

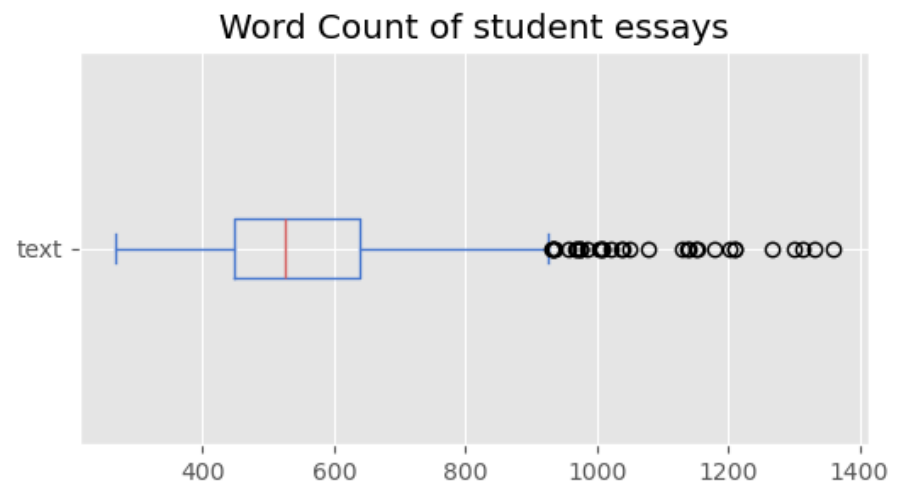
Kaggle Competition: Essay Classification Challenge

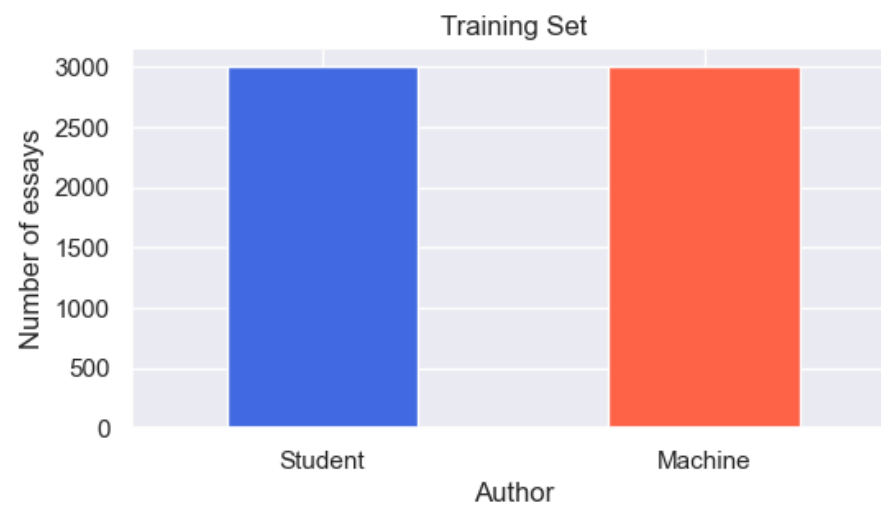
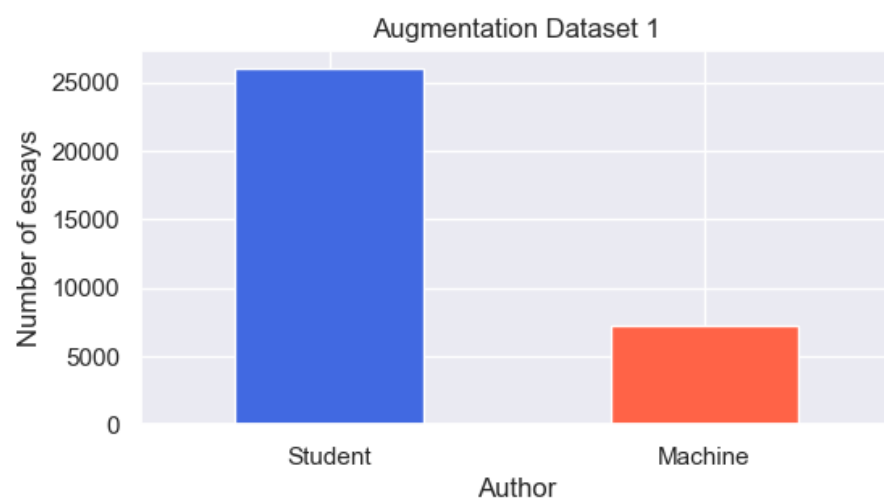
Practical Data Science - 2023

