

Johns Hopkins University, AMS Department

Risk genes to residual tumor of breast cancer patients

CMM final project

Meng Chen

Introduction

After skin cancer, breast cancer is the most common cancer diagnosed in women in the United States.¹ Breast cancer happens when cells in your breast grow and divide in an uncontrolled way, creating a mass of tissue called a tumor.² Breast cancer can spread outside the breast through blood vessels and lymph vessels. When breast cancer spreads to other parts of the body, it is said to have metastasized.³ Studies have shown that your risk for breast cancer is due to a combination of factors. The main factors that influence your risk include being a woman and getting older. Most breast cancers are found in women who are 50 years old or older.⁴ If you have a strong family history of breast cancer or inherited changes in your BRCA1 and BRCA2 genes, you may have a high risk of getting breast cancer.⁵

Residual cancer burden is estimated from routine pathologic sections of the primary breast tumor site and the regional lymph nodes.⁶ Residual cancer burden after neoadjuvant chemotherapy can accurately predict disease recurrence and survival across all breast cancer subtypes.⁷

Fine needle aspiration

Fine needle aspiration (FNA) uses a very thin, hollow needle attached to a syringe to take out a small amount of fluid and very small pieces of tissue from the tumor. The doctor can aim the needle while feeling the tumor, if it's near the surface of the body. If the tumor is deeper inside the body and can't be felt, the needle can be guided while being watched on an imaging test such as an ultrasound or CT scan.⁸

Core biopsy

Needles used in a core biopsy are slightly larger than those used in FNA. They remove a small cylinder of tissue (about 1/16 inch in diameter and 1/2 inch long). The core needle biopsy is done with local anesthesia (drugs are used to make the area numb) in the doctor's office or clinic. Like FNA, a core biopsy can sample tumors that the doctor can feel as well as smaller ones that must be seen using imaging tests.⁹

Gene samples of patients with breast cancer were taken by either FNA(Fine needle aspiration) or core biopsy from primary tumor before they were taking chemotherapy. Patients with no residual tumor in the breast and the axillary lymph node are labeled as pathological complete response (pCR). While patients with a residual disease in either the breast, axilla or both are labeled as residual disease (RD).

Our dataset includes 139 patients and 13299 genes. Values of genes were normalized and log2-scaled. There is no missing value in the dataset. The ratio of RD over pCR is approximately 11:3 (figure 1). Our goal is to find a compact set of genes which is related to the response and is able to

¹ <https://www.mayoclinic.org/diseases-conditions/breast-cancer/symptoms-causes/syc-20352470>

² <https://my.clevelandclinic.org/health/diseases/3986-breast-cancer>

³ https://www.cdc.gov/cancer/breast/basic_info/what-is-breast-cancer.htm

⁴ https://www.cdc.gov/cancer/breast/basic_info/risk_factors.htm

⁵ https://www.cdc.gov/cancer/breast/basic_info/risk_factors.htm

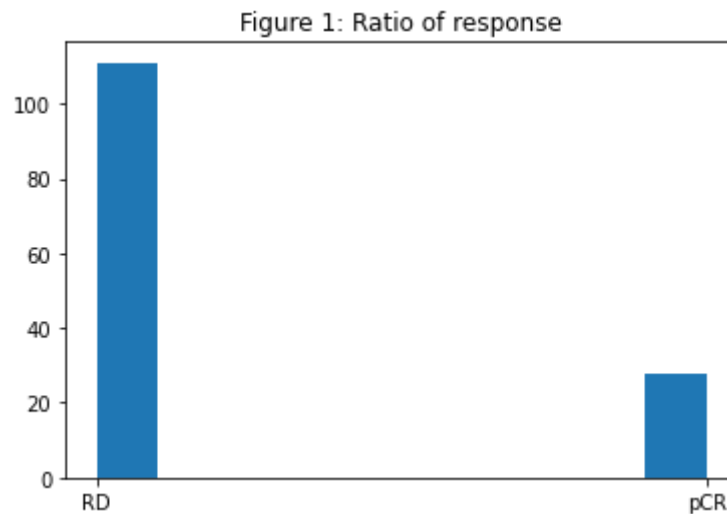
⁶ <https://ascopubs.org/doi/10.1200/jco.2007.10.6823>

