he didn't know one existed, and then following up with promises that he would indeed take them down.

He's done just that, as the SpaceX **NORP** Facebook page is now gone, after having been live earlier **today** **DATE** (as you can see from the screenshot included taken at **around 12:10 PM ET** **TIME**).

Why do NER/POS

- Information extraction
- Relation extraction

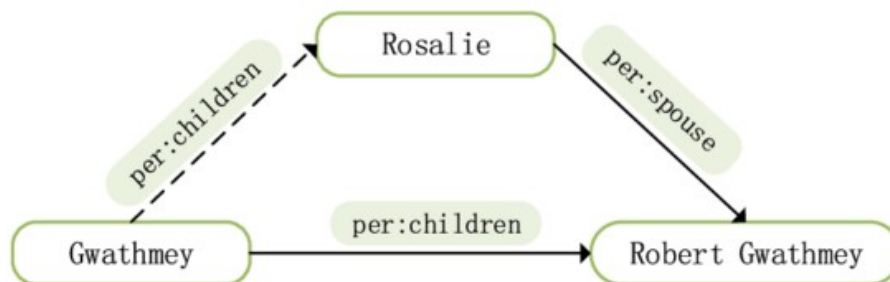


Relation Extraction

- extract relationship between each entity (nouns) in the text

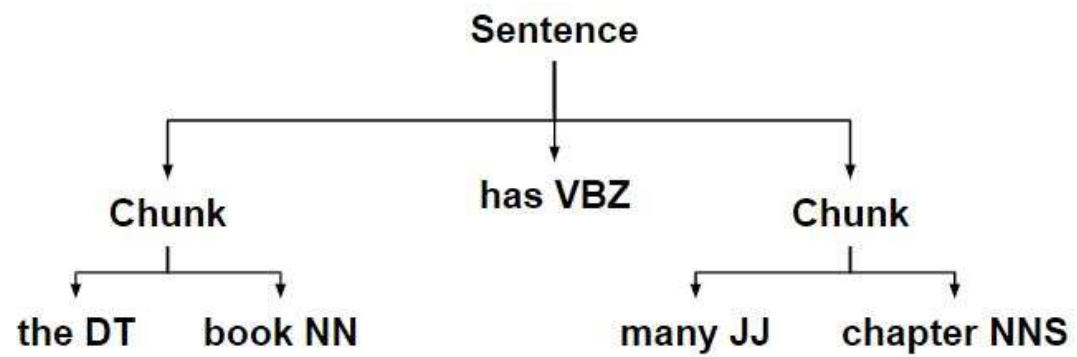
Gwathmey was born in 1938, the only child of painter Robert Gwathmey and his wife, Rosalie, a photographer.

(a)



(b)

