

# CS59000

## HOMEWORK 1

Penghao Wang

February 26, 2016

### 1 Introduction

The aim of this project is to develop a structured learning system. The data is from the University of Carnegie Mellon. They provided a fast and robust Java-based tokenizer and part-of-speech tagger for tweets, its training data of manually labeled POS annotated tweets, a web-based annotation tool, and hierarchical word clusters from unlabeled tweets. The data is divided into three part, train, dev and test part. The challenge of this project is using the traditional NLP processing tool to work with the Twitter data.

In this project, two learning method is used. First is the MEMM, which is local learning. The second is SVM, which is global learning.

### 2 Feature Sets

I believe the feature setting is the most important part in this project. Bad feature set will lead to a terrible result. For instance, when I started the project, I wanted to do it as the same way that last semester do. That is, set all the words in the training set as the features. However, I found it is impossible to do it. Also, it is not the right way to do it. There are several reasons. First, it will be come very very slow, and it is impossible to do it in time. Second, the feature is not relevant. In this project, several features are set. I will introduce some of them.

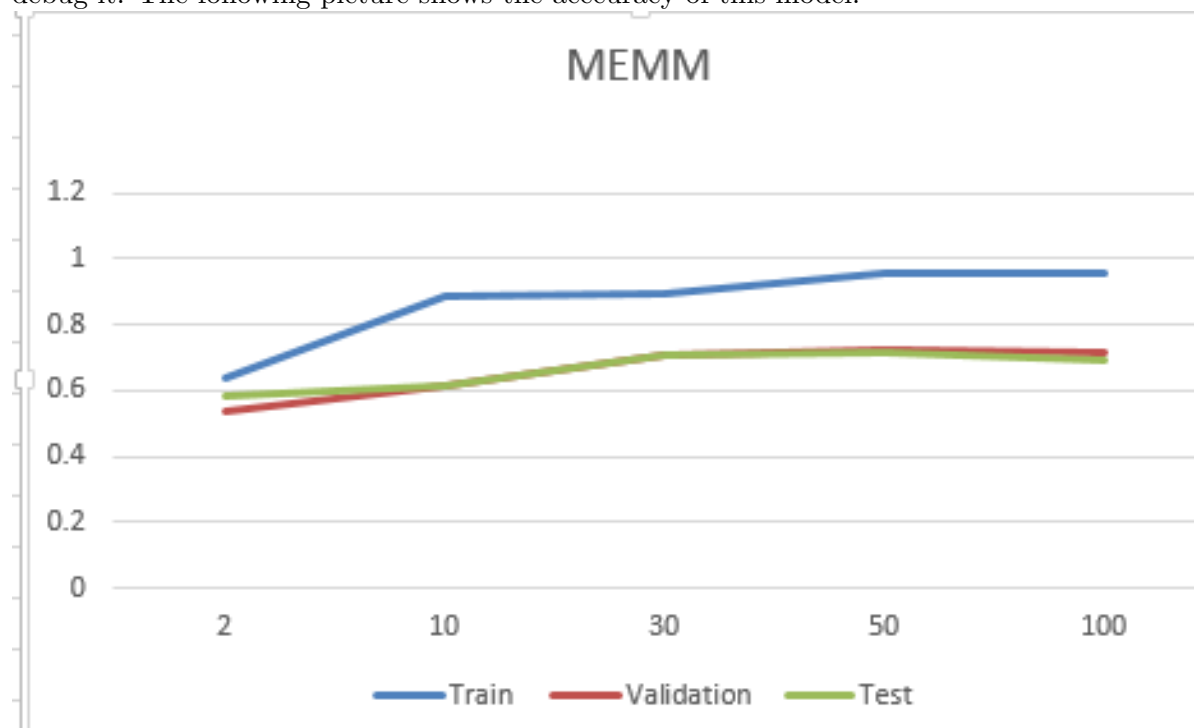
1. The word is in the startList or not. To achieve that, a startWordList is created for it. In the function, when a word come in, it will be compared with the startWordList. If it is in the list, the value should be 1. Otherwise, it will be 0.
2. Whether the word is in the first of the last. For each generation, the first word is marked as one. The other value of the other word will be 0. If the word position is equal to the length of the sentence, it means the word is the last word. The feature position that judges the last word will be one.



