

Asha Guruvayurappan | City University of London

- Construct and compare two models for a binary classification problem to predict Breast Cancer type.
- Investigate the efficiency of logistic regression and random forest models in predicting whether a cancer is benign or malignant
- Compare the results of similar implementations obtained by **Adel S. Assiri , Saima Nazir and Sergio A. Velastin in Breast Tumor Classification Using an Ensemble Machine Learning Method[2]**

- Logistic regression is the go-to method for a supervised binary classification problem (classifying cancer type to be Benign or Malignant)
- The paper **Breast Tumor Classification Using an Ensemble Machine Learning Method by Adel S. Assiri , Saima Nazir and Sergio A. Velastin[2]**, concludes that Logistic Regression is one of the most accurate machine learning models for this problem.

- RF is a more stable model when compared to LR, as noise over data drastically decrease model performance for LR.
- Hyperparameter optimization might not always have a positive impact on model for LR.
- Although the accuracy of overall model is high; recall for malignant cases(model predicting malignant cases to be benign) are low which would have consequence in real world.



- ## How it works?

- Random forest combines multiple Decision Trees to reach a single result. It uses bagging and feature importance when building individual trees and create an uncorrelated forest to obtain an accurate and stable prediction.

Advantages

- RF is a flexible and easy to use machine learning algorithm
- It can be used on both classification and regression problems
- it reduces overfitting in decision trees and helps to improve the accuracy.

Disadvantages

- Training time of RF is much more when compared to any other models as it has to train and combine numerous trees
- It requires much computational power.
- **They are unstable**, meaning that a small change in the data can lead to a large change in the structure of the optimal decision tree

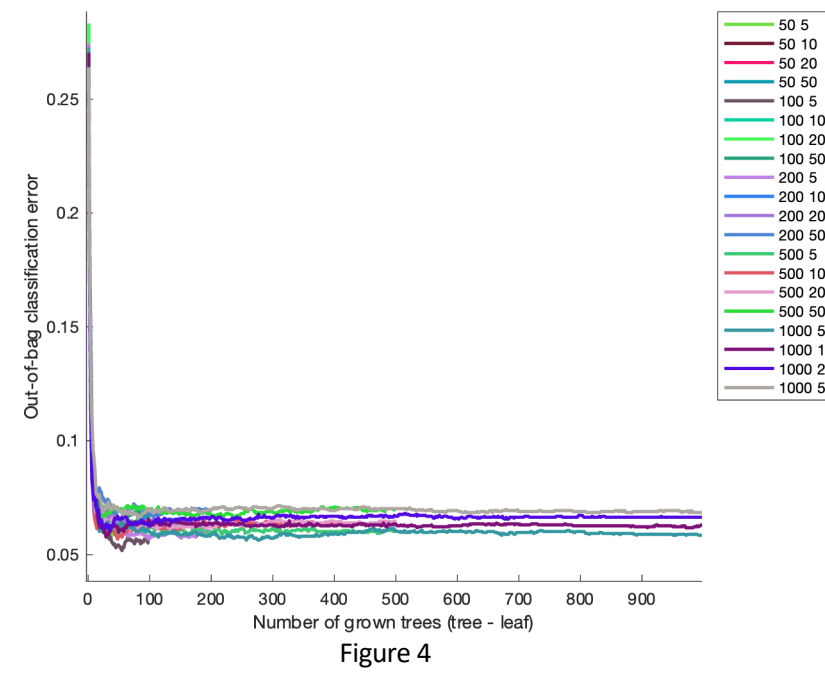
- Logistic regression will have lower training time when compared to Random Forest.
- Logistic regression is expected to perform better than RF as specified in **Breast Tumor Classification Using an Ensemble Machine Learning Method[2]**
- Feature reduction (using algorithms) should improve accuracy of both the models.

- The original dataset is split into 70% - 30% as training – testing sets resulting in 399 training and 171 testing records.
- SMOTE technique is imposed on the training set to balance and increase the training records[4]. SMOTE sampling can benefit feature selection which can be applied to reduce feature dimensions[7].
- Versus the original published study[2], Gaussian Noise is induced onto training set to avoid overfitting of models[5] and to evaluate the accuracy with noise.

