



01.112 Machine Learning, Fall 2019

Design Project

Due 10 Dec 2019, 11.59pm

This project will be graded by Richard Sun
Please submit your work to eDimension.

Please form groups for this project early, and start this project early.

Instructions

Please read the following instructions carefully before you start this project:

- This is a group project. You are allowed to form groups in any way you like, but each group must consist of either 2 or 3 people.
- You are strictly NOT allowed to use any external resources or machine learning packages. You will receive 0 for this project if you do so.
- Part 1 deadline is Friday 8 Nov 2019 5pm. **Please start working on part 1 as early as possible** as this part is done **individually**, and you do not need to form a team before you start. Annotated training and development set will be shared with you all by 12 Nov 2019.
- Each group should submit code together with a report summarizing your work, and give clear instructions on how to run your code. Please also submit your system's outputs. Your output should be in the same column format as that of the training set.

Project Summary

In an interview with Michael I. Jordan, the Pehong Chen Distinguished Professor in the Department of Electrical Engineering and Computer Science and the Department of Statistics at the University of California, Berkeley and a giant on Machine Learning, he was asked “If you got a billion dollars to spend on a huge research project that you get to lead, what would you like to do?”. He answered: “I’d use the billion dollars to build a NASA-size program focusing on natural language processing”.

Indeed, one of the most challenging problems within industry and academia is to design intelligent systems that are capable of comprehending human languages. Natural language processing (NLP) is regarded as one of the most challenging yet most promising directions within the field of artificial intelligence in the next 10 years. In this project, we will be working together on building some basic models for solving several simple NLP problems.

Many problems within the field of NLP are essentially **structured prediction problems**, among which sequence labeling is the simplest class of problems. The hidden Markov model (HMM) that we have learned

in class is a **simple structured prediction model**. In this design project, we would like to design our sequence labelling model for informal texts using the HMM that we have learned in class. We hope that your sequence labelling system for informal texts can serve as the very first step towards building a more complex, intelligent language processing system. Specifically, we will focus on building **four NLP systems** – an **English noun phrase chunking system**, a **Chinese address chunking system** as well as **two sentiment analysis systems for Tweets/Weibo**.

The files for this project are in the files `EN.zip` (the English noun phrase chunking dataset), `AL.zip` (the Chinese address chunking dataset, which is generously shared by Alibaba to use for our class), as well as `SG.zip`, `CN.zip` (the latter two will be available on 12 Nov 2019, after we all have finished part 1). For each dataset, we provide a labelled training set `train`, an unlabelled development set `dev.in`, and a labelled development set `dev.out`. The labelled data has the format of one token per line with token and tag separated by tab and a single empty line that separates sentences.

The format for the `SG` dataset (as well as `CN` dataset) can be something like the following:

```
Best O
Deal O
Chiang B-positive
mai I-positive
Tours I-positive
, O
The O
North O
of O
Thailand B-neutral
To O
Get O
special O
Promotion O
and O
free O
Transfer O
roundtrip O
. O
Contact O
: O
... O
http://t.co/sSn10BTZ O
```

where labels such as `B-positive`, `I-positive` are used to indicate **B**eginning and the **I**nside of the entities which are associated with a positive sentiment. `O` is used to indicate the **O**utside of any entity. Similarly for `B-negative`, `I-negative` and `B-neutral`, `I-neutral`, which are used to indicate entities which are associated with negative and neutral sentiment, respectively.

The format for the `EN` dataset (similarly for the `AL` dataset) can be something like the following:

```
Cant B-VP
wait I-VP
```

```
for B-PP
the B-NP
ravens I-NP
game I-NP
tomorrow B-NP
..... O
go B-VP
ray B-NP
rice I-NP
!!!!!!! O
```

where labels such as B-VP, I-VP are used to indicate **B**eginning and the **I**nside of a **V**erb Phrase. Similarly for other tags such as B-NP, I-NP and B-PP, I-PP, which are used to indicate spans which are associated with noun phrases and propositional phrases etc. O is used to indicate the **O**utside of any phrase.

