

Goal:

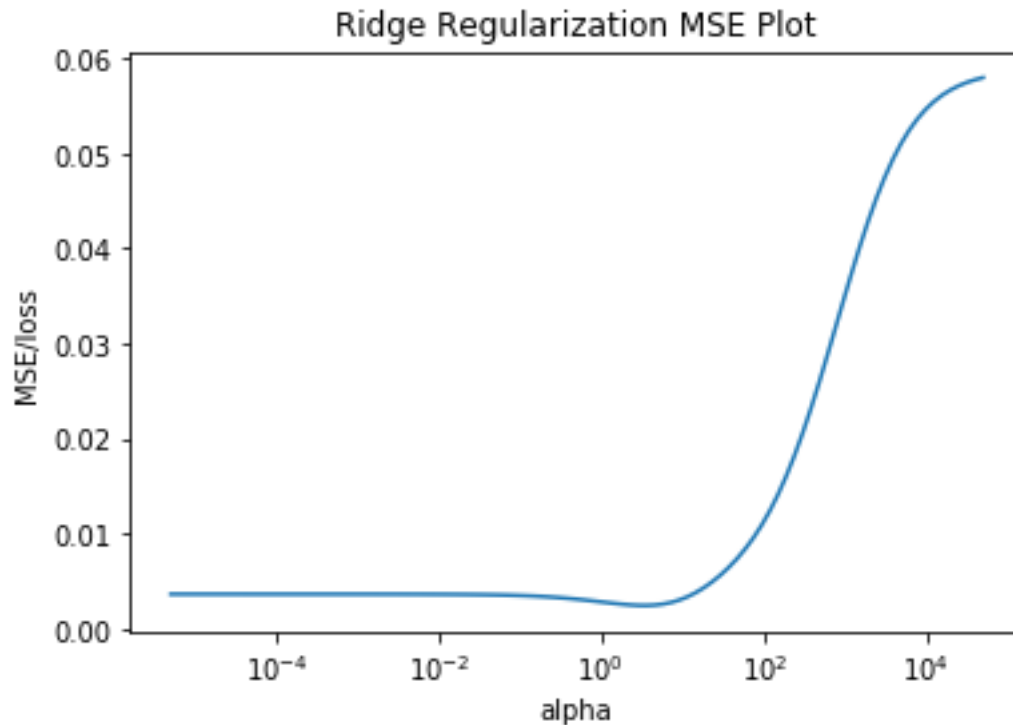
- Use regularized regression to select the informative features
- Calculate the confidence interval of the prediction using the bootstrapping method
- Create Support Vector Machine Classifier and evaluate its performance using ROC (Receiver Operating Characteristic) and PR (Precision – Recall) Curve
- Perform Principal Component Analysis and evaluate its performance using ROC and PR Curve

Problem 1

Create a predictor of the bacterial growth attribute by using only the expression of the genes as attributes. Use regularized regression technique to select genes that are informative for this task. Report the number of features that have non-zero coefficients and the 10-fold cross-validation generalization error of the technique.

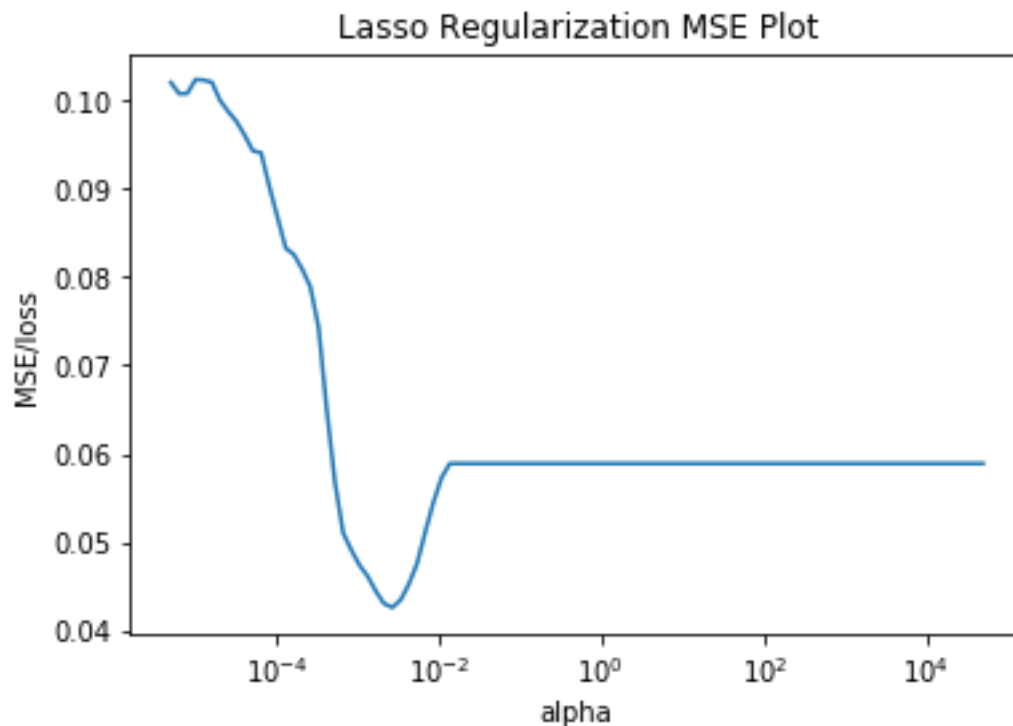
The goal of regularization is to prevent the model from overfitting the training sample. It uses the tuning parameter(λ) to change the model complexity.

For ridge regularization, there is a slight decrease in the error rate after we increase α . However, the error starts to increase dramatically after it reaches its minimum. Using 10-fold cross validation, I found that the best α is 9.15 and its MSE (generalization error) is 0.003049. The number of features that have non-zero coefficients are 4433.



For lasso regularization, there is a significant decrease in MSE as I increase α . However, the MSE starts to increase back after it reaches its minimum and remains consistent after that. Using

10-fold cross validation, I found that the best alpha is 0.00017 and its MSE (generalization error) is 0.0133. The number of features that have non-zero coefficients are 98.



After comparing the results return from the lasso and ridge regularization, I have decided to use Lasso because it reduces much more features, which is more suitable in this case because our dataset is relatively small.

Problem 2

Extend the predictor and calculate the confidence interval of the prediction by using the bootstrapping method.

Bootstrapping is a method that resample the data with replacement.

Assumptions:

- 1) Each sample is randomly selected from the population and the selection is independent and unbiased.
- 2) The sample distribution is a good approximation to the population distribution and the population is infinite.

Problem 3

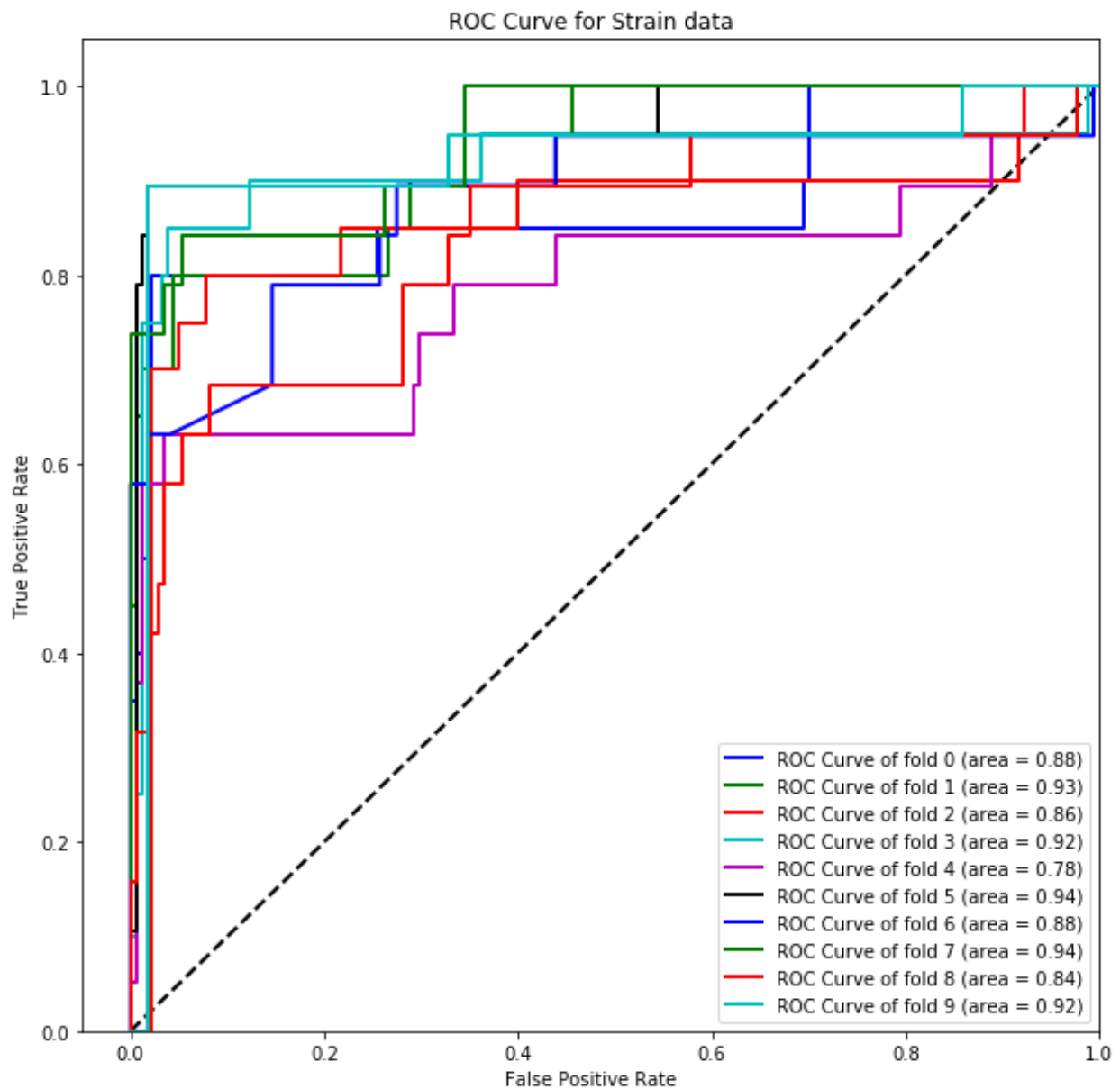
The 95% confidence interval of the predicted growth for a bacterium whose genes are expressed exactly at the mean expression value by using the bootstrapping method is 0.3942 and 0.4351.

