

The Problem:

Breast Cancer Diagnosis

The Problem:

Breast Cancer Diagnosis:
False Negatives

(minimize)

The Data:

31 Variable: Diagnosis ("M" of "B") and 30 Data Points

- a) radius (mean of distances from center to points on the perimeter)
- b) texture (standard deviation of gray-scale values)
- c) perimeter
- d) area
- e) smoothness (local variation in radius lengths)
- f) compactness ($\text{perimeter}^2 / \text{area} - 1.0$)
- g) concavity (severity of concave portions of the contour)
- h) concave points (number of concave portions of the contour)
- i) symmetry
- j) fractal dimension ("coastline approximation" - 1)

Simple Models

Simple Model Results

TrainLogReg	
True Neg 256 63.05%	False Pos 7 1.72%
False Neg 4 0.99%	True Pos 139 34.24%

TrainKmeans	
True Neg 259 63.79%	False Pos 63 15.52%
False Neg 1 0.25%	True Pos 83 20.44%

TrainNeurNet	
True Neg 260 64.04%	False Pos 11 2.71%
False Neg 0 0.00%	True Pos 135 33.25%

DevLogReg	
True Neg 48 57.14%	False Pos 3 3.57%
False Neg 2 2.38%	True Pos 31 36.90%

DevKmeans	
True Neg 50 59.52%	False Pos 13 15.48%
False Neg 0 0.00%	True Pos 21 25.00%

DevNeurNet	
True Neg 48 57.14%	False Pos 1 1.19%
False Neg 2 2.38%	True Pos 33 39.29%