⁷ <https://ascopost.com/issues/january-25-2020/residual-cancer-burden-is-prognostic-of-outcomes-across-breast-cancer-subtypes/>

⁸ <https://www.cancer.org/treatment/understanding-your-diagnosis/tests/testing-biopsy-and-cytology-specimens-for-cancer/biopsy-types.html>

⁹ <https://www.cancer.org/treatment/understanding-your-diagnosis/tests/testing-biopsy-and-cytology-specimens-for-cancer/biopsy-types.html>

predict the residual disease.

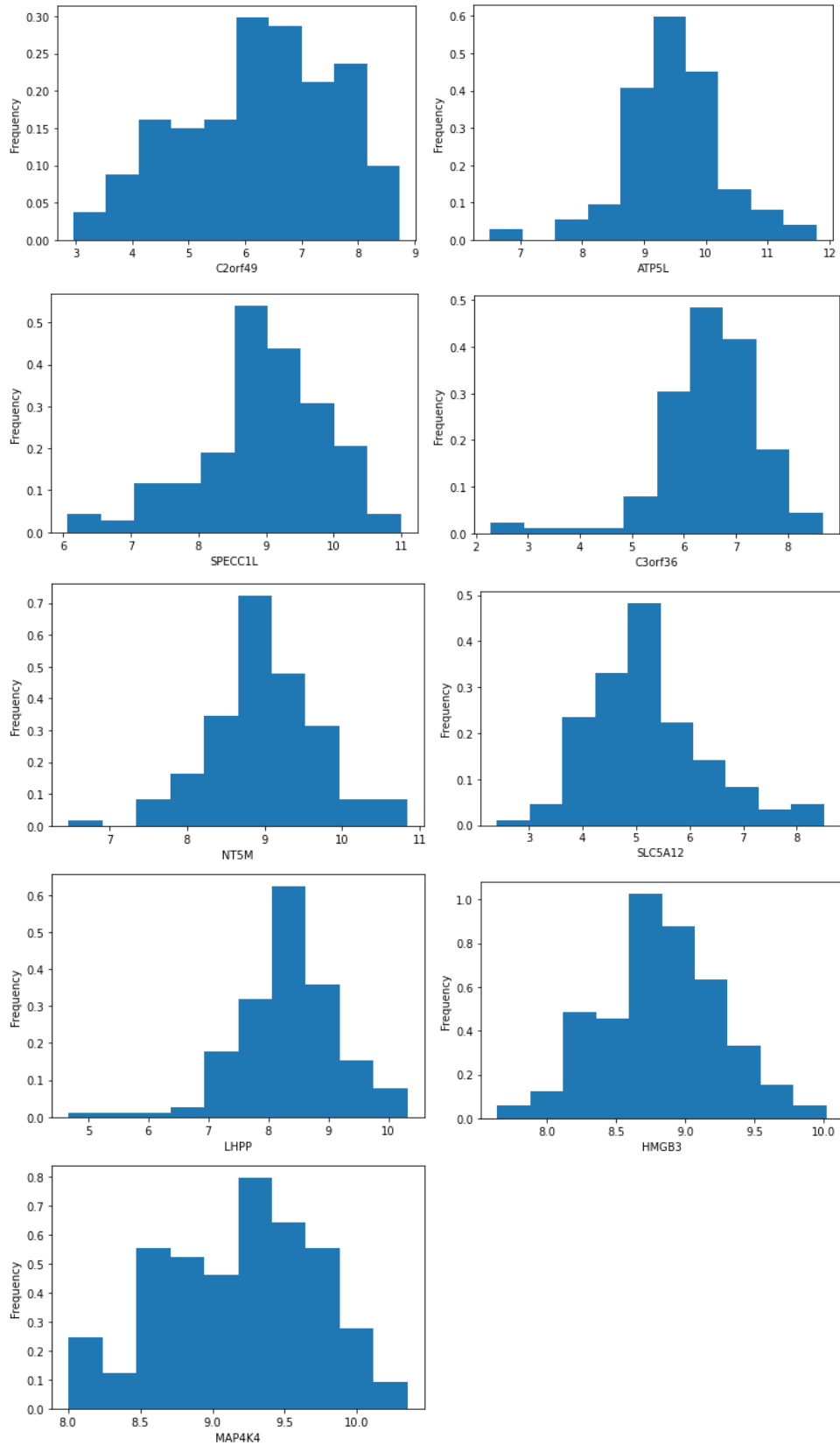


Methods

Our first step is to reduce the dimension of genes to find sets of genes which are most likely to predict the response. To achieve this, we use three methods which are joint mutual information, Pearson Correlation and Wilcoxon Rank-Sum Test. Specifically, we want to leave 20 most informative genes, so we find three sets of genes A, B and C among 13299 genes. Genes in A have 20 highest joint mutual information of genes and the response. Genes in B have 20 highest Pearson Correlation to the response. Genes in C have 20 smallest P-values of Wilcoxon Rank-Sum test with the null hypothesis values of a gene related to two labels are equally spread. In the second step, we split our data into training and test set with proportion 8:2. Notice that our data is greatly imbalanced, we keep the proportion same for two labels of response in training and test set. In addition, we can also implement weighted methods to eliminate the effect of imbalance. Next step is to use different machine learning methods, particularly, Logistics Regression, SVM, LDA, QDA, Random Forest and Artificial Neural Network, to train and predict using A, B, C three sets. We try to find the best methods and gene set based on ROC curve and AUC (Area under curve), since the data is imbalanced AUC instead of accuracy score is an appropriate criterion for the prediction.