Problem 4

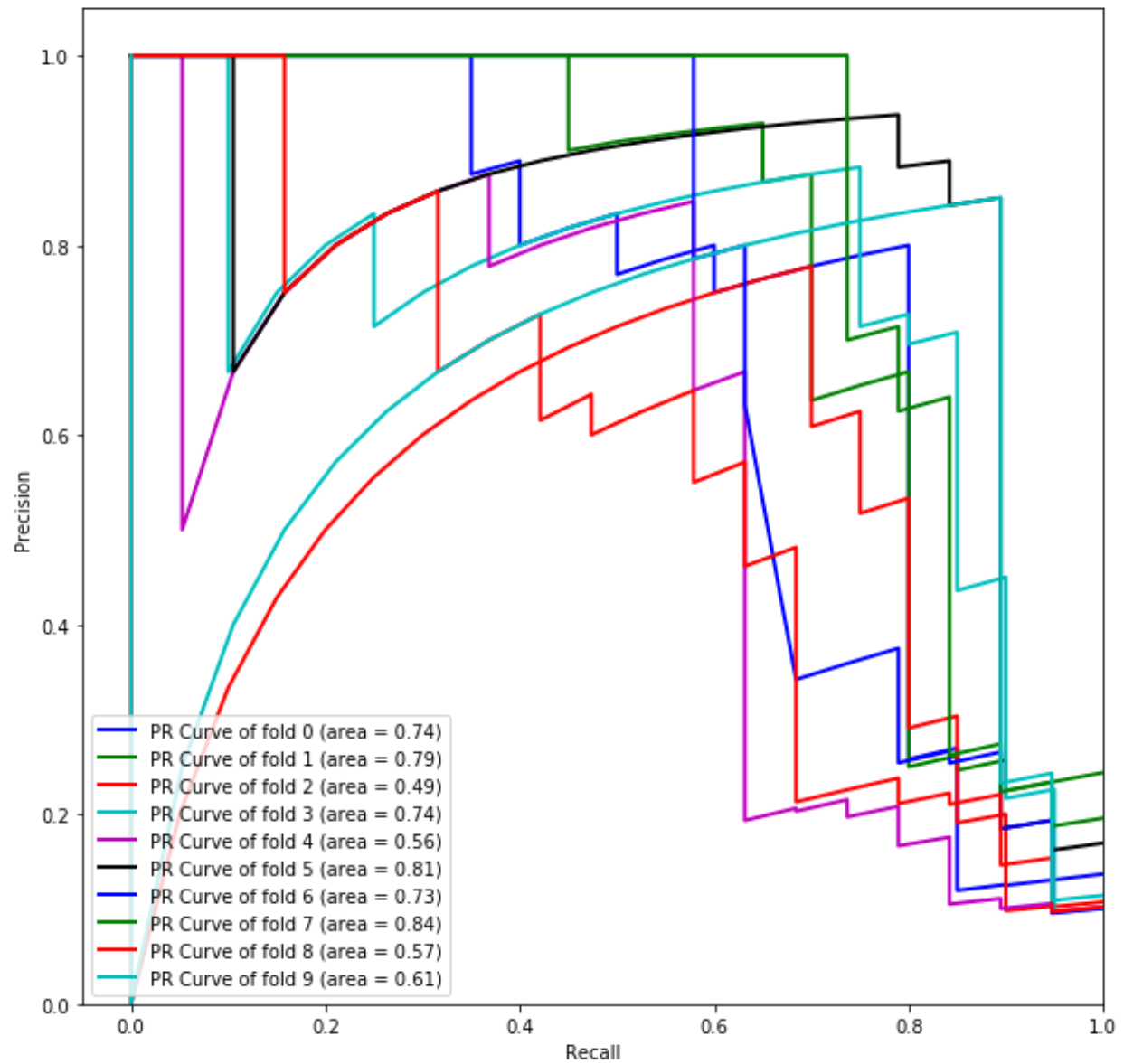
Create four separate SVM classifiers to categorize the strain type, medium type, environmental and gene perturbation, given all the gene transcriptional profiles.

The number of features for all 4 classifiers are 98 because I used the non-zero weighted features from the lasso regularization as X.

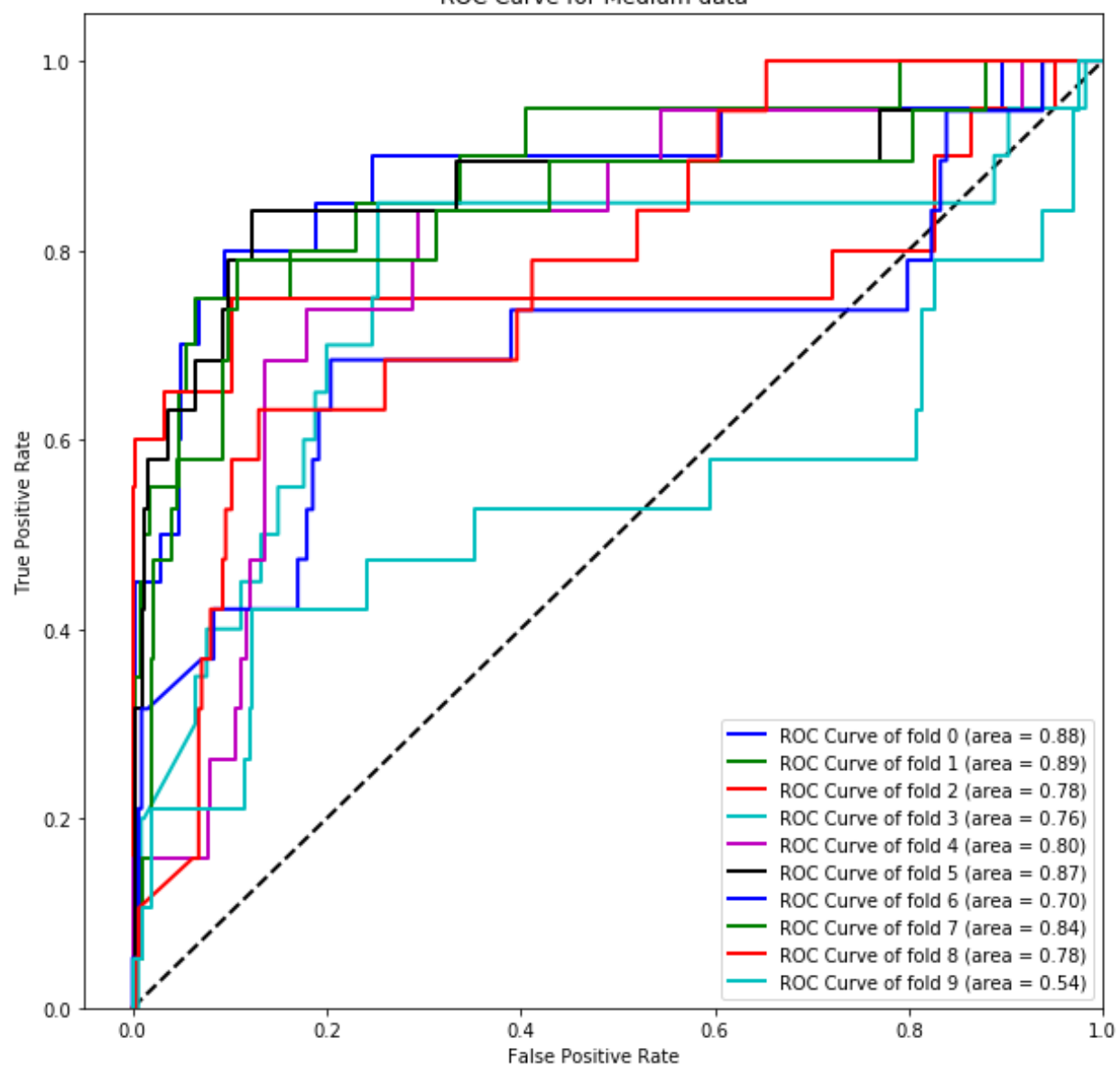
The classifiers do a pretty good job in predicting GenePerturbed, followed by Stress and Strain. However, it does not perform too well for the Medium feature.