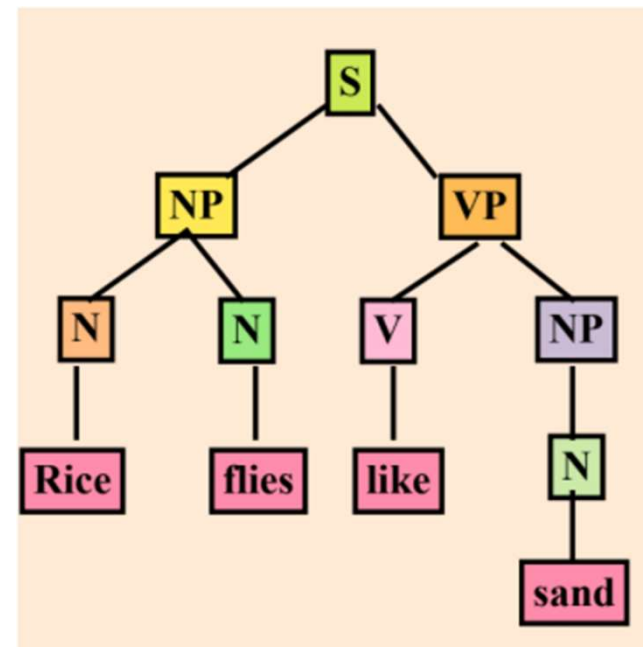
Sentence chunking



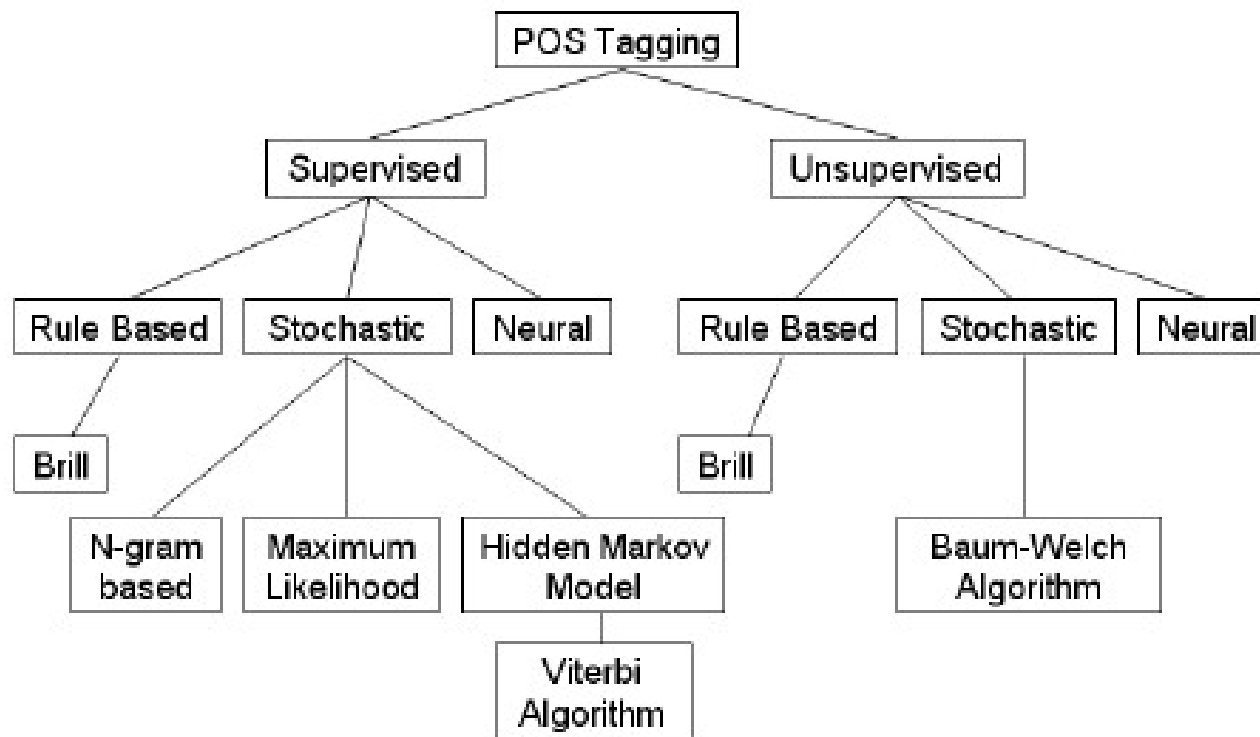
Sentence Parse Tree

Grammar

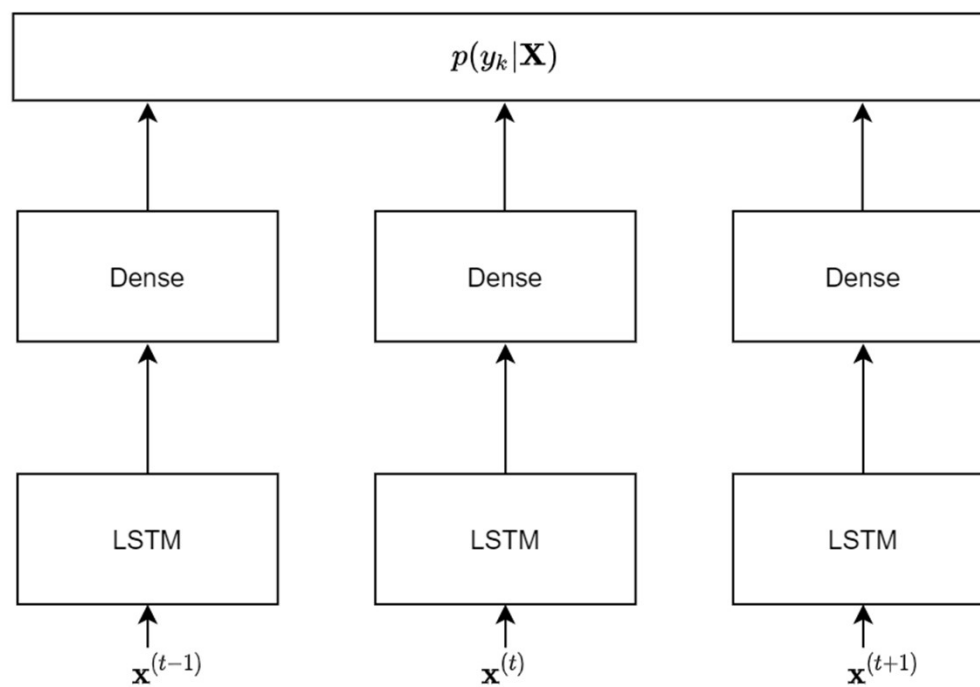
$S \rightarrow NP VP$	$N \rightarrow \text{Rice}$
$VP \rightarrow V NP$	$N \rightarrow \text{flies}$
$VP \rightarrow V$	$N \rightarrow \text{sand}$
$NP \rightarrow N N$	$V \rightarrow \text{like}$
$NP \rightarrow N$	$P \rightarrow \text{like}$
$PP \rightarrow P NP$	



