# CS 584-04: Machine Learning

1. (20 points) Describe your strategies for developing the new models.
   I have chosen to use logistic regression for the region 1 and decision tree for region 2 and 3. Since at a company, explaining why the model is predicting a certain way is important for a business to decide whether to use that model or not, I thought it would be apt to use decision tree as it's very transparent. Also, visualizing and interpreting the tree is easy and this makes it easy to explain to executives or others. Decision trees make the feature importance clear and relations can be viewed easily.
   Initially, I planned to use decision tree for all 3 regions but it so happened that the it did not work well in the 1$^{st}$ region as the class distribution in the last node at a depth of 6 was still very uniform( it contained instances of both classes almost equally rather than a majority of one) and so opted to use logistic regression for region 1 as it is simple to understand too.

2. (20 points) Show how you selected the predictors into each new model.  Your objective is to exclude variables that do not contribute to the goodness-of-fit of the new model.

   For logistic regression the features were selected and ranked using recursive feature elimination. First, the estimator is trained on the initial set of features and the importance of each feature is obtained. Then, the least important features are pruned from current set of features and this procedure is recursively repeated until the number of features specified to select is reached. I chose to use 9 features and not too many. The output of RFE were the following features:
   Vibration, Engine_Coolant_Temp, Speed_OBD, Voltage_Control_Module, Ambient_air_temp, Engine_Oil_Temp, GPS_Longitude, GPS_Latitude, Accel_Ssor_Total

   Since decision trees themselves rank the features by the information gain and then use that feature to split upon, setting a reasonable depth ensures the best features which produce low entropy nodes after the split are used.
   Features used in region 2: Trip_Distance, Ambient_air_temp, Engine_RPM, GPS_Altitude, GPS_Latitude, Voltage_Control_Module, Engine_Load, Intake_Air_Temp

   In region 3, even though I set the max depth of the tree to be 6, the leaves are pure at depth 5 and hence, it doesn't proceed further.
   Features used in region 3: GPS_Latitude, Speed_OBD, Accel_Ssor_Total, GPS_Bearing, Mass_Air_Flow_Rate, Turbo_Boost_And_Vcm_Guage

3. (10 points) List the primary model specifications and the key model results (e.g., decision tree diagram, logistic regression, parameter estimates, support vector machine hyperplane equation, etc.)
   For regions 2 and 3, decision trees with a max_depth of 6 and the criterion of "entropy" is used

Region 2: Accuracy Score 2=  0.8741258741258742

```
                                    Trip_Distance <= 294.347
                                    entropy = 0.87
                                    samples = 1772
                                    value = [1256, 516]
                                    class = 0
                         True                          False
                    Ambient_air_temp <= 8.5      entropy = 0.0
                    entropy = 0.997              samples = 663
                    samples = 1109              value = [663, 0]
                    value = [593, 516]          class = 0
                    class = 0
              Engine_RPM <= 1599.5        entropy = 0.0
              entropy = 0.998            samples = 134
              samples = 975             value = [134, 0]
              value = [459, 516]         class = 0
              class = 1
         GPS_Altitude <= 565.5       entropy = 0.0
         entropy = 0.987            samples = 67
         samples = 908             value = [67, 0]
         value = [392, 516]         class = 0
         class = 1
      GPS_Latitude <= 48.715        Engine_RPM <= 1302.25
      entropy = 0.972              entropy = 0.662
      samples = 844               samples = 64
      value = [339, 505]          value = [53, 11]
      class = 1                   class = 0
  Voltage_Control_Module <= 14.33   Engine_Load <= 9.216   Intake_Air_Temp <= 10.5   entropy = 0.0
  entropy = 0.962                   entropy = 0.391         entropy = 0.928          samples = 32
  samples = 818                    samples = 26            samples = 32             value = [32, 0]
  value = [315, 503]               value = [24, 2]         value = [21, 11]          class = 0
  class = 1                        class = 0               class = 0
entropy = 0.959  entropy = 0.0   entropy = 0.918  entropy = 0.0   entropy = 0.954  entropy = 0.337
samples = 814    samples = 4     samples = 3      samples = 23    samples = 16     samples = 16
value = [311, 503] value = [4, 0] value = [1, 2]   value = [23, 0] value = [6, 10]  value = [15, 1]
class = 1        class = 0       class = 1        class = 0       class = 1        class = 0
```
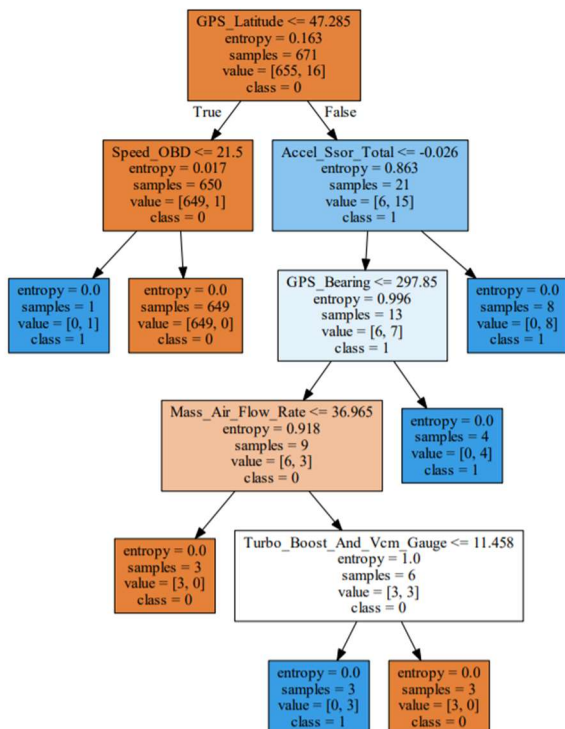
Region 3: Accuracy Score 3=  1.0

```
                         GPS_Latitude <= 47.285
                         entropy = 0.163
                         samples = 671
                         value = [655, 16]
                         class = 0
                 True                    False
            Speed_OBD <= 21.5      Accel_Ssor_Total <= -0.026
            entropy = 0.017        entropy = 0.863
            samples = 650          samples = 21
            value = [649, 1]       value = [6, 15]
            class = 0              class = 1
    entropy = 0.0   entropy = 0.0    GPS_Bearing <= 297.85   entropy = 0.0
    samples = 1     samples = 649    entropy = 0.996         samples = 8
    value = [0, 1]  value = [649, 0] samples = 13           value = [0, 8]
    class = 1       class = 0        value = [6, 7]          class = 1
                                     class = 1
                    Mass_Air_Flow_Rate <= 36.965   entropy = 0.0
                    entropy = 0.918                samples = 4
                    samples = 9                   value = [0, 4]
                    value = [6, 3]                 class = 1
                    class = 0
            entropy = 0.0    Turbo_Boost_And_Vcm_Gauge <= 11.458
            samples = 3      entropy = 1.0
            value = [3, 0]   samples = 6
            class = 0        value = [3, 3]
                             class = 0
                    entropy = 0.0   entropy = 0.0
                    samples = 3     samples = 3
                    value = [0, 3]  value = [3, 0]
                    class = 1       class = 0
```

