

Chapter 9 Context-Free Grammars for English

~~Syntax~~ $S \rightarrow NP VP$

Syntax: constituents

English - SVO

Nepali - SOV

Transitive verbs - direct objects. eg I played football.

Intransitive verbs : eg I cried.
(no direct objects) I laughed

Nepali language - Dependency grammar followed instead of CFG.
Panini Grammar

CFG also called PSG (Phrase-Structured Grammar)

→ Grammar for L0

$S \rightarrow NP VP$

$NP \rightarrow \{ \text{Pronoun} \mid \text{Proper-Noun} \mid \text{Det Nominal} \}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Nominal} \mid \text{Inoun}$

$VP \rightarrow \text{Verb} \mid \text{Verb } NP$

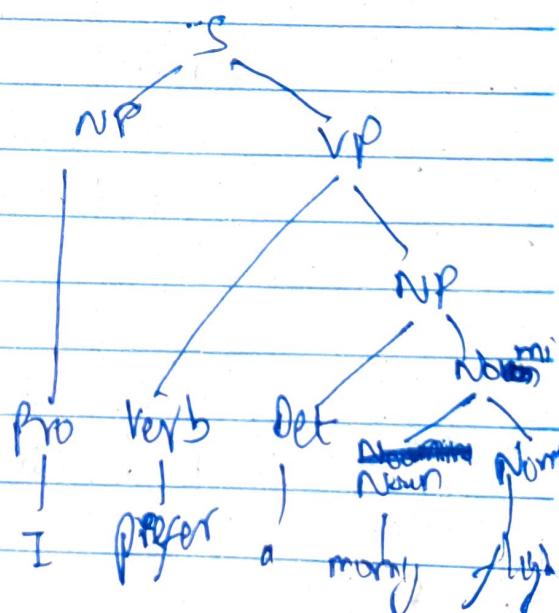
$\text{Verb } NP$

$\mid \text{Verb } NP PP$

$\mid \text{Verb } PP$

$PP \rightarrow \text{Preposition } NP$

I prefer a morning flight



6 < 5
30

- Bracket notation of parse tree

Sentence level construction

1. Imperative structure $\rightarrow S \rightarrow VP$ (No subject) eg Show the lowest fare
 2. Declarative structure - Subject NP + VP $\rightarrow NP VP$ (Subject-NP)
 3. Yes-no question structure: $S \rightarrow Aux NP VP$ (Subject-NP)
 4. Wh-question structure: $S \rightarrow wh\text{-}NP VP$ eg What airlines fly from Nepal to Thailand?
- eg Do any of these flights have stops? - Yes-no ques

Noun phrase

- head - central noun in NP
- eg a boring lecture
- prenominal and postnominal ^{head} modifiers
- determiners (optional)
- Mass nouns / Material nouns - no determiners
- pre determiners and post determiners
- all the flights



Nepali - no preposition, only postposition ~~except~~

Agreement - not confined in CFG anymore
- adds context

English - number agreement only
French / German / Nepali - ~~Gender~~ and number agreement
(pronouns as well)

1st

Gill POS tagging

PAGE NO.:
DATE: 28 April

- Parsing

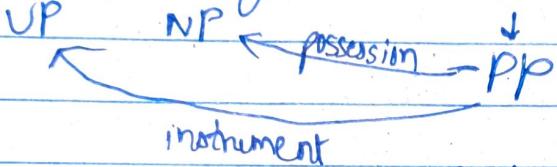
Bottom-up parsing
- CKY algorithm

Top-down
- Earley algorithm

Chart parser

→ mix of top-down and bottom-up

→ I saw the boy with the telescope



Ambiguity

→ I shot the elephant in my pajamas.

Problems with top-down parser - 1) left recursion

2) Attachment ambiguity

Coordination ambiguity

eg {old [men and women]}

{old men (and women)}



Farkley algorithm

- Predictor
- Scanner
- Completer

$O(n^3)$ n - no of words in the input

CKY Parsing (Bottom up)

CNF grammar used i.e.

