

30/8/24

Paper

⑤ → Talking about ethics → meaningless

⑥ → Ethics doesn't have a legal analogue or any universal acceptance.

'Human Rights' → by UN → group of countries

→ Ethics are different from H.R.

→ H.R. → can be enforced universally.

↳ ∴ easier to practice across borders.

→ But also a problem → huge critiques.

⑦ Snake oil → 'china lagana'

Most of AI claims → snake oil.

+ → various snake oil → various claims.

Eg: Emotion recognition.

→ snake oil → where fraudulent datasets used to train the model.

Eg: Sentiment analysis

NLTK chapter ③

Recap: - HTML parsing using 'Beautiful soup'

→ ① Writing a file

→ Path where file to be saved.

→ In terminal:-

f = open ("path" example.txt, "w")

means to write.

w → write
a → append
r → read

classmate

Date _____

Page _____

→ f.write("A fox jumped ... dog")
→ f.close()

→ file saved,

→ Now to add something again to it:-

→ f.open("path \\ example.txt", "a")
appends the data instead of overwriting it

→ f.write("what a lovely day")

→ f.close()

→ Now to add something in new line.

→ f.open(---)

→ f.write("In a quick fox")
next line sentence element

'cv' → comma separated values.

→ f.close()

→ Now to read a file

→ f.open("path", "r") → read

→ for i in f:
print(i)

→ Everything printed.

→ NLP Pipeline:-

From internet

Text / HTML

Regex

Text

Vocab

Regex/
Beautiful
Soup
preprocessing
tokenization
etc.

→ Regular Expressions for detecting word patterns:-

→ To find tense of words,
eg:- ing for present continuous.

→ • import nltk

• import re

• wordlist = [w for w in nltk.corpus.CERT if w.islower()]
↳ only lower-case

• wordlist[:50] → print first 50 words

• stuff = [w for w in wordlist

if re.search('ing\$', w)]

\$ → tell for words ending with 'ing'.

• stuff[:50]

→ Now for words starting with 'win'.

• stuff [... if re.search('^win', w)]

^ → for start

for all
english
language
words

just
change
this

→ Other operators in chapter 3-4

→ stemmers:

→ process of converting every word to its smallest form, usually removing: ing, es, s, etc. from last.

→ lemmatizer

→ similar to stemmers

→ segmentation:

- Breaking down sentence into words.
- Usually difficult for some languages like Chinese.

NLTK Chap. (4) :

- '=' → assignment ; '==' → comparison
- sequence → list, string, tuple
- C++ → object oriented prog. (class, struct etc.)
- python → 'OOP' or 'procedural and declarative'
- variable scope → local or global.
- coding discipline → usually not followed, but helps in better understanding and increases readability.

- Algorithm design → to ↑ efficiency of a code
- ↳ code consumes less resources
 - ↳ works faster.
 - ↳ less computatⁿ → less harm to environment.

Example:-

• $a = [1, 2, 4, 6, 7, 2, 3, 1, 4]$

to convert into ascending order.

• answer = [1]

• for i in range(len(a)): give size of list.

$small = a[i]$

for x in a :

⇒ Please use CHATGPT :)

→ Time complexity: ^(T.C.) → time to run a code for large iterations.

$O()$ → Big O notation → to represent TC.

eg:-

```

for i in range(n):
    for j in range(n):
        i + j;
    
```

T.C. = $O(n^2)$

NLPK chapter 5:-

→ Part of Speech (POS) tagging

→ Given a sentence / text → to find POS tag for every word
i.e. grammatical position eg. noun, verb, adjective etc.

→ Used to find position and context of each word in whole text.

→ Try out @ ~~corenlp~~ corenlp.hun

→ In nltk → ① word_tokenize(text)
② pos_tag(text)

⇒ Tagged corpora:-

↳ tuple consisting of {text: tag}

→ Taggers:-

→ Default tagging:- Tagging everything same (eg:- noun) before training.

⇒ Regular exp. tagger

ing	→	gerund
ed	→	past
es	→	present
auld	→	modal
?		

→ Not ML yet ;c

→ Look up tagger

→ From dataset find freq. of tags
eg: Noun : 60%
verb : 20%
etc.

→ Now, check for all words in your text and compare with dataset

Dataset text
Bat → noun Bat → 'noun'

→ Larger dataset → more accuracy

→ But this method → not much useful

→ We use 'Bigrams': - Bigram Tagger
i.e. noun is followed by a verb.

→ take two words → tagging them