



[Real-world Sentence Boundary DetectionUsing Multitask Learning: A Case Study on French]

This repository is a submission "Real-world Sentence Boundary DetectionUsing Multitask Learning: A Case Study on French" for NAACL 2021. You can check the data in the "data" folder both for training and testing. We serve around 50 percent of our training data because of license issues. You can test three different models, namely 'Baseline', 'multi-task with XLM-Roberta', and 'multi-task with CamemBERT' models.

1. Dependencies

This is a list of packages that required to run the codes.

- Python 3.7 interpreter
- Pytorch 1.6.0
- transformers from Huggingface

1-1. Install Pytorch

- Windows, Linux

```
conda install pytorch-cpu torchvision-cpu -c pytorch
```

- MacOS

```
conda install pytorch torchvision -c pytorch
```

1-2. Install Required packages

- Windows, Linux

```
pip install -r requirements.txt
```

2. Training and testing

2-1. Baseline with XLM-Roberta

```
python baseline.py
```

2-2. Multi-task with XLM-Roberta

```
python multi-task-roberta.py
```

2-3. Multi-task with CamemBERT

```
python multi-task-camembert.py
```

3. Performance

Metric	Value
F1-Score of the all SBD	0.7846
F1-score of the middle SBD	0.4834
POS Accuracy	0.9811
F1-Score of the NER	0.2613