

ASSIGNMENT 2: NAMED ENTITY RECOGNITION FOR PROCEDURAL EXTRACTION

Motivation: The motivation of this assignment is to get practice with deep learning models for sequence labeling tasks such as Named Entity Recognition (NER).

Scenario: Scientific text documents contain various procedures in chemistry and biology for performing lab experiments. Each procedure contains a sequence of *instructions* i.e. sentences (E.g. A chemical compound recipe). To make it simple, we extract all these *instructions* from various procedures and make one training data. Our goal is to train a Named Entity Recognition (NER) system that is able to tag each word in a given *instruction* so that relevant information can be extracted downstream.

Problem Statement: The goal of the assignment is to build an NER system for *instructions*. The input of the code will be a set of tokenized *instructions* and the output will be a label for each token in the sentence. Labels will be from 37 classes (including Other (O)) such as *Reagent, Action, Amount, Temperature, Concentration* etc. The complete list of labels is provided to you in the file named *labels.json*.

Labeled Training Data: We are sharing a training and a dev dataset of *instructions*, named `train.txt` and `dev.txt` respectively. In the train/dev file, each line contains one token of a sentence, followed by a tab ("`\t`") and its token label. Sentences are separated by the blank lines. There are 7,199 labeled sentences (examples) in `train.txt` and 2,268 ones in `dev.txt`. We also provide the sample test input and output files to you.

BIO Tagging Scheme: The scheme used for tagging tokens is BIO (Read <https://lingpipe-blog.com/2009/10/14/coding-chunkers-as-taggers-io-bio-bmewo-and-bmewo/>). Basically, every token is either the beginning (B) of a chunk, the continuity of a chunk (I) or outside the chunk (O). E.g. "*Barack Obama went to Greece today*" -> "*Barack B-PER Obama I-PER went O to O Greece B-LOC today O.*" Of course, there are other types of tagging schemes also possible like simply BO tagging, where I- is not explicitly tagged, and all contiguous tokens of the same type are combined to extract one entity. In such a schedule, the tagging will be "*Barack PER Obama PER went O to O Greece LOC today O.*" You may read on other possible schemes like BIOUE scheme.

The Task: You need to write a sequence tagger that labels the given *instructions* in a tokenized test file. The tokenized test file follows the same format as training data except that it contains only a sentence token in each line immediately followed by newline character '\n' (i.e. no token labels in test file will be given at inference time). You should label the test file in the same format as the training data. The format of your output file will be the label (of string type) for each corresponding token of the test file in one line immediately followed by '\n' (do not output the token itself in output file but only the labels). So, the output will have the same no. of lines as the text file with matching blanks marking the end of sentence. A few points to note are following-

1. You need to use deep learning for this assignment. The allowed models are CNN, LSTM, Transformers (or any extension of these, including attention or custom-built models). **However, you are NOT allowed to use any pretrained Language Models such as BERT, ELMO, GPT.** You must train all models from scratch. You are allowed to use pre-trained word vectors from word2vec, Glove or FastText. If you wish to use any other pre-trained information, you should ask on Piazza.
2. You may like to create additional features for each token, e.g. whether the token is capitalized or not, whether it's a number or not etc. You may also try features from lower level syntactic processing like POS tagging or shallow chunking. You may use off-the-shelf POS taggers with proper reference/citation in the writeup file. However, such taggers (or other codes) cannot also use any pre-trained language models.
3. Defining task-specific features such as specific regular expressions indicative of specific types could be useful.
4. You may like to create your own tokenizer for the task. Please read the data carefully to design one.
5. You can use any off-the-shelf tool/code, but only for feature engineering (but not for making the NER model itself), as long as it does not use any pre-trained language model.
6. You are welcome to use probabilistic models like CRF on top of deep learning models. Example, read up on BiLSTM-CRF or Transformer-CRF models for the task of sequence labeling.

Submission Format:

The submission deadline for the assignment is 17th October.

We will follow a similar format as in A1, with changes in output format etc. for this problem. Please read below:

1. Submit a zip file on moodle with name <kerberos_id.zip> (E.g. `csz198394.zip`). Unzipping this should generate a directory with name <kerberos_id> (E.g. `csz198394`) having 3 files – `install_requirements.sh`, `run_model.sh` and `writeup.txt`.
2. The command for training is - `bash run_model.sh train <train_file_path> <val_file_path>`. This should generate a tagging model named `_model` (E.g. `csz198394_model`) along with possibly other files needed for inference.
3. Command for inference - `bash run_model.sh test <test_file_path> outputfile.txt`. This should generate `outputfile.txt` having predicted label in string format (`B-Action`, `I-Action` etc.) followed by `\n` in each line for the token in the corresponding line of test file.

Note:- Please see the sample input and output files provided to you. We will follow the exact same format for test input and expected output.

Evaluation Criteria:

1. This assignment is worth 100 points.
2. **We will take the mean of micro-F1 and macro-F1 scores based on all 36 NER labels excluding Other.**
2. Bonus points awarded for outstanding performers

What is allowed? What is not?

1. The assignment is to be done individually.
2. You should use Python 3.7 and PyTorch for this assignment.
3. You must not discuss this assignment with anyone outside the class. **Make sure you mention the names in your write-up in case you discuss with anyone from within the class.** Please read academic integrity guidelines on the course home page and follow them carefully.
4. Feel free to search the Web for papers or other websites describing how to build named entity recognizers. Cite the references in your writeup.

5. As mentioned earlier, you are NOT allowed to use any pretrained LMs such as BERT, ELMO, GPT etc. or any software based on these.
5. We will run plagiarism detection software (amongst submissions + commonly available public repositories for NER). Any team found guilty will be awarded a suitable penalty as per IIT rules.
6. Your code will be automatically evaluated. You will get a significant penalty if it does not conform to output guidelines.