

# Cosi 134 - Project 2: Neural Discourse Relation Classification

Due: October 25, 2019

## 1. What is Discourse Relation Classification?

A typical text consists of sentences that are glued together in a systematic way to form a coherent discourse. Discourse Relation Classification is the task of classifying the discourse relation between two adjacent or non-adjacent discourse units. This project is motivated by the CoNLL 2016 Shared Task: Shallow Discourse Parsing. More information can be found here: <http://www.cs.brandeis.edu/~clp/conll16st/index.html> (<http://www.cs.brandeis.edu/~clp/conll16st/index.html>)

### 1.1 Discourse Analysis in the Penn Discourse Treebank

We use the Penn Discourse Treebank as the data set. The PDTB annotates a text with a set of discourse relations. A discourse relation is composed of:

- a discourse connective, which can be a coordinating conjunction (e.g., "and", "but"), subordinating conjunction (e.g. "if", "because"), or a discourse adverbial (e.g., "however", "also"). In an implicit discourse relation, a discourse connective is omitted;
- two Arguments of the discourse connective, Arg1 and Arg2, which are typically text spans the size of clauses or sentences;
- the sense of the discourse connective, which characterizes the nature of the relationship between the two arguments of the connective (e.g., contrast, instantiation, temporal precedence).

### 1.2 Examples of discourse relations

Here is a paragraph taken from the document wsj 1000 in the PDTB. A discourse relation classifier will output a discourse relation for each pair of arguments, as visualized below. Arg1 is shown in red, Arg2 is shown in blue, and the discourse connective is underlined.

### Explicit Discourse relations:

According to Lawrence Eckenfelder, a securities industry analyst at Prudential-Bache Securities Inc., "Kemper is the first firm to make a major statement with program trading." He added that "having just one firm do this isn't going to mean a hill of beans. But this prompts others to consider the same thing, then it may become much more important."

The discourse connective is "but", and the sense is Comparison.Concession.

### Implicit Discourse Relations

According to Lawrence Eckenfelder, a securities industry analyst at Prudential-Bache Securities Inc., "Kemper is the first firm to make a major statement with program trading." He added that "having just one firm do this isn't going to mean a hill of beans. But if this prompts others to consider the same thing, then it may become much more important."

The omitted discourse connective is "however", and the sense is Comparison.Contrast.

## 2. What you need to do for this project

### 2.1 Data

We use Penn Discourse TreeBank (PDTB) 2.0, a 1-million-word Wall Street Journal corpus, for our experiments. The training data are from Sections 2-21, the development set is from Section 22, and Section 23 is used as the evaluation/test set.

#### relations.json

This file is from The Penn Discourse Treebank (PDTB), with gold standard annotation. Each line in the file is a json line. In Python, you can turn it into a dictionary.

```
In [10]: import json
pdtb_file = open("data/train/rerelations.json", encoding='utf-8')
relations = [json.loads(x) for x in pdtb_file]
example_relation = relations[10]
example_relation
```

```
Out[10]: {'Arg1': {'CharacterSpanList': [[9, 50]],
  'RawText': 'Solo woodwind players have to be creative',
  'TokenList': [[9, 13, 0, 0, 0],
    [14, 22, 1, 0, 1],
    [23, 30, 2, 0, 2],
    [31, 35, 3, 0, 3],
    [36, 38, 4, 0, 4],
    [39, 41, 5, 0, 5],
    [42, 50, 6, 0, 6]]},
  'Arg2': {'CharacterSpanList': [[54, 77]],
  'RawText': 'they want to work a lot',
  'TokenList': [[54, 58, 8, 0, 8],
    [59, 63, 9, 0, 9],
    [64, 66, 10, 0, 10],
    [67, 71, 11, 0, 11],
    [72, 73, 12, 0, 12],
    [74, 77, 13, 0, 13]]},
  'Connective': {'CharacterSpanList': [[51, 53]],
  'RawText': 'if',
  'TokenList': [[51, 53, 7, 0, 7]]},
  'DocID': 'wsj_0207',
  'ID': 3183,
  'Sense': ['Contingency.Condition'],
  'Type': 'Explicit'}
```

The dictionary describes the following component of a relation:

- Arg1 : the text span of Arg1 of the relation
- Arg2 : the text span of Arg2 of the relation
- Connective : the text span of the connective of the relation
- DocID : document id where the relation is in.
- ID : the relation id, which is unique across training, dev, and test sets.
- Sense : the sense of the relation
- Type : the type of relation (Explicit, Implicit, Entrel, AltLex, or NoRel)

The text span is in the same format for Arg1, Arg2, and Connective. A text span has the following fields:

- CharacterSpanList : the list of character offsets (beginning, end) in the raw untokenized data file.
- RawText : the raw untokenized text of the span
- TokenList : the list of the addresses of the tokens in the form of (character offset begin, character offset end, token offset within the document, sentence offset, token offset within the sentence)

**raw/ws\***

These files contain the raw text data from PDTB, corresponding to character/token/sentence offsets in relations.json.

## 2.2 Project Requirements:

1. Implement data processing code that (1) reads in training/develop/test data in the format as described in Section 2.1, and (2) outputs sense classifications in the format as described in Section 2.3.
2. Design and implement a Neural Network model for discourse relation classification with Tensorflow Keras.
3. Experiment with different hyper-parameters, feature representations, and neural network architectures. Experiment with both a feedforward and a convolutional Network.
4. Evaluate your system on explicit, implicit, and all relations, using the evaluation script described in Section 2.4.
5. Write a project report:
  - Describe your feature representation, neural network architecture, and hyper-parameters design.
  - Describe your code structure, and how to run/test your system.
  - Describe your experiments, and report experimental results and analysis.
  - Limit your description within 3 pages.

**Extra credit experiments:** Code your feedforward neural network from scratch, using the backpropagation code you wrote for Homework 1. Extra points will be awarded if you can show that with the same hyperparameters, your feedforward is competitive against a Keras-based model.

## 2.3 What should the system output look like?

The system output must be in json format. It is very similar to the training set except for the TokenList field. The TokenList field is now a list of document level token indices. If the relation is not explicit, Connective field must still be there, and its TokenList must be an empty list. You may however add whatever field into json to help yourself debug or develop the system. Below is an example of a relation given by a system.

```
Out[13]: {'Arg1': {'TokenList': [275, 276, 277, 278, 279, 280, 281, 282]},
          'Arg2': {'TokenList': [329, 330, 331, 332, 333, 334, 335, 336, 337]},
          'Connective': {'TokenList': []},
          'DocID': 'wsj_1000',
          'Sense': ['Expansion.Conjunction'],
          'Type': 'Implicit'}
```

## 2.4 Evaluation

Precision, recall, and F of the discourse relation classification are used as the evaluation metrics. To evaluate, run: `python scorer.py /relations.json /output.json`

The scorer should not be modified. If you run into issues, contact me or the TAs