Overall, our goal is to build a sequence labelling system from such training data and then use the system to predict tag sequences for new sentences. Specifically:

- We will be building two sentiment analysis systems for two different languages from scratch, using our own annotations.
- We will be building yet another two NLP systems (one for noun phrase chunking and the other for address chunking) for two different languages using annotations provided by others.

1 Part 1 (15 points, due 8 Nov 2019 at 5pm. Please budget your time well.)

The first and most important step towards building a supervised machine learning system is to get annotated data. This is also often one of the most difficult and most challenging steps in building a practical machine learning system, as we will see in this project. To allow each of us to have a full end-to-end experience on how challenging it is to build a practical supervised machine learning system, in the first part of this project, we will work together to get annotated data for performing sentiment analysis from social media data (some of them are collected from local social media users). You will receive 10 points if you complete the annotation. An additional 5 points will typically be awarded to you too unless we found the quality of your annotation is unacceptable. We will use an automatic approach to assess the quality of your annotations. Your annotations will be compiled and distributed to your fellow students in order to complete Part 2 and Part 3 of this project.

Please visit the following site for the annotation interface:

http://ml-project.statnlp.org/annotation.php?id=/**YOUR_ID_HERE**/

You then need to key in your student ID to proceed to the next step for annotation. For example, if your student ID is 1001949, please visit the following link for annotation:

<http://ml-project.statnlp.org/annotation.php?id=1001949>

Alternatively, visit the following site and type in your student ID:

<http://ml-project.statnlp.org/annotation.php>

You need to log in to start the annotation process. The user name and password are both your student ID. We also provide a sample collection of annotated data, which is available here:

<http://ml-project.statnlp.org/annotation.php?id=1000000>

Essentially, we are interested in annotating all major entities together with their sentiment information. Detailed instructions on the annotation can be found at <http://ml-project.statnlp.org/annotation.pdf>. The annotation interface is straightforward to use, but if you have questions there is a manual here:

<http://brat.nlplab.org/manual.html>

Disclaimer: to grant us the right to re-distribute your annotated data to your other fellow classmates and for potential future usage, by submitting your annotations online, you agree that your annotated data will be in public domain unless otherwise stated. Please contact Thilini Cooray if you have questions or doubts on this.

2 Part 2 (25 points)

Recall that the HMM discussed in class is defined as follows:

$$p(x_1, \dots, x_n, y_1, \dots, y_n) = \prod_{i=1}^{n+1} q(y_i | y_{i-1}) \cdot \prod_{i=1}^n e(x_i | y_i) \quad (1)$$

where $y_0 = \text{START}$ and $y_{n+1} = \text{STOP}$. Here q are transition probabilities, and e are emission parameters. In this project, x 's are the natural language words, and y 's are the tags (such as O, B-positive).

- Write a function that estimates the emission parameters from the training set using MLE (maximum likelihood estimation):

$$e(x|y) = \frac{\text{Count}(y \rightarrow x)}{\text{Count}(y)}$$

(5 points)

$k = 3$

- One problem with estimating the emission parameters is that some words that appear in the test set do not appear in the training set. One simple idea to handle this issue is as follows. First, replace those words that appear less than k times in the training set with a special word token #UNK# before training. This leads to a “modified training set”. We then use such a modified training set to train our model.

To be researched

During the testing phase, if a word does not appear in the “modified training set”, we replace that word with #UNK# as well.

This method is related to the idea called **smoothing**. The resulting emission parameters are called **smoothed emission parameters**. Let us assume such a smoothing method for the emission parameters is always used in the subsequent questions of this project (you may also do so in part 5 if you choose to).

(10 points)

- Implement a simple sentiment analysis system that produces the tag

$$y^* = \arg \max_y e(x|y)$$

for each word x in the sequence.

For all the four datasets EN, AL, CN, and SG, learn these parameters with `train`, and evaluate your system on the development set `dev.in` for each of the dataset. Write your output to `dev.p2.out` for the four datasets respectively. Compare your outputs and the gold-standard outputs in `dev.out` and report the precision, recall and F scores of such a baseline system for each dataset.