Preprocess

We randomly select 9 genes in the total set and plot the distribution of them. We can see that they are all approximately normal shaped and roughly range in the interval from about 5 to 10. This could be the result of log-scale.



The genes in A, B and C three sets are showed bellowed. We can see every set is entirely exclusive to each other, which means there is no agreement made by three criterions. Although, joint mutual information, Pearson Correlation and Wilcoxon Rank-Sum Test all measure the relationship

between genes and the response, the imbalanced proportion of sample size and genes size results from the vast amount of genes leads to great deviation of results measured by them. To examine which set of genes best separate and predict the response, we need to resort to other methods.

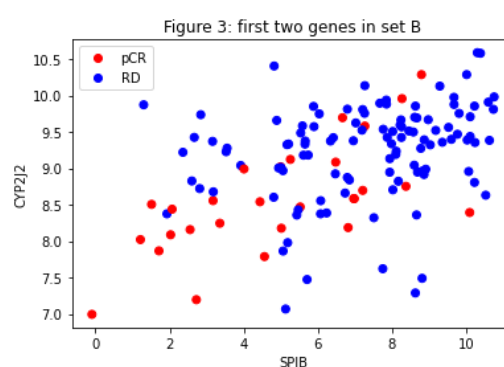
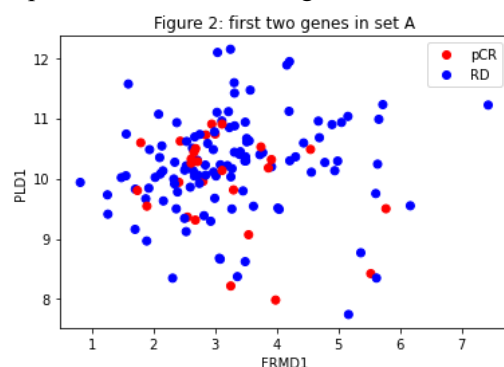
12414	FRMD1	4046	SPIB	9777	TRIM44
18	PLD1	3369	CYP2J2	5060	LOC101928635///ALDH1A2
39	FAU	13108	TJAP1	8188	MISP
42	RPL18	3498	SURF2	9778	ADI1
383	HIF1A	10733	RNF25	8203	NCAPH
386	IP07	3902	SDS	556	NUDCD3
570	SDC1	10386	ENY2	564	LOC101926921///DAB2
762	TMX4	10836	COL5A3	8190	TSPYL4
769	ANXA2	713	PFDN1	8260	NTAN1
1098	MX1	3058	ATP2B2	12028	PAQR5
1159	VEZF1	245	IQGAP1	8180	THAP11
1285	TRIM2	2967	ABCG1	11823	MCUR1
1389	AFG3L2	3502	RBMS2	8175	LRRC47
1586	CASP3	1794	TTC37	8227	MZT2A///MZT2B///PHGDH
1813	MYO1E	2100	STX7	9788	NDUFA4
1891	TFAM	436	SNU13	8194	MCF2L
1926	CLK2	1658	CNPY2	2312	BAK1
2087	EML3	10295	CARHSP1	9746	EIF3L
2102	VPS16	7951	SETD2	9767	COPG1
2204	DYNC1LI2	1407	TNFAIP2	568	APEH