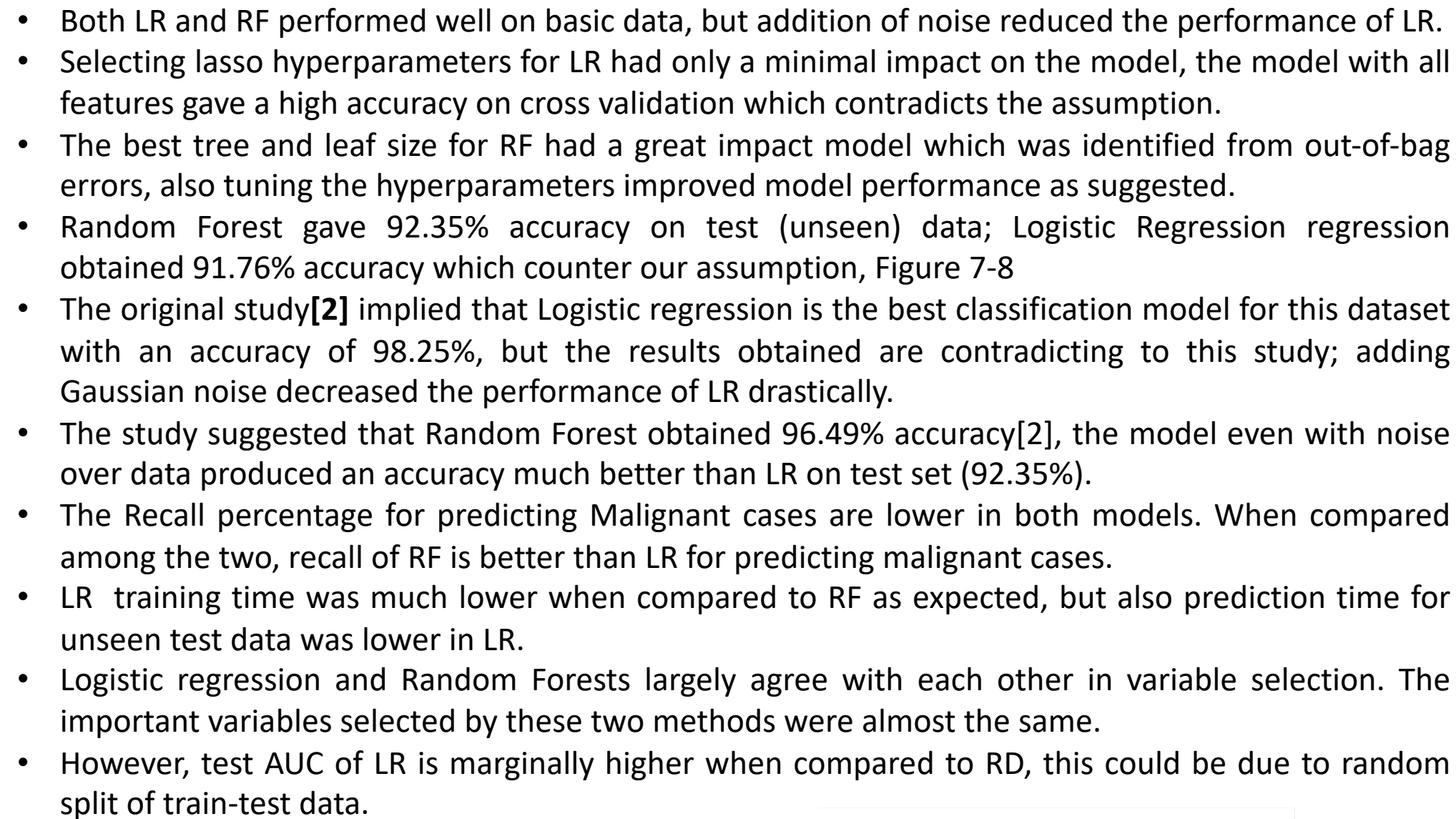
- A Binomially distributed model created for LR was Cross validated by 10-fold
- The goal is to represent the target class's posterior probabilities as linear functions of the predictors ensuring that they sum to one and stay within the range of 0 to 1. Setting a threshold value allows us to determine the expected target class from these probabilities (e.g., 0.5).
- Average accuracy(Cross validated by 10-fold) after performing Lasso Regularization for feature did not have a huge impact on the model when compared to all features.
- Although the AUC for train and test data was high, the percentage of predicted malignant cases were low in tested data.