$A \rightarrow BC$ or 2 non terminals or
 $A \rightarrow w$ single terminal

Conversion to CNF

$S \rightarrow ABC$

$S \rightarrow XBC$

$X \rightarrow AB$

Operations (CKY)

1. Recognition
2. Parsing

→ POS tagging in Nepali texts
AB tags

→ Penn Treebank tags etc - popular
112 forms of single verb

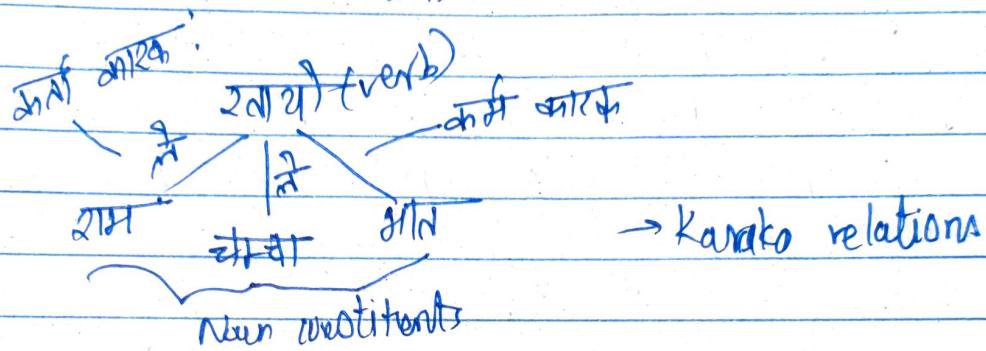
→ Chunker, grammar analyzer, POS tagger - Nepali

Grammar analyzer

- Dependency grammar formalism
- free word order language - Nepali

रामले भात रामा)

भात रामले रामा)



Statistical Parsing

Probabilistic / Stochastic CFG (PCFG).

- 4 parameters (N, Σ, R, S)

N - non terminal symbols

Σ - set of terminal symbols disjoint from N

R - set of rules or productions, ^{each} of the form

A \rightarrow B [p]; diff from CFG

A - non terminal

B - string of symbols from a set of strings ($\Sigma \cup N$)^{*}

p no betw 0 & 1 expressing $P(B|A)$

S - start symbol

$$p = P(A \rightarrow B)$$

$$= P(A \rightarrow B | A)$$

$$= P(RHS | LHS)$$

consider all possible expansions of a non-terminal

$$\sum p(A \rightarrow B) = 1$$

Prob of parse tree T is

$$P(T, S) = \prod_{i=1}^n P(RHS_i | LHS_i)$$

$$P(T, S) = P(T) P(S|T) = P(T) \quad (\because P(S|T) = 1)$$

Prob of sentence = sum of prob of parse trees

→ Most probable derivation (parse tree) for a sentence

Yield

Consider all possible parse trees for a sentence s .

The string of words s - yield of parse tree over s .

Disambiguation algo finds the parse tree which ~~is most~~ is most probable given s

$$\hat{T}(s) = \underset{T \text{ s.t. } S = \text{yield}(T)}{\operatorname{argmax}} P(T|S)$$

$$\hat{T}(s) = \underset{T \text{ s.t. } S = \text{yield}(T)}{\operatorname{argmax}} \frac{P(T, S)}{P(S)}$$

≈ approximate

$$\hat{T}(s) = \underset{T \text{ s.t. } S = \text{yield}(T)}{\operatorname{argmax}} P(T)$$

Prob of next word w_i in a sentence given all words, ^{seen} so far

$$P(w_i | w_1, w_2, \dots, w_{i-1}) =$$

10-9-10 12-3 10-

PAGE NO. :
DATE:

- Probabilistic CKY
- learning PCFG Rule Prob
- Head words

ground truth
gold standard

Unification grammar / Feature - Structure Grammar

Semantic Level

Chapter 14 Representing Meaning

Lexical semantics - slide

meaning conveyed by a word

act

→ wound

→ share

...

:

Natural language understanding (NLU) - current focus
Natural language generation (NLG)

JHSN JHN

world knowledge - KU is located in Dhulikhel

C

Lexical Semantics

WordNet - largest lexical database

- contains synonyms set (synset)

Relation among lexemes & their senses

1. Homonym - lexemes that share a form. (phonological, orthographic)
eg Bat - wooden same word diff meaning
Bat - animal

Homophones - sound same homograph

It
Eat

2. Synonym

3. Polysemy - different senses of same word
has floppy ears
has a good ear for jokes.

4. Hyponym / Hypernym

Lemma (stem / root - interchangeable

L I F E N E

OKAY?
OKAY!