The precision score is defined as follows:

$$\text{Precision} = \frac{\text{Total number of correctly predicted entities}}{\text{Total number of predicted entities}}$$

The recall score is defined as follows:

$$\text{Recall} = \frac{\text{Total number of correctly predicted entities}}{\text{Total number of gold entities}}$$

where a gold entity is a true entity that is annotated in the reference output file, and a predicted entity is regarded as correct if and only if it matches exactly the gold entity (*i.e.*, both their *boundaries* and *sentiment* are exactly the same).

Finally the F score is defined as follows:

$$F = \frac{2}{1/\text{Precision} + 1/\text{Recall}}$$

Note: in some cases, you might have an output sequence that consists of a transition from O to I-negative (rather than B-negative). For example, “O I-negative I-negative O”. In this case, the second and third words should be regarded as one entity with negative sentiment.

You can use the evaluation script shared with you to calculate such scores. However it is strongly encouraged that you understand how the scores are calculated.

(10 points)

3 Part 3 (20 points)

- Write a function that estimates the transition parameters from the training set using MLE (maximum likelihood estimation):

$$q(y_i|y_{i-1}) = \frac{\text{Count}(y_{i-1}, y_i)}{\text{Count}(y_{i-1})}$$

Please make sure the following special cases are also considered: $q(\text{STOP}|y_n)$ and $q(y_1|\text{START})$.

(5 points)

- Use the estimated transition and emission parameters, implement the Viterbi algorithm to compute the following (for a sentence with n words):

$$y_1^*, \dots, y_n^* = \arg \max_{y_1, \dots, y_n} p(x_1, \dots, x_n, y_1, \dots, y_n)$$

For *all datasets*, learn the model parameters with `train`. Run the Viterbi algorithm on the development set `dev.in` using the learned models, write your output to `dev.p3.out` for the four datasets respectively. Report the precision, recall and F scores of all systems.

Note: in case you encounter potential numerical underflow issue, think of a way to address such an issue in your implementation.

(15 points)

4 Part 4 (20 points)

- Use the estimated transition and emission parameters, implement an algorithm to find the 7-th best output sequences. Clearly describe the steps of your algorithm in your report.

Run the algorithm on the development sets `EN/dev.in` and `AL/dev.in` only. Write the outputs to `EN/dev.p4.out` and `AL/dev.p4.out`. Report the precision, recall and F scores for the outputs for both languages.

Hint: find the top-7 best sequences using dynamic programming by modifying the original Viterbi algorithm.

(20 points)

5 Part 5 – Design Challenge (20 points)

- Now, based on the training and development set, think of a better design for developing an improved sentiment analysis system for tweets using any model you like. Please explain clearly the model/method that you used for designing the new system. We will check your code and may call you for an interview if we have questions about your code. Please run your system on the development set `EN/dev.in` and `AL/dev.in`. Write your outputs to `EN/dev.p5.out` and `AL/dev.p5.out`. Report the precision, recall and F scores of your new systems for these two languages.

(10 points)

- We will evaluate your system's performance on two held out test sets `EN/test.in` and `AL/test.in`. The test sets will only be released on 8 Dec 2019 at 11.59pm (48 hours before the deadline). Use your new system to generate the outputs. Write your outputs to `EN/test.p5.out` and `AL/test.p5.out`.

The system that achieves the overall highest F score on the test sets will be announced as the winner.

(10 points)

Hints: Can we handle the new words in a better way? Are there better ways to model the transition and emission probabilities? Or can we use a discriminative approach instead of the generative approach?

Perhaps using Perceptron?¹. Any other creative ideas? Note that you are allowed to look into the scientific literature for ideas.

Items To Be Submitted

Upload to eDimension a single ZIP file containing the following: (Please make sure you have only one submission from each team only.)

- A report detailing the approaches and results
- Source code (.py files) with README (instructions on how to run the code)
- Output files
 - EN/
 - 1. dev.p2.out
 - 2. dev.p3.out
 - 3. dev.p4.out
 - 4. dev.p5.out
 - 5. test.p5.out
 - AL/
 - 1. dev.p2.out
 - 2. dev.p3.out
 - 3. dev.p4.out
 - 4. dev.p5.out
 - 5. test.p5.out
 - CN/
 - 1. dev.p2.out
 - 2. dev.p3.out
 - SG/
 - 1. dev.p2.out
 - 2. dev.p3.out

¹<http://www.aclweb.org/anthology/W02-1001>