

CS 334 — Homework 4 Solutions

Jack Anstey

Problem 1: Preprocessing + Model Assessment

a) See q1.py

b) See q1.py

c) See q1.py

d) See q1.py

Problem 3: Spam Detection using Naive Bayes and Logistic Regression

- a) For this problem, I found that using MultinomialNB was the best option for both the binary and count datasets as the number of mistakes were either to or less than that of BernoulliNB. With that, the number of mistakes I found were:

```
Number of mistakes when using the binary dataset and multinomial naive bayes: 70  
Number of mistakes when using the count dataset and multinomial naive bayes: 546
```

When using the count dataset, the sklearn implementation matched the number of mistakes that my own implementation had when using an optimal epoch. Surprisingly, my own implementation, when using an optimal epoch, was slightly better than the sklearn implementation (65 mistakes vs 70). I believe this is due to sklearn not using the exact optimal epoch, which I think is fine as their implementation runs significantly faster and provides nearly the same result in regards to accuracy.

- b) When using logistic regression, the results were even better.

```
Number of mistakes when using the binary dataset and logistic regression: 38  
Number of mistakes when using the count dataset and logistic regression: 546
```

The number of mistakes again stayed the same when using the count dataset, but the number of mistakes when using the binary dataset were nearly cut in half - only 38 mistakes were made. Overall, it seems that logistic regression and the binary dataset performs the best for this classification problem.