## 4 Evaluation

### 4.1 MEMM

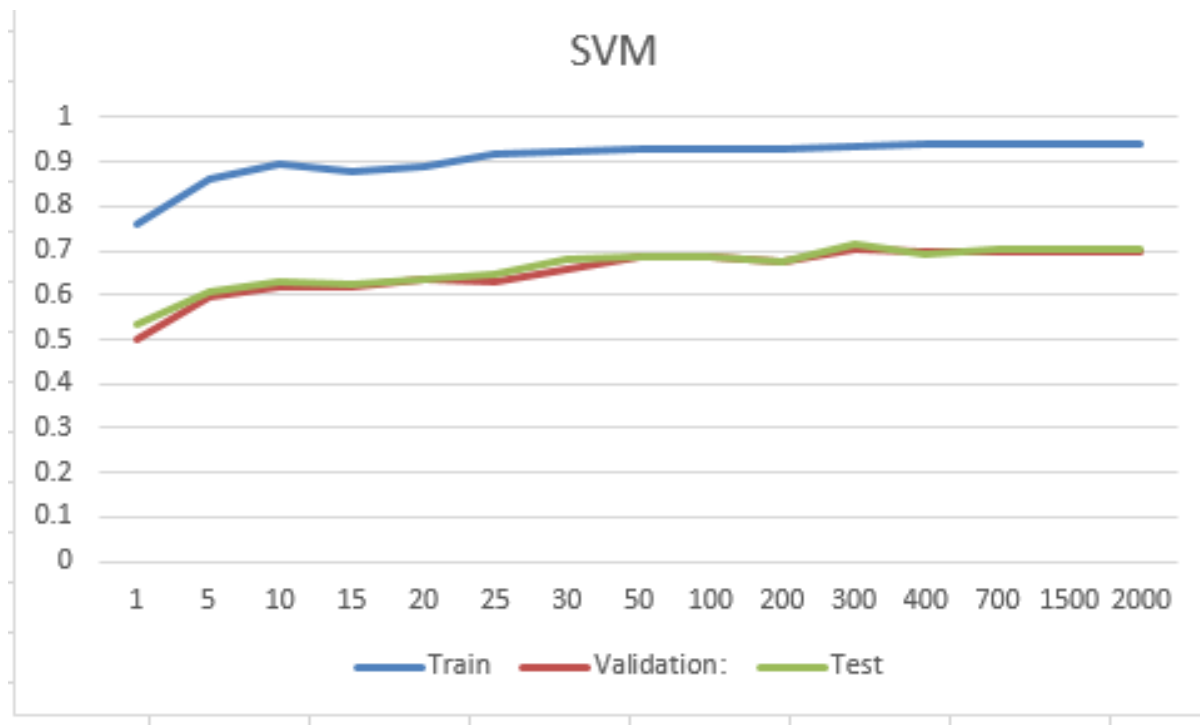
MEMM is a local learning model. Memm is combination of hidden Markov models and maximum entropy models. I use Stochastic gradient descent in this project. Then use the Viterbi to solve it. This Viterbi is different from the Viterbi to solve the HMM. It took me a lot of time to debug it. The following picture shows the accuracy of this model.



Because of my computer, it took very long time to train the model. I don't make to much iterations. But I believe it will converge soon after the 100 iterations. As the iteration increase, the accuracy of the training increase. However, the accuracy of validation and the test set decrease after 30 iterations. It could be overfitting. The best iteration should be 50.

### 4.2 SVM

In this project, I use SVM as the global training method. The following picture shows the accuracy of the model.



The blue line the the training accuracy, the red and green line is stand for the accuracy of the validation and the test. The accuracy of the training increase as the iteration goes up. These are the accuracy of the tag but not the entire tweets. However, the accuracy of the test is going down. It means overfitting. The best iteration should be 300 iteration. The overall accuracy is OK for the training part when compared with the original source.

## 5 Future work

There are several things should be improved. The first one is the feature set. There should be more feature that I create. However, because of the time limit, I don't have the time to do it. The second is the space of the program. It is kind of slow to train the model. More than that, I should do analysis about what tags that are classified wrong the most and improve the program.