Ideal Case

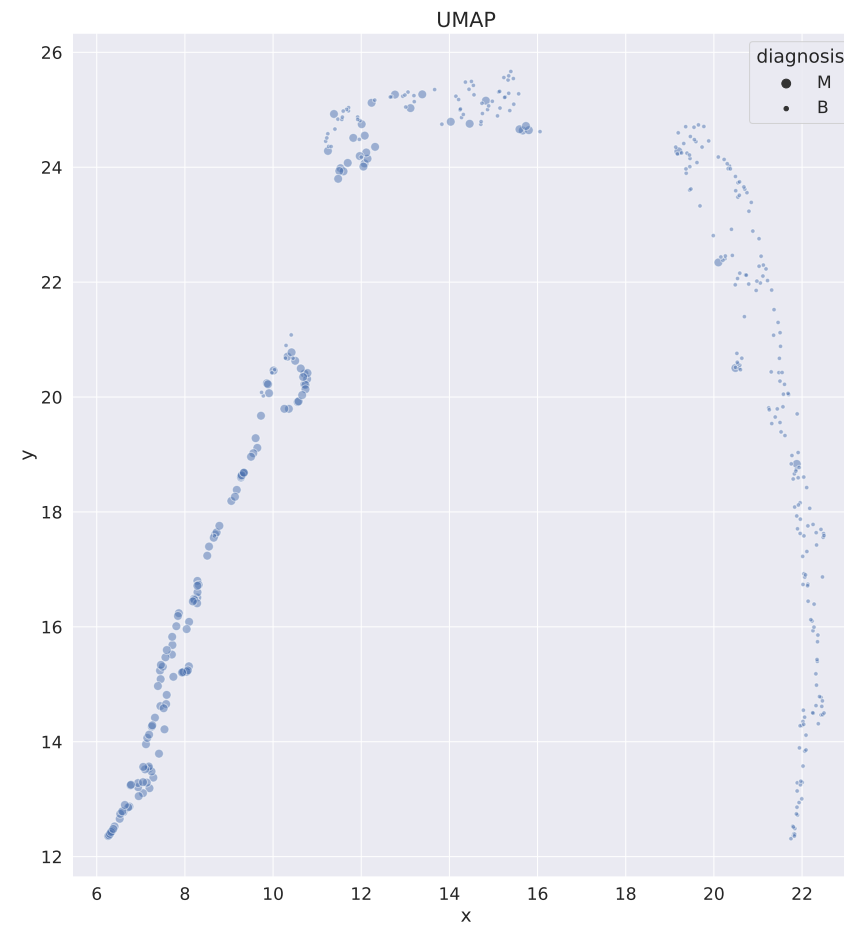
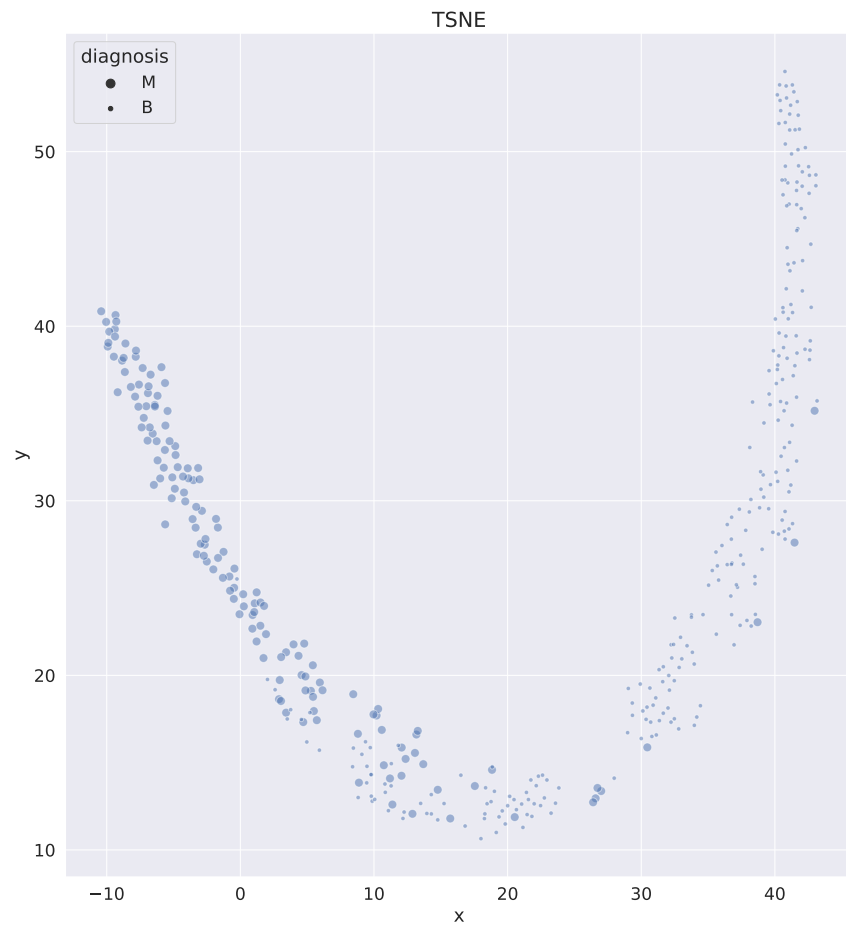
(Whiteboard)

