... of any point of time (depends how many communities we want)

→ SNA can also be based on Context Analysis. (till now → structural analysis)

→ Ex: can we infer something from a coll" of tweets

→ Here, we're looking at English level. But it can be extended to other languages

→ NLP :-

POS Tagging [Part of Speech]   → small part of NLP

→ To find meaning of context : we' we to look at words, coll" of words → sentences, coll" of sentences and so on.

→ Looking at word → Noun / pronoun / Adjective /..

⤷ Have to classify in automated version

φ° Most algo are statistical in nature with accuracy ~ 97%

→ How to input the context & issues faced while doing POS Tagging & Context Analysis.

φ° In 100 B.C., they defined 8 classes :

1. Noun
2. Verb
3. Pronoun
4. Prepositions          Basis of most of
5. Adverb          → European language
6. Conjunction
7. Participle
8. Article.

Objective of this : to understand word, we need to know its class.

→ Now, we have almost 45 tags (these 8 classes are further classified)

Jagging Algos

* 1. Rule based : People have read article & labelled the words manually & classify accordingly. They've made their own rules.

Dictionary → Tells class of word.

2. How you treat a word in that sentence is decided by set of rules (disambiguate)

○ No learning : learning is done by human being

✓ 2. statistical based :

Have stats of how this word is used previously.

Ex. A student

↳ most probably an article — determiner → ① telling no.
Stats: This classifier will tell it to be a determiner

● learn from past history in a statistical way & disambiguate words accordingly.

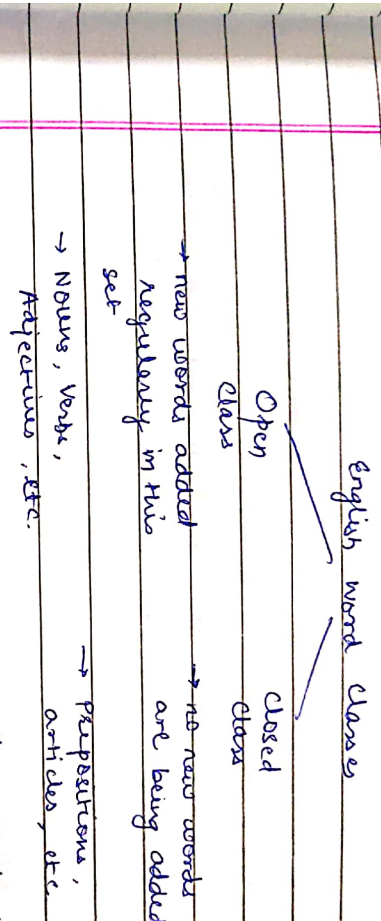1.) → Not very popular (we've very good ML models) (in 2)

2.) → we've HMM

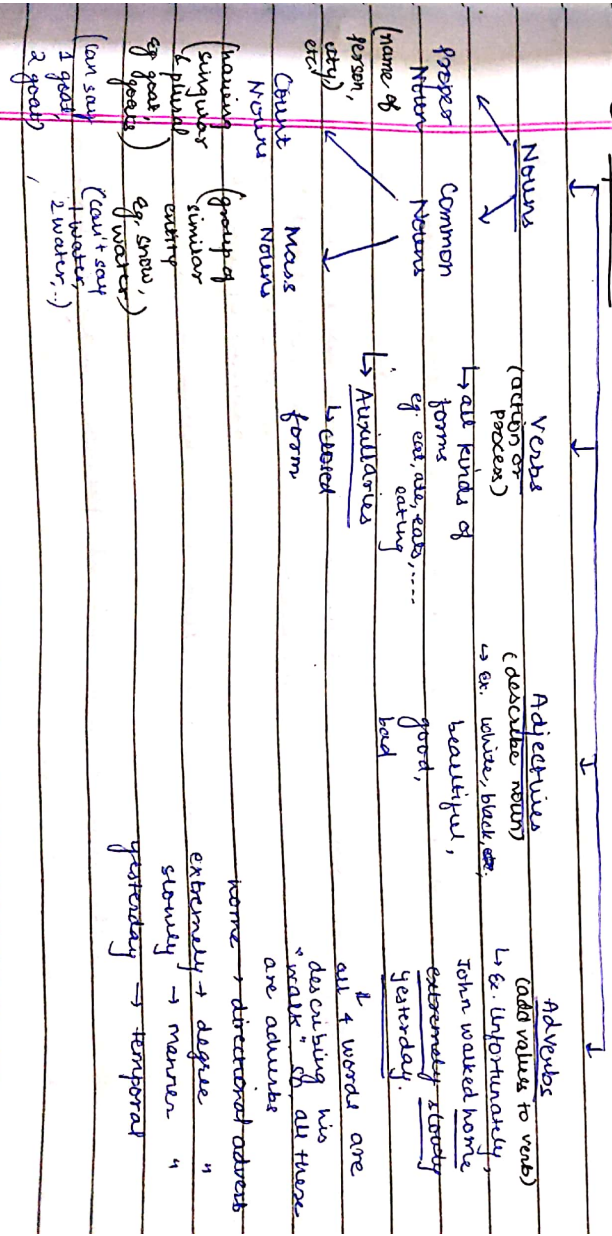He'll be looking at Rule-based approach.