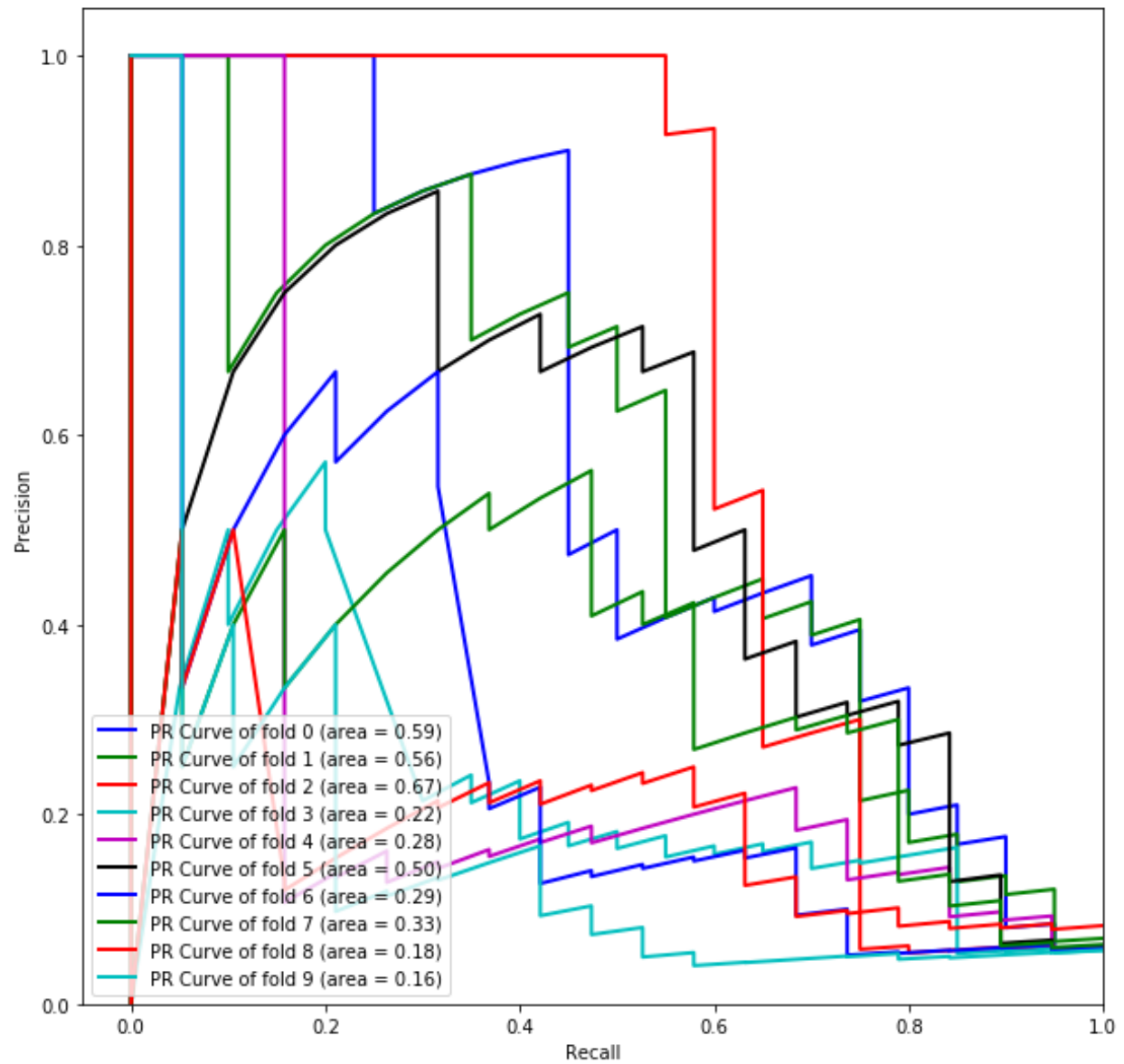
PR Curve for Strain data

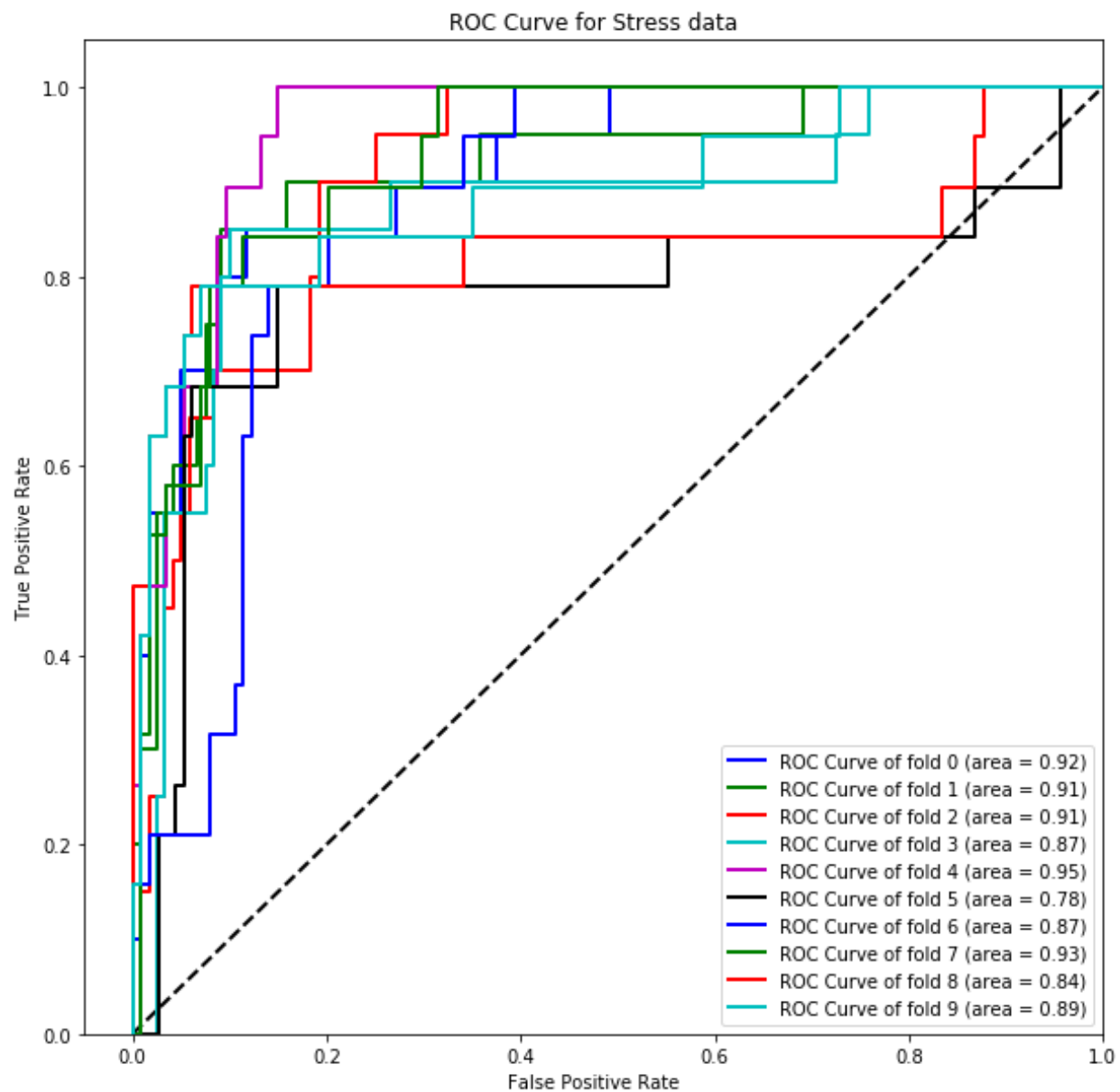


ROC Curve for Medium data

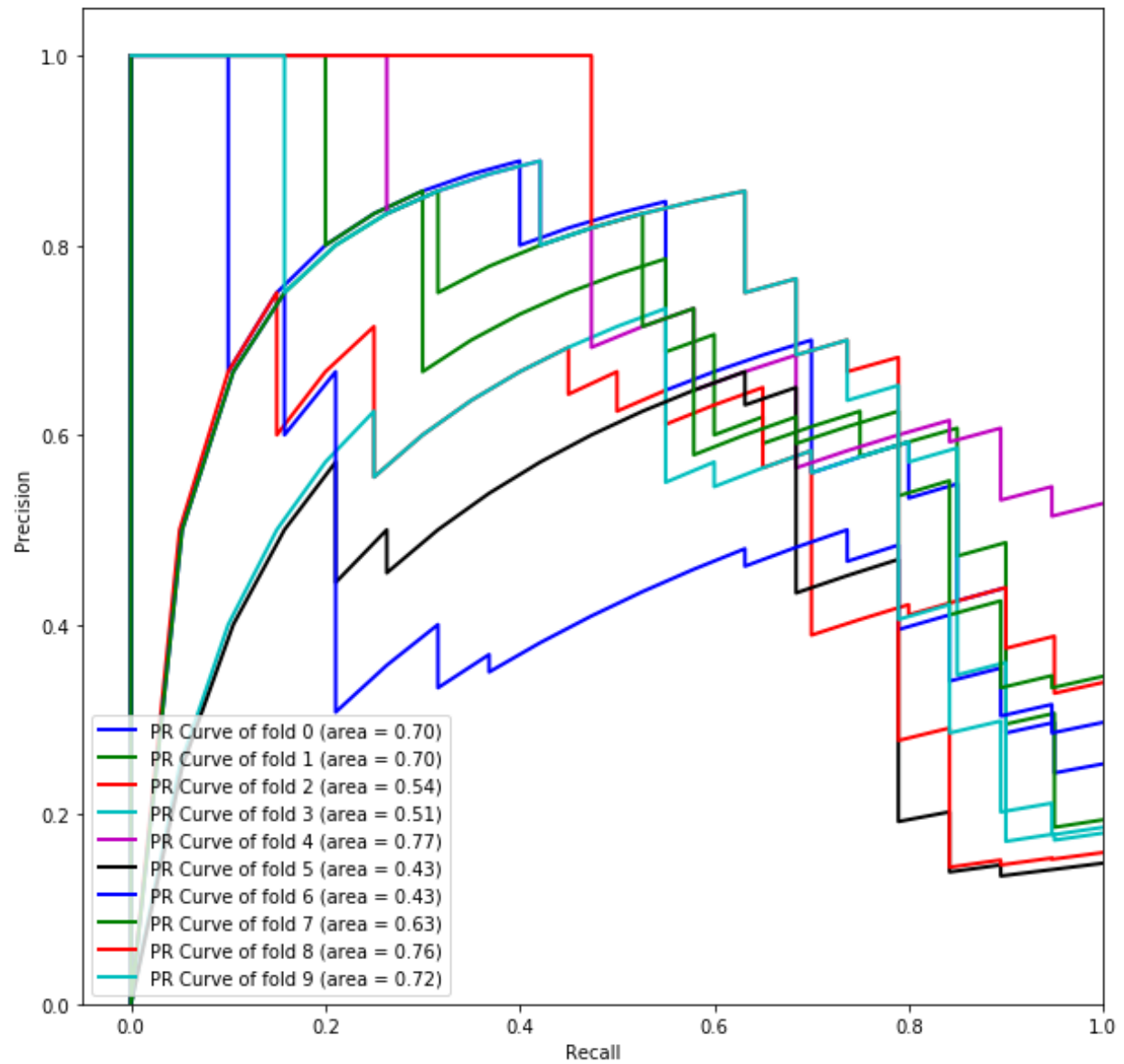