The Ensemble

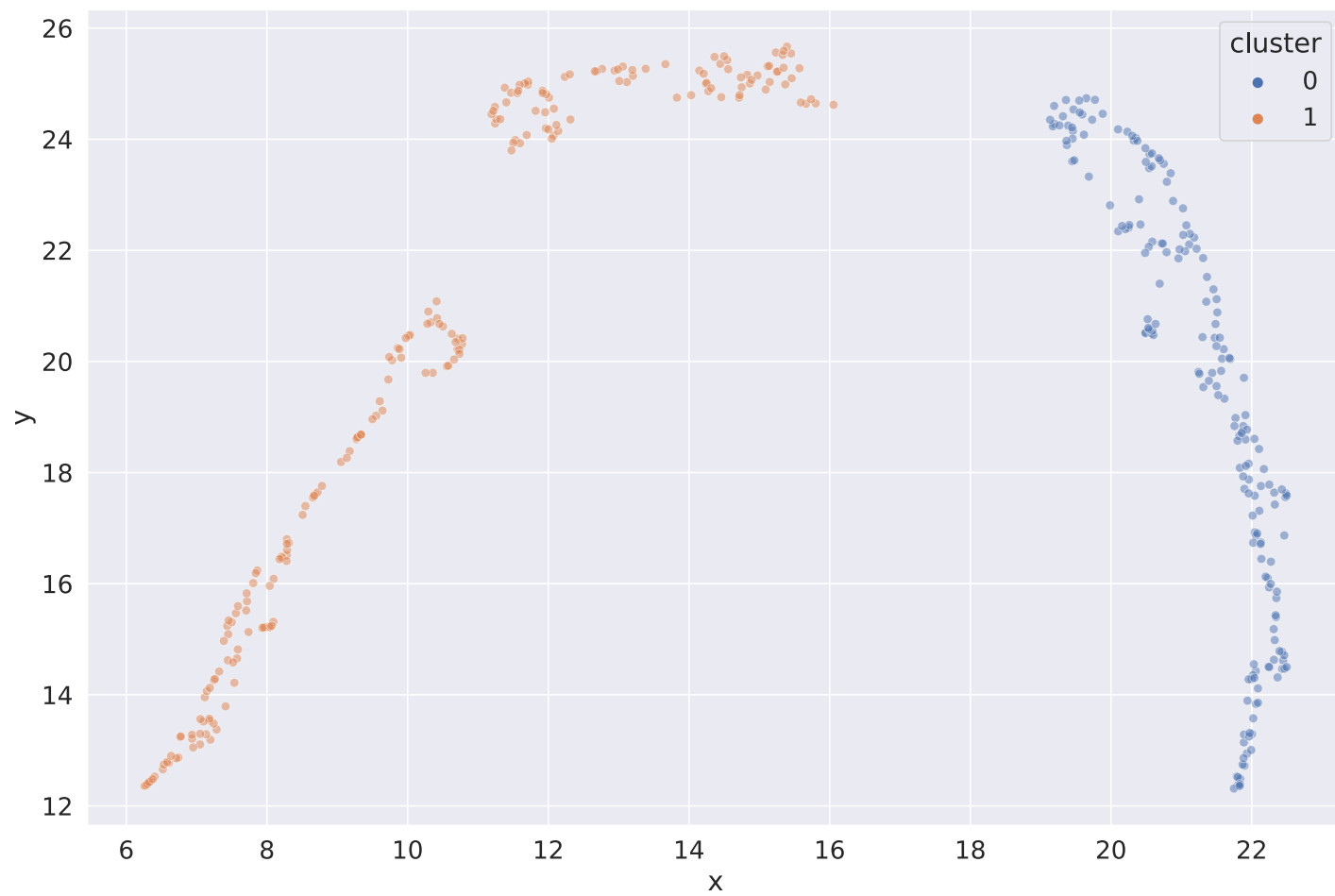
(Whiteboard)