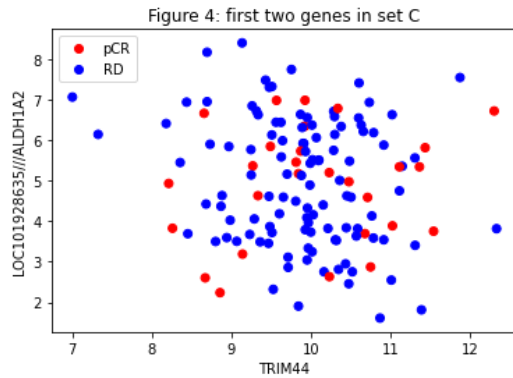
Gene set A

Gene set B

Gene set C

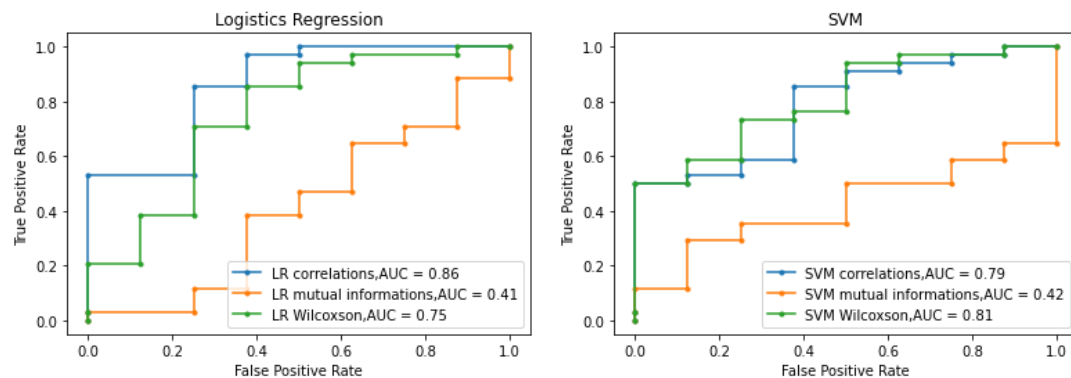
Next we will implement them separately to 6 machine learning methods. But first we can obtain some rough idea based on the scatter plot for the first 2 genes in each set with respect to two labels(Fig 2,3,4). We notice that clear separation can only be observed in set B, while scatters in other two plots are similarly spread out, which implies that Pearson Correlation might be the best separation criteria for this gene sets.





