

CL Seminar 1

Pattern-Based Context-Free Grammars for Machine Translation
by Koichi Takeda

Abhinav S Menon

Overview

- ▶ Context-Free Grammars
- ▶ Pattern-Based CFGs
- ▶ The Translation Algorithm
- ▶ Features and Agreement
- ▶ An Example

Context-Free Grammars

Context-Free Grammars (CFGs) are a common formalism for describing the syntax of natural (and programming) languages. For example,

$$S \rightarrow NP VP$$
$$NP \rightarrow A N$$
$$VP \rightarrow V NP$$
$$N \rightarrow \text{cat} \mid \text{dog} \mid \text{mouse}$$
$$A \rightarrow \text{the} \mid \text{a}$$
$$V \rightarrow \text{chased} \mid \text{hugged}$$

Pattern-Based CFGs: Patterns

Pattern-Based CFGs or PCFGs consist of a set of patterns – analogous to the set of rules that make up a CFG.

Each pattern describes how to translate a certain type of phrase from the source to the target language.

An example pattern is

$$\begin{aligned} \text{NP}_1 [\text{miss}] \text{V}_2 \text{NP}_3 &\leftarrow \text{S}_2 \\ \text{S}_2 &\rightarrow \text{NP}_3 [\text{manquer}] \text{V}_2 \text{à NP}_1 \end{aligned}$$

Pattern-Based CFGs: Skeletons

The two rules (one in the source language and one in the target language) that make up the pattern are called its *skeleton*.

For example, the skeleton of the example above

$$\begin{aligned} \text{NP}_1 [\text{miss}] \text{V}_2 \text{NP}_3 &\leftarrow \text{S}_2 \\ \text{S}_2 &\rightarrow \text{NP}_3 [\text{manquer}] \text{V}_2 \text{à NP}_1 \end{aligned}$$

is

$$\begin{aligned} \text{NP V NP} &\leftarrow \text{S} \\ \text{S} &\rightarrow \text{NP V à NP} \end{aligned}$$

Pattern-Based CFGs: Constraints

$$\text{NP}_1 [\text{miss}] \text{V}_2 \text{NP}_3 \leftarrow \text{S}_2$$
$$\text{S}_2 \rightarrow \text{NP}_3 [\text{manquer}] \text{V}_2 \text{à NP}_1$$

There are two types of constraints in the above pattern:

- ▶ Link constraints (subscripts): specifying the correlation between the constituents of each part.
- ▶ Head constraints (enclosed in []): specifying lexeme- or phrase-level conditions on the applicability of the pattern.

The latter makes it possible to translate phrases that are constructed differently in the target language, as in the above example.

The Translation Algorithm: Essentials

The basic algorithm consists of three simple steps:

- ▶ Parse the input using the source language CFG skeletons.
- ▶ Propagate link and head constraints from the source to the target language and build a target language derivation sequence.
- ▶ Generate the output using the target sequence.

The Translation Algorithm: Parsing

Parsing of CFGs is a well-studied area. Many algorithms are known for this task – we will consider the one described in Earley (1970).

This algorithm consists of stepping through the input string one at a time, and trying to parse the upcoming symbols given the incomplete parses until the current point.

It generates all possible derivations for the input string (sometimes called a *parse forest*) in the given CFG.

The Translation Algorithm: Priorities

Under some conditions, it is very easy for the number of target skeletons to increase exponentially for the same number of source skeletons.

Therefore, when two or more patterns share the same source skeleton, we must decide which pattern is more likely to yield a correct translation.

The Translation Algorithm: Priorities

The priority order is as follows:

- ▶ Prefer a pattern with more head constraints (a more *specific* pattern).
- ▶ Prefer a pattern with more terminal symbols.
- ▶ Prefer the shortest derivation sequence.

These preferences can be expressed as numeric *costs* for patterns.

Features and Agreement

The constraints can be extended to lay restrictions on the features of words or phrases as well, to increase its descriptive power.

For example, in French, pronouns as direct objects come *before* the verb. This can be expressed as

$$\begin{aligned} V_1 NP_2 &\leftarrow VP_1\{+OBJ\} \\ VP_1\{+OBJ\} &\rightarrow NP_2\{+PRO\} V_1 \end{aligned}$$

Integration of Bilingual Corpora

Given a bilingual corpus, we can translate all the sentences and compare them.

If the translation is correct, nothing needs to be done.

If there was a derivation sequence that resulted in an incorrect translation, use the vocabulary as head constraints.

If there was no paired derivation sequence, add specific patterns for idioms or collocations (if any).

An Example

$$\text{NP}_1 \text{ VP}_1 \leftarrow S_1$$
$$S_1 \rightarrow \text{NP}_1 \text{ VP}_1$$
$$\text{VP}_1 \text{ ADV}_2 \leftarrow \text{VP}_1$$
$$\text{VP}_1 \rightarrow \text{VP}_1 \text{ ADV}_2$$
$$[\text{know}] \text{VP}_1 \{+ \text{OBJ}\} \text{ well} \leftarrow \text{VP}_1$$
$$\text{VP}_1 \rightarrow [\text{connaître}] \text{VP}_1 \{+ \text{OBJ}\} \text{ bien}$$
$$\text{V}_1 \text{ NP}_2 \leftarrow \text{VP}_1 \{+ \text{OBJ}\}$$
$$\text{VP}_1 \{+ \text{OBJ}\} \rightarrow \text{V}_1 \text{ NP}_2 \{- \text{PRO}\}$$
$$\text{V}_1 \text{ NP}_2 \leftarrow \text{VP}_1 \{+ \text{OBJ}\}$$
$$\text{VP}_1 \{+ \text{OBJ}\} \rightarrow \text{NP}_2 \{+ \text{PRO}\} \text{ V}_1$$

An Example

he \leftarrow NP{+PRO + NOMI + 3RD + SG}
NP{+PRO + NOMI + 3RD + SG} \rightarrow il

me \leftarrow NP{+PRO + CAUS - 3RD + SG}
NP{+PRO + CAUS - 3RD + SG} \rightarrow me

knows \leftarrow V{+FIN + 3SG}
V{+FIN + 3SG} \rightarrow sait

knows \leftarrow V{+FIN + 3SG}
V{+FIN + 3SG} \rightarrow connait

well \leftarrow ADVP
ADVP \rightarrow bien

well \leftarrow ADVP
ADVP \rightarrow beaucoup

An Example

According to these rules, the sentence He knows me well is parsed.

Fragment	Translation
He knows me well	S
He	il
knows me well	VP
well	bien
knows me	VP
knows	connait
me	me

We get the correct output: Il me connait bien.

References

Takeda, Koichi. “Pattern-based context-free grammars for machine translation.” arXiv preprint [cmp-lg/9607011](https://arxiv.org/abs/cmp-lg/9607011) (1996).