- Bootstrap-aggregated (*bagged*) decision trees combine the results of many decision trees, which reduces the effects of overfitting and improves generalization
- Best results were found at the following hyperparameters: 1000 trees and minimum leaf size of 5 by iterating through different combinations of tree and leaf size. Figure 4
- Identifying important predictors above threshold of 0.5 had a slight improvement on the model.
- Feature selection, tree size and leaf size were obtained by validating out-of-bag records in the training data.

Figure 4



- Both training and test data had no missing values, identify how the models would perform in case of presence of missing values. If the performance reduces; what can be done further to improve it.
- Try and train with a large number of records as the current dataset had only 569 records, also try to improve the features of dataset by integrating more patient information like age, sex.
- Try to improve prediction of malignant cases(recall)



	Logistic Regression		Random Forest	
	B	M	B	M
Precision	0.89	0.98	0.91	0.94
Recall	0.99	0.79	0.97	0.84
F1 Score	0.93	0.87	0.94	0.89
AUC Test	0.991		0.978	
Training Acc.	94.03		95.61	
Test Acc.	91.76		92.35	
Training Time (Avg.)	0.043s		31.28s	
Testing Time(Avg.)	0.02s		0.77s	

Figure 5

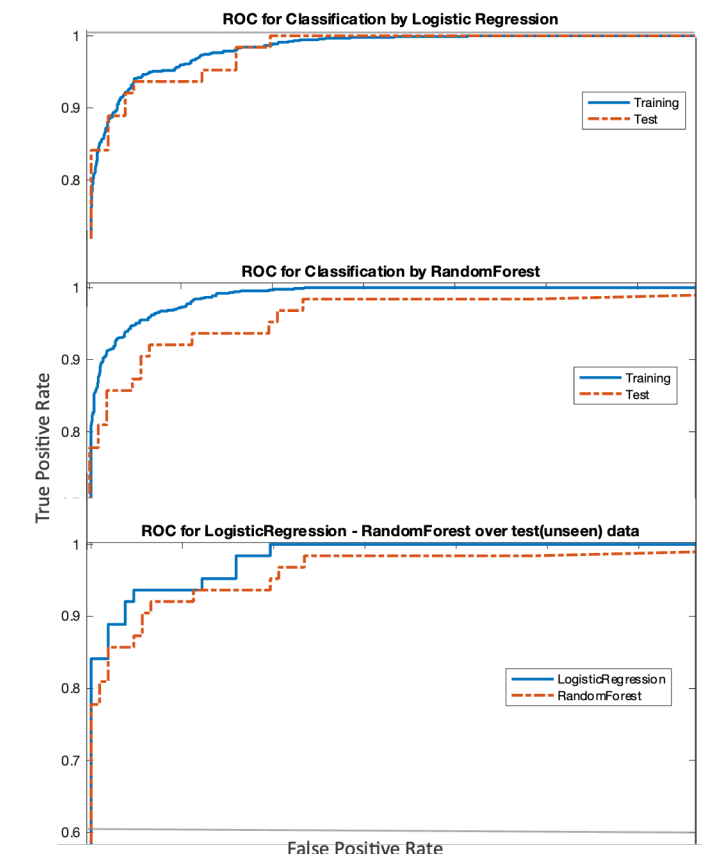


Figure 1

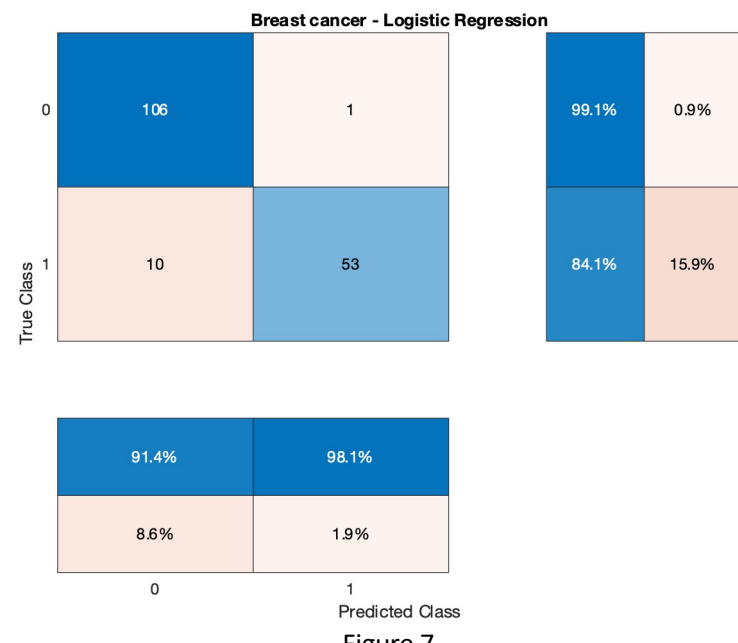


Figure 2

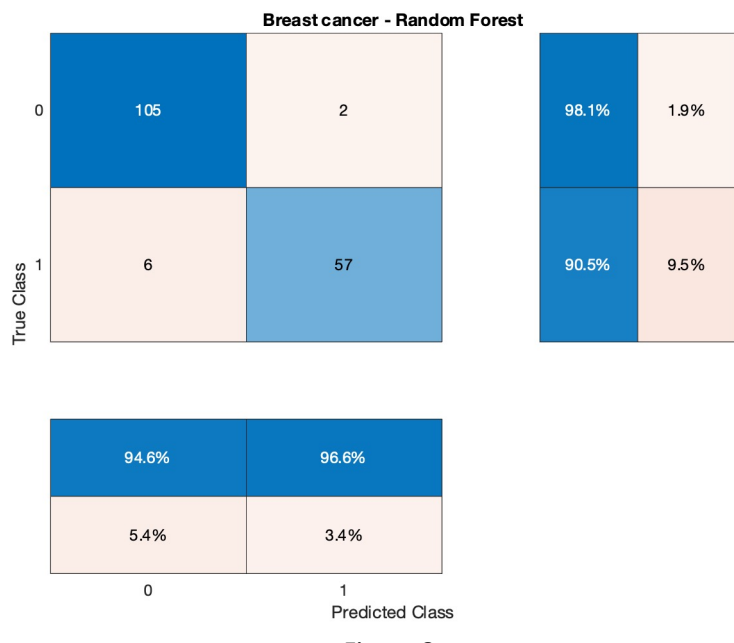


Figure 8

References