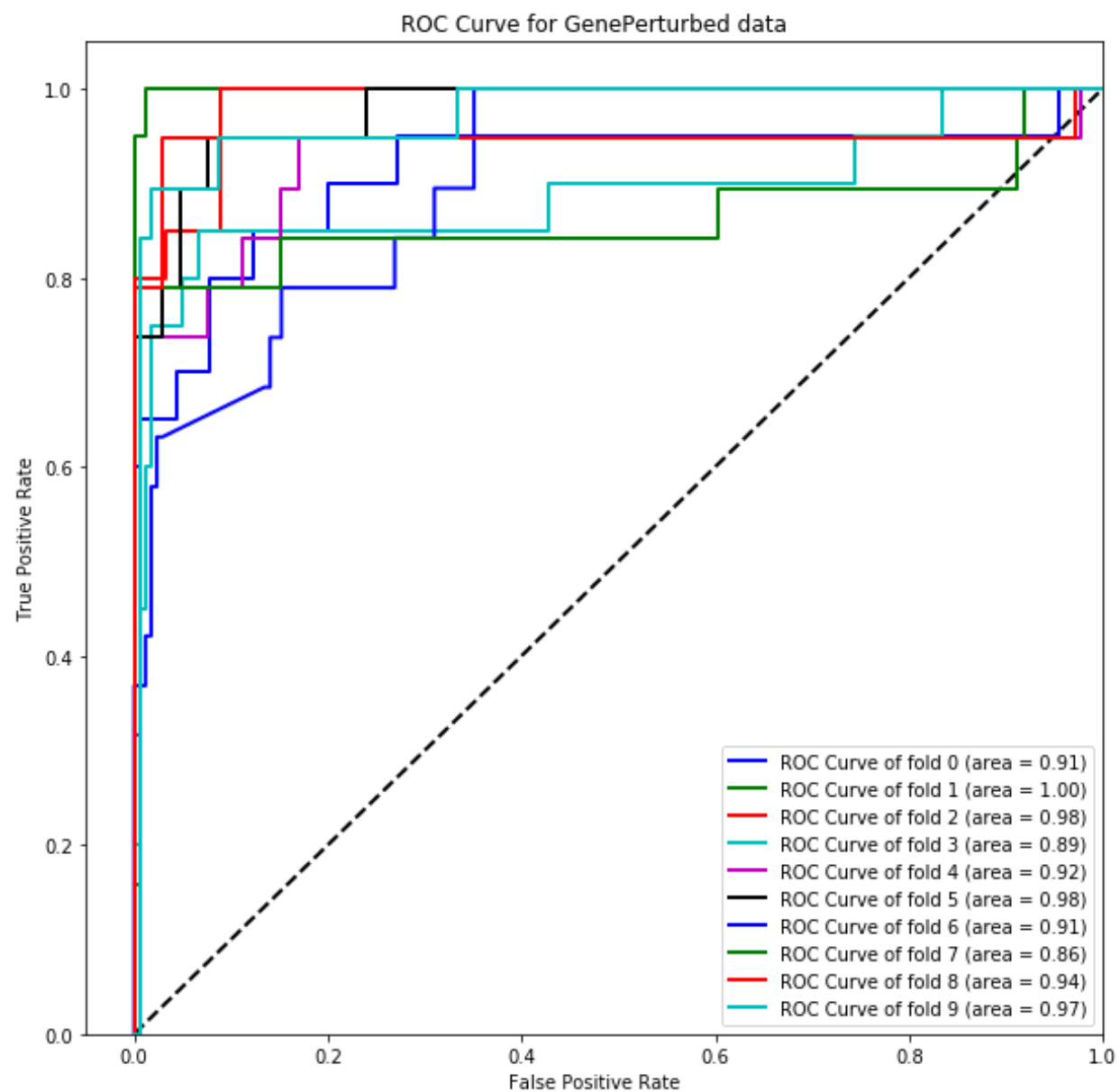
PR Curve for Medium data

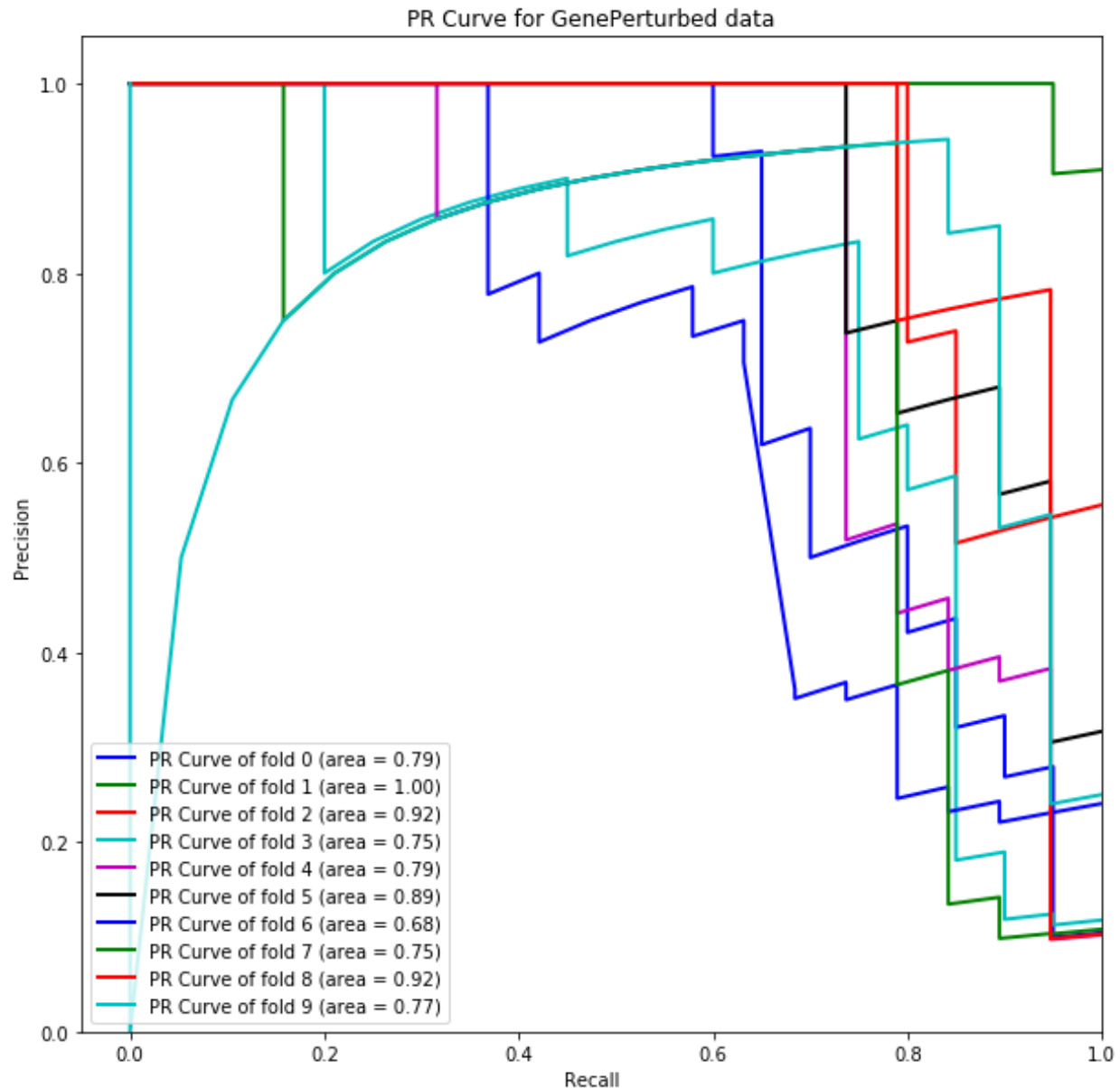




PR Curve for Stress data



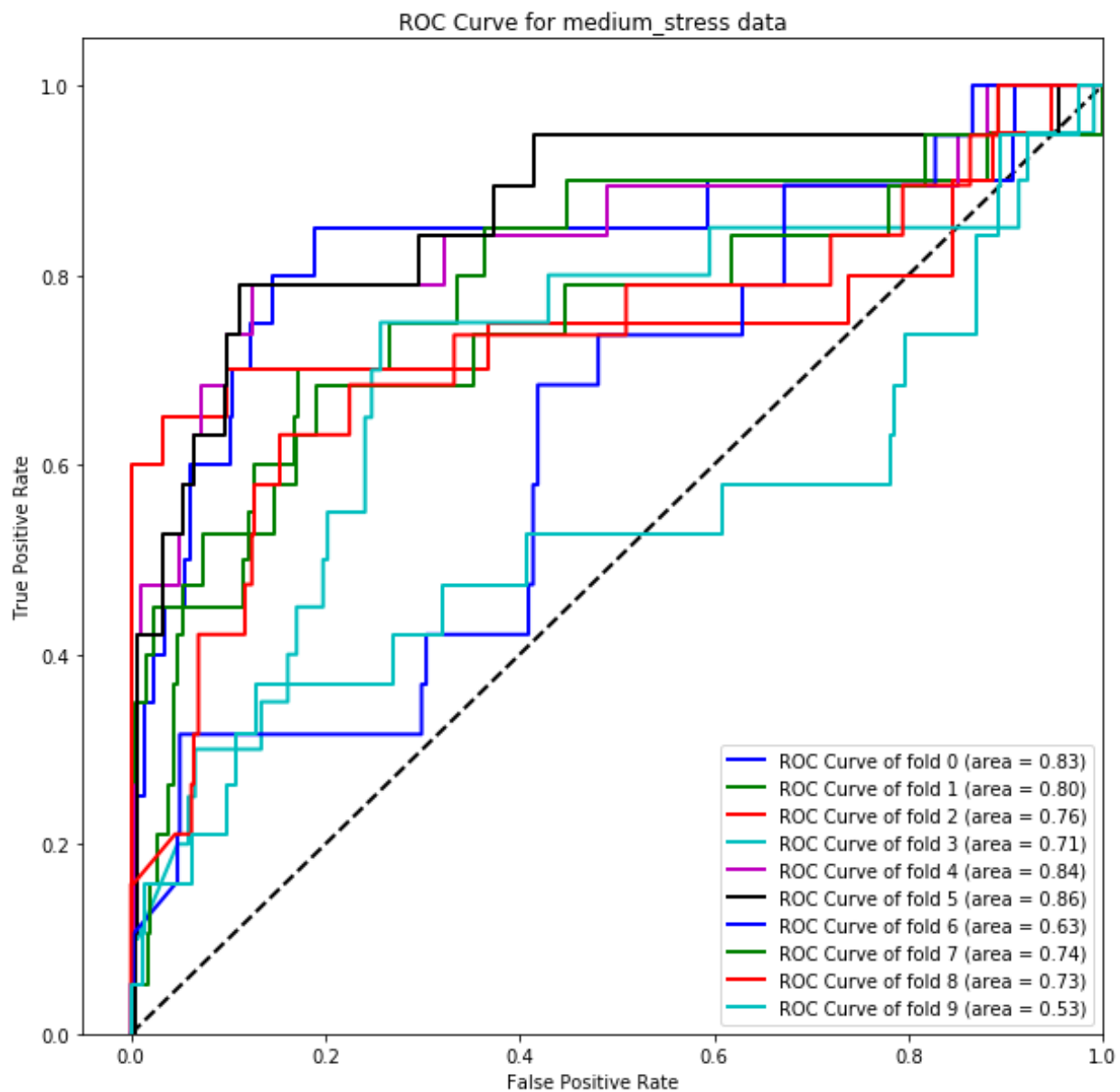