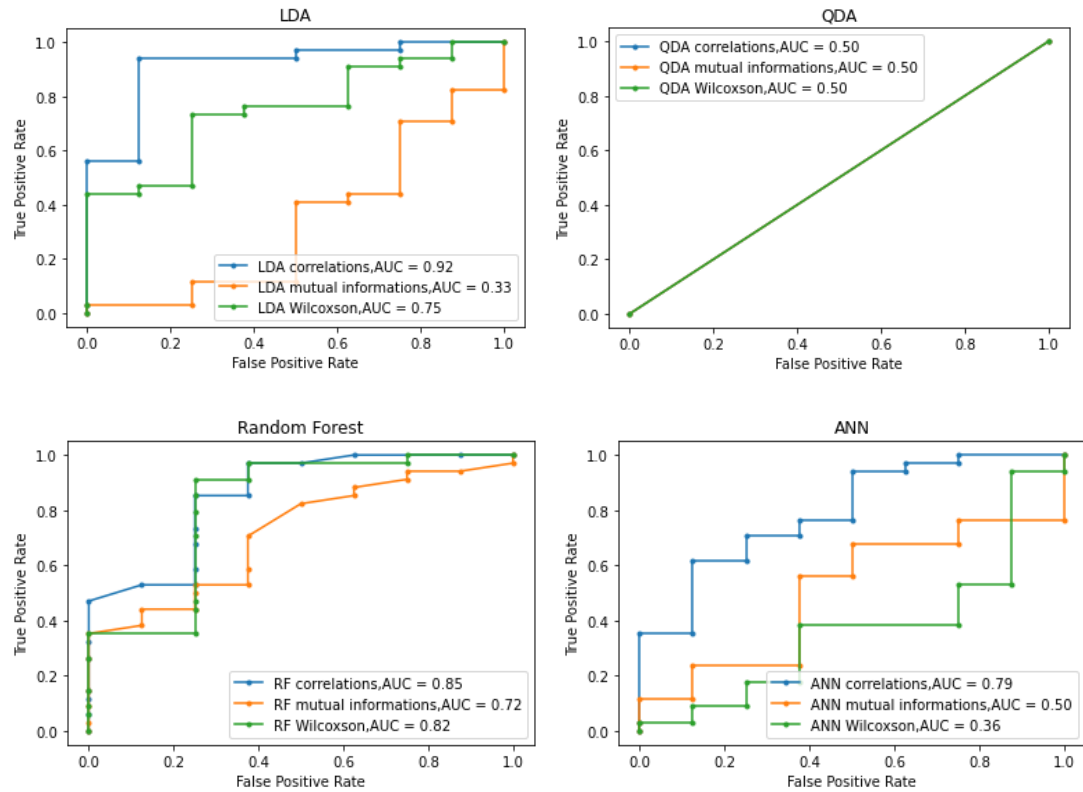
Results

The corresponding ROC curve and AUC are showed below. We find that the average AUC for prediction with set A is 0.476, set B is 0.842 and set C is 0.698. Therefore the best set for prediction is set B which is 20 genes with highest Pearson Correlation, which conforms to our speculation. We notice that joint mutual information doesn't seem to be a good criteria since mutual information might lead to selection of redundant and irrelevant features.¹⁰ Among all machine learning methods, excluding set A, Random Forest and LDA has highest average AUC 0.835, then is Logistics Regression 0.805, SVM 0.8, ANN 0.575. We notice that QDA fails to predict because of some genes are collinearly. The computation in QDA involves the inverse of design matrix like linear regression. Thus the collinearity leads to an inaccuracy in computing inverse of design matrix. The inaccuracy of ANN to the prediction with set C might attribute to the inappropriate design of network structures which is unavailable.¹¹



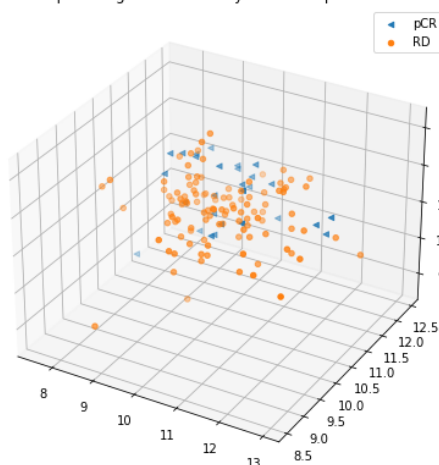
¹⁰ <https://www.sciencedirect.com/science/article/pii/S0957417415004674>