→ there ... language in which we don't have any adjectives

Challenges behind building Rule based approach

→ Need to know English structure

English word class
- Open Class
- Closed Class

**Open Class**
→ new words added regularly in this set
→ Nouns, Verbs, Adjectives, etc.

**Closed Class**
→ no new words are being added
→ Prepositions, articles, etc.

(1) Open Class → some of class if are made to understand contexts
→ effeciently (contrib. from English language)

Nouns

Verbs
(action or process)
↳ all kinds of forms
eg. eat, ate, eats, eating ...
↳ Auxiliaries
↳ ebed
bad

Adjectives
(describe noun)
→ eg. white, black etc.
beautiful,
good,

Adverbs
(add values to verb)
↳ eg: Unfortunately, John walked home, extremely slowly yesterday.

all 4 words are describing his "walk" go all these are adverbs

home → directional adverb
extremely → degree "
slowly → manner "
yesterday → temporal

**Nouns**
- Proper Noun (name of person, city, etc.)
- Common Noun
  - Count Noun (having singular & plural entry) eg goat 1 goat, 2 goats
  - Mass Noun (group of similar entity) eg snow, water (can't say 1 water, 2 water...)

④ __Closed Classes__

1) Prepositions — on, in, ∅, at, by

2) Determiners — Articles determining something
   a, an, the, this, that

3) Pronouns — he, she, who, I

4) Conjunctive — and, but, or (join 2 sentences)

5) Auxillary verbs — can, may, should

* 6) Particles — up, down, on, off, in, out (give a little bit
                                              of ambiguity)

7) Numerals — one, two, first, second

__Note__

→ Preposition always occur before nouns

→ Particles are used in combin" with verbs

   Ex. turned down ──→ reject
                   meaning
        ↓        ↓
       verb    particle

   but when given to parser → verb + particle → should be
                                                considered as
                                                1 meaning
                                                word)

   Ex. rule out → eliminate  ⎫ also. has to look at
                               ⎬ both words (1, prev.
       find out → discover  ⎭  word is verb → find
                               meaning of comb" of words)

Ex, Rule
based.  → Whenever we see down, look prev. words & find meaning
          accordingly
          ( In stats based → it'll look at previous usage of word &
            classify accordingly )

          → Conjunctions
             ↓                    ↓
          Co-ordinating        sub-ordinating
                                  ↓
          → and, or             → that
            joining 2 sentences    ↓ verb
                                Ex, I thought that this course
                                will be easy
                                'that' is telling about the verb.

→ Pronouns

Personal | Possessive | wh - pronouns

→ He, she, I, me

→ mine, yours, his, her → who, what, whom, whoever

→ Auxillary Verbs : be, do, can, have, had, are

→ some other set of words :

* Interjection : oh, ah, uh, alas

* Negation : no, not

* Politeness marker : please, thankyou

→ There've various algo available to label. we'll see → Rule based algo

→ People have classified set of Tags for POS Tags [to maintain std.]

A popular tag is,

Tag sets :

(i) 45 - tag - Penn Tree Bank tag set (1993) → more popular

(2) 61 - tag C 5 Tagset (1997)

→ found from 87 - tag Brown Corpus (1960-70), done in Brown University

→ It's observed only 14% of words are ambiguous but they come most frequently in sentences.

So, they give 14% ambiguity in sentences.

Q) why reduced 87 → 45 ?

Algo is more efficient when it works with less no. of tags.

Wed → Speech & Language Processing
by Jurafsky & Martin → 2nd edition
(5th chapter → 1st & 2 sub-section)

→ In some analysis, we can't simply remove '$', '(', ')', etc
" " ") So, we've tags for these also.

↳ Have to check whether these add values or lessen values
of contexts

Rule based approach :

Eng CG → One of the algos
↓
constructed
grammar

lexicon → dictionary, classifying
step 1 → Finding all possibilities of all the words
step 2 → applies large set of constraints to input sentence to rule out
incorrect POS (from list we obtained earlier)

End term →

Before mid-sem : 33%
After           :  66%    Today's → very less

15 minutes video      Report → 28th
                      Video → 10th may