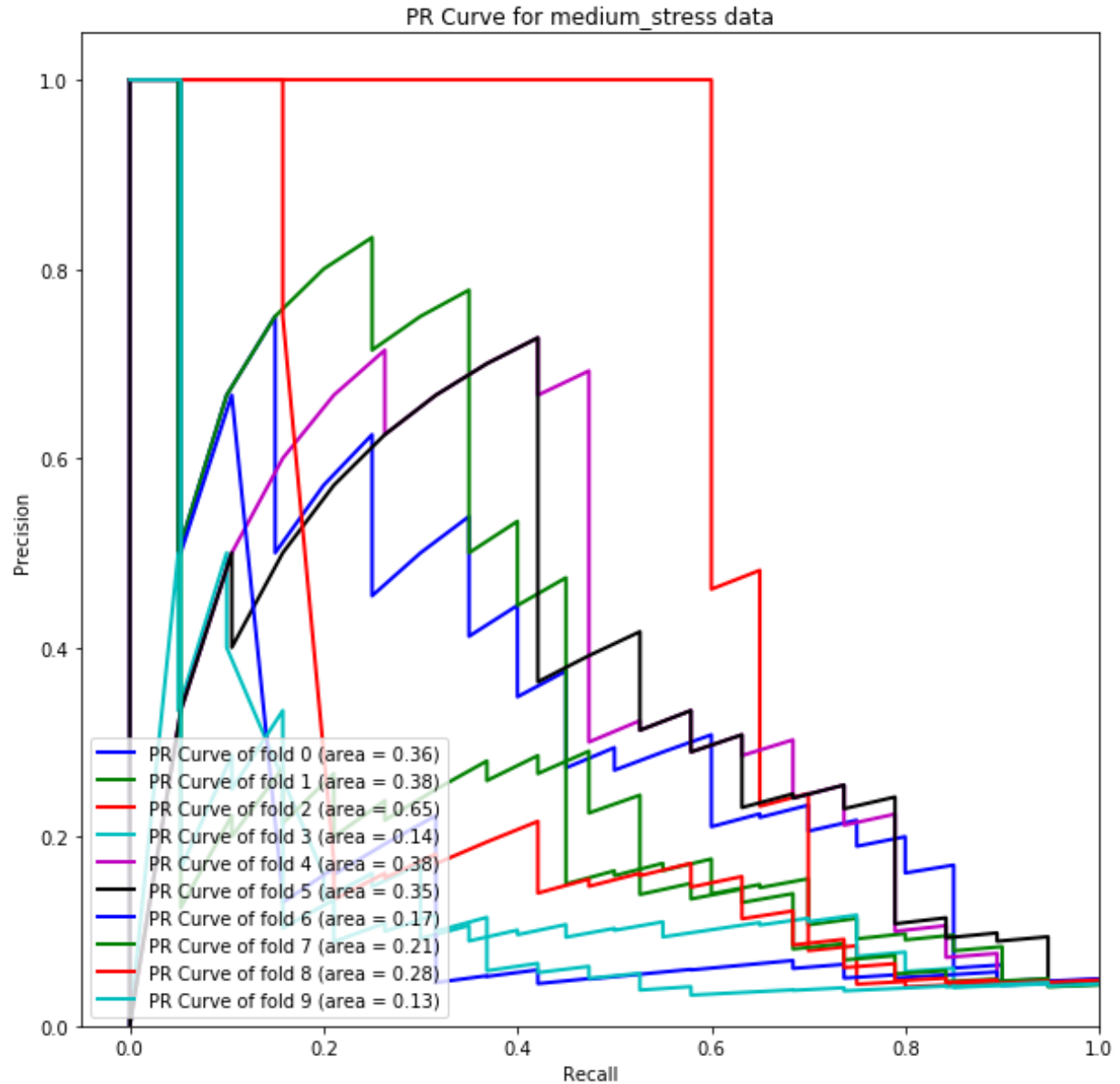


Problem 5

Create one composite SVM classifier to simultaneously predict medium and environmental perturbations and report the 10-fold cross-validation AUC/AUPRC value.

