

## Assignment 12

import re, string

Loading the data

```
filename = '/content/novel.txt'
```

```
file = open(filename, 'rt', encoding = 'utf-8')
```

```
text = file.read()
```

```
file.close()
```

Split the data into words by white space.

```
words = text.split()
```

```
print(words[:120])
```

['The', 'one', 'morning', 'John', 'George', 'Sims', 'Warner', 'from', 'trouble', 'dreams', 'he', 'found', 'himself', 'transformed', 'into', 'and', 'into', 'a']

Here we are seeing that the punctuation is preserved (E.g. armour-like and wasn't) which is nice. We can also see that end of sentence punctuation is kept with the last word (E.g. thought) which is not great.

So this time let's try to split the words at non-word character.

```
words = re.split('[^\\w+]', text)
```

```
print(words[:120])
```

['The', 'one', 'morning', 'John', 'George', 'Sims', 'Warner', 'from', 'trouble', 'dreams', 'he', 'found', 'himself', 'transformed', 'into', 'and', 'into', 'a']

Here we are seeing that the words would like through. Have been converted into thought. But the problem is that the word

like wasn't are converted into two words like and +, we need to fix it

In python, we can use `string.punctuation` to get bunch of punctuation at once. we will use that to remove punctuations from our text

```
print(string.punctuation)
```

```
!"#$%&'()*+,-./:;<=>?@[ ]^_`{|}~
```

So now we will split the words by white space and then remove all the punctuations which have been recorded in the data

```
words = text.split()
```

```
re_punc = re.compile('[%s]' % re.escape(string.punctuation))
```

```
stripped = [re_punc.sub('', word) for word in words]
```

```
print(stripped[:120])
```

```
['one', 'morning', 'when', 'Gregor', 'saw', 'wore', 'from', 'troubled', 'dreams', 'he', 'found', 'himself',  
'transformed', 'in', 'his', 'bed', 'into', '-----']
```

Here we can see that we don't have the words like thought, but we have words like wasn't

Sometimes that text also contains the characters which are not printable. we need to filter those out too. To do this, we can use python `string.printable` which gives us bunch of characters that can be printed. So we will remove the characters which are not present in this

```
re_print = re.compile('[^%s]' % re.escape(string.printable))
```

```
result = [re_print.sub('', word) for word in stripped]
```

print(result[:120])

['one', 'morning', 'when', 'Gregor', 'Samsa', 'woke', 'from', 'troubled', 'dreams', 'he', 'found', 'transformed', 'in', 'his']

Now let's make all the words into lowercase. This will reduce our vocabulary. But this has some disadvantages also. After doing this, two words such as Apple as a company and apple as a fruit will be considered a same entity

Result = [word.lower() for word in result]

print(result[:120])

['one', 'morning', 'when', 'gregor', 'samsa', 'woke', 'from', 'troubled', 'dreams', 'he', 'found', 'transformed', 'in', 'his']

Also, words with one character won't contribute to most of the NLP tasks. So we will be removing these too.

result = [word for word in result if len(word) > 1]

print(result[:120])

['one', 'morning', 'when', 'gregor', 'samsa', 'woke', 'from', 'troubled', 'dreams', 'he', 'found', 'transformed', 'in', 'his']