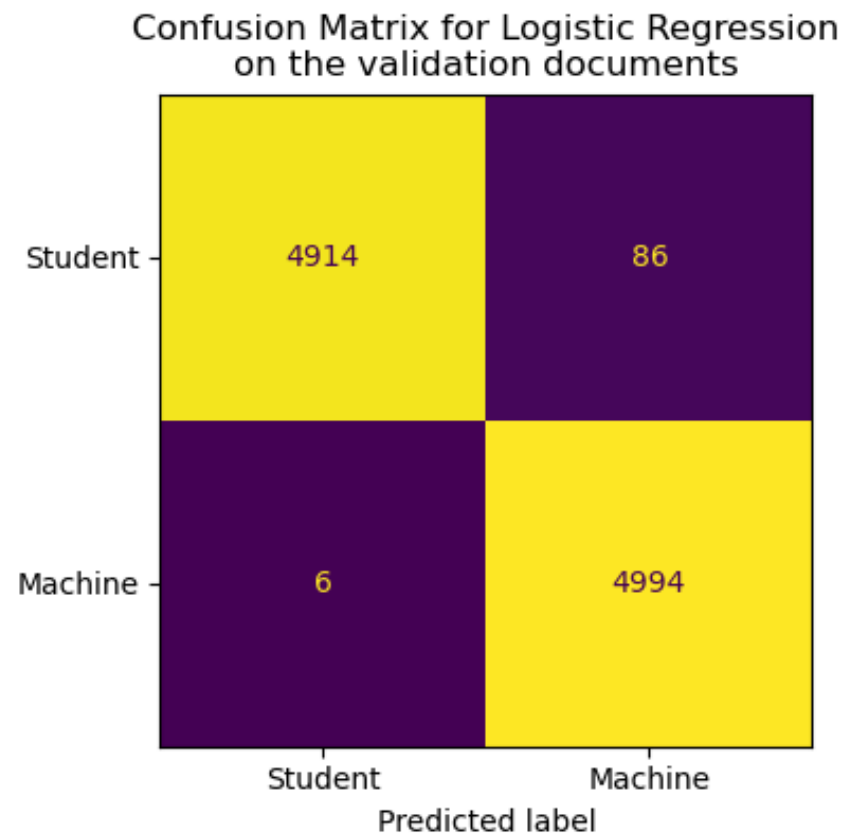
Phevos A. Margonis

Athens University of Economics and Business

