

Micah Fadrigio (Student ID #42923836), Cecilia Nguyen (Student ID #44328584)
Professor Xie
CS 273
24 March 2023

Final Project Report

- Describe the problem you chose
 - The methods you used to address the problem
 - Which model(s) you tried
 - How you trained them
 - How you selected any parameters they might require
 - How they performed on the test data
 - Tables of performance of different approaches
 - Plots of performance used to perform model selection (i.e., parameters that control complexity).
 - Describe to the best of your ability who was responsible for which aspects (which learners, etc.), and how the team as a whole put the ideas together.
 - We both contributed to initial data preprocessing and EDA. Cecilia utilized Support Vector Machines to train her model. Micah trained models using Random Forest and Adaboost. We made sure to keep our formats similar and explored different cross-validation methods to find optimal hyperparameters in terms of performance.
1. **Data exploration.** Visualize and analyze the data in different ways. It is really important that you are familiar with it. You should describe how you made various design choices, based on the dataset exploration.
 2. **Model exploration.** Exploration of at least one or two advanced techniques that you learned about in class. For example, using neural networks, support vector machines, or random forests are great ideas; but linear models and KNN classifiers are too basic. You should always explore in some depth the various options available to you for parameterizing the model, controlling complexity, etc. (This should involve more than simply varying a parameter and showing a plot of results.)
 3. **Data preprocessing and feature design.** You may have to design features, or optimize your models to deal with special aspects of the data (missing features, too many features, large numbers of zeros in the data; possible outlier data; etc.). Your report should describe what aspects you chose to focus on.
 4. **Performance validation.** You should practice good form and use validation or cross-validation to assess your models' performance, do model selection, combine models, etc. You should not simply try a few variations and assume you are done.
 5. **Adaptation to under- and over-fitting.** Machine learning is not very "one size fits all"; it is impossible to know for sure what model to choose, what features to give it, or how to set the parameters until you see how it does on the data. Therefore, much of machine learning revolves around assessing performance (e.g., is my poor performance due to underfitting, or overfitting?) and deciding how to modify your techniques in response. Your report should describe how, during your process, you decided how to adapt your models and why.

Dataset: Breast Cancer Wisconsin Diagnostic Data Set

Source: <https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

Introduction

Cancer is a disease in which cells in the human body grow abnormally as the damaged cells grow and multiply in the absence of biological signals. These cells may disrupt vital organ functions through tumors and destroy healthy cells, resulting in death. Cancer is the second-most common death, after heart attacks. Specifically, breast cancer has affected many women, being the second most common cause of death due to cancer after lung cancer. In the US, there are high annual incidence rates of breast cancer: 128.6 per 100,000 whites and 112.6 per 100,000 among African Americans. In regards to the Wisconsin data set we will use, we expect bias-the demographic to be more White-if there's an application of the model to patients of different ethnicities.

Breast cancer cells may form tumors, lumps of tissue which can be classified as either malignant or benign, affecting a patients' diagnosis and further treatment. A more accurate and timely diagnosis on patient outcomes with machine learning can revolutionize the healthcare industry. Utilizing the Wisconsin Breast Cancer Diagnosis Data Set, we will be classifying malignant and benign tumors. The feature measurements were computed from digital images of breast mass, where the image describes the characteristics of the cell nuclei, such as the mean radius, concavity, perimeter, texture.

EDA and Feature Preprocessing

The dataset initially contains 18,208 data points, with 33 features. We will be transforming our dataset to be machine learning ready. Since one of the features is all null values, we exclude it from the dataset. We will also be label encoding the diagnosis features into 0/1. Since our dataset is large, we will also be performing feature selection to improve model performance, reduce computational complexity, and reduce noise.

Noticing mean features as having high correlations with other features, we selected all the mean features from the dataset. Since mean features are already built for standard error features and worst features. Similar information with standard error and worst features would cause multicollinearity and increase computational complexity. The mean features will also be part of the standardization of the variables to allow for our model performance to not be biased due to extreme magnitudes between features. Standardization using a standard scalar would transform the features to follow standard normal distribution with zero mean and variance 1, and therefore make the features be comparable. If necessary, if model performance on test data is not a high percentage (~90%), we will add features most important.