Vector Space (Distributional) Lexical Semantics

Information Retrieval System

The Vector Space Model

Graphic Representation

Term Weights : Term frequency

TF-IDF Weighting

Term frequency

Inverse Document frequency

Similarity Measure

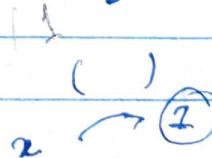
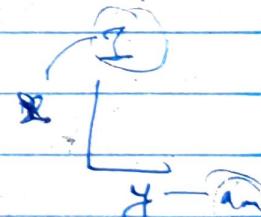
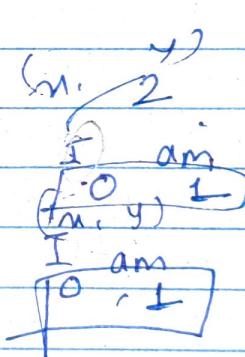
Cosine similarity

Gold standard

Ground truth data

Bench-marked datasets

SemEval competition - NLP popular



Ques

Discourse (highest level of text analysis)

?
Pragmatics (contextual meaning) ?

?
Semantics (word level)

?
Syntax

?
Morphology

} Document

→ Pragmatics
context

- i) Linguistic context (cotext)
- ii) Physical context

acts-locutionary, illocutionary, ~~perlocutionary~~, perlocutionary

Discourse and Discourse analysis

13th May

LAA

000 to 025
delivery
diagnosis
Sequence
by
Sequence

Discourse

Discourse and Discourse Analysis

To make logical
decisions for

OK

RESULT

Text linguistics and discourse analysis

Discourse - both text and speech

Information Retrieval - return a complete document based on query.

Information Extraction - extract part of information about a document eg model name.

Reference Resolution

Anaphora Resolution

Co-reference Resolution

Re Referents \longleftrightarrow Named-entity

Referring expressions.

Cataphora vs Anaphora

He . . . John . . . - Cataphora
 ↑
 referent

* - not well formed sentence. (wrong sentence)

Recently

Salient \rightarrow more important / pertinent

Sentiment Analysis

16th May

- Opinion mining
- Polarity :- +ve, -ve, neutral

Product Reviews

- Movie
- Book
- Any other product

Opinion extraction / mining

Sentiment mining

Subjectivity analysis



sentiment analysis

subject to change

Google Trends

- Polarity detection / Classification
- Degree / Intensity of

Baseline Algorithm

- Tokenization
- Feature extraction - (Adjective, Adverb) - Potential features
- Classification using (Naive Bayes, MaxEnt, SVM)

→ Sentiment bearing words
may be misleading

So, all words need to be considered
eg * didn't like this movie
liked this movie

Naive Bayes (Boolean Multinomial)

- For sentiment analysis,
word occurrence more important than word frequency.
eg worst movie - -ve
good . . good . . worst - -ve

So, duplicate words removed

Why reviews hard to classify?

- Implicit opinions - (not expressed explicitly)
^{opinions}

MPC A Subjectivity Cues Lexicon
Bing Liu Opinion lexicon
Sentimentnet

too [adj]

too bad - -ve

too good + +ve

too good to be true - -ve (didn't work)

- 1 I bought a mobile phone last week.
- 2 It fits in my pocket. : Size : small
- 3 It is affordable. : Price : affordable
- 4 It is cool. : Texture : cool
- 5 Battery life is just a few hrs. Battery life : short

Finding sentiment of a sent

- important for finding aspects or attributes.
eg size in sent 2 : size = small
- sent 3: price

→ Detection of friendliness

Information Retrieval search engine (google)

Precision: fraction of retrieved documents that are relevant to user's information need

$$= \frac{TP}{TP + FP}$$

Recall: fraction of relevant documents in collection that are retrieved.

$$\frac{TP}{TP + FN}$$

Rare words more important in collection of documents

Frequent words more important inside a document.

- Boolean search
- Rank retrieval search

Information extraction and Named Entity Recognition

- POS tagging - 1st step
- Phrase level (Named entity tagging)

Unstructured text → structured info

Question Answering

Gold leveled - true

What is the capital of Nepal?

Kathmandu - gold leveled answer.

2nd Text summarization