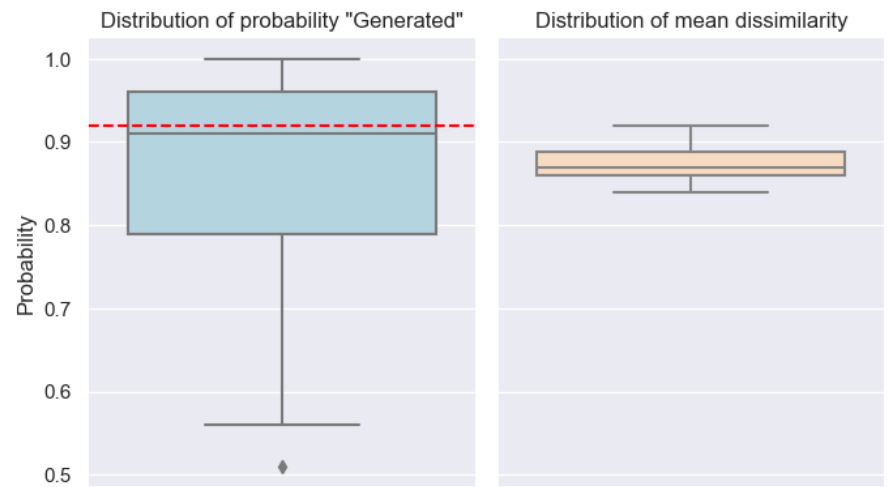
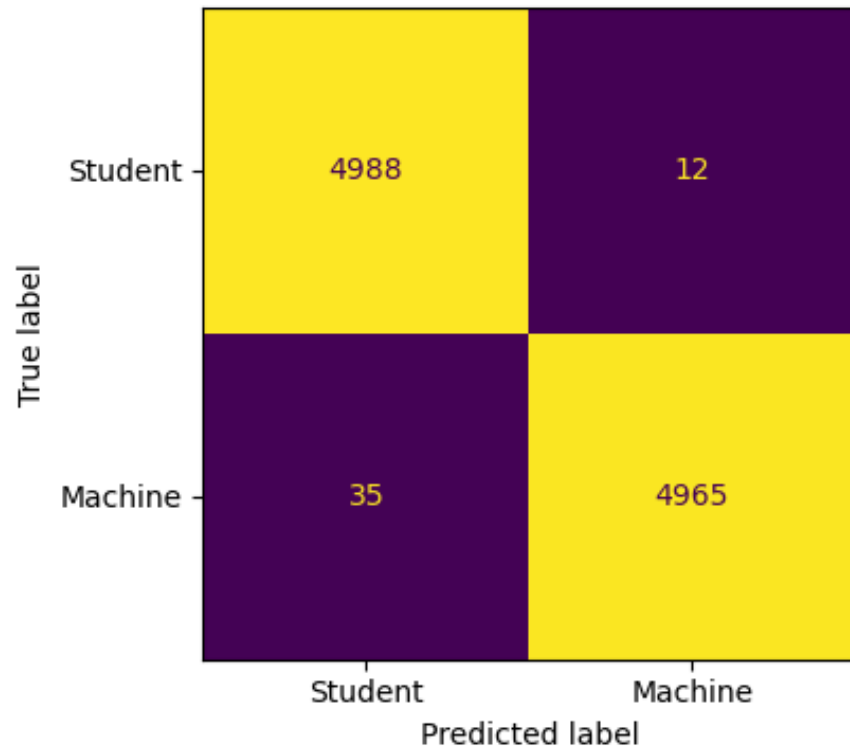
December 14, 2023