TSNE vs UMAP

TSNE vs UMAP

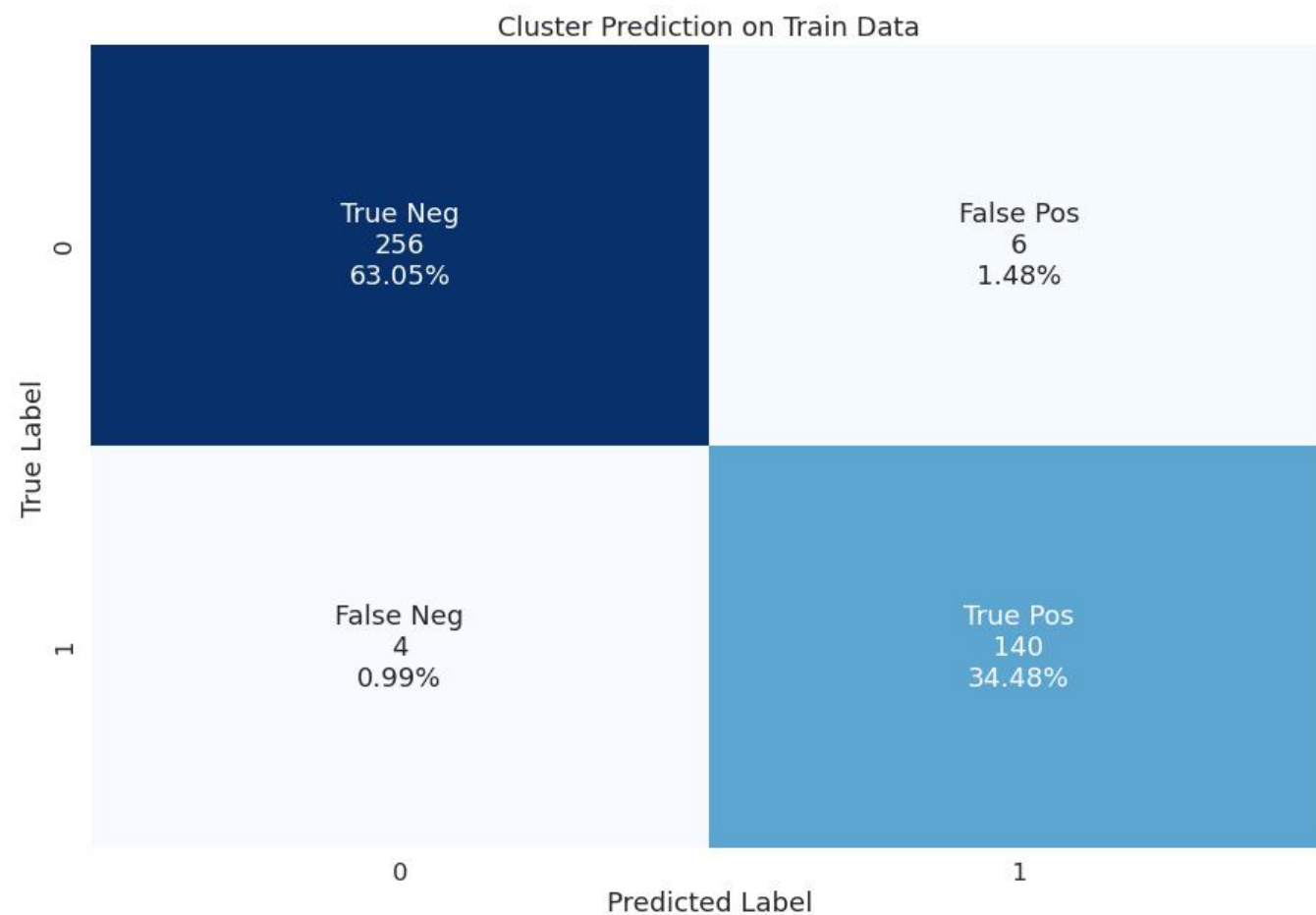


Train Kmeans



Train NN on each Cluster

To predict classification



Problems with UMAP

(Streamlit)

UMAP Predict Model

(Whiteboard)

Test Loss Functions:

