

Concept Tagging for the Movie Domain

Giovanni De Toni (197814)

University of Trento

Via Sommarive, 9, 38123 Povo, Trento TN

giovanni.detoni@studenti.unitn.it

Abstract

This work focuses on the well-known task of concept tagging phrase elements. This presents a relatively important challenge in the NLP field since it is the starting point for more complex applications. Weighted Finite State Transducer and statistical methods were used and an analysis of their performances is also provided.

1 Introduction

One of the first steps of some NLP task is to analyze a given phrase to understand its underlying concept. For instance, imagine we are using a vocal application to order something to eat (e.g., "I would like a tortel di patate delivered at Piazza Trento 1, please"). In order for a machine to understand correctly what we want, firstly it needs to convert our utterances to a word representation, secondly it needs to assign a concept to each of the words it heard (e.g., "tortel di patate" equals to requested food and "Piazza Trento 1" equals to the delivery address). This basic operation is of utmost importance for all the applications which can be built upon this. For instance, the quality of your food assistant is dependent on the quality of this concept tagger. Imagine what could happen if the machine were to swap the delivery address with the requested food! The scope of this project is to provide a simple concept-tagger by developing a WSTF (Weighted Finite State Transducer) applied to the movie domain. The report is detailed as follows: in the first section we define more formally the problem statement, then we proceed with an analysis of the given dataset and with the description of the models employed. Ultimately, we discuss the obtained results while underlying their strengths and limitations.

2 Problem Statement

Given a sequence of tokens $\langle t_1, \dots, t_n \rangle$ and given a pool of concepts $\langle c_1, \dots, c_m \rangle$, we want to find the most likely assignment $\langle t_i, c_i \rangle$ such that it maximizes the following probability:

$$c_1, \dots, c_n = \arg \max_{c_1, \dots, c_n} P(c_1, \dots, c_n | t_1, \dots, t_n)$$

Thanks to the Markov assumption (the probability of the i -th concept depends only on the $(i-1)$ -th concept and the probability of the i -th token depends only on the i -th concept), we can simplify the previous formula. Moreover, we can also compute its component by MLE (Maximum Likelihood Estimation):

$$c_1, \dots, c_n = \arg \max_{c_1, \dots, c_n} P(t_i | c_i) P(c_i | c_{i-1})$$

In the previous formula, $P(t_i | c_i) = \frac{C(c_i, t_i)}{C(c_i)}$ and $P(c_i | c_{i-1}) = \frac{C(c_{i-1}, c_i)}{C(c_i)}$ (where $C(x)$ counts the occurrences of x inside the given dataset).

3 Data Analysis

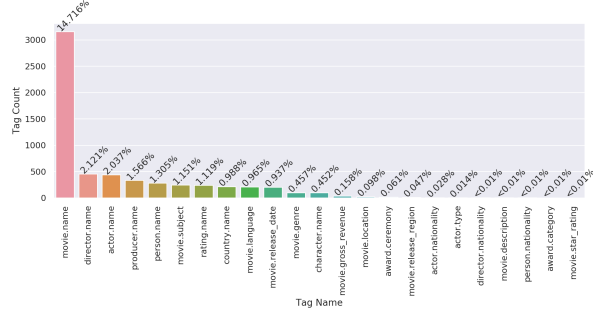
The dataset used is called NL2SparQL4NLU. This work used only the "*.conll.txt" files (more specifically, one for training the language model and one for testing it). Each file is written using the token-per-line CONLL format with tokens and NLU concept tags. An analysis of the content of the two files follows:

NL2SparQL4NLU.trai.conll.txt	
# of lines	24791
# of sentences	3338
# of unique tokens	1728
# of unique concepts (without the prefix)	24

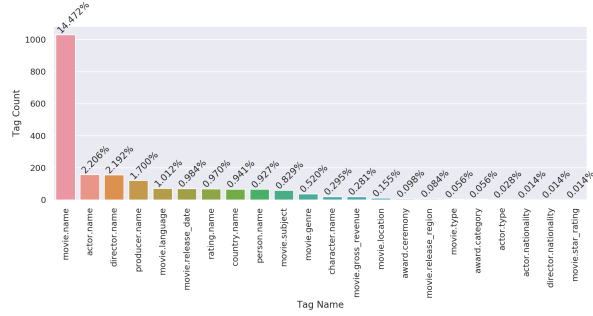
Table 1: Description of the content of the train dataset.

NL2SparQL4NLU.test.conll.txt	
# of lines	8201
# of sentences	1084
# of unique tokens	1039
# of unique concepts (without the prefix)	23

Table 2: Description of the content of the test dataset.



(a) Train dataset.



(b) Test dataset.

Figure 1: Distributions of the various concepts in the train and test datasets. The *O* was removed to make it easier to understand the weight of the others.

The distribution of the tokens was analyzed. It was found to behave accordingly to the Zipf's Law, however, for brevity, it was not reported here since it was not showing anything interesting. The OOV rate for the tokens is around 23.68%. There are also some tokens which are the results of misspelling (e.g., "dispay", "Scorcese"). The distribution of the various concepts was also analyzed (both for the training and the test datasets). The distribution of the various concepts is shown in Figure X.

4 Models

5 Experiments

6 Results

7 Credits