# CENG 463 - Introduction to Natural Language Processing
## Fall 2023-2024
## Programming Assignment 2

KARADAYI, Meriç

e2448553

January 14, 2024

---

# 1 Dependency Parsing

## 1.1 Implementation Details

As also task mentions I have used spacy library. After a little documentation research, I found out that the CLI documentations of spacy is much more detailed than the python file documentations. Therefore I decided to use the command line functionalities spacy provides. My implementation can work properly when spacy and spacy-transformers are installed on the environment. My implementation consists of 3 main commands.

- Preparing datasets

- Creating configuration file

- Training model

### 1.1.1 Preparing datasets

The training, validation and testing datasets are provided in .conllu formats. Spacy does not support files in conllu format directly. It takes data in.spacy file format. However, it is providing a function that converts the conllu file into spacy file directly. Moreover, I have divided my sentences into documents with 10 sentences as spacy suggesting.

```
#Example command
python –m spacy convert tr_imst−ud−train.conllu . —converter conllu —n−sents 10
```

### 1.1.2 Creating configuration file

Spacy wants a config file which is storing all settings and hyperparameters that model is going to use while training. The components which are going to used and trained in pipeline also is added on this part. I have preferred to use transformers rather than tok2vec embedders. Even though tok2vec would work more efficiently, I thought that transformers would outperform the tok2vec in the result, so I have used transformers. As I said pipeline components are being specified in the config file. Since we are implementing a dependency parser, I have only added a parser to my pipeline. After these adjustment are added as flag, spacy automatically generates a base cfg file.

```
python –m spacy init config config.cfg —lang tr —pipeline parser
—optimize accuracy —gpu  —force
```

### 1.1.3 Training model

Base configuration file is created but some values inside the cfg file still need to be changed such as train and validation set' paths. These parameters can also be passed on this part. Furthermore all the settings in config file can be modified while calling the train command. Therefore, I have changed some parameters such batch size and number of max epoch. The default value of batch size is 128 but we do not have a too large dataset, so I think use of batches with 128 data is not going to help as we wanted. I modified it to 32. Spacy does not limit the number of epochs when it is training. There are some methods to stop like early stopping etc. but I observed after 100 epochs the learning of the model is starts to slow down.

```
python –m spacy train config.cfg —output . —gpu−id 0 — —paths.train tr_imst−ud−train.spacy
—paths.dev tr_imst−ud−dev.spacy —training.max_epochs 100 —nlp.batch_size 32
```

## 1.2 Evaluation

The benchmark of the model are as figures below.

```
================================= Training pipeline =================================
ℹ Pipeline: ['transformer', 'parser']
ℹ Initial learn rate: 0.0
E    #        LOSS TRANS...   LOSS PARSER   DEP_UAS   DEP_LAS   SENTS_F   SCORE
---  ------  -------------   -----------   -------   -------   -------   ------
  0     0         2374.51       1889.85     22.49      1.51      0.00     0.12
  9   200       339833.86     324056.37     57.18     45.79     73.76     0.51
 19   400       137150.82     150584.63     67.69     58.73     96.29     0.63
 29   600        60106.04      76717.77     68.74     59.64     97.09     0.64
 38   800        25792.08      43701.57     68.70     60.18     97.43     0.64
 48  1000        12368.92      31995.54     68.95     60.54     96.78     0.65
 58  1200         7998.24      28033.78     69.84     61.35     96.74     0.66
 68  1400         5620.77      26917.62     69.19     60.70     97.34     0.65
 78  1600         5298.43      26951.66     68.72     60.18     96.88     0.64
 87  1800         4802.05      26429.69     69.43     61.24     96.93     0.65
 97  2000         3892.71      25529.23     69.47     61.18     96.13     0.65
✓ Saved pipeline to output directory
batch/model-last
```

Figure 1: Evaluation metrics during training on validation set

```
============================== Results ==============================

TOK      97.53
UAS      72.24
LAS      63.76
SENT P   97.17
SENT R   98.26
SENT F   97.71
SPEED    1722
```

Figure 2: Evaluation metrics on testing

- Unlabeled attachment score: 0.7224

- Labeled attachment score: 0.6376

Labeled attachment score (LAS) and unlabeled attachment score (UAS) are both metrics used to evaluate the performance of dependency parsers. They analyze the grammatical structure of sentences by determining the relationships between words. Labeled Attachment Score evaluates both the correct attachment of words and the correct identification of the relationships between the words. However Unlabeled Attachment Score focuses only on the correct attachment of words. It does not consider the relationships between words.

$$\text{LAS} = \frac{\text{\# of Correct Attachments and Labels}}{\text{Total Number of Words}}$$

$$\text{UAS} = \frac{\text{\# of Correct Attachments}}{\text{Total Number of Words}}$$