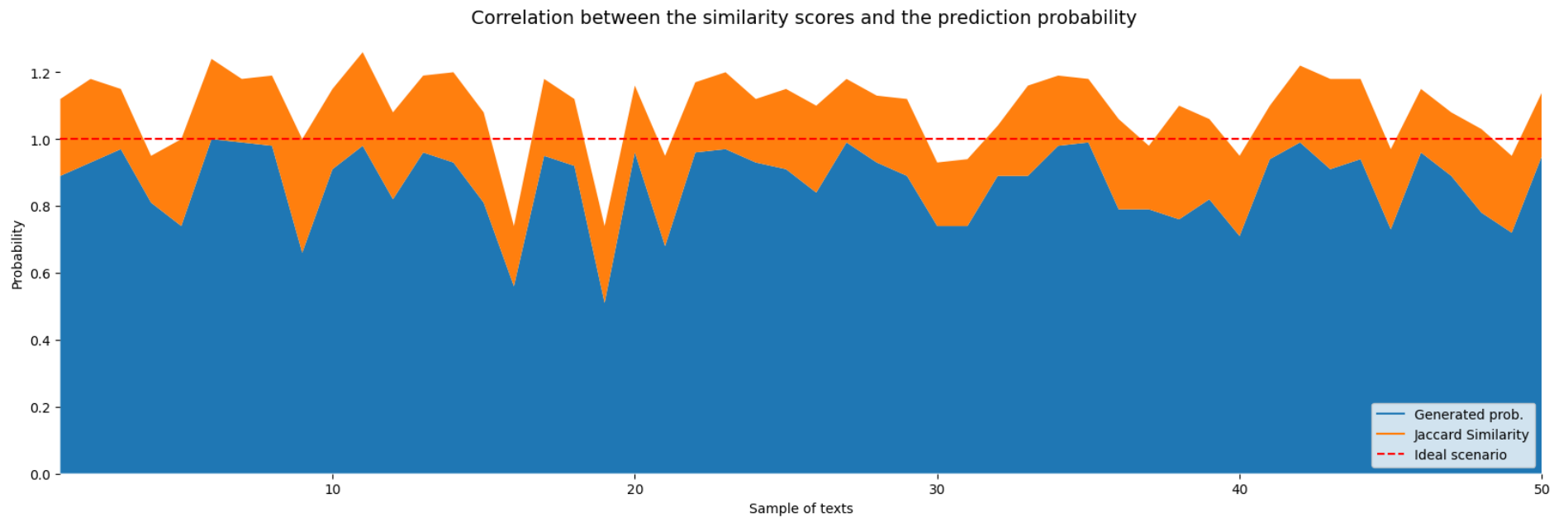




Confusion Matrix for Logistic Regression
on the validation documents

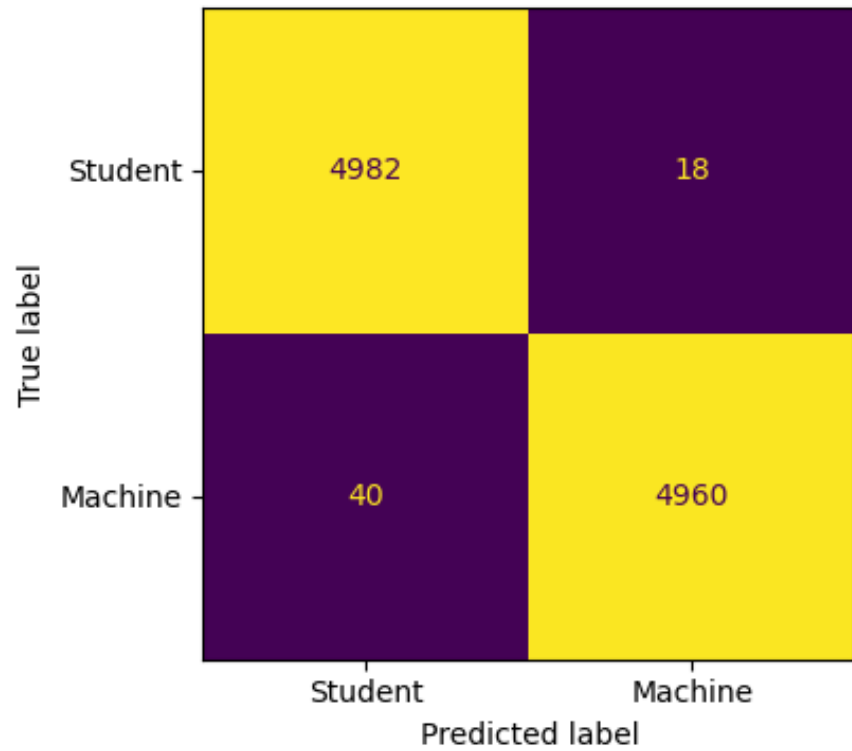


Thesis: The more dissimilar a generated text is from an original, the easier it is should be to detect.
By removing the texts with probability greater than the max average dissimilarity, we force the model to focus on the nuances of each text.

