

# ICA 4 POS Tagging

(September 30th)

In this activity, we will work on POS Tagging

**(Same as ICA 1) From NLTK import gutenberg corpus and open up the book *Alice in Wonderland* by Lewis Carol.**

**No need for preprocessing the text, in fact we want it to be un-preprocessed.**

Here is a resource that will help:

<http://www.ling.helsinki.fi/kit/2009s/clt231/NLTK/book/ch05-CategorizingAndTaggingWords.html>

Step 2: Run the sentences in the first three paragraphs of *Alice in Wonderland* through POS tagger using the below command.

```
nltk.pos_tag(text)
```

Below is a list of tags and what they represent.

Tag	Meaning	Examples
ADJ	adjective	<i>new, good, high, special, big, local</i>
ADV	adverb	<i>really, already, still, early, now</i>
CNJ	conjunction	<i>and, or, but, if, while, although</i>
DET	determiner	<i>the, a, some, most, every, no</i>
EX	existential	<i>there, there's</i>
FW	foreign word	<i>dolce, ersatz, esprit, quo, maitre</i>
MOD	modal verb	<i>will, can, would, may, must, should</i>
N	noun	<i>year, home, costs, time, education</i>
NP	proper noun	<i>Alison, Africa, April, Washington</i>
NUM	number	<i>twenty-four, fourth, 1991, 14:24</i>
PRO	pronoun	<i>he, their, her, its, my, I, us</i>
P	preposition	<i>on, of, at, with, by, into, under</i>
TO	the word <i>to</i>	<i>to</i>
UH	interjection	<i>ah, bang, ha, whee, hmpf, oops</i>
V	verb	<i>is, has, get, do, make, see, run</i>
VD	past tense	<i>said, took, told, made, asked</i>
VG	present participle	<i>making, going, playing, working</i>
VN	past participle	<i>given, taken, begun, sung</i>
WH	<i>wh</i> determiner	<i>who, which, when, what, where, how</i>

Step 3: Using the output from the previous step, list tags in order of decreasing frequency. Which are the top 5 most frequent tags?

Step 4: The POS tagging task can be accomplished using regular expressions. For example, we can assume that any word ending in “-ed” is a past participle of a verb.

Using a regular expression based POS tagger, with the following three regexes, determine at least three of the words from the output of Step 2 above which would be incorrectly tagged.

```

(r'.*ing$', 'VBG'),      # gerunds
... (r'.*ed$', 'VBD'),    # simple past
... (r'.*s$', 'NNS'),     # plural nouns

```