1 SVM

SVM is a supervised algorithm that can be utilized for classification. SVM increases the dimensionality of data in order to find a hyperplane to separate the two classes through finding a hyperplane. The support vector classifier is the decision boundary between the two classes, while the maximized margins (distance between line and points on the margin, called support vectors) allow for misclassifications to reduce bias and make the model more robust to outliers, in the form of a soft SVM. To determine how to increase dimensionality, a kernel function is used to assign new coordinates of the data. The regularization parameter C controls the softness of the margins and decision boundary broadness. The gamma parameter is used for the rbf kernel, computing the similarity between input data points.

1.1 Preprocessing Phase

Since SVM can be sensitive to highly correlated data, our mean data features were standardized to decorrelate the features. We also included the regularization parameter values when creating our SVM classifier, C, in our model, to discourage the model from relying too heavily on any feature. Additionally, we utilized cross-validation methods to select the most appropriate kernel function to mitigate the effects of correlation.

1.1 Training and classification

Imbalanced data also had to be dealt with, where 63% of observations belong to the benign class and 37% belong to the malignant class, the weighted class SVM was implemented in order to assign higher misclassification penalties to training instances for our smaller malignant class.

To find optimal parameters and perform hyperparameter tuning, we utilized cross validation methods. We began with first splitting the model into train and test sets, then validation and training sets (Used for hyperparameter tuning). We then defined a parameter grid with the following hyperparameters. In our first method included we defined a search space with a range of kernel, gamma, and C values. The parameters will give us the best model given our parameters chosen, and this cross validation will help in also determining the dimensions required.

```
kernels = ['linear', 'poly', 'rbf']
C_values = [0.001, 0.01, 0.1, 1, 10, 100]
gamma_values = [0.001, 0.01, 0.1, 1, 10, 100]
```

General Cross Validation: Define parameters to tune

Finding that the best model from the first method included an rbf kernel function, we used the parameter in our model and implemented Bayesian optimization for cross validation, which runs in a reasonable time compared to the other methods with long runtime, such as Grid Search and Random Search.

```
param_bounds = {'C': (0.1, 100), 'gamma': (0.01, 10)}
```

Bayesian Optimization: Define parameters to tune

1.1 Experimental results

After conducting General Cross Validation and Bayesian optimization on the train data, it determined the optimal values for the hyperparameters (shown below). Our two models utilized the default cross validation parameter of 5 folds, which was then applied on the test set to measure the model's performance using the optimal hyperparameter choices.

```
Best parameters: SVC(C=1, class_weight='balanced', gamma=0.1, probability=True)
Kernel: rbf
```

Optimal Parameters Generated from General Cross Validation Method

```
Total time taken by 5-fold Cross Validation with SVM is 7.17618862700002 seconds
Classification Report
      precision    recall  f1-score   support

     0       0.96      0.95      0.95        75
     1       0.90      0.92      0.91        39

 accuracy          0.93
 macro avg          0.93
weighted avg          0.94

Recall Score: 1.00
Precision Score: 0.95
Accuracy: 0.98
```

Overall Performance General Cross Validation

```
Best parameters: {'C': 1.7827503100600173, 'gamma': 0.01}
```

Optional Parameters Generated From Bayesian Optimization Method

```
Total time taken by 5-fold Cross Validation with SVM is 34.198644588 seconds
Average Performance Measure of SVM Model using Bayesian Optimization
Classification Report
      precision    recall  f1-score   support

     0       0.94      0.99      0.96        75
     1       0.97      0.87      0.92        39

 accuracy          0.95
 macro avg          0.93
weighted avg          0.95

Recall Score: 0.87
Precision Score: 0.97
Accuracy: 0.98
```

