## 3.   NLP – Replacing and Correcting words

## Contents

### 3. NLP – Replacing and correcting words

## 1. Text conversions

✓  We can convert the text from lower to upper and upper to lower case

---

| | |
|---|---|
| **Program Name** | Converting lower case to upper case<br>demo1.py<br><br>text ="hello good morning"<br>print(text.upper()) |
| **Output** | <br>HELLO GOOD MORNING |

---

| | |
|---|---|
| **Program Name** | Converting upper case to lower case<br>demo2.py<br><br>text ="HELLO GOOD MORNING"<br>print(text.lower()) |
| **Output** | <br>hello good morning |

## 2. Removing numbers

✓ By using regular expression we can remove numbers from the text

---

| | |
|---|---|
| Program Name | Removing numbers from the text<br>demo3.py<br><br>import re<br><br>myString = 'Box A has 4 red and 6 white balls, while Box B has 3 red and 5 blue balls.'<br><br>output = re.sub(r'\d+', '', myString)<br>print(output) |
| Output | Box A has red and white balls, while Box B has red and blue balls. |

---

**3. Removing punctuations**

✓ By using regular expression we can remove the punctuations from the text

| | |
|---|---|
| **Program Name** | Removing the punctuations from the text<br>demo4.py<br><br>import re<br><br>text = "Hello $@#$# Good !@#!@# morning #*#@&@#"<br><br>print("Text is:", text)<br><br>res = re.sub(r'[^\w\s]', '', text )<br><br>print("After punctuations:", res) |
| **Output** | <br><br>Text is: Hello $@#$# Good !@#!@# morning #*#@&@#<br>After punctuations: Hello  Good  morning |

**4. Removing whitespaces**

✓ We can remove the whitespaces in string by using strip() method.

| | |
|---|---|
| **Program Name** | Removing whitespaces from text<br>demo5.py<br><br>text = "          a sample string        "<br><br>print(text)<br>res = text.strip()<br>print(res) |
| **Output** | <br><br>          a sample string<br>a sample string |

## 5. Part of Speech Tagging (POS)

- ✓ The goal of POS is to assign the various parts of a speech to every word of the provided text like nouns, adjectives, verbs, etc.
- ✓ This is normally done based on the definition and the context.
- ✓ Install textblob library,
  - ○ pip install textblob

| | |
|---|---|
| Program Name | Removing whitespaces from text<br>demo6.py<br><br>```python<br>from textblob import TextBlob<br>import nltk<br><br>nltk.download('averaged_perceptron_tagger')<br><br>myString = "Parts of speech: an article, to run, fascinating, quickly, and, of"<br><br>output = TextBlob(myString)<br>print(output.tags)<br>``` |
| Output | [('Parts', 'NNS'), ('of', 'IN'), ('speech', 'NN'), ('an', 'DT'), ('article', 'NN'), ('to', 'TO'), ('run', 'VB'), ('fascinating', 'VBG'), ('quickly', 'RB'), ('and', 'CC'), ('of', 'IN')] |

Some examples are as below:

| Abbreviation | Meaning |
| --- | --- |
| CC | coordinating conjunction |
| CD | cardinal digit |
| DT | determiner |
| EX | existential there |
| FW | foreign word |
| IN | preposition/subordinating conjunction |
| JJ | adjective (large) |
| JJR | adjective, comparative (larger) |
| JJS | adjective, superlative (largest) |
| LS | list market |
| MD | modal (could, will) |
| NN | noun, singular (cat, tree) |
| NNS | noun plural (desks) |
| NNP | proper noun, singular (sarah) |
| NNPS | proper noun, plural (indians or americans) |
| PDT | predeterminer (all, both, half) |
| POS | possessive ending (parent\ 's) |
| PRP | personal pronoun (hers, herself, him,himself) |
| PRP$ | possessive pronoun (her, his, mine, my, our ) |
| RB | adverb (occasionally, swiftly) |
| RBR | adverb, comparative (greater) |
| RBS | adverb, superlative (biggest) |
| RP | particle (about) |
| TO | infinite marker (to) |
| UH | interjection (goodbye) |

| VB | verb (ask) |
| --- | --- |
| VBG | verb gerund (judging) |
| VBD | verb past tense (pleaded) |
| VBN | verb past participle (reunified) |
| VBP | verb, present tense not 3rd person singular(wrap) |
| VBZ | verb, present tense with 3rd person singular (bases) |
| WDT | wh-determiner (that, what) |
| WP | wh- pronoun (who) |
| WRB | wh- adverb (how) |

| | |
| --- | --- |
| **Program Name** | pos example<br>demo7.py<br><br>```from nltk.corpus import wordnet

syn = wordnet.synsets('hello')[0]
print("Syn tag : ", syn.pos())