MSLE

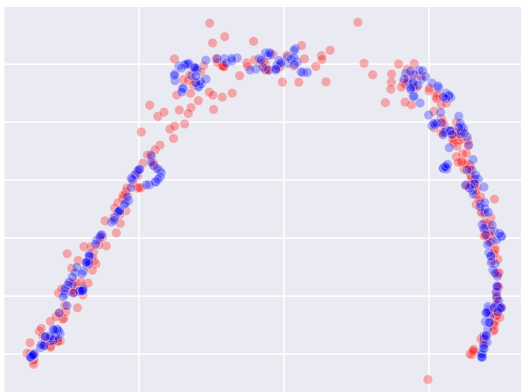
MAPE

Huber

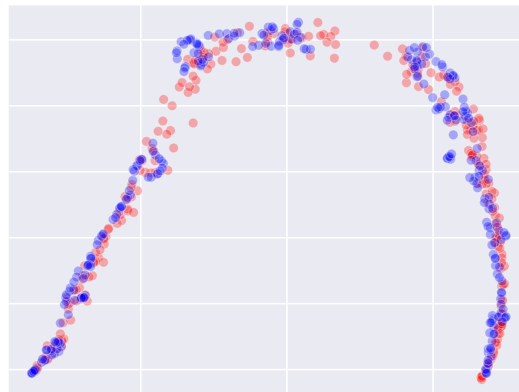
UMAP Parameters

Neural Net on UMAP Results

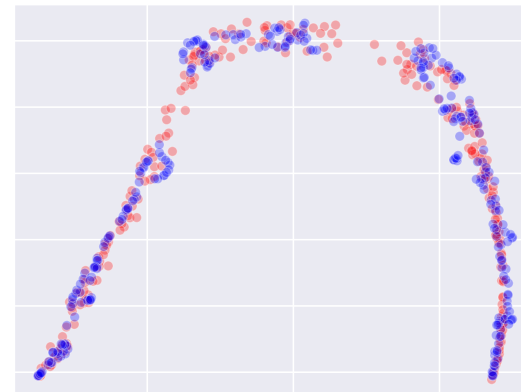
MSLE - Train



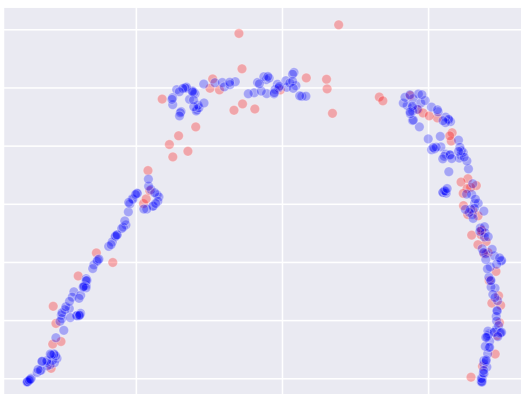
MAPE - Train



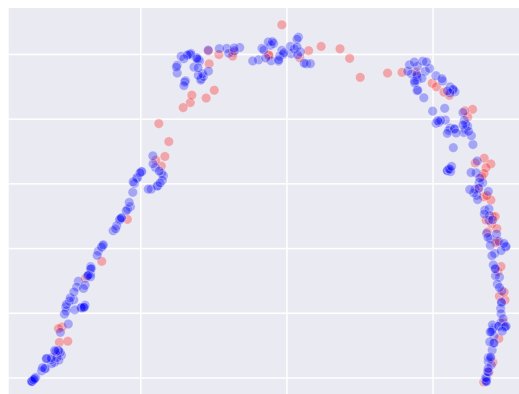
Huber - Train



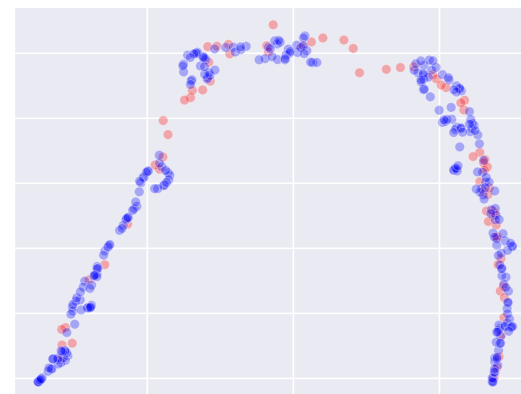
MSLE - Dev



MAPE - Dev



Huber - Dev

