

Linguistic terminologies

~~False~~ Morphology (structured) (word)

↓
Syntax (phrase/sentence)

↓
Semantics (meaning)

↓
Pragmatics (contextual ^{meaning} information)

↓
Discourse (higher than the paragraph level)

- Turing test - 3 actors (2 humans & 1 machine)

- Eliza

False positive (Type 1) - Precision / Accuracy
False negative (Type 2) - Recall / Coverage

- Antagonistic efforts

1. minimizing FP
2. minimizing FN

So, we measure F1-score (HM of Precision & Recall)

Every word is token but every token is not a word
San Francisco - 2 tokens, but 1 word

→ UNIX for poets

Tokenization - 1 word per line

state-of-the-art → 1 word

Edit Distance

Minimum edit distance

I	N	T	E	*	N	T	I	O	N
*	E	X	E	C	U	T	I	O	N
1	2	2	-	1	2	-	-	-	-
(Deletion)	(Substitution)			(Insertion)					

Total cost = 8 units

Alignment in Computational Biology

Named entity extraction - extract proper nouns from text

How to find Min Edit Distance?

Word type

Tokens

Whatever that goes into vocabulary

- all the occurrences of a word in a text

Minimum Edit as Search

for 2 strings

X of length n

Y of length m

We define $D(i, j)$

Dynamic Programming

→ Tabular computation of $D(n, m)$

→ solving problems by combining solutions to subproblems

Bottom - up

→ Compute $D(i, j)$ for small i, j

→ And compute larger $D(i, j)$

Levenshtein

→ Initialization

$$D(i, 0) = i$$

$$D(0, j) = j$$

→ Recurrence Relation

for each $i = 1 \dots M$

for each $j = 1 \dots N$

$$D(i, j) = \min D(i-1, j) + 1$$

$$D(i, j-1) + 1$$

$$D(i-1, j-1) + 2;$$

$$0;$$

$$x(i) \neq y(j)$$

$$x(i) = y(j)$$

Page 107

Assignment - 1 Deadline : Friday

1) Fill in the table with min edit distance values + back trace pointers (arrows)

2) Explain the process (also back trace pointers.)

Performance

Time $O(nm)$

Space $O(nm)$

Back trace $O(n+m)$

Weighted Edit Distance

→ Some words are more likely to be misspelled.

Unicode and Multilingual computing

Software Localization

Locale : Date
Currency
Days of a week

Internationalization

Localization

Globalization

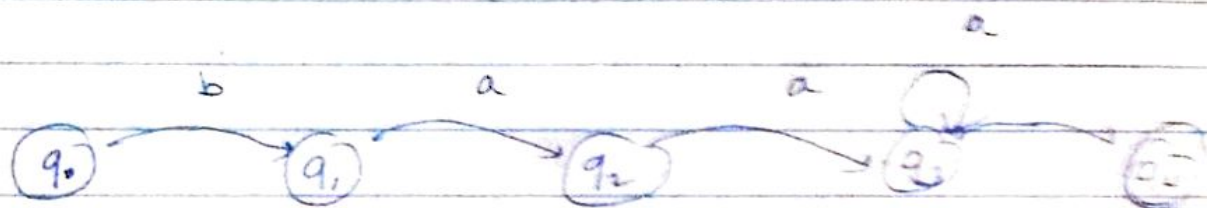
Text processing not possible - Breeti, Kantipur, Annapurna

Glyph - graphical representation of a character

Nepali NLP

Information Retrieval
Extraction

Finite State Automata



DFSA

determined path for an input

NFSA

multiple paths for an input

Morphology and Finite State Transducers

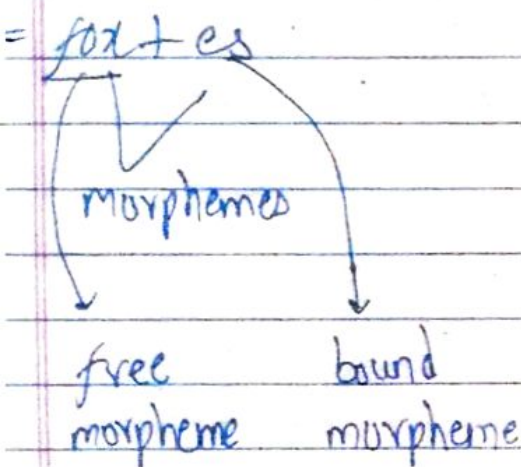
Study of word forms

Surface form cat

Lexical form cat

Smallest form of a word is morpheme

foxes



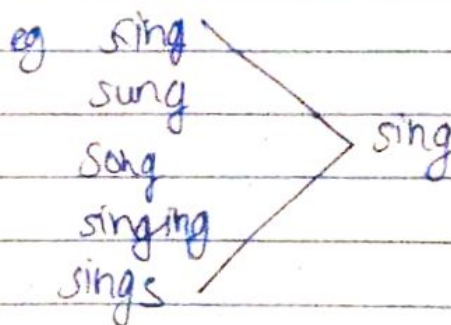
Morphological parsing

- break down a particular word into free morpheme and bound morpheme

vs stemming

- only concerned about the stem (not about others)
- breaking down a particular word into its corresponding stem or root word

foxes
→ fox + es



- IR
- Machine translation
- Spell checking

stem - free/main morpheme

affixes - bound morpheme

prefix, suffix, infix, circumfix

1. Inflectional morphology

N
foxes
↓
foxes
N

→ Noun to noun
→ no change in class / parts of speech

2. Derivational morphology

→ result in different class

V
compute + ing
↓
computing (N)

→ eg. Verb to noun

GNP - (Nepali)
Gender Number Person

English - NP only
He eats
She eats

HMM Hidden Markov Model POS tagging

maximize $P(\hat{r}_1 | \hat{w}_1)$

Window size

NN VB

Sequence Labelling -

POS tagging approaches

Classification learning

and
most frequent
word in English
Language

Chapter 8 Word Classes and Part of Speech Tagging

↳ lowest level of syntactic analysis

Parts of speech (word class / morphological class
lexical tag)

English Word Classes

1. Closed class
- ~~not~~ relatively fixed membership eg. prepositions
2. Open class
noun, verb, adjective, adverbs

Parts of speech tagging

- 1) Rule-based taggers
- 2) Stochastic taggers
- used a training corpus to count the probability

Brill tagger - hybrid of both