This classifier performs worse than the two individual classifiers together for these predictions. The null hypothesis is that the combined classifier performs better or equal to the two individual classifiers. Our result proves that our null hypothesis is incorrect.

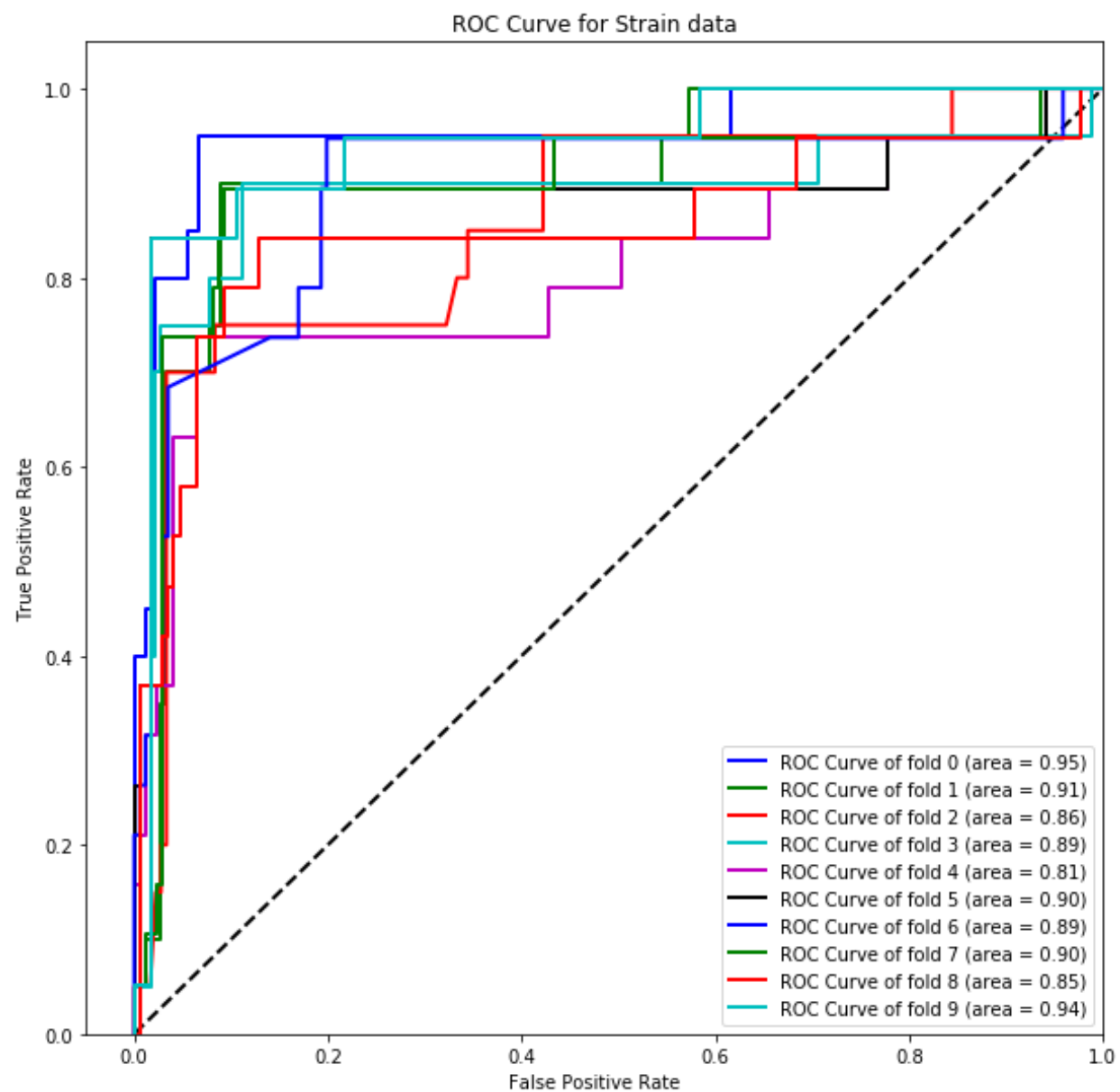




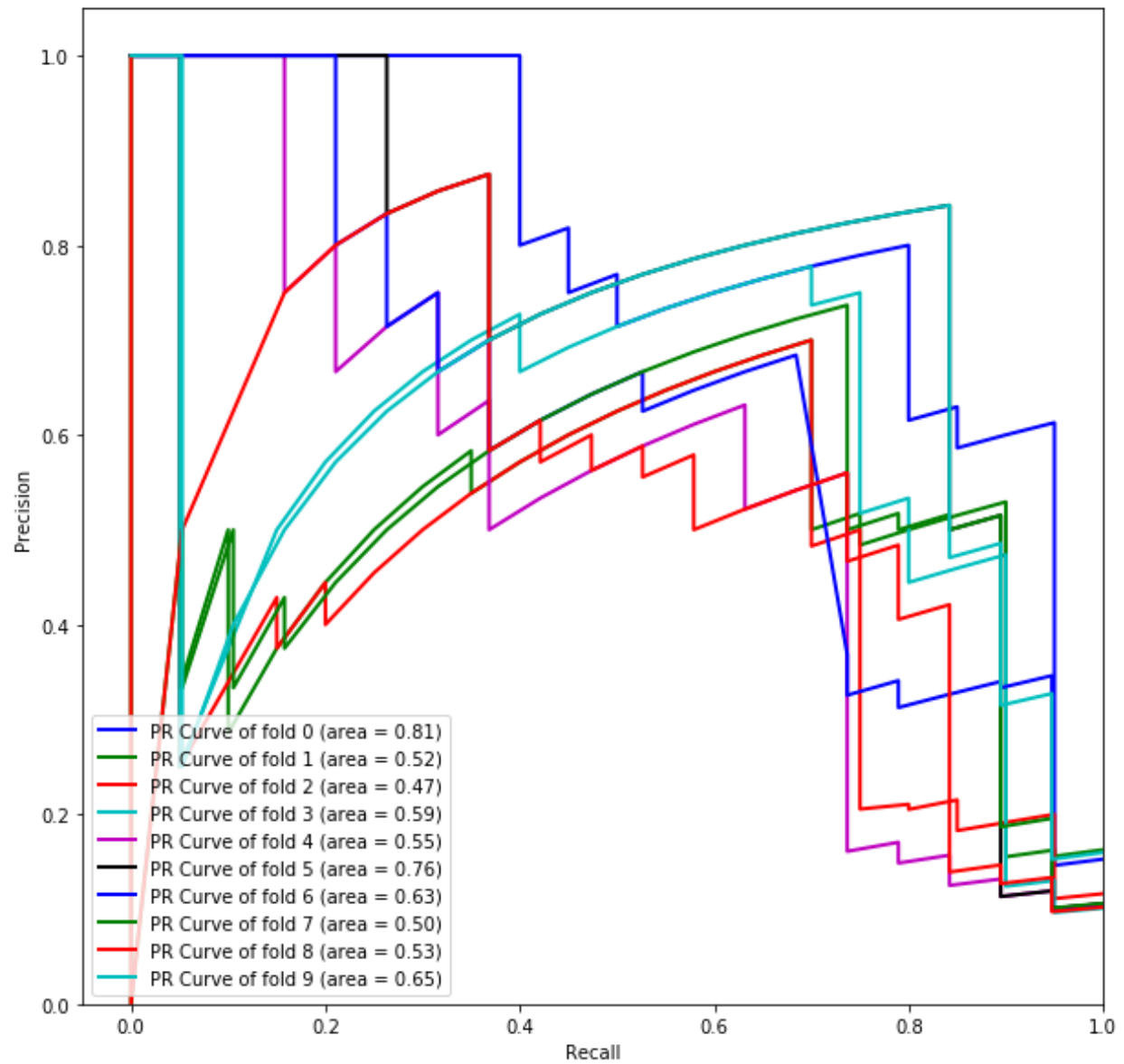
Problem 6

Perform Principal Component Analysis, keeping only the 3 Principal Components (PCs) as features for the SVM classifier.