syn = wordnet.synsets('doing')[0]
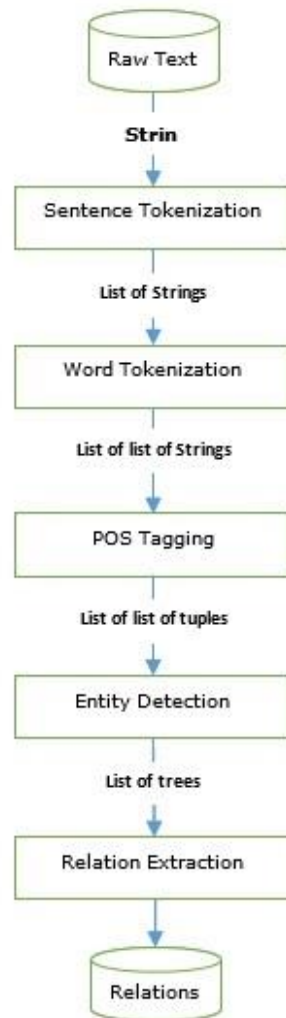print("Syn tag : ", syn.pos())

syn = wordnet.synsets('beautiful')[0]
print("Syn tag : ", syn.pos())``` |
| **Output** | ```Syn tag :  n
Syn tag :  v
Syn tag :  a
Syn tag :  r``` |

**6. Information Extraction**

- ✓ We need to understand the tags and parsers to build information extraction engine.
- ✓ Let us see a basic information extraction pipeline

**7. Information extraction has many applications including**

- ✓ Business intelligence
- ✓ Resume harvesting
- ✓ Media analysis
- ✓ Sentiment detection
- ✓ Patent search
- ✓ Email scanning

## 8. Collocations: Bigrams and Trigrams

### What is Collocations?

- ✓ Collocations are the pairs of words occurring together many times in paragraphs.
- ✓ It is calculated by the number of those pair occurring together to the overall word count of the paragraph.
- ✓ We can say that finding collocations requires calculating the frequencies of words and their appearance in the context of other words.

### Bigrams and Trigrams

- ✓ Collocation can be categorized into two types,
  - o Bigrams combination of two words
  - o Trigrams combination of three words
- ✓ Bigrams and Trigrams provide more meaningful and useful features for the feature extraction stage.
- ✓ These are especially useful in text-based sentimental analysis.

| | |
|---|---|
| **Program Name** | Bigram example<br>demo8.py<br><br>```import nltk```<br><br>```text = "Data Science is a totally new kind of learning experience."```<br>```Tokens = nltk.word_tokenize(text)```<br>```output = list(nltk.bigrams(Tokens))```<br><br>```print(output)``` |
| **Output** | [('Data', 'Science'), ('Science', 'is'), ('is', 'a'), ('a', 'totally'), ('totally', 'new'), ('new', 'kind'), ('kind', 'of'), ('of', 'learning'), ('learning', 'experience'), ('experience', '.')] |

Program
Name

Trigram example
demo9.py

```
import nltk

text = "Data Science is a totally new kind of learning experience."
Tokens = nltk.word_tokenize(text)
output = list(nltk.trigrams(Tokens))

print(output)
```

Output

```
[('Data', 'Science', 'is'), ('Science', 'is', 'a'), ('is', 'a', 'totally'), ('a', 'totally',
'new'), ('totally', 'new', 'kind'), ('new', 'kind', 'of'), ('kind', 'of', 'learning'), ('of',
'learning', 'experience'), ('learning', 'experience', '.')]
```

## 9. Wordnet

- ✓ Wordnet is an NLTK lexical database for English.
- ✓ It can be used to find the meaning of words, synonym or antonym.

### synset

- ✓ Synset is a special kind of a simple interface that is present in NLTK to look up words in Wordnet.
- ✓ Synset instances are the groupings of synonymous words that express the same concept.

| | |
|---|---|
| Program Name | wordnet example<br>demo10.py<br><br>```from nltk.corpus import wordnet<br><br>syn = wordnet.synsets('hello')[0]<br><br>print ("Synset name :", syn.name())<br>print ("Synset meaning :", syn.definition())<br>print ("Synset example :", syn.examples())``` |
| Output | Synset name : hello.n.01<br>Synset meaning : an expression of greeting<br>Synset example : ['every morning they exchanged polite hellos'] |

<table>
<tr><td>Program Name</td><td>wordnet example<br>demo11.py<br><br>from nltk.corpus import wordnet<br><br>syn = wordnet.synsets('boy')[0]<br><br>print ("Synset name :", syn.name())<br>print ("Synset meaning :", syn.definition())<br>print ("Synset example :", syn.examples())</td></tr>
<tr><td>Output</td><td>Synset name : male_child.n.01<br>Synset meaning : a youthful male person<br>Synset example : ['the baby was a boy', 'she made the boy brush his teeth every night', 'most soldiers are only boys in uniform']</td></tr>
</table>

<table>
<tr><td>Program Name</td><td>wordnet example<br>demo12.py<br><br>from nltk.corpus import wordnet<br><br>syn = wordnet.synsets('good')[0]<br><br>print ("Synset name :", syn.name())<br>print ("Synset meaning :", syn.definition())<br>print ("Synset example :", syn.examples())</td></tr>
<tr><td>Output</td><td>Synset name : good.n.01<br>Synset meaning : benefit<br>Synset example: ['for your own good', "what's the good of worrying?"]</td></tr>
</table>