Overall Performance Bayesian Optimization Method

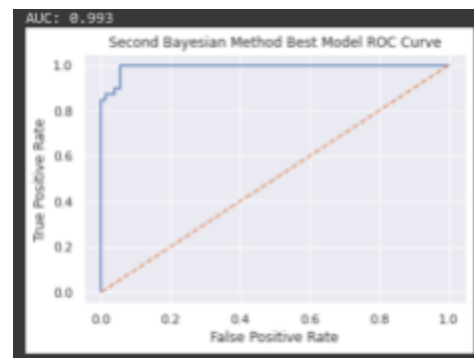
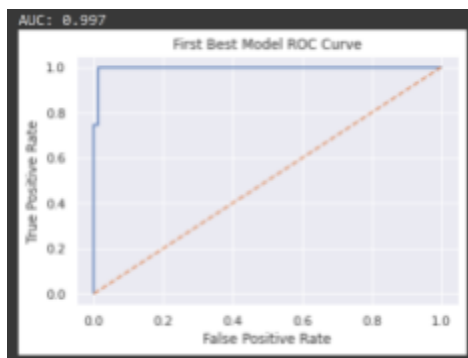
Referring to the performance metric (classification report) above for test data, for medical diagnosis, we prioritize high precision (ability to avoid false positives) and high recall (ability to avoid false negatives). F1 score is a model performance measure that captures both, in which an F1 score of 1 indicates perfect precision and recall. The classification reports are given above.

For General Cross Validation SVM, the model performed well on unseen data, with high precision and high recall, as indicated by the F1 score of 0.91 for malignant and 98% testing accuracy. For Bayesian Optimization, the model also performed well on unseen data, with high precision and high recall, as indicated by the F1 score of 0.92 for malignant and 98% testing accuracy. However, for the General Cross Validation SVM, a high recall score of 100 would mean that the model had correctly identified all patients who have the disease, but it can also classify healthy patients as having a tumor, leading to false positives that can lead to unnecessary/harmful treatment for healthy patients.

Classification Report (Train)				
	precision	recall	f1-score	support
0	0.96	0.97	0.96	282
1	0.95	0.94	0.94	173
accuracy			0.96	455
macro avg	0.95	0.95	0.95	455
weighted avg	0.96	0.96	0.96	455

Classification Report (Train)				
	precision	recall	f1-score	support
0	0.98	0.99	0.99	282
1	0.99	0.97	0.98	173
accuracy			0.98	455
macro avg	0.98	0.98	0.98	455
weighted avg	0.98	0.98	0.98	455

To check for possible overfitting we can compare the F1 scores of training and test data. If both F1 scores are high, it means that the model is probably not overfitting and able to generalize well to new data. For both models, they both have high F1 values for train/test sets, meaning that the classifiers can generalize well to new data.



Another alternative to visualize the performance of our classification model would be to look at an ROC/AUC curve, where our y-axis is the True Positive Rate (recall) and our x-axis is the False Positive Rate (1-precision). Since we have an imbalance in the classification of tumors, we would like to utilize ROC/AUC since it is a more robust evaluation, and provides a more comprehensive measure. The diagonal line represents a random classifier which makes random guesses of the classes in a 50-50 manner. Our ROC curve is around 1, indicating a high classification ability. Looking at the AUC's (Area under the ROC curve), it is also high (~99). Both of these measurements indicate that our classifier model has a high classification ability to separate benign/malignant tumors.

2 Random Forest

Random Forest is an ensemble learning algorithm that combines multiple decision trees, each trained on a random subset of features and the data. The final prediction is made by aggregating the predictions of all individual trees. These characteristics allow random forest to be able to prevent overfitting, generalize well to unseen data, and reduce the effect of correlated features.

2.1 Preprocessing Phase

Since random forest reduces the effect of correlated features to some extent, all of the “mean” type features were used and highly correlated features were not removed.

2.2 Training & Classification

Imbalanced data can reduce the effectiveness of random forest by leading to biased selection of features that are more prevalent in the majority class, poor tuning of hyperparameters, and misclassification of the minority class. To address class imbalance, we used Stratified K-fold cross validation, and GridSearchCV combined. Stratified K-fold cross validation is essentially K-fold cross validation, except that it ensures that within each fold, the proportion of each class is roughly equal. We also used GridSearchCV for hyperparameter tuning which will search over a range of hyperparameters to find the combination that yields the best performance on the validation set. Using these two methods together will allow us to make accurate hyperparameter choices.