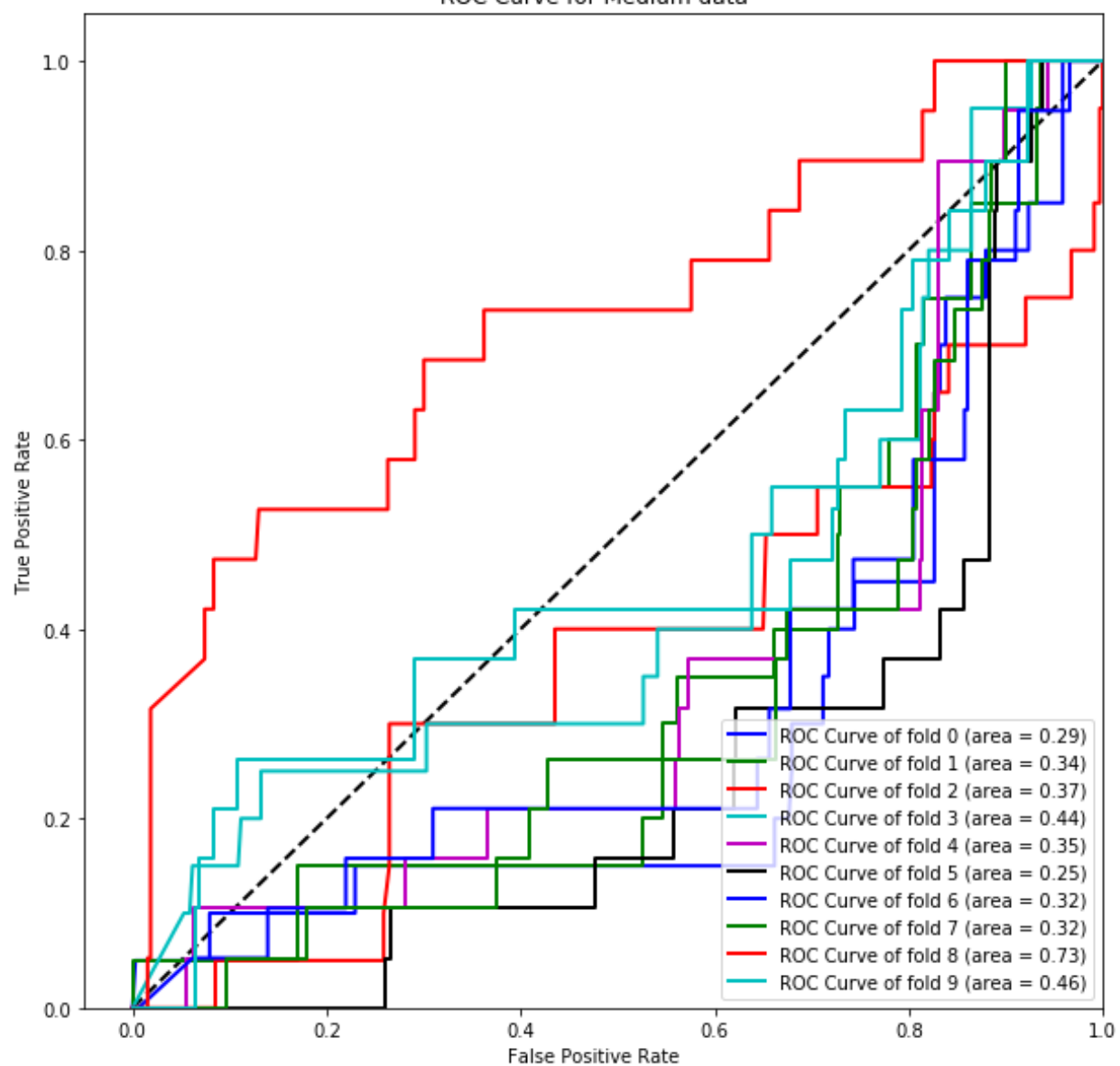
The PCs retain most of the classification performance while reducing the dimensionality because AUC/AUPRC value for Strain, Stress and GenePerturbed are still relatively high given we only have three principal components. However, AUC/AUPRC value for Medium drops significantly.



