

DAG transducer

This project is the two-part system for our ACL2018 submission.

1 Dataset

There are three directories in directory `data/export`.

- `wsj00a` is training data sample.
- `wsj20a` is dev dataset sample
- `wsj21a` is test dataset sample

In directory `output`, there are four files. They are the generation result of our complete experiment.

- `output/test-pred.txt` is the final sentences
- `output/test-output.txt` is the final lemma sequences
- `output/test-sentences.txt` is the standard sentences
- `output/test-lemmas.txt` is the standard lemma sequences

2 Requirements

- Python3 3.5 or newer
- CMake 2.8 or newer
- Modern C++ compiler that supports C++14
- C++ Boost Library 1.58 or newer
- Git
- Torch (use for OpenNMT)

Run following command in project root to install python packages

```
1 pip install -r requirements.txt > /dev/null
```

3 Rule Extraction

Run following command in project root:

```
1 bash scripts/extract_rules
```

Or run following commands step by step:

```
1 # Dump gzip data files to binary format
2 python3 rule_extraction/reader.py
3 # Generate intermediate graphs for train data
4 python3 rule_extraction/data_stream.py
5 # Extract induced-rules
6 python3 rule_extraction/extract_compiled_rules.py
7 # Generate extended-rules
8 python3 rule_extraction/extend_rules.py
```

Several files will be generated in directory `data`. `data/rules.txt` is the rules file for reading and `data/compiled-rules.txt` is the rules file for transducer.

Every rule in `rules.txt` is like:

```
{NP_VP:ARG1(a0, a1)} => _v_1 => {NP:ARG2(b0), NP:ARG1(c0)}
a0 = c0
a1 = "to" + L + b0
```

- Variables are directly put after state labels.
- Direction of state is not displayed since we can induce it from statements.

4 DAG transducer

Run following command in project root to build and train DAG transducer:

```
1 bash scripts/train_transducer
```

Or run following commands step by step:

```
1 mkdir build && cd build
2 cmake ../transducer -DCMAKE_BUILD_TYPE=Release && make
3 cd ..
4 # Train the main perception for selecting rules
5 build/Release/transducer -c scripts/train.config -m data_driven
6 # Train the perception for generator dynamic-rules
```

```
7 build/Release/transducer -c scripts/train.config -m data_driven.generator
8 # Transduce test dataset
9 build/Release/transducer -c scripts/transduce.config -m data_driven
```

Three files will be generated in directory `data`.

- Two feature files: `data/train.feats` and `data/dynamic.feats`
- Transduction result for test dataset: `data/test-output.txt`

5 Morphology Generator

First follow link [Install Torch](#) to install Torch. Make sure executable 'th' is in PATH. Then install OpenNMT by running following commands in project root:

```
1 luarocks install tds
2 luarocks install bit32 # if using LuaJIT
3 git clone https://github.com/OpenNMT/OpenNMT
```

Then run following command in project root:

```
1 bash scripts/opennmt_run
```

Or run following commands step by step:

```
1 # Generate standard sequence of lemmas for training
2 python3 rule_extraction/process_sentence.py
3 bash scripts/opennmt_prepare_data
4 bash scripts/opennmt_train
```

The final model file will be `data/final-model_epoch*.t7`. Run following command to generate final sentences:

```
1 bash scripts/opennmt_infer data/final-model_epoch*.t7
```

- `data/test-pred.txt` is the output for test dataset
- `data/dev-pred.txt` is the output for dev dataset