¹¹ <https://www.asquero.com/article/advantages-and-disadvantages-of-artificial-neural-networks/>

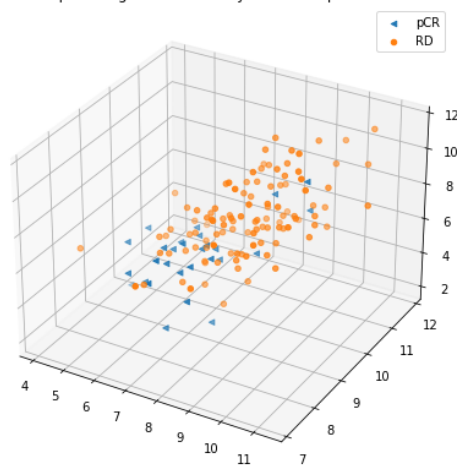


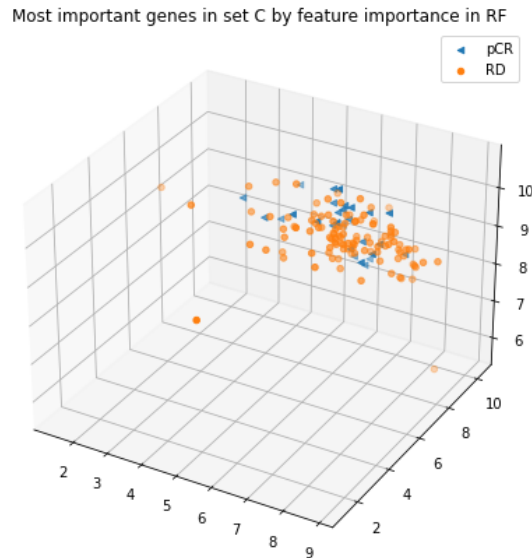
We can also notice from the above result that Random Forest is the only way with fair results in all three sets. There is a criterion in Random Forest called feature importance which measures the percentage importance of features. For the three gene sets A, B and C, the top 3 most important genes in set A are: IPO7, MYO1E and CLK2, with importance 0.154, 0.106 and 0.064. Top 3 most important genes in set B are: ATP2B2, RNF25 and SDS, with importance 0.096, 0.090 and 0.084. Top 3 most important genes in set C are: THAP11, NUDCD3 and ADI1, with importance 0.111, 0.104 and 0.068. Those top 3 genes are not intersecting with the pairs of genes we obtained from sorting in A, B and C. One thing we notice is that importance for genes in set B are more equally spread, however, in other two sets are converging to the top 2 genes. We also plot three similar but in 3-D scatter plots. We can see that still, genes in set B are most separable respect to labels in response. But we can see a little separation in set A as well as dots labeled as ‘RD’ are smaller in one axis.

Most important genes in set A by feature importance in RF



Most important genes in set B by feature importance in RF





Conclusion

From our results, we conclude that LDA trained by 20 genes with highest Pearson Correlations to residual disease can do the best prediction whose AUC reaches 0.92. This could be the result of normality of genes and they also have approximately same variance as showed in the Chapter Preprocess, which perfectly fit the assumptions of LDA. There is also one contradiction with respect to the inaccuracy prediction by QDA. The warning of Python indicates "Variables are collinear", however, LDA and Logistic Regression which also involve inverse of design matrix don't pop out this warning. This might be the result of underlying functions of Python in those methods, however no exact explanation was found. There is also limitation in this analysis. CNN instead of ANN could be implemented to this data, where CNN tends to be a more powerful and accurate way of solving classification problems.¹²In addition, by using pretrained CNNs, we don't need to adjust the structure of network by trials and they are more sophisticated built.

¹² <https://viso.ai/deep-learning/ann-and-cnn-analyzing-differences-and-similarities/>