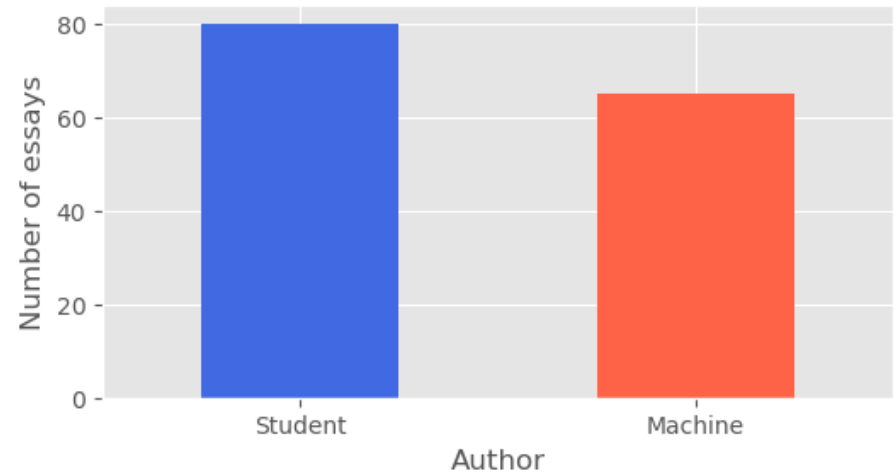


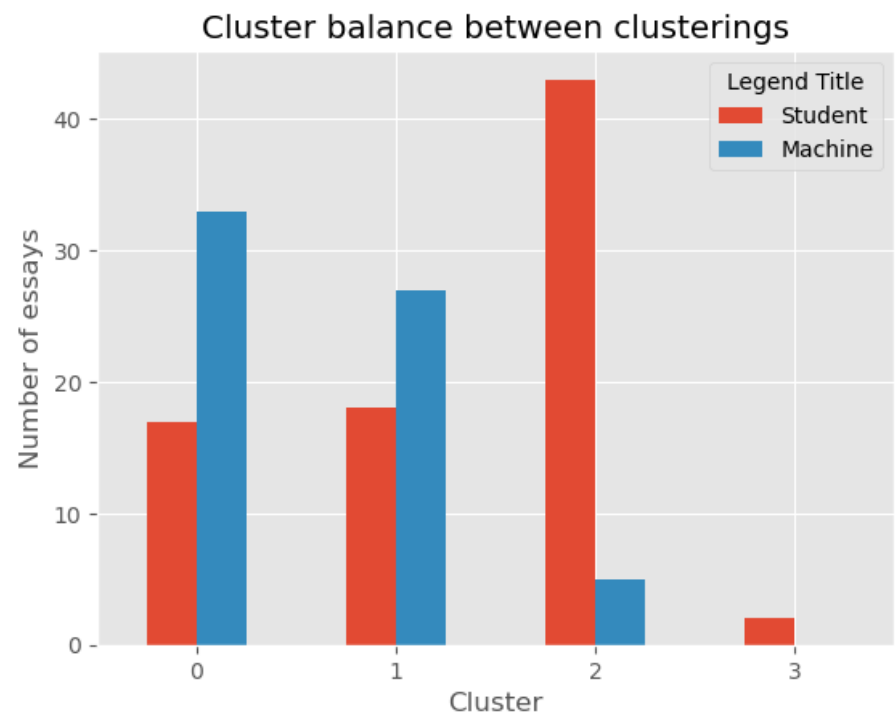
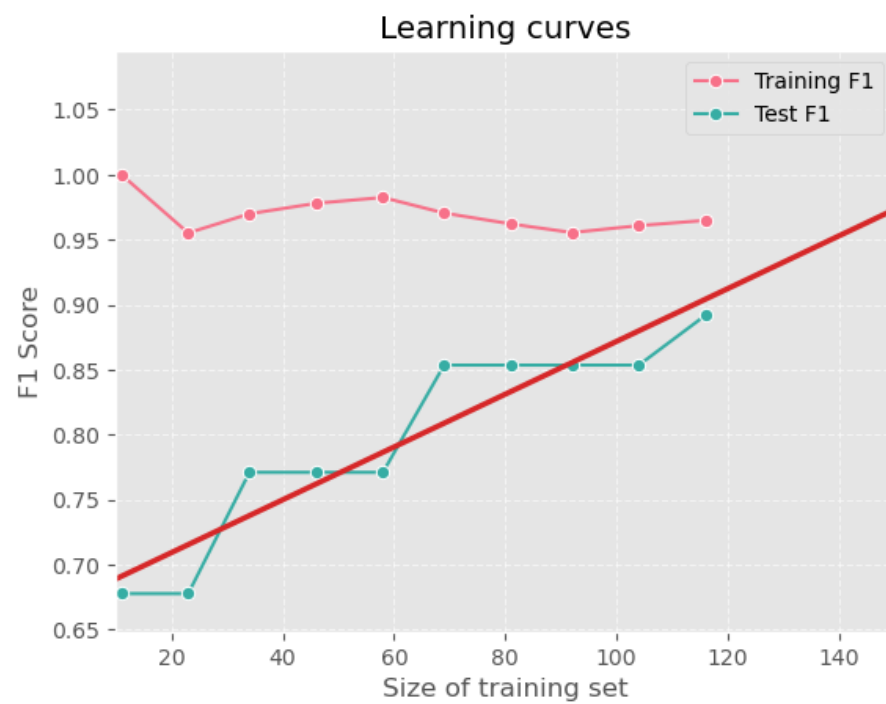
Thesis: The more similar pairs of student-LLM texts are, the harder it should be for the classifier to distinguish the True class.
 Following that logic, the correlation between the probability that a text is generated, should be inversely proportional to the maximum similarity of that text to all the rest.
 In the plot above, the ideal scenario is captured at points $x=5$, $x=9$, where the probabilities sum to one.

