

Question 1) The CS 464 Case

First of all, we can say the following:

$$P(S_M | H) = 0.87, P(S_M | L) = 0.21, P(S_M | F) = 0.04$$

and

$$P(H) = 0.64, P(L) = 0.24, P(F) = 0.12$$

Question 1.1)

$P(S_M)$ can be calculated as the share of motivated students for each grade type which is the equation:

$$P(S_M) = P(H) * P(S_M | H) + P(L) * P(S_M | L) + P(F) * P(S_M | F)$$

or

$$P(S_M) = 0.64 * 0.87 + 0.24 * 0.21 + 0.12 * 0.04 = \mathbf{0.612}$$

Question 1.2)

$$P(H | S_M) = \frac{P(H) * P(S_M | H)}{P(S_M)} \quad \text{or} \quad P(H | S_M) = \frac{0.64 * 0.87}{0.612} = \mathbf{0.9098}$$

Question 1.3)

$$P(S_U) = 1 - P(S_M) = 0.388 \quad \text{and} \quad P(S_U | H) = 1 - P(S_M | H) = 0.13$$

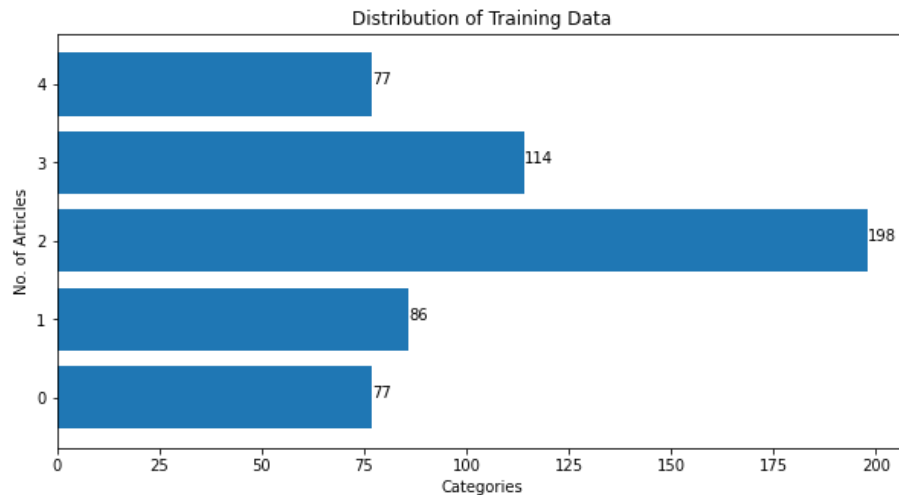
Finally,

$$P(H | S_U) = \frac{P(H) * P(S_U | H)}{P(S_U)} = \frac{0.64 * 0.13}{0.388} = \mathbf{0.2144}$$

Question 2)

Question 2.1.1)

The class distribution can be seen in the following plot:



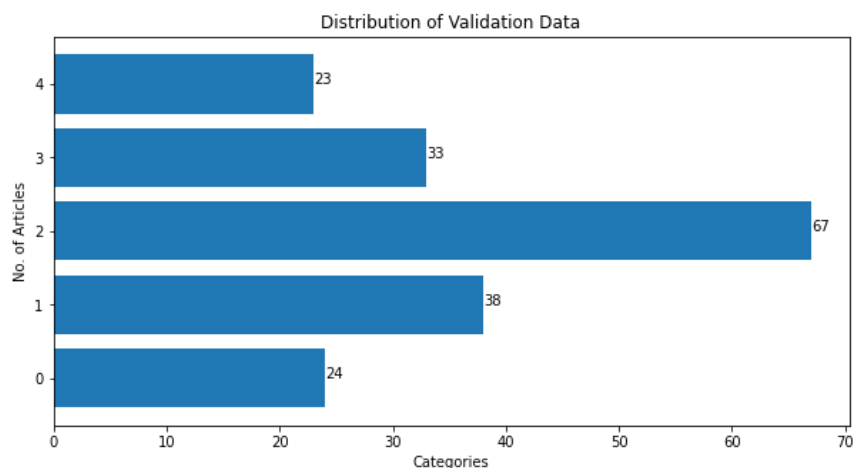
(Figure 1: Distribution of Training Data)

Question 2.1.2)

As can be seen from the graph above, the class 2 labeled articles form one thirds of all training data, the data is skewed through football articles. I believe this creates a bias towards the class 2, due to higher likelihood. To prevent this a dataset that is uniformly distributed must be used.

Question 2.1.3)

The validation set has a similar distribution, and it is skewed through class 2. The class likelihood is misleading since the skewed class label has a higher probability and the overall probability is higher for skewed classes.



