

» Dealing with html

» Dealing with html

```
from bs4 import BeautifulSoup
```

```
from bs4 import BeautifulSoup
```

```
raw = BeautifulSoup(html, ...)
```

30/8/24

Paper 7 : AT snake oil
meaning fraud

NLTK chapter 3

- Reading local files

for local files
f = open('path of folder', 'w') // example.txt at the end

```
f = open('path of folder', 'w')
f.write("a quick brown fox to jump over the lazy dog")
f.close()
```

f. dose

» [example named test file will be created with a quote written in it]

$f = \text{open}(\text{same as earlier})$

f. write ("a new line")

f. close

[earlier text will get deleted]

» write 'a' instead of ('w') to append it and not erase the earlier text.

» "\n a new line" with 'a' — append with new line.

» 'r' used to read sentences from text file.

- NLP Pipeline

HTML \rightarrow ASCII $\xrightarrow{\text{tokenize the texts}}$ text \rightarrow Vocab

• Regular expressions for detecting word patterns

→ ~~#~~ to find patterns from text using python

en → english

\$ → word ending with this pattern

wislower → w is of lower index

^ → word starting with this pattern

+aa → word having — aa in it

aa+ → of pattern aam, aali, etc.

• Stemmers: part of the word which gets modified into new word
eg: li

• Process of stemming is lemmatisation.

• part of speech tagging (pos-tagging)

• tagged corpora

• dictionaries in python

then

earlier