

ASSIGNMENT 1: SENTIMENT MINING and DOMAIN ADAPTATION

Motivation: The motivation of this assignment is to get practice with text categorization using Machine Learning algorithms. The second part of the assignment is to make efforts on domain adaptation within sentiment mining.

PART I

Problem Statement: The goal of the assignment is to build a sentiment categorization system for tweets. The input of the code will be a set of tweets and the output will be a prediction for each tweet – positive or negative.

Training Data: We are sharing a [training dataset](#) of positive and negative tweets. The format of each line in the training dataset is <"label", "tweet">. It has a total of 1.6 million tweets. A label of 0 means negative sentiment and a label of 4 means positive sentiment.

This dataset has been labeled not manually but automatically. In this automatic labeling process tweets with positive smileys were searched and labeled positive and tweets with negative smileys were labeled negative. The smileys have been stripped off from the tweets to make the classification task real.

The Task: You need to write a classifier that predicts for a given new tweet (which does have a positive or a negative sentiment) its sentiment polarity. To accomplish this, as a baseline algorithm, you could use all words as features and learn a classifier (naïve Bayes or logistic regression). To improve your system performance over the baseline I list a few ideas below.

1. Try changing the classifiers. Try SVMs, random forests or even multilayer perceptrons.
2. If you use logistic regression try playing with regularizers – try L1 instead of L2. You could also implement a different feature selection procedure.
3. Try to work with the features. You could lemmatize. You could get rid of stop words and highly infrequent words. You could use tfidf-based weighting.
4. Try to use existing sentiment resources discussed in class. Examples: SentiWordnet or General Inquirer.
5. You could work with bigrams (in addition to unigrams).
6. You could do tweet normalization on words where letters are repeated for intensity, e.g., haaaaaaate and loooooooooove. You can even do more normalization using some internet slang dictionary.

7. You could define new features, like you could pos-tag each word and use the tagged word as a feature instead of the original word. You can use presence of capitalization or all caps as features. You could use this code for POS tagging of tweets:
https://github.com/aritter/twitter_nlp
8. You could look at data quality – if you find that data is not consistently good, you could label some tweets by hand, and use semi-supervised techniques to correct the data errors (or even remove parts of the data).
9. Your idea here...

Methodology and Experiments:

Option 1: Remove 10% of data randomly as test set and 10% as development set. Train on 80% of the data. If needed, do parameter fitting on the devset and finally test on the test set.

Option 2: Do a 10-fold cross validation. Train on 8 folds, use 9th as devset and 10th as test set. Repeat this 10 times with different folds as test set. This is a more robust option.

As you work on improving your baseline system document its performance. Perform error analysis on a subset of data and think about what additional knowledge could help the classifier the most. That will guide you in picking the next feature to add.

Test Format:

Your final program will take input a set of tweets. Each line will have one tweet (without double quotes). Your program will output predictions (0 or 4) one per line – matching one prediction per tweet. [Here](#) is a sample test file and the [sample](#) output.

PART II

Problem Statement: The goal of the assignment is to build a sentiment categorization system for tweets *for a specific 'person'*, e.g., “Arvind Kejriwal” or “Salman Khan”. You will be provided a set of tweets on the person. First you will label a set of tweets manually (50 per student). These will be aggregated from all students to create a domain-specific training set. Second, you will use the training dataset from Part I and this additional training set to create the sentiment categorization system for the specific person.

Training Data: We are sharing 100 tweets to each student, one tweet per row. Write down 0, 2, or 4 below every tweet as its label and return the file so that we can aggregate it.

The Task: You need to write a classifier that predicts for a given new tweet its sentiment polarity with respect to the specific person. Note that this is a three-way classification: 0 for negative, 2 for neutral,

and 4 for positive. As a baseline system you can set thresholds based on your system from Part I (out-of-domain training) to make these predictions. To improve your system performance, use the additional hand-annotated training data (in-domain data) and retune the parameters. Some ideas below:

1. Simply combine both datasets and relearn.
2. Upweight the in-domain training data.
3. Use ideas from bootstrapping as discussed in class to induce domain-specific features.
4. Use PMI of top words in out-of-domain training data with those in in-domain-training data to upweight the features in in-domain data.
5. Your idea here... (feel free to search for ideas on the Web).

What to submit?

1. Submit your tagging for Part II by Wednesday, February 17th. Submit in a text file titled tagging-Entryno.txt. Late submission is not allowed for this.
2. Submit your best code for Part I (best if trained on all training data and not just on a subset) by Thursday, 25th February 2016, 11:55 PM. The code should **not** need to train again. You should submit only the testing code, after the models have been trained. That is, you should not need to access the training data anymore.

Submit your code in a .zip file named in the format **<EntryNo>.zip**. Make sure that when we run “unzip yourfile.zip” the following files are produced in the present working directory:

compile.sh
run.sh

You will be penalized if your submission does not conform to this requirement.

Your code will be run as ./run.sh inputfile.txt outputfile.txt. The outputfile.txt should only contain a sequence of numbers (0 or 4), one per line, with total number of lines matching inputfile.txt. We are providing a format checker. Make sure your code passes format checker before final submission.

Your code should work on our Baadal servers with 2GB memory. You will be penalized for any submissions that do not conform to this requirement.

3. Submit your best code for Part II (best if trained on all training data and not just on a subset) by Tuesday, 8th March 2016, 11:55 PM. The code should **not** need to train again. You should submit only the testing code, after the models have been trained. That is, you should not need to access the training data anymore.

Submit your code in a .zip file named in the format **<EntryNo>.zip**. Make sure that when we run “unzip yourfile.zip” the following files are produced in the present working directory:

compile.sh
run.sh

writeup.txt

You will be penalized if your submission does not conform to this requirement.

Your code should work on our Baadal servers with 2GB memory. You will be penalized for any submissions that do not conform to this requirement.

4. The writeup.txt should have two lines as follows

First line should be just a number between 1 and 3. Number 1 means C++. Number 2 means Java and Number 3 means Python.

Second line should mention names of all students you discussed/collaborated with (see guidelines on collaboration vs. cheating on the course home page). If you never discussed the assignment with anyone else say None.

After these first two lines you are welcome to write something about your code, though this is not necessary.

Evaluation Criteria

- (1) 20 points are for completion of tagging for Part II
- (2) 70 points for performance of your code in Part I
- (3) 60 points for performance of your code in Part II.
- (4) Bonus given to outstanding performers.

What is allowed? What is not?

1. The assignment is to be done individually.
2. You may use Java, or C++ or Python 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 outside your team.** 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 a tweet sentiment classifier. However, you should not use (or read) other people's sentiment mining code.
5. You can use any pre-existing ML softwares for your code. Popular examples include Weka (<http://www.weka.net.nz/>), Python Scikit (<http://scikit-learn.org/stable/>). However, if you include the other code, use a different directory and don't mix your code with pre-existing code.
6. We will run plagiarism detection software. Any team found guilty will be awarded a suitable penalty as per IIT rules.
7. Your code will be automatically evaluated. You get significant penalty if it does not conform to output guidelines. Make sure it satisfies the format checker before you submit.