

CSE 5525 Homework 3: Tagging

Alan Ritter

In this assignment you will implement the structured perceptron and Viterbi algorithms for part-of-speech tagging. Then you will experiment with your trained models on a small dataset of tweets annotated with parts-of-speech and named entities. These experiments will explore the question of how well part of speech taggers perform when applied to domains other than for which they were trained. For example, how well does a Wall Street Journal trained part-of-speech tagger perform when applied on Twitter?

We provide you with starter Python code to help read in the data and evaluate the results of your model's predictions. You are *strongly* encouraged to make use of the provided code. If you really prefer to implement everything from scratch for some reason, please talk to the instructor first. Your submitted code should run on the command line in a unix-like environment (e.g. Linux, OSX, Cygwin).

The experiments required to complete the assignment will take some time to run, so it is highly recommended to start early. We recommend you read through this entire document and run the sample code before getting started.

We have provided an evaluation script to train your models and generate predictions on the test data as follows:

```
> #Trains a part-of-speech tagger on the provided Twitter training set
> bash eval.sh twitter
> #Trains a part-of-speech tagger on the provided penn treebank data
> bash eval.sh ptb
#Trains a part-of-speech tagger on the provided IRC chat data
> bash eval.sh nps
```

Your model's predictions on the test data will be output into the directory, `eval/`. You can check the accuracy of your model on the test data by using the provided scripts `accuracy.py` for POS-tagging and the Perl script `conlleval.pl` for named entity recognition.

When you first run the starter code, the tagger will always predict every word is a noun. You will need to implement the viterbi algorithm for decoding in addition to parameter updates analogous to the perceptron algorithm in Homework #2.

Word's Most Frequent Tag Baseline (2 points)

Before getting started with implementing the perceptron tagger, write a simple program to implement the following baseline. First count the number of times each word occurs with each tag in the training dataset (`data/twitter_train_universal.txt`). Now generate an output file (using the same format as the training data) that predicts the tag of each word using the following heuristic: simply tag each word in the test dataset (`data/twitter_test_universal.txt`) with the tag that appears most frequently in the training dataset (breaking ties arbitrarily). Tag all the unseen words in the test set as nouns. Report your accuracy on the test data using the provided script like so:

```
> python accuracy.py mft_baseline.out data/twitter_test_universal.txt
```

Viterbi Algorithm (6 points)

Implement the Viterbi Algorithm for a bigram perceptron tagger. The provided code in `Data.py` will read in the provided training data. You should make use of log-scores (unnormalized log probabilities) - each multiplication in Viterbi should be replaced with addition, and unnormalized probabilities are simply dot products of feature vectors and weights. Note: you will need to complete the next part of the assignment before you can test if your implementation is properly working.

The methods you will need to implement is `ViterbiTagger.Viterbi`. Before you do this, the classifier always predicts 'Noun'.

Structured Perceptron (4 points)

Next, implement the structured perceptron algorithm. For this you will need to modify `ViterbiTagger.Train`. Include parameter averaging as in Homework #2. Report your performance (accuracy) training and testing on the Twitter data.

Cross-Domain Experiments (2 points)

Next, try training your POS tagger in each of the following scenarios and report accuracy:

- Train on the provided penn-treebank data and test on Twitter.
- Train on the provided IRC-chat data and test on Twitter.
- Train on all the data (irc + ptb + twitter) and test on Twitter.

What can you say about the performance of part-of-speech taggers when they are applied on text outside their training domain?

Extra Credit: Named Entity Recognition (2 points)

Train your tagger on the provided named entity recognition dataset `twitter_ner_train.txt` and report precision, recall and F_1 on `twitter_ner_test.txt` using the provided script `conlleval.pl`. Next, add additional features to the tagger specifically for the named entity recognition task, and report performance. Commonly used features for named entity recognition include lists of first and last names.¹ Lists of companies, products, etc... can be scraped from various places on the web, such as Wikipedia.²

¹http://www.census.gov/topics/population/genealogy/data/1990_census/1990_census_namefiles.html

²https://en.wikipedia.org/wiki/List_of_companies_of_the_United_States