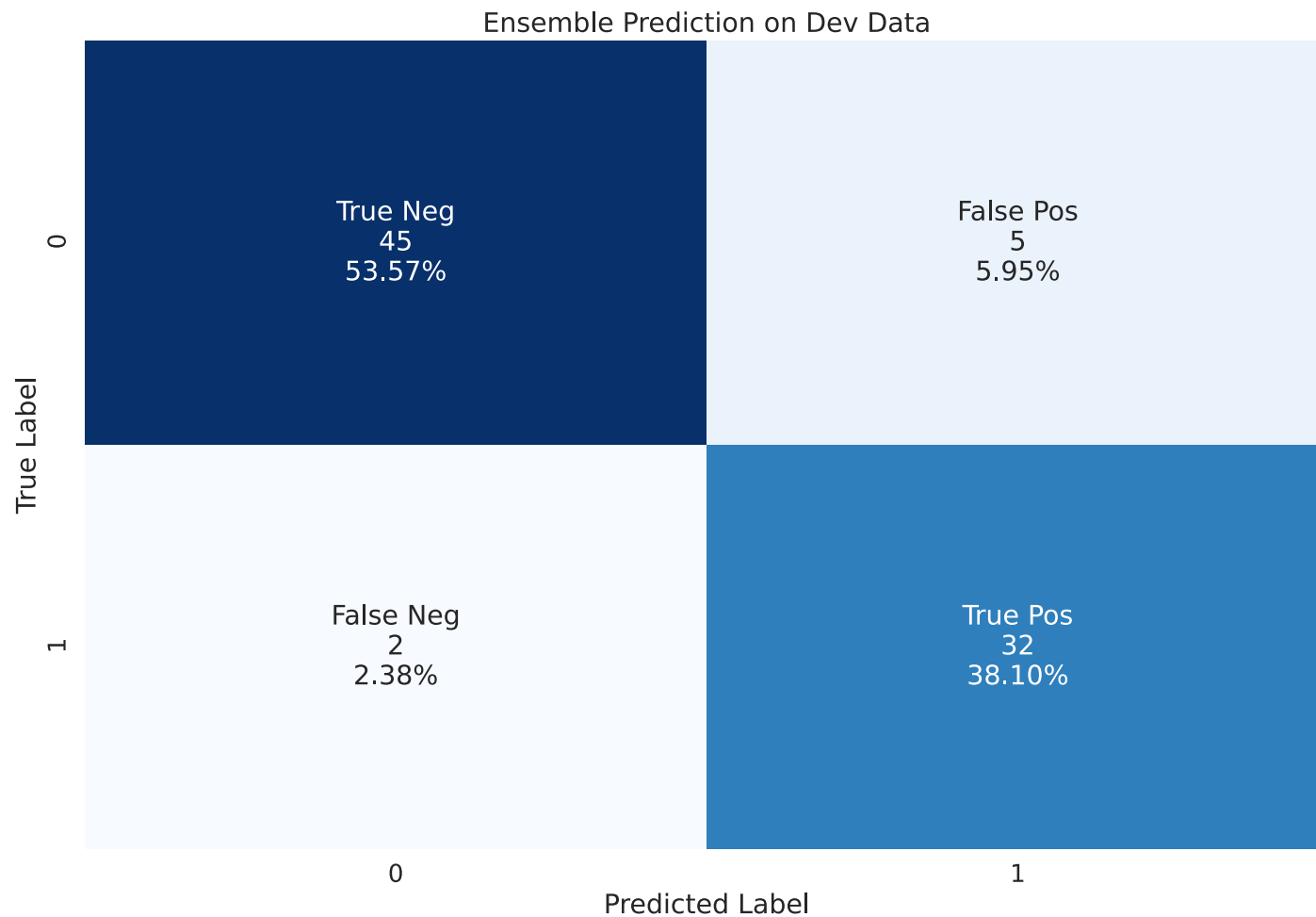


UMAP Parameters

the KMeans model was able to predict 66 correct memberships with the MSLE loss function the Dev set
the KMeans model was able to predict 67 correct memberships with the MAPE loss function the Dev set
the KMeans model was able to predict 67 correct memberships with the Huber loss function the Dev set
Total observations in Dev set is: 70