- [1] Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries; by *Hyuna Sung, MD; Jacques Ferlay, MSc, MEd; Rebecca L Siegel, MPH; Mathieu Laversanne, MSc; Isabelle Soerjomataram, MD, MSc, PhD; Ahmehin Jemal, DMV, PhD; Freddie Bray, BSc, MSc, PhD*
- [2] Breast Tumor Classification Using an Ensemble Machine Learning Method; by *Adel S. Assiri, Saima Nazir and Sergio A. Velastin*; Published: 29 May 2020
- [3] Breast Cancer Classification Using Machine Learning Techniques: A Review; by *Srwa Hasham Abdulla, Ali Mokki Saqheer, Haidi Veisi*; Published online: 19 August 2021
- [4] Fernandez, A.; Garcia, S.; Herrera, F.; Chawla, N.V. *SMOTE for Learning from Imbalanced Data: Progress and Challenges, Marking the 15-year Anniversary*; Published: 2018
- [5] Reducing Noise and Bias in Neural Network Models; by *Sovrit Ranjan Rathi*; Published February 3, 2020
- [6] *Mathematics*; United Kingdom: <https://www.mdpi.com/help/stats/feature-selection>, <https://www.mdpi.com/help/stats/feature-selection>, <https://www.mdpi.com/help/stats/feature-selection>
- [7] Solanki, V.; Chakrabarti, P.; Jasinski, M.; Leonowicz, Z.; Bolshev, V.; Vinogradov, A.; Jasinska, E.; Gono, R.; Nami, M. A Hybrid Supervised Machine Learning Classifier System for Breast Cancer Prognosis Using Feature Selection and Data Imbalance Handling Approaches. *Electronics* 2021, 10, 699
- [8] Personalized analysis of breast cancer using sample-specific networks; by *Ke Zhu, Cong Pian, Qiong Xiang, Xin Liu, Yuan Yuan Chen*; published on May 15, 2020
- [9] Discovering the shades of Feature Selection Methods by *Sweetha Manoj*—April 23, 2021