

Homework 2)

1.1

$$NLL(w) = - \sum_{i=1}^n (1 - y_i) * \log(1 - \sigma(\mathbf{w}^T \mathbf{x})) + y_i * \log(\sigma(\mathbf{w}^T \mathbf{x}))$$

$$-NLL(w) = \sum_{i=1}^n (1 - y_i) * \log(1 - \sigma(\mathbf{w}^T \mathbf{x})) + y_i * \log(\sigma(\mathbf{w}^T \mathbf{x}))$$

Taking the derivative of each piece of the summation separately (with respect to \mathbf{w}), we get:

$$d/dw [(1 - y_i) * \log(1 - \sigma(\mathbf{w}^T \mathbf{x}))] = [(1 - y_i) * 1 / (1 - \sigma(\mathbf{w}^T \mathbf{x}))] * \sigma(\mathbf{w}^T \mathbf{x}) * (1 - \sigma(\mathbf{w}^T \mathbf{x})) * x * -1$$

$$= (1 - y_i) * \sigma(\mathbf{w}^T \mathbf{x}) * x * -1$$

$$= x * (y_i - 1) * \sigma(\mathbf{w}^T \mathbf{x})$$

$$= xy_i \sigma(\mathbf{w}^T \mathbf{x}) - x \sigma(\mathbf{w}^T \mathbf{x})$$

Now, taking the derivative of the other piece:

$$d/dw [y_i * \log(\sigma(\mathbf{w}^T \mathbf{x}))] = y_i * (1 / \sigma(\mathbf{w}^T \mathbf{x})) * \sigma(\mathbf{w}^T \mathbf{x}) * (1 - \sigma(\mathbf{w}^T \mathbf{x})) * x$$

$$= y_i * (1 - \sigma(\mathbf{w}^T \mathbf{x})) * x$$

$$= xy_i - xy_i \sigma(\mathbf{w}^T \mathbf{x})$$

Now, adding this derivative with the above result yields:

$$xy_i \sigma(\mathbf{w}^T \mathbf{x}) - x \sigma(\mathbf{w}^T \mathbf{x}) + xy_i - xy_i \sigma(\mathbf{w}^T \mathbf{x}) = xy_i - x \sigma(\mathbf{w}^T \mathbf{x})$$

$$= x(y_i - \sigma(\mathbf{w}^T \mathbf{x}))$$

So the final result is:

$$d/dw [-NLL(w)] = \sum_{i=1}^n x(y_i - \sigma(\mathbf{w}^T \mathbf{x}))$$

$$\text{Therefore, } d/dw [NLL(w)] = - \sum_{i=1}^n x(y_i - \sigma(\mathbf{w}^T \mathbf{x}))$$

1.2

a.

$$\sum_{i=1}^n (1 - y_i) * \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) + y_i * \log(\sigma(\mathbf{w}^T \mathbf{x}_i))$$

b. Update \mathbf{w}_{t-1} by:

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta * (\mathbf{y}_t - \sigma(\mathbf{w}_{t-1}^T \mathbf{x}_t)) * \mathbf{x}_t$$

c. $O(k * \log(d))$, where k is the number of nonzero elements of \mathbf{x} and d is the total number of features in \mathbf{x}

d.

Using a large value of η (learning rate) may cause the optimal solution to be skipped over as the iterations occur. A small value of η will increase the time of convergence, and thus, cause too many iterations to result in the optimal solution.

e.

Update each w_j in \mathbf{w}_t by the following rule:

$\mathbf{w}_{\text{updated_j}} = w_j + \eta * ((y_t - \sigma(\mathbf{w}^T \mathbf{x}_t)) * x_j - (2 * \mu * w_j))$ where \mathbf{x}_t is the corresponding vector of feature values for the t^{th} patient, y_t is the label for the t^{th} patient, x_j is the feature value associated with weight w_j , and η is the learning rate

The time complexity is $O(n)$ where n is the number of features.

f.