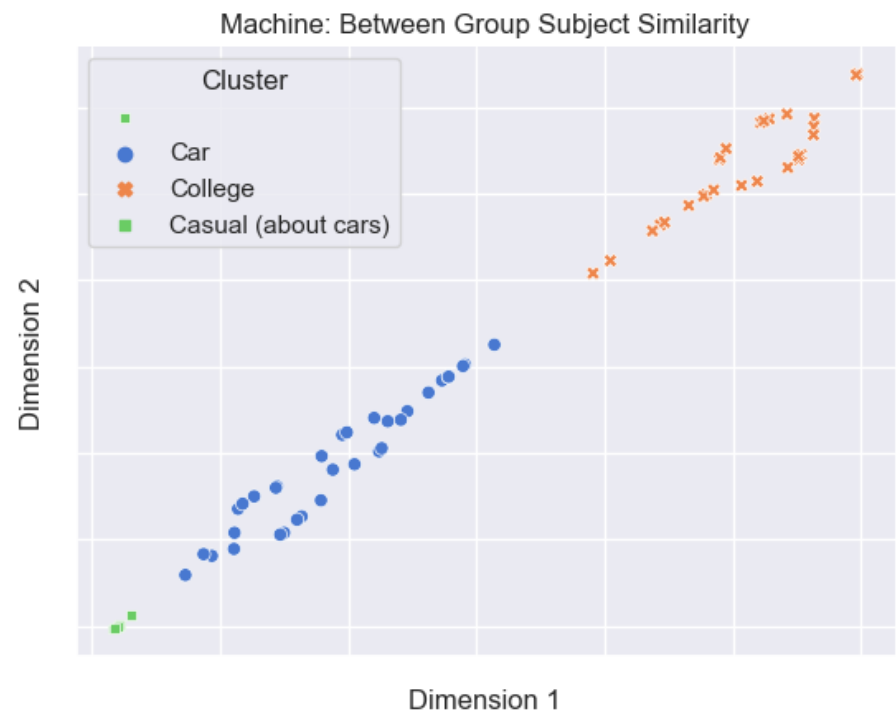
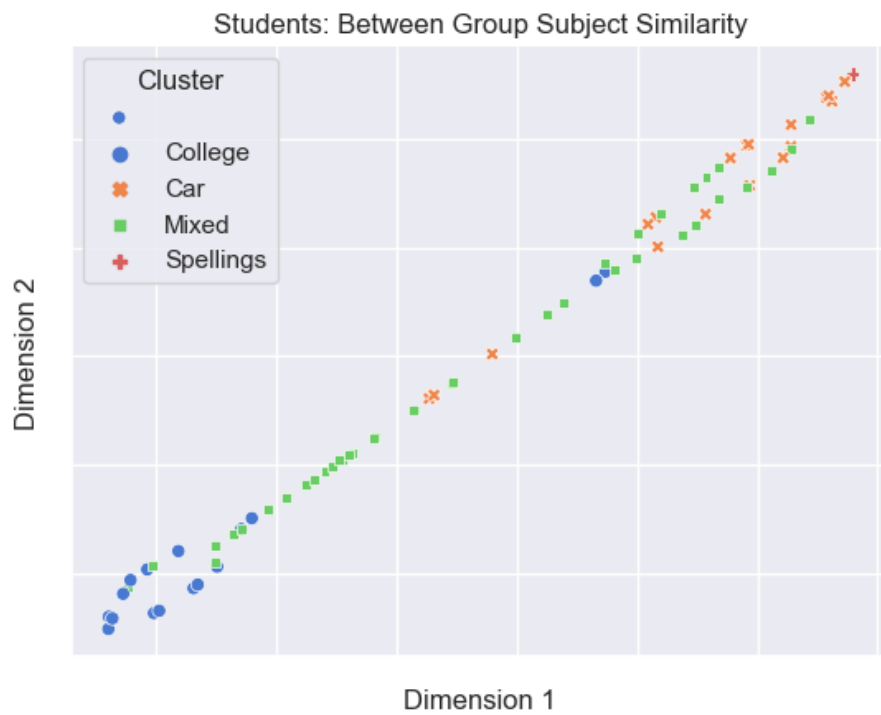
Confusion Matrix for Logistic Regression
on the validation documents



Training set







Similarities Between Clusterings