Region 1 : Logistic Regression;  Accuracy Score =  0.8443067389620449
The coefficient values are :
Vibration -0.23324671
Engine_Coolant_Temp 0.26943419
Speed_OBD -0.04039195
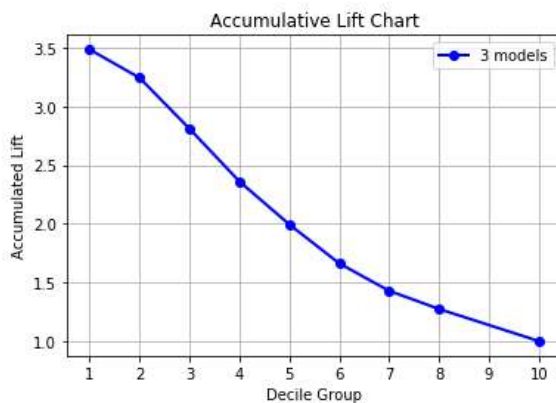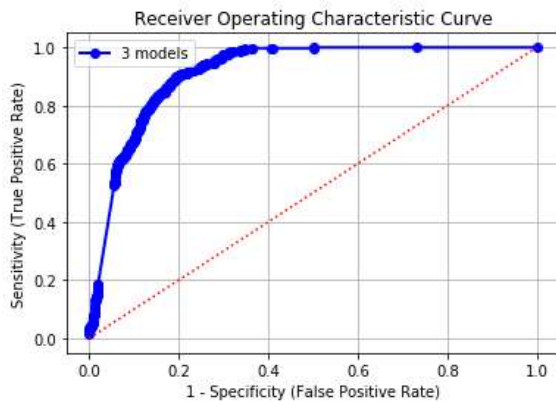Voltage_Control_Module  2.43774278
Ambient_air_temp -0.660558
Engine_Oil_Temp 0.0393996
GPS_Longitude -1.52197286
GPS_Latitude 0.31678666
Accel_Ssor_Total -0.26052172

4.  (30 points) Show the model comparison results and list all supporting tables and charts (e.g., Area Under Curve, Root Average Squared Error, Misclassification Rate, ROC curve, Lift or Accumulated Lift curves)



Area Under Curve = 0.915815
Root Average Squared Error = 0.307828
Misclassification Rate = 0.179514

5. (20 points) Argue that you have actually found a better model than the current model. In the current model, the trucks are sent for maintenance only based on their region and the number of days elapsed, both of which are not strong indicators if the truck needs maintenance. However, the models I have used considers various parameters of the truck, learns this data to uses it to predict when maintenance is required. The intuition behind this is that some of the parameters indicate wear of different features of the truck and hence my model is learning when the truck is worn enough to warrant maintenance. Also my model significantly increases the area under the curve metric and lowers the RASE, while reducing the Misclassification rate on the test data by a wide margin.