UAS is approximately 10 percent more than the LAS. When the formulas is considered, it makes totally sense. Because the correct predictions LAS is marking is basically a subset of the correct prediction UAS is marking. LAS does not count the words whose attachment is correct but label is incorrect unlike UAS does.
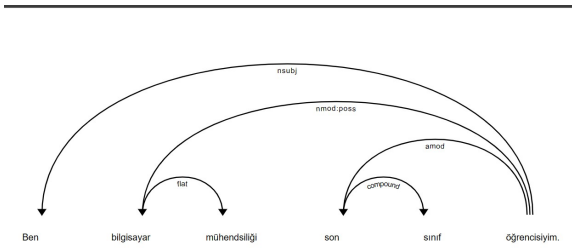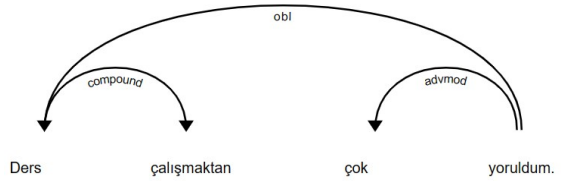
## 1.3 Experiment

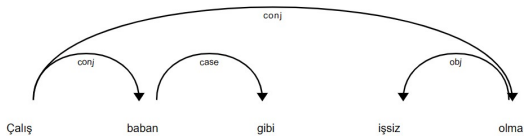Figure 3: Correct parsing

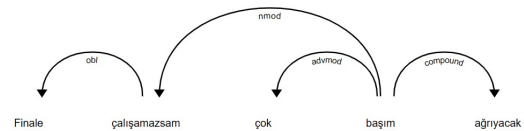Figure 4: Correct parsing

Figure 5: Incorrect parsing

Figure 6: Incorrect parsing

The correct parsing results are sentences which are grammatically totally correct, so it performed a better result. However, turkish is a flexible language and words on sentences can be in inverse order, they do not need to have a strict order, so this caused our model to fail. In addition soma sentence can be derived to multiple meanings which people can understand from the context, this also a weak side which cause false predictions.

# 2 Extractive Summarization

## 2.1 Implementation Details

My implementation consist of phases.

- Take file name as input

- Create vectorizer

- Getting most informative sentences

### 2.1.1 Creating vectorizer

1. Reads all the file's name ending with ".xml"

2. Parses XML files and extract the sentences from each file

3. Creates a mega document which each element is merge of sentence of each cases

4. Create a IDF vectorizer according to megadoc.

### 2.1.2 Getting most informative sentences

1. Create a TF-IDF matrix specific the file user entered using the created vectorizer

2. Calculate cos similarities and return the 5 sentences which has the largest cos similarity value on average. These are the most informative sentences
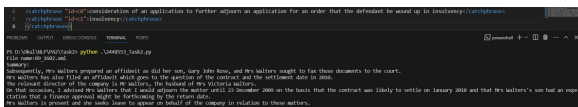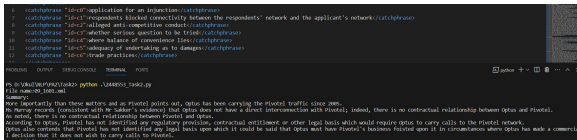
## 2.2 Experiment



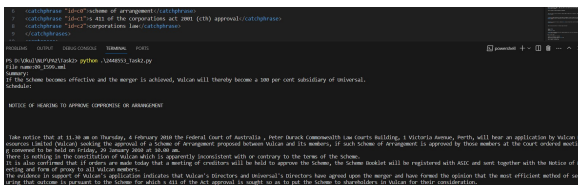Figure 7: Experiment-1
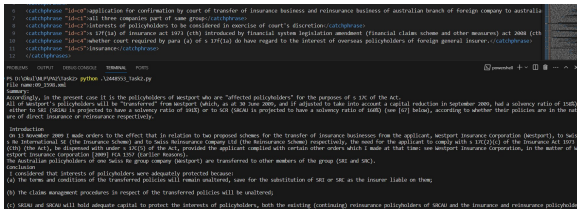


Figure 8: Experiment-2

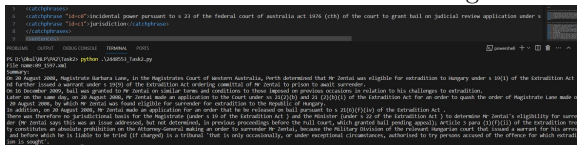

Figure 9: Experiment-3



Figure 10: Experiment-4



Figure 11: Experiment-5

## 2.3 Enhancement

We have used TF-IDF embedding for this task, however instead of relying on TF-IDF, consider using more advanced sentence embedding techniques can provide better performance. For example, we have used transformers on Task-1, if we have used transformers also in this this would help our model to perform better. Furthermore we did not add any semantic relationship or meanings to our words during training, adding such features can enhance the performance. We have compared our sentences with cosine similarities but we did not perform any clustering. If we use clustering algorithms, we could provide better summary option with our model.

## 2.4 Evaluation

According to my research on how to evaluate a summarization algorithm. ROUGE (Recall-Oriented Understudy for Gisting Evaluation) which is a set of metricss can be used primarily is such tasks. The main point of ROUGE is that it measures overlapping n-grams. There is a library called rouge in Python. That libraries is evaluating the generated summary with respect to a reference summary. The case files have also catchphrase information, we could use merge this catchphrases and take it as reference summary to evaluate our model with rogue.