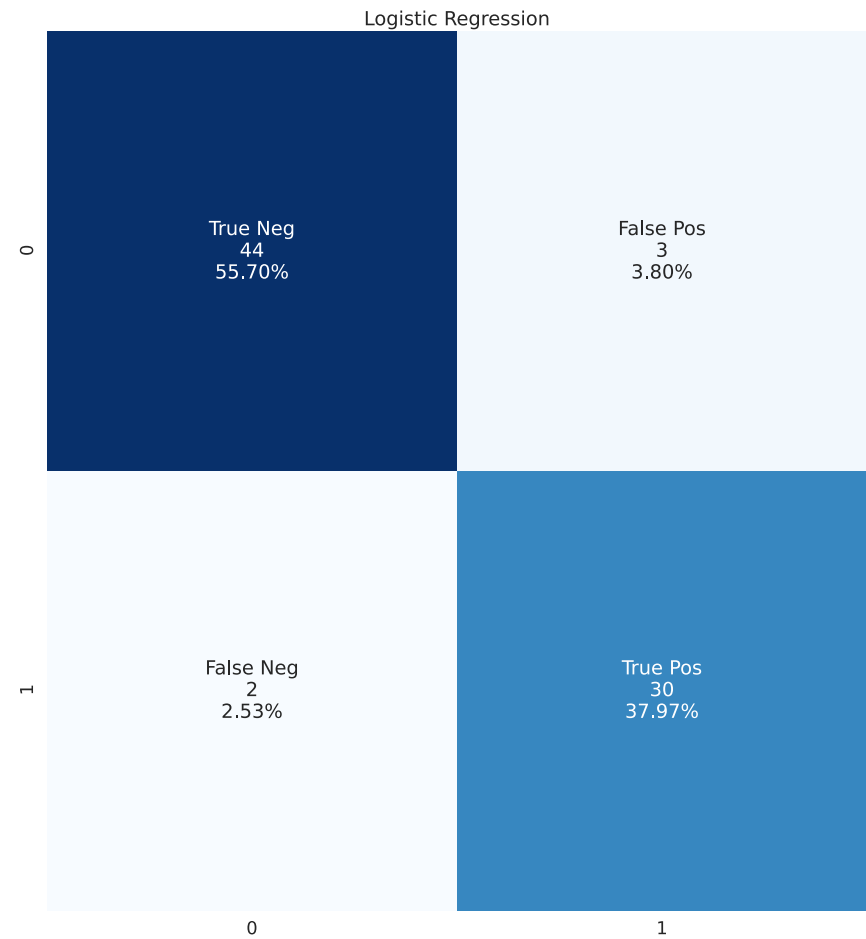
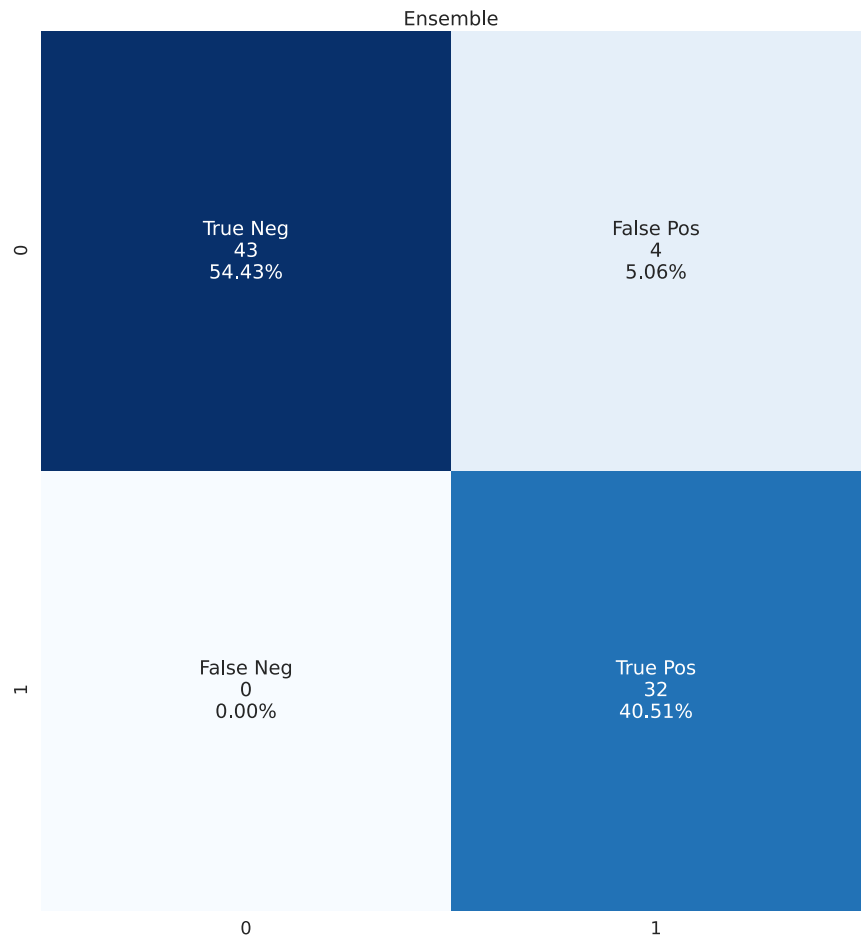
the KMeans model was able to predict 54 correct memberships with the MSLE loss function the Test set
the KMeans model was able to predict 55 correct memberships with the MAPE loss function the Test set
the KMeans model was able to predict 56 correct memberships with the Huber loss function the Test set
Total observations in Test set is: 56

Ensemble Model Results



Ensemble Model Results

Test Data Results



Conclusion

- I would not rely on this model.
- Does not give consistent results; any good results could be because of data split
- Does not significantly outperform a simple Logistic Regression Model.
- **UMAP** uses a **stochastic method** of reducing dimensionality. Because our **data was normalized**, we **assumed data** was **homogeneous**. If data is not homogeneous, any stochastic method would be inappropriate.
- Have less than **600 data points; need more data**.