27/Aug                                    Project Gutenberg

1830's → Slavery → citizen → spike
Bhashini project → under Ind Govt.
→ corpus → english & corpura - Latin.

30/Aug:-   Sept 28 → midterm presentation
          Midterm Quiz → Reasoning Questions.
                            No coding
                            Conceptual Technological Questions
                                 ( HL part )

          (Open Laptop/Notes)
          (Tut before midterm)
          ( duration → 1+ hr )
          ( handwritten )

Papers:-   (4, 5, 6, 7) in document
                (Set of slides) - AI snake oil
                                        (Fraud)

                        claim itself      dataset based
                        is fraud          on fradulent

NLTK - chapter 3:-    urllib , utf-8
                            decoder
   for html → proser ⟹ Beautiful soup.              f.open
* Reading local files        → 'w' (to write) → f. write (" - ")
                                  (override)        f. close
                          → 'a' → append

* NLP Pipeline.

→ fxn:- wordpunct_tokenize } takes care of punctuations also.

* strings:-

* regex:- import re

↳ nltk.corpus.words.words } list of words

(ending → ing $ )    [ ] → one of set match
 begin → ^ win       a + } → ≥ 1 of a
                     a* } → ≥ 0 of a → Klune closure.

Normalizing text ——→ everything lower case

Stemming → finding smallest possible word

Lemmatization ——→ process of stemming.

regex to remove various suffixes of a word.

Segmentation ✓  Tokenisation / word segmentation.

(EX:- chinese → no spaces)

Ch_4:-

Generator Expr → X

→ Questions of Style:- ⟨ procedural → lot of functions
                        ⟨ declarative → lot of variables & iterations.

→ variable scope  &  parameter type.
                   ( X ) → checking not allowed.

Program devp:

↳ Error management → generally no coding discipline
                      like C++ / Java
                              ↑
                             gdb

Algorithmic design:-

↳ (less resources / time)
        ↓                ↓
  (grandious)        faster → small

        (mergesort vs bubblesort)

Ch-5:- → Is task of ML in lang :- POS Tagger
                                    classifier learning task.

→ (position & context)  ⇆  other
   of 1 word              words   words- nltk.word_tokenize ("sentence")
                                   nltk.pos_tag (words).

                                                          context
Tagged corpora:-
POS Tagger ——→ (Default Tagger → "Noun")              (A^s) ?
              { Regex Tagger ——→ ting" → gerunds...etc)
| (a followed by b) }{ lookup Tagger → A's model size } → (acy. → 94%)
                                                           accuracy.

bigram → a followed by b

unigram tagger ≅ lookup tagger.

→ ( Using NLTK, make a Pos Tagger" ) ⟶