The data was split into train and test sets and a parameter grid was defined, which consisted of different values for the following hyperparameters: number of decision trees, the maximum number of features randomly selected for each base learner, the maximum depth of each base learner, the minimum number of samples required to be in a leaf node, and the criterion to evaluate the quality of a split. After conducting GridSearchCV on the training set, it determined the optimal values for the hyperparameters (shown in first table). Stratified 10 fold cross validation was then applied on the training set to score the selected combination of hyperparameters (second table).

Name of Hyperparameter	Hyperparameter Values	Optimal Values Determined By Grid Search CV
n_estimators	[50, 100, 200]	200
max_features	range(1, 11)	4
max_depth	range(2, 15)	7
min_samples_leaf	range(1, 3)	1
criterion	['gini'], ['entropy']	'gini'

Precision	Recall	F1 Score	Accuracy
0.913	0.944	0.925	0.941

Performance Measures (via 10-fold CV) of Best Combination of Hyperparameter Value

2.3 Experimental Results

Precision	Recall	F1 Score	Accuracy
0.947	0.923	0.935	0.956

Performance of Random Forest on Test Data

In the context of medical diagnosis, it is important to prioritize high precision (ability to avoid false positives) and high recall (ability to avoid false negatives). F1 score is a model performance measure that captures both, in which an F1 score of 1 indicates perfect precision and recall. The random forest model with hyperparameter choices generated from GridSearchCV is expected to perform well on unseen data, with high precision and high recall, as indicated by the F1 score of 0.93 and 96% accuracy (seen above).

3 Gradient Boosting (AdaBoost)

AdaBoost is a gradient boosting algorithm that, similar to random forest, averages over many weak base learners to produce a strong classifier. AdaBoost works by iteratively training weak learners on different subsets of the training data and adjusting the weights of the training examples to emphasize the examples that are misclassified by the current set of weak learners. AdaBoost is suitable for our problem because it handles imbalanced data (to some degree) by assigning higher weights to misclassified samples (the minority class) in order to improve overall accuracy. However, a drawback is that when the dataset contains highly correlated features, the weak learners in AdaBoost may place too much emphasis on those features which contribute similar information to the model. This can lead to overfitting and reduced generalization ability.

3.1 Preprocessing Phase

For the dataset belonging to the **AdaBoost classifier 1**, highly correlated features were removed from the “mean” type features in which `perimeter_mean`, `area_mean`, `concavity_mean`, and `concave_points_mean` were removed. Six features in total were used for the first **AdaBoost classifier 1**. In an attempt to improve this model, these features were later scaled prior to conducting Lasso cross validation, however, no features were eliminated. Six features in total were used for the **AdaBoost classifier 1**.

For the dataset belonging to the **AdaBoost classifier 2**, highly correlated features were removed from all type features (“mean”, “standard error”, and “worst”). The features for **AdaBoost classifier 1** are a subset of the features for **AdaBoost classifier 2**, and the 4 additional features for **AdaBoost classifier 2** are `'symmetry_worst'`, `'symmetry_se'`, `'smoothness_se'` and `'texture_se'`. In an attempt to improve this model, these features were later scaled prior to conducting Lasso cross validation, however, no features were eliminated.. Ten features in total were used for the **AdaBoost classifier 2**.