How to do POS?

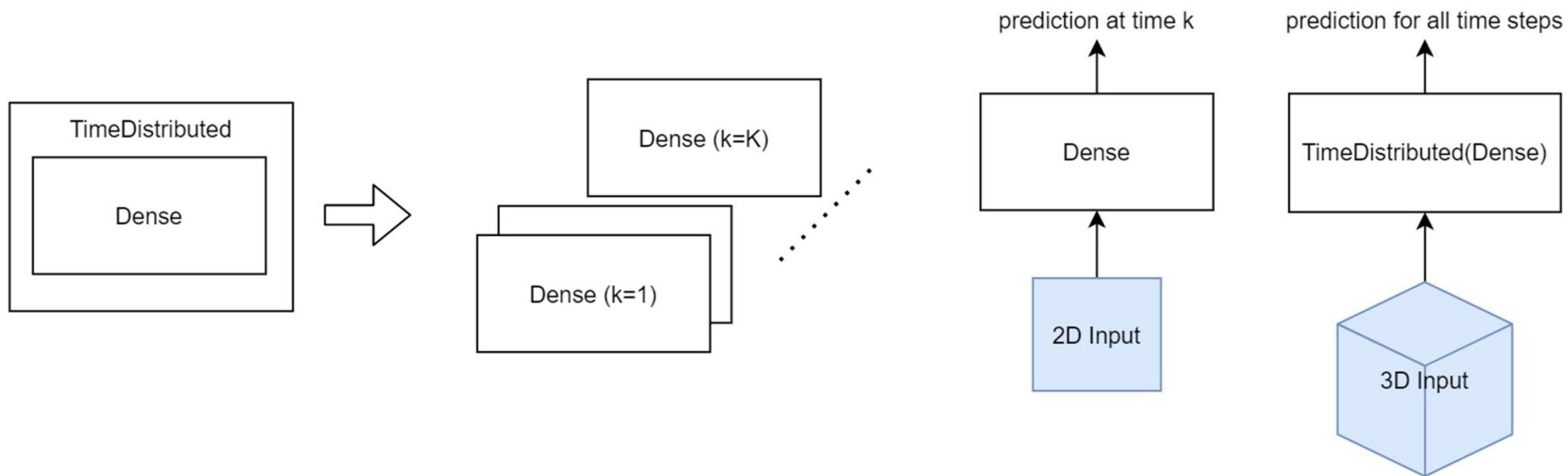


One way to tag tokens



TimeDistributed

TimeDistributed(Dense(n_tags, activation="softmax"))



Time Distributed

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 50)]	0
embedding (Embedding)	(None, 50, 50)	1758950
dropout (Dropout)	(None, 50, 50)	0
bidirectional (Bidirectional)	(None, 50, 200)	120800
time_distributed (TimeDistributed)	(None, 50, 17)	3417

note: the number of classes = 17
and sequence length = 50

Loss of sequence tagging

- $L(t)$ = loss at time t
- Sequence tagging is doing classification for each time step
- The loss $L(t)$ is just regular cross-entropy like classification