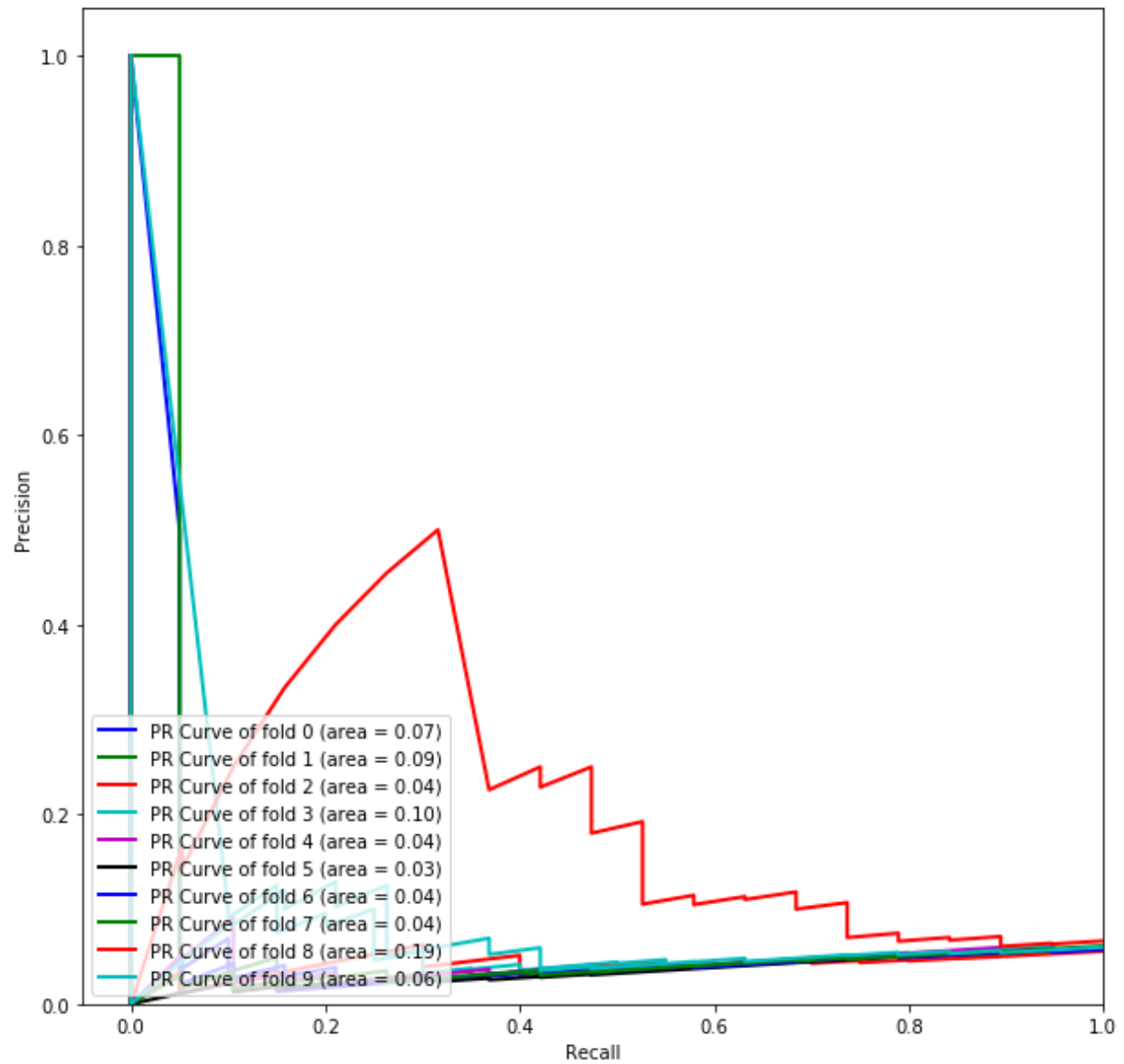
PR Curve for Strain data

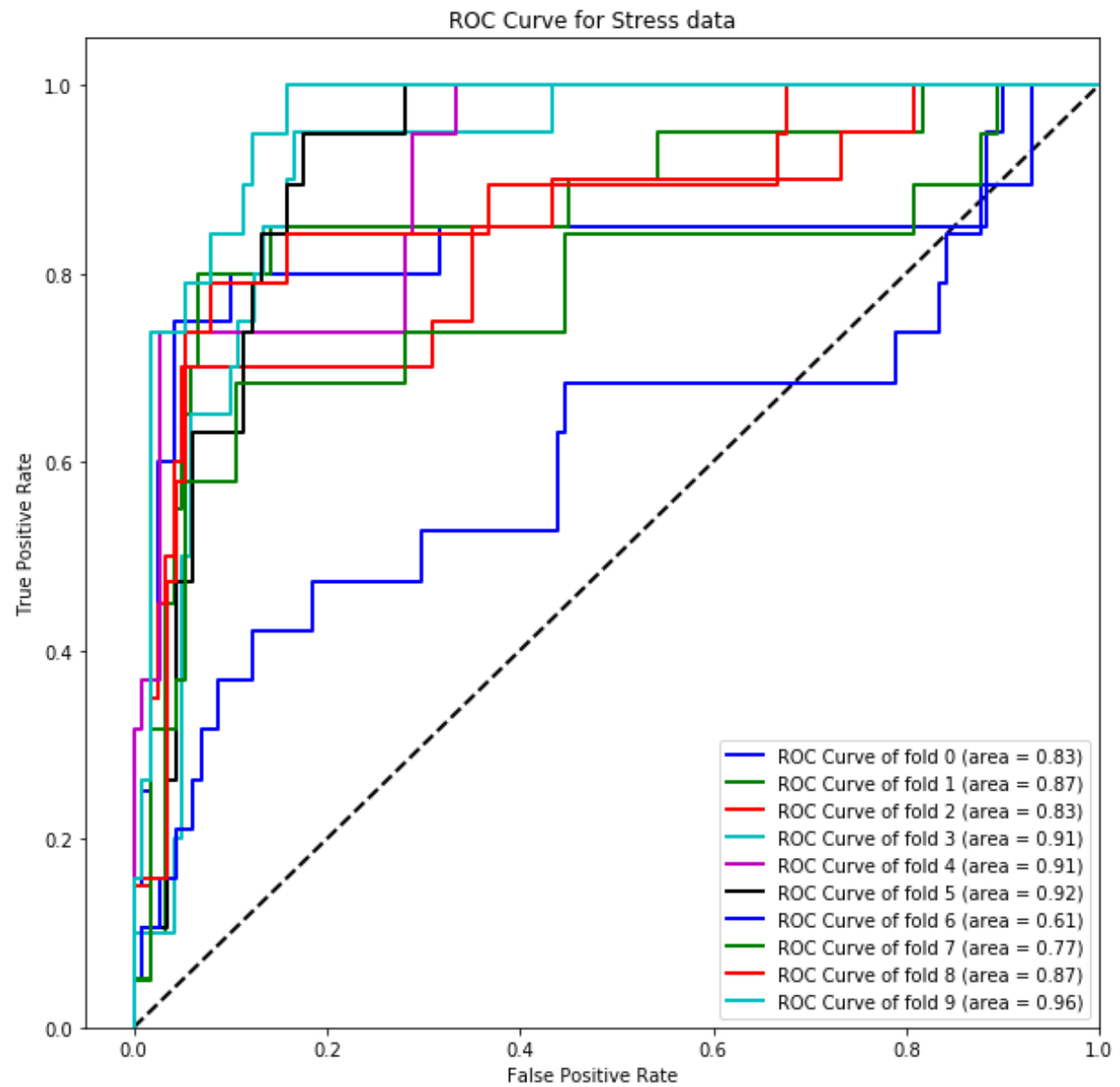


ROC Curve for Medium data

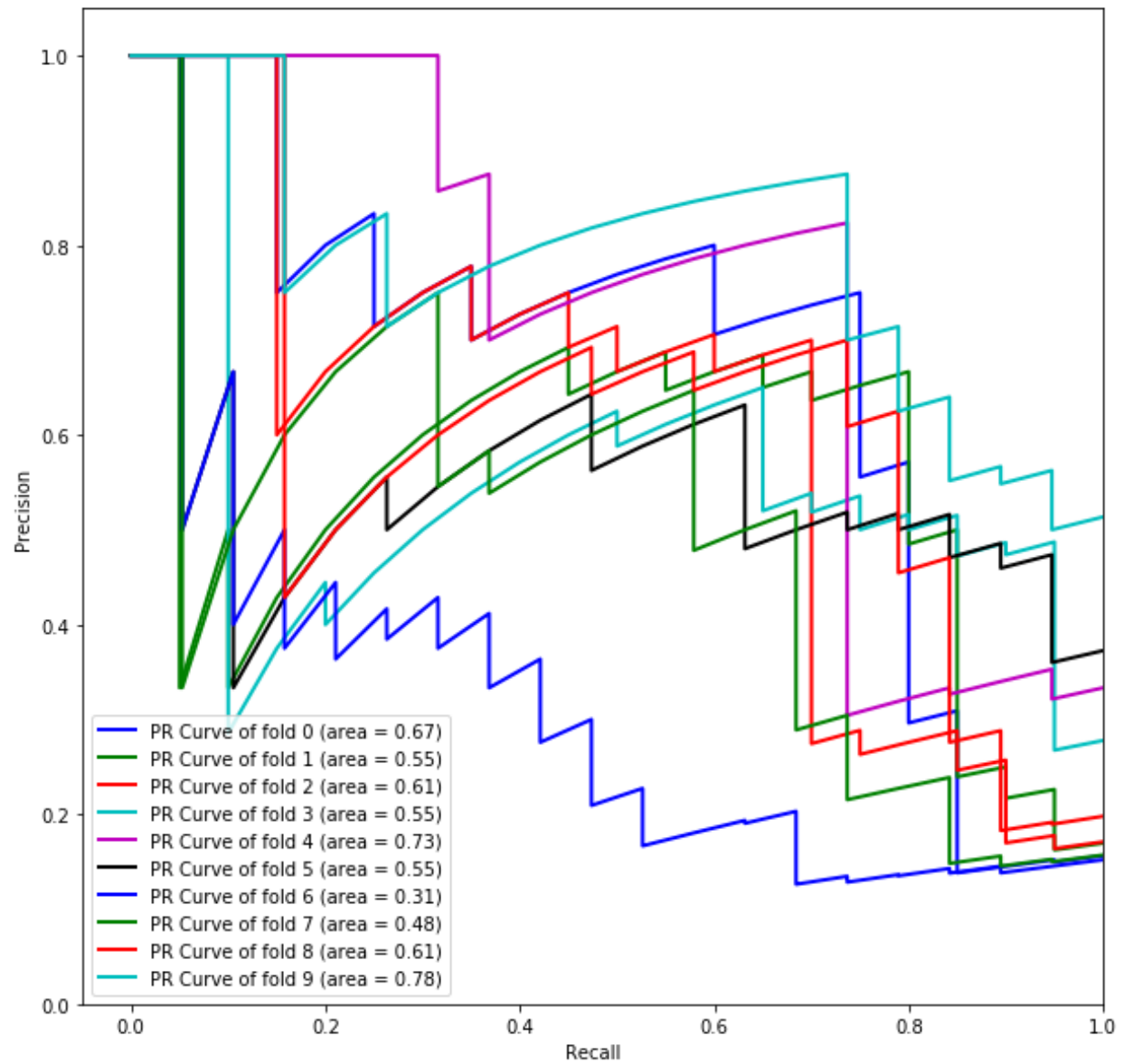


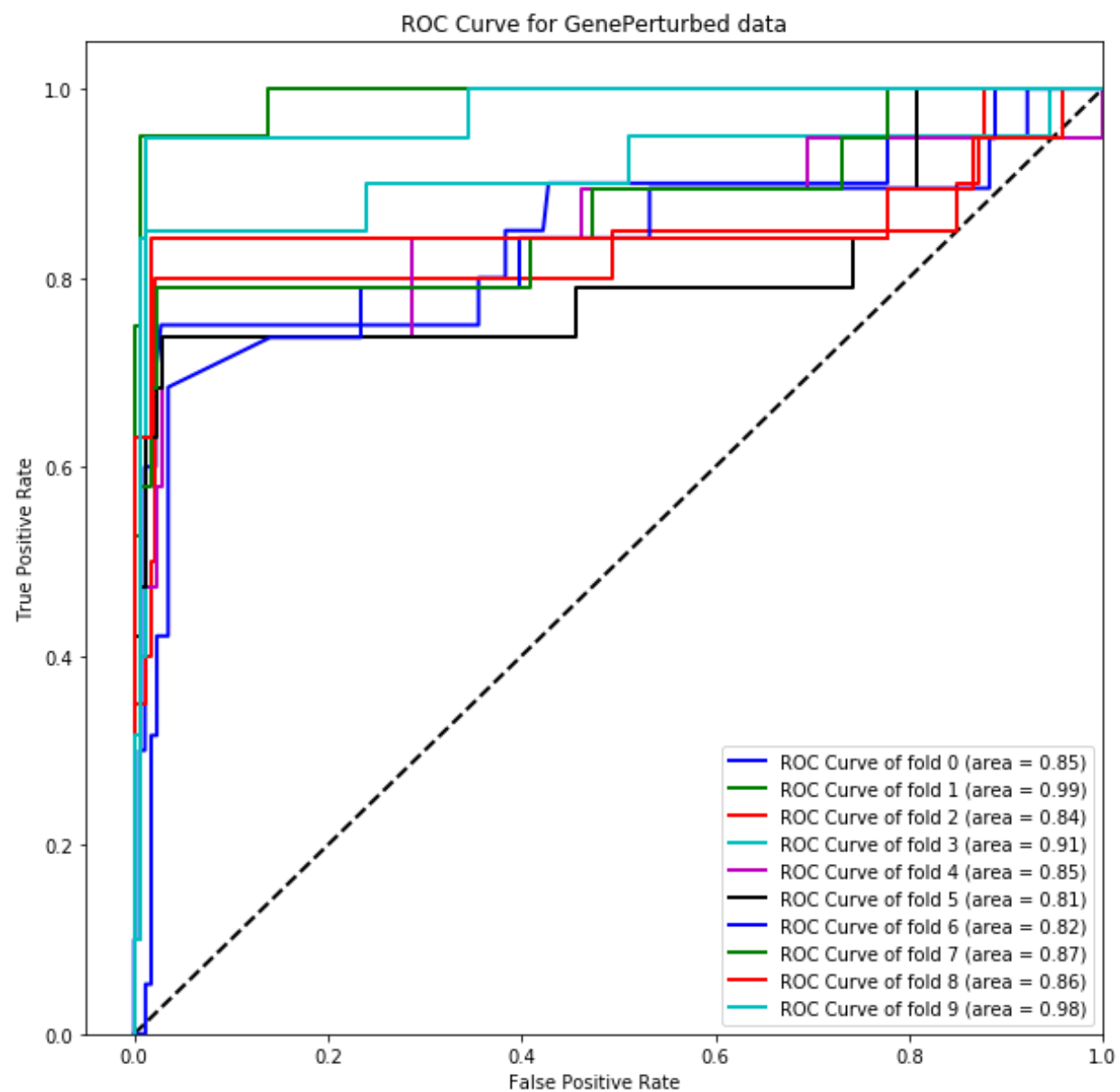
PR Curve for Medium data





PR Curve for Stress data





PR Curve for GenePerturbed data