For x_{ti} in \mathbf{x}_t :

 If x_{ti} not zero:

 Update \mathbf{w}_{ti} by:

$$\mathbf{w}_{\text{updated_ti}} = \mathbf{w}_{ti} + \eta * ((y_t - \sigma(\mathbf{w}^T \mathbf{x}_t)) * x_{ti} - (2 * \mu * \mathbf{w}_{ti}))$$

 Else

 Do Nothing

2.2

Metric	Alive Patients	Deceased Patients
Event Count		
Average	682.6741	1029.059
Max	12627	16829
Min	1	2
Encounter Count		
Average	18.6694	24.861
Max	391	375
Min	1	1
Record Length		
Average	194.6541	151.397
Max	3103	2601
Min	0	0

Common Diagnosis

Alive DIAG

DIAG320128	1019
DIAG319835	721
DIAG317576	719
DIAG42872402	674
DIAG313217	641

Dead DIAG

DIAG320128	415
DIAG319835	413
DIAG313217	374
DIAG197320	346
DIAG132797	297

Common Laboratory Test

Alive LAB

LAB3009542	66910
LAB3000963	57733
LAB3023103	56967
LAB3018572	54667
LAB3007461	53548

Dead LAB

LAB3009542	32747
LAB3023103	28376
LAB3000963	28288
LAB3018572	27364
LAB3016723	27041

Common Medication

Alive DRUG

DRUG19095164	12452
DRUG43012825	10388
DRUG19049105	9329
