

CSE/LIN 467/567: Extra assignment report

Jui Kate (50288204)

Vakul Bhatia (50290967)

Abstract

The spanish language has different dependency structure, grammar, syntax from English language. The spanish corpus in penn-treebank format was prepared to generate dependency structure. The dependency structure for English language was known and given. We found these structures for Spanish as well. All executable codes are available at

<https://github.com/JuiKate/Computational-Linguistics/tree/master/extra-assignment>

1 Introduction

We have taken Spanish language to work on. The objective of this project is to make required changes in java files of Stanford parser, build the project, create executable jar file, modifying generate-data-corrected.sh file to use the updated jar file to generate dependency structure for the given input Spanish penn-treebank file.

2 Converting rules

We made changes in the following java files of package edu.stanford.nlp.trees

1. **EnglishPatterns.java** – We identified words which are used in Spanish and which will be detected in the dependency conversion. Copular verbs like "ser", "estar" were added which are frequently and commonly used in Spanish. Time related words were added like the days of the week, months

of the year, relative time like

today/tomorrow/morning/evening, seasons of the year etc.

2. **SemanticHeadFinder.java** – We added header rules so that it identifies the head nodes and then the dependency structure is understood. We added head rules for tokens like CL/CONP/PRED. Without adding these rules the parser gives error since it cannot identify head rules for these tags.
3. **UniversalSemanticHeadFinder.java** – Rules used in SemanticHeadFinder.java were similarly added in this java file.
4. **EnglishGrammaticalRelations.java** – We made few changes in this file since Spanish has a different grammatical structure compared to English so few grammatical rules were updated as per Spanish grammar.

3 Parsing results

Sample dependency parsing results:

1	Zedillo	—	N	N	—	2	nsubj	—	—
2	asegura	—	V	V	—	0	root	—	—
3	en	—	PREP	PREP	—	2	prep	—	—
4	California	—	—	N	N	—	3	dep	—
5	que	—	C	C	—	13	dep	—	—
6	las	—	ART	ART	—	8	dep	—	—
7	reformas	—	—	N	N	—	8	dep	—
8	políticas	—	—	ADJ	ADJ	—	13	dep	—
9	en	—	PREP	PREP	—	8	prep	—	—
10	México	—	N	N	—	9	dep	—	—
11	no	—	ADV	ADV	—	13	advmod	—	—
12	tienen	—	V	V	—	13	dep	—	—
13	marcha_atrás	—	—	N	N	—	2	dep	—
14	.	—	PUNCT	PUNCT	—	2	dep	—	—
1	EEUU	—	N	N	—	5	nsubj	—	—
2	concederá	—	—	V	V	—	5	dep	—
3	la	—	ART	ART	—	5	dep	—	—
4	residencia	—	—	N	N	—	5	dep	—
5	legal	—	ADJ	ADJ	—	0	root	—	—
6	a	—	PREP	PREP	—	5	prep	—	—
7	medio	—	Q	Q	—	6	dep	—	—
8	millón	—	N	N	—	7	dep	—	—
9	de	—	PREP	PREP	—	7	prep	—	—
10	centroamericanos	—	—	N	N	—	—	9	dep
11	.	—	PUNCT	PUNCT	—	5	dep	—	—

4 Discussion and conclusion

Due to the difference between grammatical structure of English and Spanish languages respective changes are required to form dependency relation. We made the changes as per our understanding and more changes can be made for improving the quality of Spanish parser.

References

- de Marneffe, M.-C. and Manning, C. D. (2008). The Stanford Typed Dependencies Representation.
In
Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation
, pages 1–8, Manchester, UK. Coling 2008 Organizing Committee.
- Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., and McClosky, D. (2014). The Stanford CoreNLP Natural Language Processing Toolkit. In Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations
, pages 55–60, Baltimore, Maryland. Association for Computational Linguistics.