# NLP HW2 Report

## Part 1

S = "I always like foreign films"

P(pos | S) = P(pos) * P(I | pos) * P(always | pos) * P(like | pos) * P( foreign | pos) * P(films | pos)
= 0.4 * 0.09 * 0.07 * 0.29 * 0.04 * 0.08 = 2.33856 * 10^-6

P(neg | S) = P(neg) * P(I | neg) * P(always | neg) * P(like | neg) * P(foreign | neg)* P(films | neg)
= 0.6 * 0.16 * 0.06 * 0.06 * 0.15 * 0.11 = 5.7024 * 10^-6

P(neg | S) > P(pos | S), therefore, the prediction is negative.

## Part 2

a) Done. Instructions are in the README.md file, the files are named as you've requested in the assignment, and I've added extensive comments to outline the logic of my code.

b) And done. I wrote my code to run the same way for the smaller corpus as for the larger corpus and structured the folders in both corpuses to look the same so all I'd have to input for my methods would be the directory of the corpus.

c) Test your classifier on the new document below: {fast, couple, shoot, fly}.

Compute the most likely class. Report the probabilities for each class. [5 points]

Prediction: action

Log probability of action: -12.509775004326938

Log probability of comedy: -13.736965594166206

When comparing log probabilities, the maximum value is the one predicted by our model.

d) My model using BOW features ran test cases on all 25,000 vectors in the test data set. Out of those 25,000 vectors, my model predicted the class of 20,372 of them correctly, **leading to a 81.488% accuracy score.** To raise my accuracy, I looked for words such as "the" which shared a high frequency and similar probability in both classes. I tried to remove the word "the" from the data set and frankly, my accuracy score went down to 74.2%. I'm confident that it was my implementation of removing the one word, and perhaps removing such a common word, yet not exactly equal on both classes, skewed the probabilities. However, I believe there is merit in removing certain words with high frequency and relatively equal probabilities on both sides to help the model reduce bias on neutral words.