DRUG19122121	7586
DRUG956874	7294

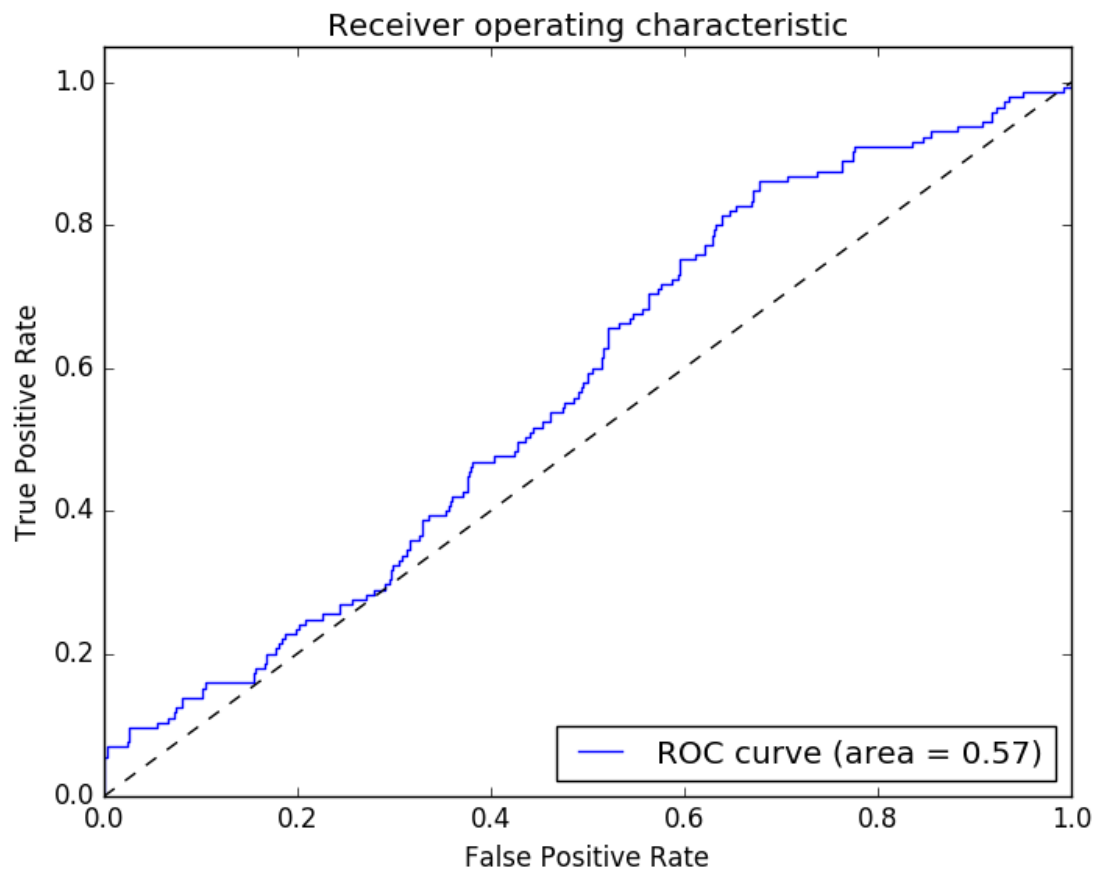
Dead DRUG

DRUG19095164	6394
DRUG43012825	5446
DRUG19049105	4323
DRUG956874	3962
DRUG19122121	3908

2.3

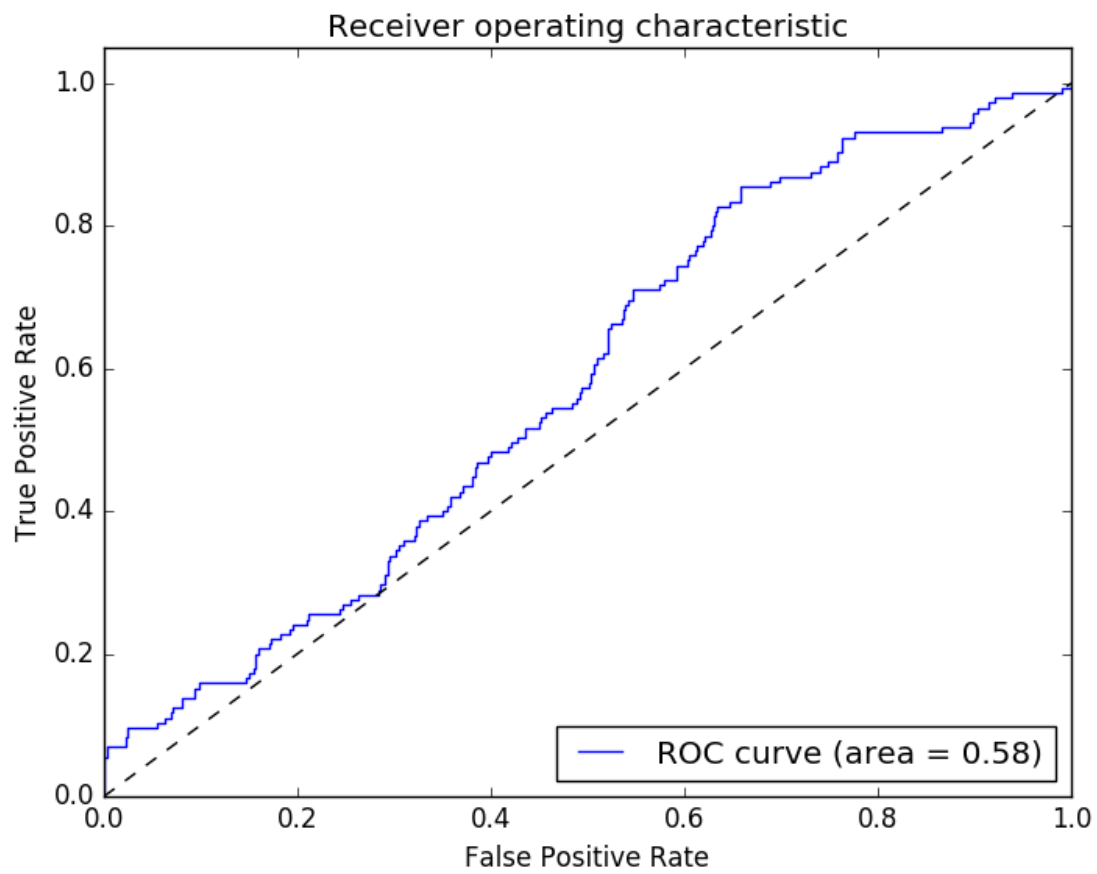
Learning Rate: .01

$\mu = .01$



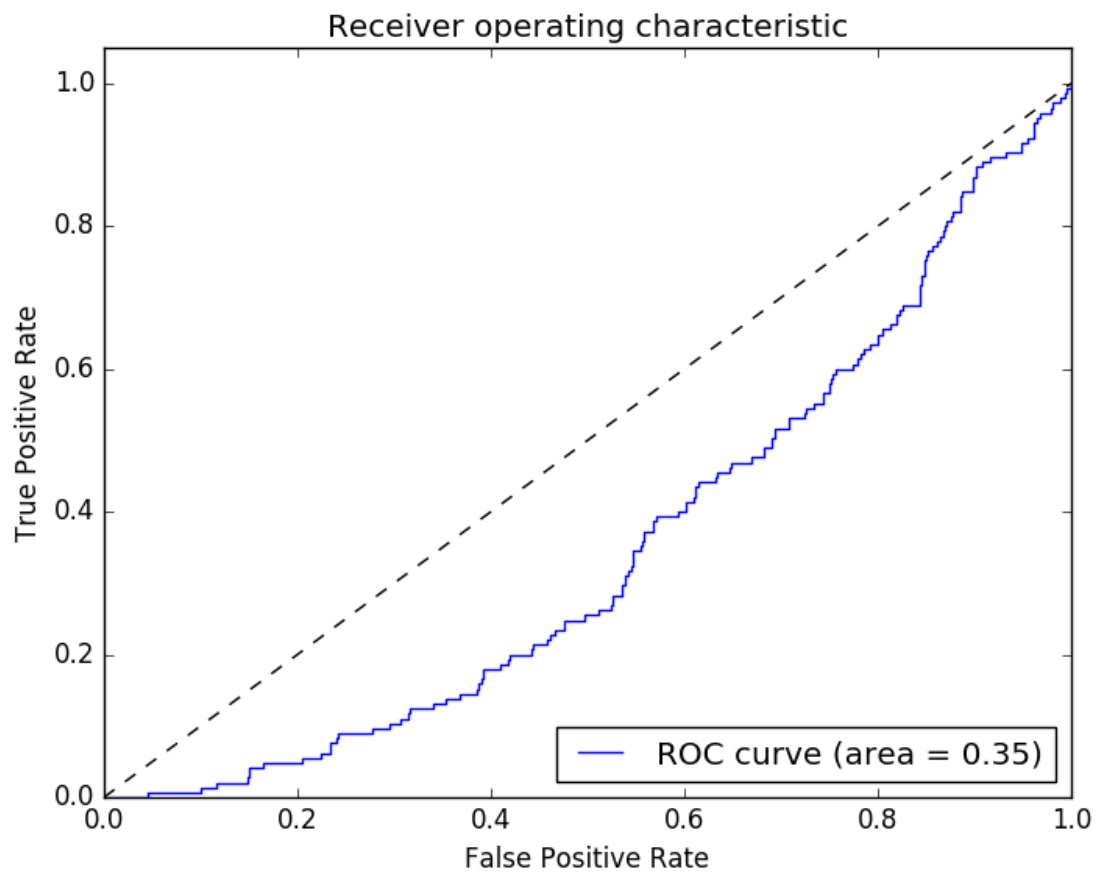
Learning Rate: .01

$\mu = .0001$



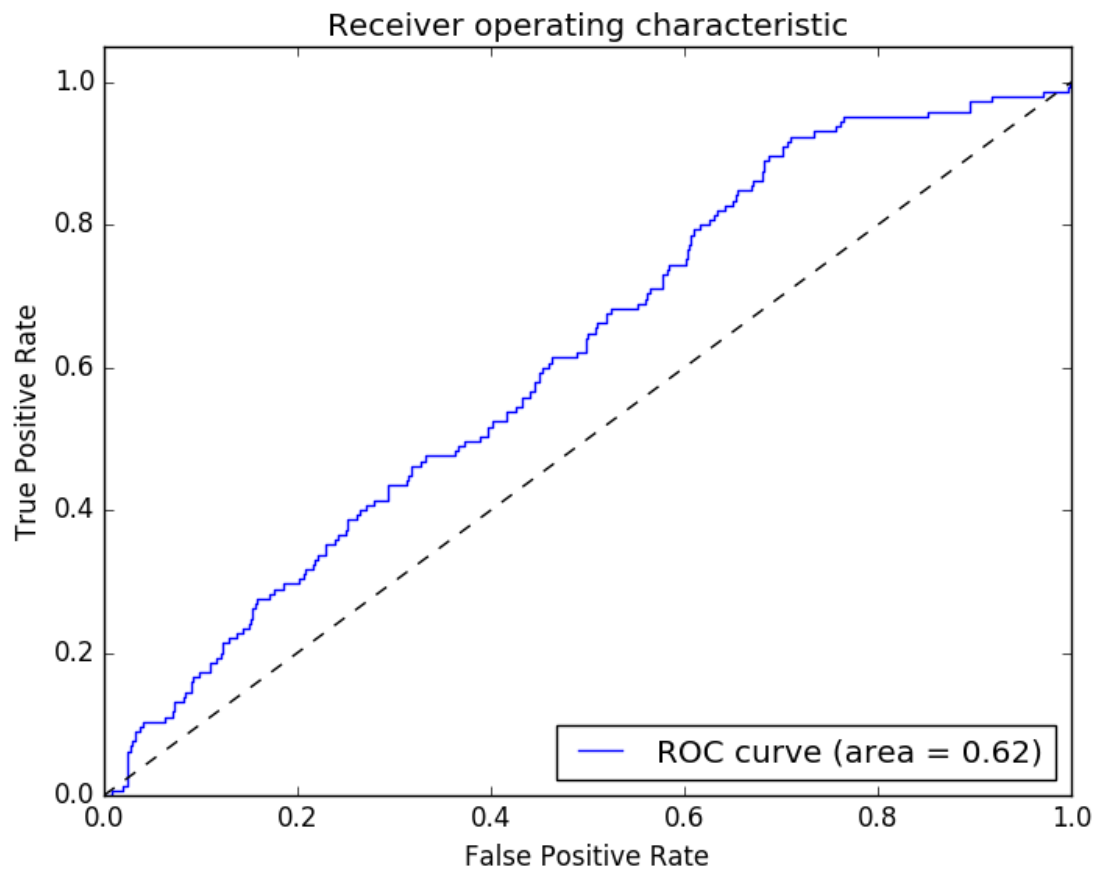
Learning Rate: .50

$\mu = .10$



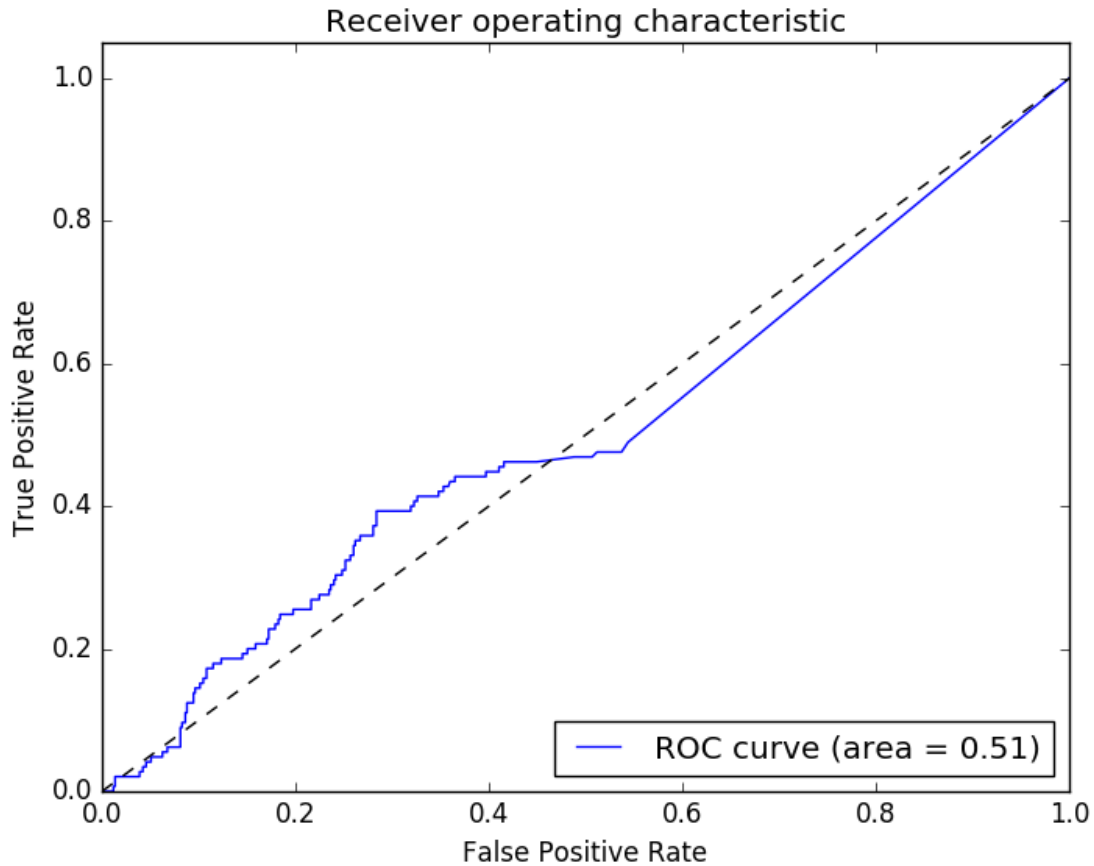
Learning Rate: 0.25

$\mu = .0001$



