

merged \longleftrightarrow generation

cat
N
re

Cats

→ Surface form is lexical form.

Step 1: Identify the morphemes

Step 2: Analyse morphemes.

Orthography changes

ಪೆಡೆ (from ಪಾಯ)

ಪೆಡೆ (from ಪಾಯ)

exceptions: ಪಾಯ \rightarrow ಪಾಯ್
~~ಪೆಡೆ~~

Refer:
 Regularity, analogy
 in histo-logy

Machine learning of morph rules

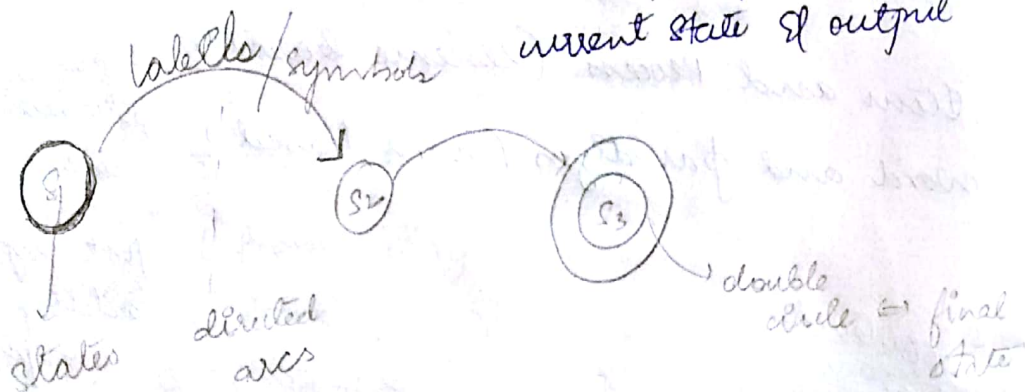
Based on large data, orthographic, suppletion, assimilation, vowel harmony etc are learnt, but we need to provide a list of all exceptions which are not regular.

CELEX \rightarrow lexical database.

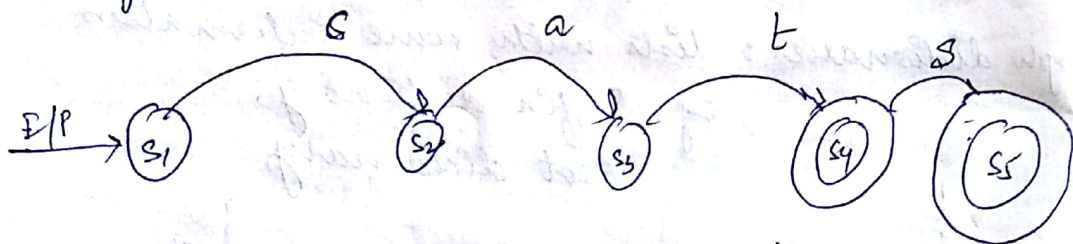
lemma, wordform and other
 (headword) details

FSA: Machine consisting of

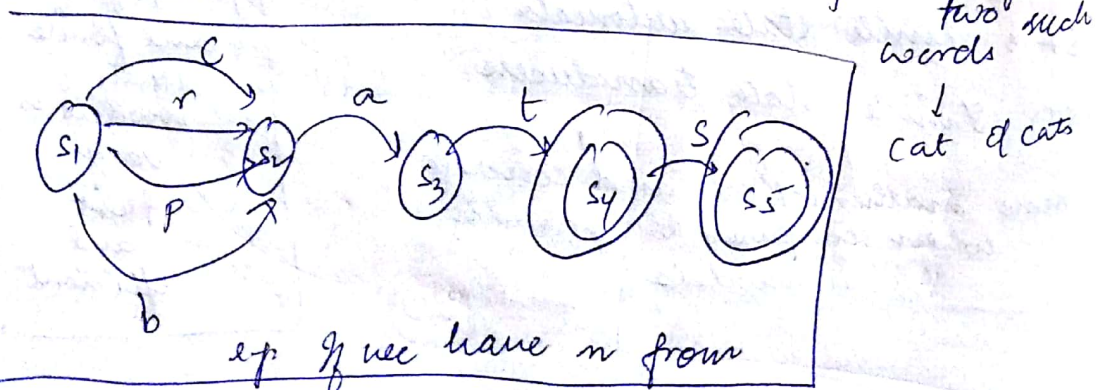
- ① Input tape
- ② Finite number of states, with one initial and one or more accepting states
- ③ Actions in terms of transition from one ~~expression~~ state to another, depending on current state & output



e.g. CATS.



