

MMAI 5400 Assignment 2 -- Sentiment Classification

For this assignment you will build a sentiment classifier to classify reviews from Trustpilot.com.

Submission

This assignment should be submitted as Python 3 code and uploaded to Canvas. The submission should be a single `.py` file, and **not** a Jupyter Notebook. The due date is on June 16 at 8:30 am.

The code will be tested and should produce the output specified below.

You will be provided with a file (`reviews.csv`) that contains close to 2000 reviews. Values in `reviews.csv` are tab separated instead of comma separated as is often the case in `csv` files. To read `reviews.csv` with `pandas` it might be good to specify tab as the delimiter, as follows:

```
data = pd.read_csv('reviews.csv', delimiter='\t')
```

or

```
data = pd.read_csv('reviews.csv', sep='\t')
```

The two are equivalent.

Task

Your task is to do sentiment analysis of the reviews. You will do this by training and testing a BoW text classifier on the review texts using `RatingValue` as labels. The ratings should be binned into negative (ratings `1` & `2`), neutral (rating `3`) and positive (ratings `4` & `5`) sentiment. The binned ratings should be coded with negative as `0`, neutral as `1` and positive as `2`.

However, the ratings are very unbalanced, there are many more positive (`2`) ratings than negative (`0`) ratings. You will need to drop positive ratings in order to balance the data so that you have approximately equal numbers of negative, neutral and positive ratings.

The resulting table should look as follows:

	Sentiment	Review
	1	"It shows ..."
	0	"Disgusting..."
	2	"Yummy..."
...

Once this is done, you should split the data into training and validation sets and save them as `training.csv` and `valid.csv`. The validation data should be used for model selection and evaluation.

For this task it is important that you report performance carefully using both accuracy, F1-score and a confusion matrix.

Deliverable

You need to submit a single Python (`PY` **NOT** `IPYNB`), that does the following:

- Loads the `reviews.csv` , preprocesses the data, splits it and saves the files as `training.csv` and `valid.csv` .
- Loads `training.csv` and trains the model.
- Loads the validation data (`valid.csv`) and prints the performance metrics on the validation set.
 - I.e. it should print something like the following:

```
accuracy: "accuracy on the test set"

F1_score: "f1-score on the test set"

Confusion_matrix:
      negative neutral positive
negative  a      b      c
neutral   d      e      f
positive  g      h      i
```

Grading

For grading the validation data (`valid.csv`) will be replaced with unseen test data (`test.csv`) and for full marks your model needs to perform above chance.

Good luck!