Learning Rate: 5.0

$\mu = 10^{-6}$

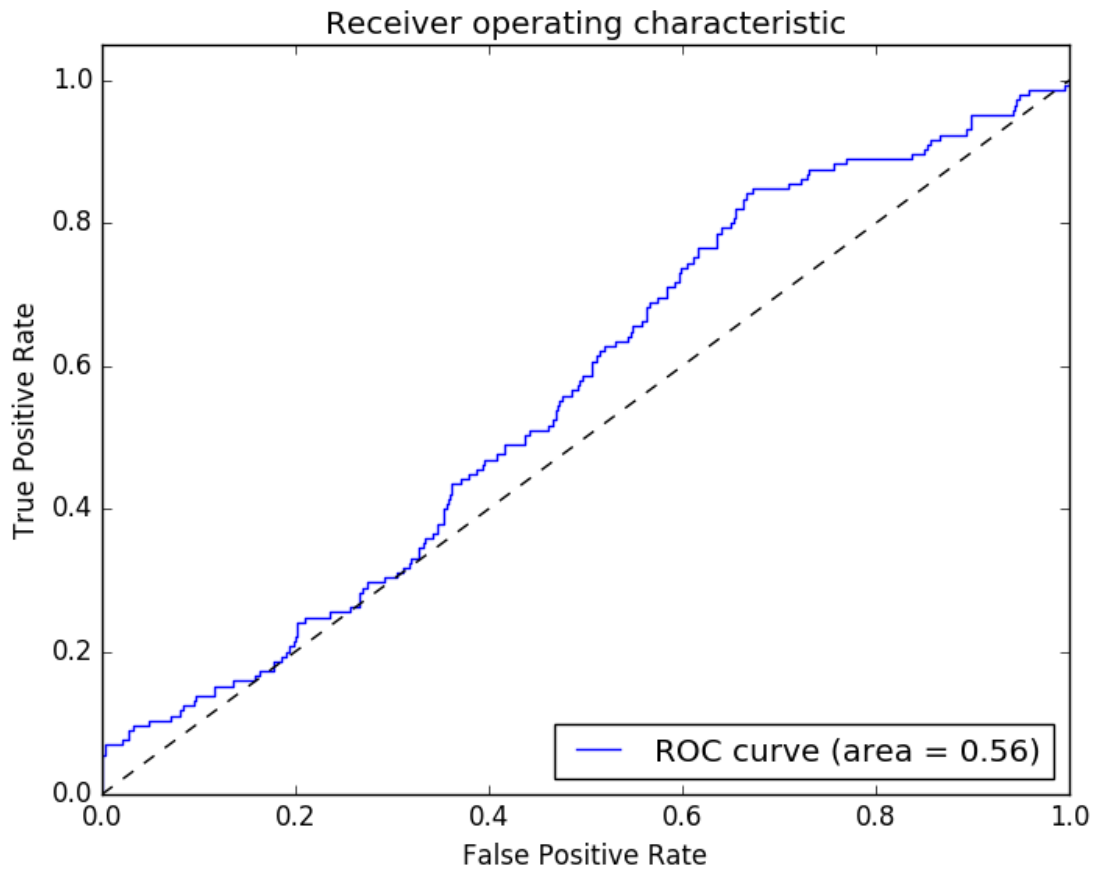


The results in the ROC charts above show the need to tune the learning parameter, η . A more aggressive (greater) value for the learning parameter is leading to worse results, which is most likely caused by the fact that the optimal solution for the logistic regression problem is being passed over by choosing too high of a value. In general, lower values of μ result in better AUC values. This is due to the fact that lower μ values result in less of an L2 norm penalty.

2.4

b.

ROC curve using five models.



c.

The ensemble methods vary depending upon the number of models used. A larger number of models results in smaller subsets of the data being trained, which results in lower AUC performance than the single logistic regression model used initially. The AUC from the ROC curve chart above results in roughly the same range as the initial runs of the logistic regression model.