To extend, we do this



new, it goes S_1 to S_2 and ~~also~~ check for new as root word.

Shallow Parsing (Chunking and LWG)

→ why morph analysis?

gives word group and the headword's relⁿ with other words.

Two levels of information stored in a sentence:

- ① Lexical level
- ② Grammar (relative arrangement of words which gives some meaning not explicit in any of the words)

→ we need morph analysis for languages which store a lot of information lexically.

e.g. Dravidian languages.

Whereas for English, we don't require morph analysis, if the main goal is just parsing.

e.g. Mohan is leaving for class now.

మోహన్ క్లాస్ కు వెళ్తున్నాడు.

→ No one-one correspondence between form and meaning: leads to ambiguity

→ In some cases, where we need to process a sentence (a little more work than just parsing) we would need morph analysis to extract features such as the roots, suffixes, etc.

e.g. Mohan is leaving for class now.

Mohan is leaving for class now.
Mohan is leaving for class now.

आ रक्ष ३ - Analytical language.
VM VAVX VADY

संज्ञासूचक - Synthetic / Fusional
+ ~~100%~~ inflections

Agglutinative - just concatenates morphemes & words

Synthetic / Fusional - changes occurs at morpheme boundaries (e.g. phonological changes such as vowel harmony)

२५
(राम ने) (गौहन को) (किताब दी) | → Tense, aspect is occurring separately (not at lexical level)

राम के भाई ने
acc to motion

"pada" → Since nouns and verbs have to be inflected before they can make a sentence together (either nominal or finite), each inflected form is a pada.
(Sanskrit)
Singular masculine

THE
PROCESS IS
CALLED
LOCAL WORD
GROUPING

- POS tagging, morph analysis and chunking are usually the three main steps followed prior to parsing. 27/9/18

eg. (My friend) (visited) (the new shop) (in the outskirts)
 (of our city)

$\begin{matrix} \text{NP} & & \text{VP} & & \text{NP} & & \text{NP} \\ \text{NN} & & \text{VBD} & & \text{NN} & & \text{PP} \\ \text{NP} & & & & & & \\ \text{NN} & & & & & & \end{matrix}$

Chunk Tags for Indian languages

NP	
JTP	adjective chunk
RBP	adverb chunk
NEGP	chunk for negatives (eg: बिना)
CCP	conjunction chunks (और is an eg)
BLK	miscellaneous

↳ words which can't be grouped under any head.

- chunks have to be kept as flat as possible.

WX notation

e e E(ee) → Telugu
 క క క

Chunk
 minimal
 phrase
 with
 internal
 dependencies
 not
 distorted

eg. (Ram) (aur) (Shyam) → the doubt of which is the head remains
 (Ram aur wah ladha)

Word compounding and Multiword Expressions 4/10/19

(MWE)

- Different conventions are followed.

eg: white spaces, hyphens or nothing

black board, blackboard, black-board.

main elementary problems in NLP:

- ① Multiword Expressions (MWE)
- ② Word Sense Disambiguation (WSD)
- ③ Named Entity Recognition (NER)

main examples of MWE:

- ① Compounds (Nominal & Verbal)
- ② Named Entities.
- ③ Idioms
- ④ (Some common prepositional /
postpositional / adverbial phrases) ?

NST- A word which is classified
as a noun, but could be
an adverb or something

eg: NSTs.

मेज के ऊपर (vs) मेज के पर
|
NST

Compounds

Nominal

Table cover
Tube light

Verbal

गिर पड़ना
कह देना
ले आना

↓
Noun
-Verb

बैठ जाऊँ कहाँ

Named

→ proper
e.g. the

(Sentence)
Speed

• Identifying compounds computationally?

• MWE - a multiword unit or a collocation of words that occur together statistically more than chance



co-occurrence is one of the main criteria

e.g. जल व पानी are synonyms,

but गंगाजल is often used whereas गंगा पानी is not!

so गंगाजल, strong coffee etc. are compounds (multiword expressions)

Compositional

vs Non Compositional

black bag

white ants

(they're termites. Neither white nor ants)

kick the bucket

(neither kicking or bucket)

* There is never a binary division of compositional and non compositional

It could be gradual. Even if binary, depends on aspect

Named Entities

→ Proper nouns

eg How do you POS tag NIT?

(International Institute of Information Technology)
NP NP NP NP