(Figure 2: Distribution of Validation Data)

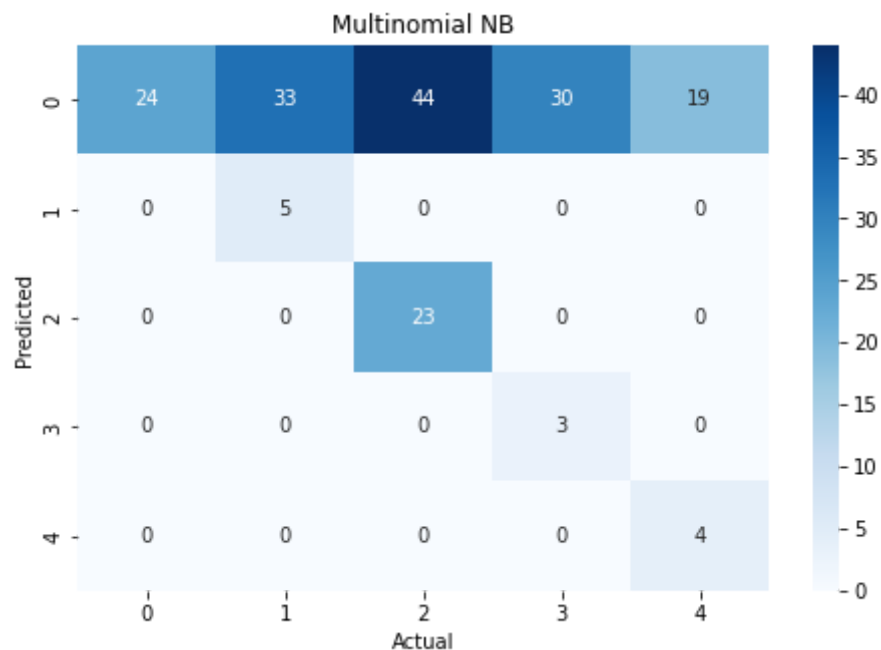
Question 2.1.4)

It makes accuracy less reliable. For example in our case, nearly 1/3 of classes is football articles. Even if our model predicts none of other classes correctly, and predicts everything as football articles, the model still will have nearly 33.3% accuracy.

Question 2.2)

I trained a model using MLE estimator and chose $\log(0)$ as $-\inf$ (`np.nan_to_num(-np.inf)`). The model had a poor accuracy because of the zero probabilities that were resulting in $-\infty$ that make the model predict the lowest class in most cases. The accuracy was low in this part which is 31.89%. The following confusion matrix shows the bias towards class 0.

Accuracy for Multinomial Naive Bayes without MAP 0.31891891891891894

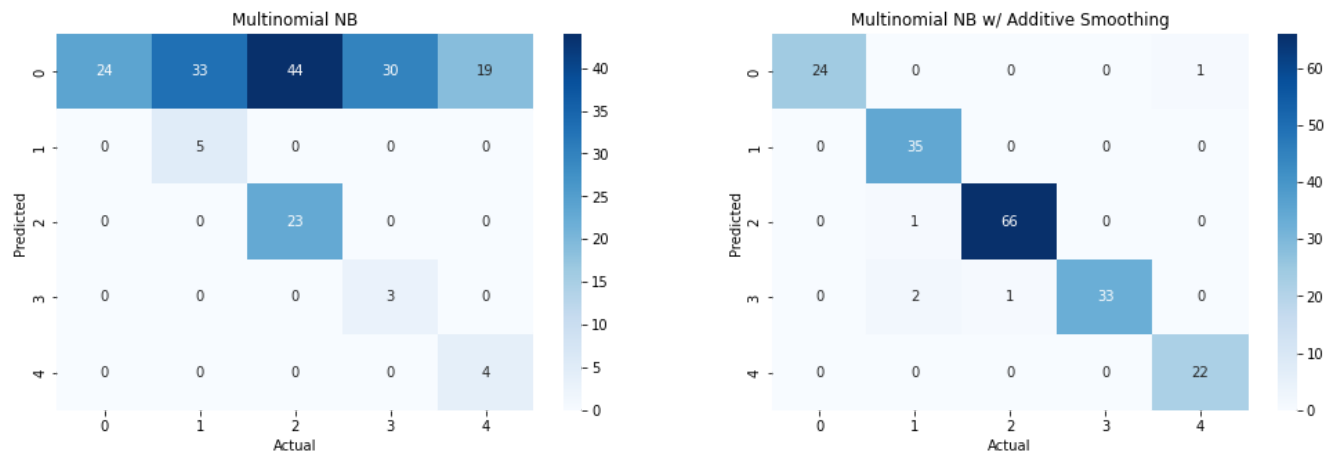


(Figure 3: Confusion Matrix of Multinomial NB Without Smoothing)

Question 2.3)

Using the same training and validation data, I trained a multinomial naive Bayes model with additive smoothing where Dirichlet prior $\alpha = 1$. Since this time there were no zero probabilities thanks to additive smoothing, the accuracy was pretty high which is 97.29%. The model has predicted 180 out of 185 articles correctly. Following confusion matrices show the difference between two models.

Accuracy for Multinomial Naive Bayes with Additional Smoothing
0.972972972972973



(Figure 4: Confusion Matrices of Two Models)

Comparing the two models you trained, how does the Dirichlet prior α effects your model? Also, interpret the structure of the dataset. Given that the dataset does not include stop words, why are the two models different? Explain by giving references to your results. You can also benefit from the statistical structure of the feature matrix.

Question 2.4)

The Dirichlet prior causes a dramatic increase in accuracy (from 31% to 97%). It reduces wrong predictions from 126 to 5. This addition means that we assume that there are already one of each words in all articles. This is meaningful because there are some words in dataset that never occurs in specific types of articles (which result in zero probabilities). The Dirichlet prior eliminates all zero possibilities in dataset and prevent us from facing -inf value that makes model predict wrong classes. This is why two models are drastically different, one includes -inf values and having a disaster because of these values, and the other one does not deal with these values and properly predicts.