With the rapid development of telecommunication industry, the service providers are inclined more towards expansion of the subscriber base. To meet the need of surviving in the competitive environment, the retention of existing customers has become a huge challenge. It is stated that the cost of acquiring a new customer is far more than that for retaining the existing one. Therefore, it is imperative for the telecom industries to use advanced analytics to understand consumer behavior and in-turn predict the association of the customers as whether or not they will leave the company. Consider the telecom_churn.csv data file posted on Blackboard (under the In-Class 2 assignment link). This data set contains customer level information for a telecom company. Various attributes related to the services used are recorded for each customer. In Python, answer the following:

- 1. (3 points) Using the pandas library, read the csv data file and create a data-frame called churn_data.
- 2. (4 points) Using AccountWeeks, ContractRenewal, CustServCalls, MonthlyCharge, and DayMins as the predictor variables, and Churn is the target variable, split the data into two data-frames (taking into account the proportion of 0s and 1s): train (80%) and test (20%).
- 3. (4 points) Create a synthetic training dataset, called them X_over and Y_over, by running over-sampling on the train dataset.
- 4. (8 points) Using X_over and Y_over datasets, build a random forest classification model with 500 trees and the maximum depth of each tree equal to 3. Estimate the cutoff value that makes the random forest classification model the closest to the perfect model based on the ROC curve. Using the optimal cutoff value, create the classification report.
- 5. (8 points) Using X_over and Y_over datasets, build a ada-boost classification model with 500 trees, the maximum depth of each tree equal to 3, and learning rate equal to 0.01. Estimate the cutoff value that makes the random forest classification model the closest to the perfect model based on the ROC curve. Using the optimal cutoff value, create the classification report.
- 6. (3 points) Using the results from part 4 and 5, what model would use to predict customer churn? Be specific.