Feature Importance (Lasso CV) of AdaBoost Classifier 1

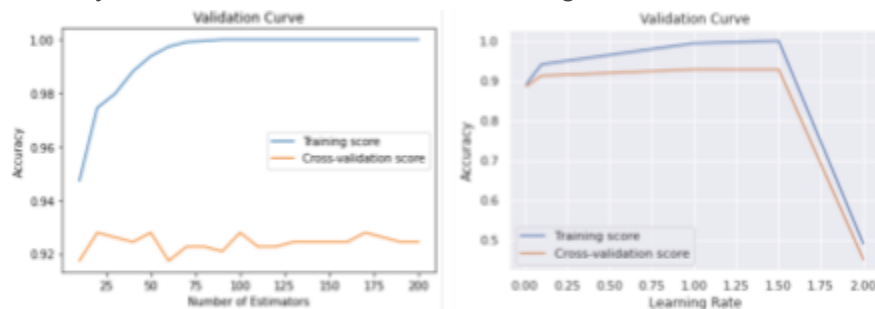


Feature Importance (Lasso CV) of AdaBoost Classifier 2

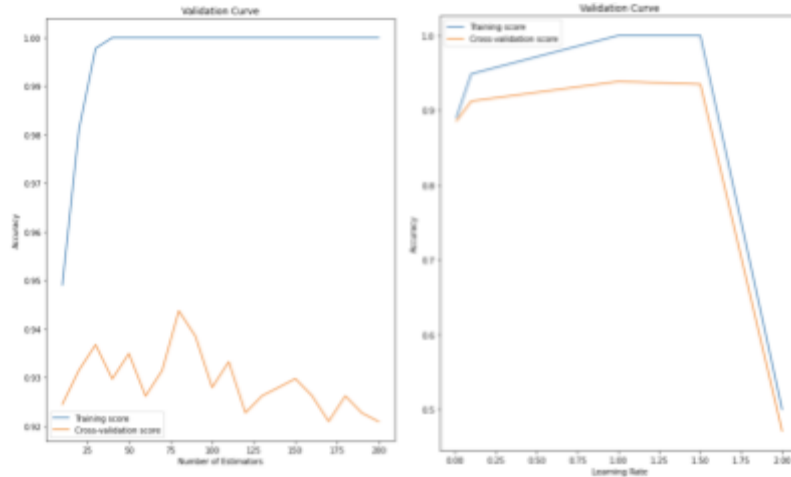
3.2 Training & Classification

The number of base learners (`n_estimators`) is the number of decision stumps and the learning rate (`learning_rate`) is the factor that scales the contribution of each base learner before adding it to the final ensemble. A small learning rate means that the classifier will take smaller steps towards the optimal solution and may require more iterations to converge, however it may result in high accuracy and low variance. A large learning rate means that the classifier will take bigger steps towards the optimal solution and may require fewer iterations to converge, however it may result in low accuracy and high variance. It is necessary to tune both hyperparameters to balance the bias-variance tradeoff.

Below are two plots that show the effect of model performance on training and validation sets as a function of `n_estimators` and `learning_rate` values, using the dataset belonging to **AdaBoost classifier 1**. When `n_estimators` = 100 and `learning_rate` = 1.5, accuracy on the validation set is maximized. Anything beyond a learning rate of 1.5 causes a significant decrease in both training and validation accuracy. These estimates are also verified using GridSearchCV.



Below are similar plots for **AdaBoost classifier 2**. When `n_estimators` = 80 and `learning_rate` = 1.5, accuracy on the validation set is maximized. Anything beyond a learning rate of 1.5 causes a significant decrease in both training and validation accuracy. These estimates are also verified using GridSearchCV.



3.3 Experimental Results

Below are the average performance metrics using 10-fold cross validation.

Model	Precision	Recall	F1 Score	Accuracy
AdaBoost Classifier 1 (n_estimators = 100, lr = 1.5, maxDepth = 1)	0.923	0.914	0.916	0.922
AdaBoost Classifier 2 (n_estimators = 80, lr = 1.5, maxDepth = 1)	0.929	0.920	0.922	0.928
AdaBoost Classifier 2 (n_estimators = 80, lr = 1.5, maxDepth = 4)	0.938	0.932	0.933	0.938

Model	Total Time
AdaBoost Classifier 1 (n_estimators = 100, lr = 1.5, maxDepth = 1)	2.952
AdaBoost Classifier 2 (n_estimators = 80, lr = 1.5, maxDepth = 1)	2.413
AdaBoost Classifier 2 (n_estimators = 80, lr = 1.5, maxDepth = 4)	4.593

Because **AdaBoost classifier 2** outperformed **AdaBoost classifier 1** on all performance measures, we wanted to improve upon **AdaBoost classifier 2**, specifically by making its base learners more complex. Keeping `n_estimators` = 80 and `learning_rate` = 1.5, we tested `max_depth` values ranging from one to five using GridSearchCV. When the maximum depth was set to 4, all performance measures increased by approximately 0.01, however, the training time nearly doubled.