$$L = \frac{1}{T} \sum_{t=0}^T L(t)$$

Example - POS tagging

- 05 - LSTM sequence tagging.ipynb

Orchid Dataset

```
1 <corpus>
2 <document TPublisher="ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ, กระทรวงวิทยาศาสตร์ เทคโนโลยีและการพลังงาน" EPublisher="National
3 <paragraph id="1" line_num="12">
4 <sentence id="1" line_num = "13" raw_txt = "การประชุมทางวิชาการ ครั้งที่ 1">
5 <word surface="การ" pos="FIXN"/>
6 <word surface="ประชุม" pos="VACT"/>
7 <word surface="ทาง" pos="NCMN"/>
8 <word surface="วิชาการ" pos="NCMN"/>
9 <word surface="&lt;space&gt;" pos="PUNC"/>
10 <word surface="ครั้ง" pos="CFQC"/>
11 <word surface="ที่ 1" pos="DONM"/>
12 </sentence>
13 <sentence id="2" line_num = "23" raw_txt = "โครงการวิจัยและพัฒนาอิเล็กทรอนิกส์และคอมพิวเตอร์">
14 <word surface="โครงการวิจัยและพัฒนา" pos="NCMN"/>
15 <word surface="อิเล็กทรอนิกส์" pos="NCMN"/>
16 <word surface="และ" pos="JCRG"/>
17 <word surface="คอมพิวเตอร์" pos="NCMN"/>
18 </sentence>
19 <sentence id="3" line_num = "30" raw_txt = "ปีงบประมาณ 2531">
20 <word surface="ปีงบประมาณ" pos="NCMN"/>
21 <word surface="&lt;space&gt;" pos="PUNC"/>
22 <word surface="2531" pos="NCNM"/>
23 </sentence>
24 <sentence id="4" line_num = "36" raw_txt = "เล่ม 1">
25 <word surface="เล่ม" pos="CNIT"/>
26 <word surface="&lt;space&gt;" pos="PUNC"/>
27 <word surface="1" pos="DONM"/>
28 </sentence>
29 </paragraph>
```

Virach Sornlertlamvanich, Naoto Takahashi and Hitoshi Isahara. Building a Thai Part-Of-Speech Tagged Corpus (ORCHID). The Journal of the Acoustical Society of Japan (E), Vol.20, No.3, pp 189-198, May 1999./p>

Orchid POS system

The following table shows POS tags as used in ORCHID:

Abbreviation	Part-of-Speech tag	Examples
NPRP	Proper noun	วินโดวส์ 95, โคโรนา, โค้ก
NCNM	Cardinal number	หนึ่ง, สอง, สาม, 1, 2, 10
NONM	Ordinal number	ที่หนึ่ง, ที่สอง, ที่สาม, ที่1, ที่2
NLBL	Label noun	1, 2, 3, 4, ก, ข, a, b
NCMN	Common noun	หนังสือ, อาหาร, อาคาร, คน
NTTL	Title noun	ครู, พลเอก
PPRS	Personal pronoun	คุณ, เขา, ฉัน
PDMN	Demonstrative pronoun	นี้, นั้น, ที่นั่น, ที่นี่
PNTR	Interrogative pronoun	ใคร, อะไร, อย่างไร
PREL	Relative pronoun	ที่, ซึ่ง, อัน, ผู้
VACT	Active verb	ทำงาน, ร้องเพลง, กิน
VSTA	Stative verb	เห็น, รู้, คือ

Homework 5

- Parse the Orchid data (xmlchid.xml) using Python's xml library
<https://docs.python.org/3/library/xml.etree.elementtree.html>
- Extract the POS tag of each word
- Store data in a list like below
- Save the list to a .pickle file

```
<sentence id="1" line_num = "13" raw_txt = "การประชุมทางวิชาการ ครั้งที่ 1">
<word surface="การ" pos="FIXN"/>
<word surface="ประชุม" pos="VACT"/>
<word surface="ทาง" pos="NCMN"/>
<word surface="วิชาการ" pos="NCMN"/>
<word surface="&lt;space&gt;" pos="PUNC"/>
<word surface="ครั้ง" pos="CFQC"/>
<word surface="ที่ 1" pos="DONM"/>
</sentence>
```

```
[
(["การ", "ประชุม", "ทาง", "วิชาการ", " ", "ครั้ง", "ที่ 1"], ["FIXN", "VACT", "NCMN", "NCMN", "PUNC", "CFQC", "DONM"]),
(next sentence),
...
]
```