

SemEval-2018 Task 7 Subtask 1: Semantic Relation Classification in Scientific Papers

Hui-Qiang Jiang Da-Wei Lee
Peking University
{1801210840,1701210963}@pku.edu.cn

April 9, 2019

Abstract

First, we test the performance with some naive ideas to transform the problem into a trainable form. Then sifting out the unfeasible one and bad performance one. Doing some research on previous competitors' paper. Then we decide to follow the LightRel model. We looking for the most important factors that affect the performance the most. Finally, focusing on them and aim to gain more improvement on the Macro F1 score.

We achieve a 56.3% F1 score on Task 1.1. And it beats the best performance of the LightRel 46.4%.

1 Introduction

2 Related Work

We spent much time working on how to transform the problem into a trainable form. We found that it is also very hard for a human (ourselves) to classify this task. (And we have drawn on the experience of the other paper.)

2.1 Naive Idea on Relation

We first doing some research on semantic relation typology. There are six major relation types of this competition, USAGE, RESULT, MODEL,

PART_WHOLE, TOPIC, and COMPARISON. And each of them may have a REVERSE relation (except COMPARISON).

The idea is putting connection words between two entities. And calculate the probability of whether it matches which relation type the most.

There are the combinations we've made:

- USAGE
 - used by
 - used for
 - applied to
 - performed on
- RESULT
 - affects
 - problem
 - yields
- MODEL
 - of and observed
 - associated to
- PART_WHOLE
 - composed of
 - extracted from
 - found in
- TOPIC
 - presents
 - of
- COMPARISON
 - compared to

2.2 Improved Relation

We try to use the words between two entities as the basis of predicting relation. And we found that LightRel has done a similar thing.

LightRel let all the sentences with the same length (i.e. same dimension feature). But the sentence in the training and test set won't always be the same. Thus they have padded some dummy words into the sentence to fill the empty space.

We think this approach is much more reasonable than the previous one. Because the relation of any two words can easily tell by the connecting words or sentence.

2.3 Feature Engineering

We have also observed some of the training samples. And conclude some possible pattern such that if the first word is end with "’s", then it possibly has maybe a sort of belonging relation between them.

In LightRel, they have done many features, they called "word-shape feature", as well.[reference] But by our experiment, we found that. Using these features will only lower the performance. So we only use enable this as comparison purpose.

- any character is capitalized
- a comma is present
- the first character is capitalized and the word is the first in the relation
- the first character is lower-case
- there is an underscore present (representing a multi-word entity)
- if quotes are present in the token

2.4 Embedding

We found that the most significant improvement of the performance is related to embedding. There are two subjects: What corpus we choose to calculate the embedding. And what tools to form the embedding. We have made several tests.

2.4.1 Corpus

We selecting the candidate corpus on the Citation Network Dataset. This dataset provides the scientific paper of recent years. LightRel chooses DBLP v5 combined with ACM v9. And extract out only the "Abstract" part of the paper using the regular expression.

We have done some combination of selection like testing by using DBLP v5 only or using DBLP v10 combined with ACM v9 etc.

2.4.2 Embedding Model

There are several tools that we have tried. Such as word2vec by Google (also used in LightRel), fastText by Facebook and BERT by Google.

We build the embedding with 300 dimension vector and skip the words with appearance less than 5 times. We use the mean of all other embeddings to deal with the out-of-vocabulary problem (i.e. the `UNK` token).

There are some tests among these three model. The usage of the `word2vec` and the `fastText` is very similar. The BERT, because there are too many parameters, need to fine-tune to fit the problem, but we haven't found a good solution, thus the performance is not as good as the previous tools. Finally, we found that `fastText` perform better than others.

3 Model and Framework

The first attempt is continuing using the model used by `LightRel`, the logistic regression model offered by `LibLinear`. But to form the required format that fit `LibLinear` is quite handy because it needs to be file format with a kind of sparse matrix representation. That become complicated when we need to do k-fold cross-validation on the training set.

Then we tried to use `Scikit Learn` library. Because the API of `Scikit Learn` are similar so we can switch from multiple models and test the performance easily. We have tried `LinearSVC`, `LogisticRegression` and `Decision TreeClassifier`. And we found the `LogisticRegression` has still had the best performance.

In the test of `LightGBM`, the performance is not quite ideal. In the first case, we think that because the dimension of features is too large, using a tree-based algorithm may not be a good idea.

Finally, we have tried the simple version `TextCNN` model, the performance is not too bad but still lower than the statistics machine learning model. So we have deprecated it later on.

4 Experiment

In the previous section, we mentioned many experiments on the different corpus, embedding tools and model. We will list the performance of each attempt and highlight the best performance on both tasks.

Because the TOPIC class in subtask 1.1 is unbalanced. So we can see the phenomenon that subtask 1.2 will overall be higher than subtask 1.1.

Task		Subtask 1.1		Subtask 1.2	
Corpus	Embedding Model	Macro-F1	Micro-F1	Macro-F1	Micro-F1
Pre-trained DBLP v5	word2vec	88.99	78.99	64.99	38.99
ACM v9 + DBLP v10	word2vec	87.99	78.99	64.99	38.99
DBLP v5	word2vec	87.99	78.99	64.99	38.99
DBLP v10	word2vec	87.99	78.99	64.99	38.99
ACM v9 + DBLP v10	fastText	87.99	78.99	64.99	38.99
DBLP v10	fastText	87.99	78.99	64.99	38.99
ACM v9 + DBLP v10	BERT	87.99	78.99	64.99	38.99

Table 1: Comparison between different corpus using LibLinear logistic regression model (in %)

Task		Subtask 1.1		Subtask 1.2	
Model	Feature	Macro-F1	Micro-F1	Macro-F1	Micro-F1
LibLinear					
Logistic Regression	embedding only	88.99	78.99	64.99	38.99
LibLinear					
Logistic Regression	+ one-hot	88.99	78.99	64.99	38.99
LibLinear					
Logistic Regression	+ shape	88.99	78.99	64.99	38.99
LibLinear					
Logistic Regression	+ cluster	88.99	78.99	64.99	38.99
LibLinear					
Logistic Regression	+ e1, e2	88.99	78.99	64.99	38.99
Scikit Learn					
Logistic Regression	embedding only	88.99	78.99	64.99	38.99
Scikit Learn					
Logistic Regression	+ one-hot	88.99	78.99	64.99	38.99
Scikit Learn					
Logistic Regression	+ shape	88.99	78.99	64.99	38.99
Scikit Learn					
Logistic Regression	+ cluster	88.99	78.99	64.99	38.99
Scikit Learn					
Logistic Regression	+ e1, e2	88.99	78.99	64.99	38.99

Table 2: Comparison between differen feature based on LibLinear and Scikit Learn Logistic Regression model (in %)

4.1 Different Corpus and Embedding Model

4.2 Different Feature

4.3 Different Model

Task	Subtask 1.1		Subtask 1.2	
Model Name	Macro-F1	Micro-F1	Macro-F1	Micro-F1
LibLinear				
Logistic Regression	88.99	78.99	64.99	38.99
Scikit Learn				
Logistic Regression	87.99	78.99	64.99	38.99
TextCNN	87.99	78.99	64.99	38.99

Table 3: Comparison between differen machine learning model using fastText embedding model on DBLP v10 corpus (in %)

5 Conclusion and Future Work