

Explanation and Analysis of Models Chosen

HARVEY KWONG and JACOB DEROSA, University at Buffalo, USA

1 Introduction

Following thorough data cleaning and exploratory data analysis, we identified 15 key predictor features. Our objective is to utilize these predictors to classify whether a given sample of *Neisseria gonorrhoeae** exhibits super resistance to specific antibiotics. In this study, resistance to azithromycin serves as the target label for our classification models. Since the target variable is binary (either true or false), we implemented and evaluated six different types of classifiers:

- K-Nearest Neighbors
- Naive Bayes
- Logistic Regression
- Support Vector Machine (SVM)
- Neural Network - (Not from Class) Harvey Kwong’s custom architecture
- Extreme Gradient Boosting (XGBoost) - (Not from class)

K - Nearest Neighbors

We selected K-Nearest Neighbors as one of the models for our analysis. KNN is a good fit for our problem due to the relatively low dimensionality of the dataset, where we focused on only the most significant features identified during our exploratory data analysis process. By reducing the number of features, KNN’s performance improves, as it works best when our data is not at a very high dimension.

For this model, we chose to use 5 neighbors, as this configuration provided the most balanced results. Fewer neighbors made the model more sensitive to noise and overfitting, while more neighbors reduced its ability to capture important variations. Overall, KNN performed strongly, achieving approximately 96% accuracy in predicting azithromycin resistance. The extended metrics for this model are as follows:

Metric	Score
Accuracy	0.9647
Precision	0.9714
Recall	0.9647
F1 Score	0.9667

Table 1. Performance Metrics for K-Nearest Neighbors

Authors’ Contact Information: Harvey Kwong, harveykw@buffalo.edu; Jacob DeRosa, jderosa3@buffalo.edu, University at Buffalo, Buffalo, New York, USA.

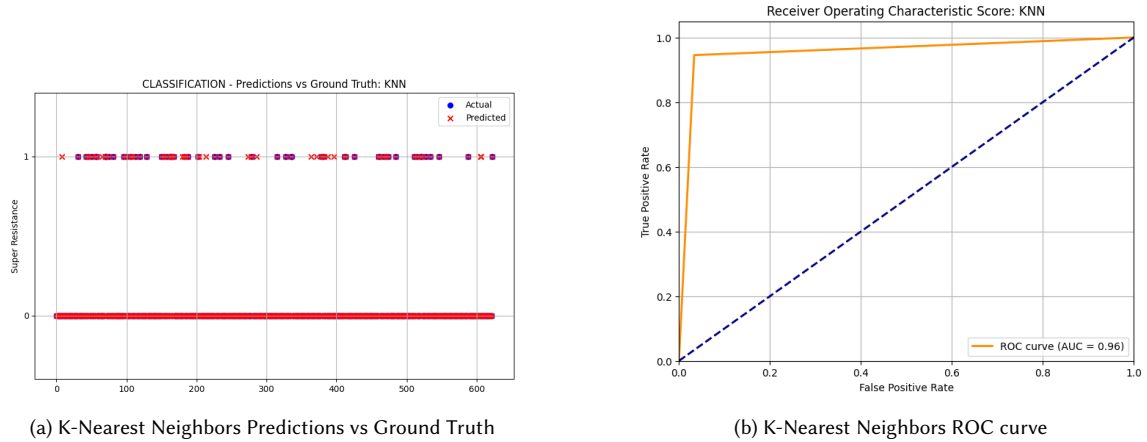


Fig. 1. K-Nearest Neighbors

Naive Bayes

We chose to implement Naive Bayes as one of our classification models. Naive Bayes is well-suited to our problem due to its efficiency and simplicity, especially when dealing with binary classification tasks like predicting antibiotic resistance. The model assumes that the features are independent, which, while a simplification because many things are interconnected in biology – performed surprisingly well in our modelling. Given the relatively reduced number of features in our dataset, Naive Bayes was an appropriate choice as it tends to work well even with limited data.

Additionally, Naive Bayes requires minimal computational resources, which makes it advantageous for fast predictions. Although it does not capture feature interactions, it still correctly predicted azithromycin resistance approximately 91% of the time. Our extended metrics for this model are as follows:

Metric	Score
Accuracy	0.9197
Precision	0.9545
Recall	0.9197
F1 Score	0.9298

Table 2. Performance Metrics for Naive Bayes

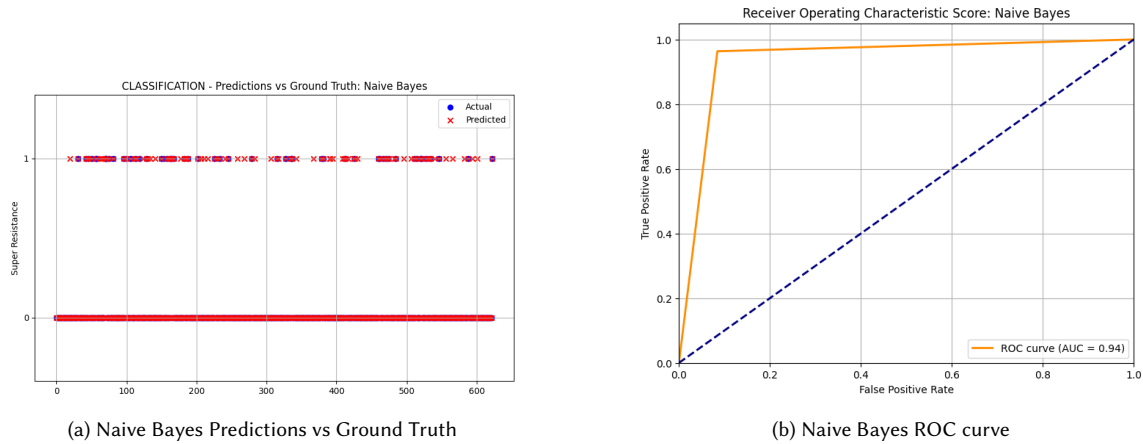


Fig. 2. Naive Bayes

Logistic Regression

We chose to include Logistic Regression as one of the models for our classification task. Logistic Regression is a commonly used algorithm for binary classification problems like ours, where the goal is to predict azithromycin resistance. One of the main advantages of Logistic Regression is its simplicity and interpretability, as it provides clear probabilistic outputs. This makes it particularly useful for understanding the relationship between the predictor variables and the target label. Though in our model, we thresholded the predicted values. Values above 0.5 went to positive, and those below went to negative.

Since our dataset contains well-processed and balanced features, Logistic Regression was able to achieve a solid performance without much additional complex tuning. Despite using a max iteration limit of 10, the model demonstrated reliable results, with an accuracy of around 95%. Our extended metrics for this model are as follows:

Metric	Score
Accuracy	0.9502
Precision	0.9681
Recall	0.9502
F1 Score	0.9550

Table 3. Performance Metrics for Logistic Regression

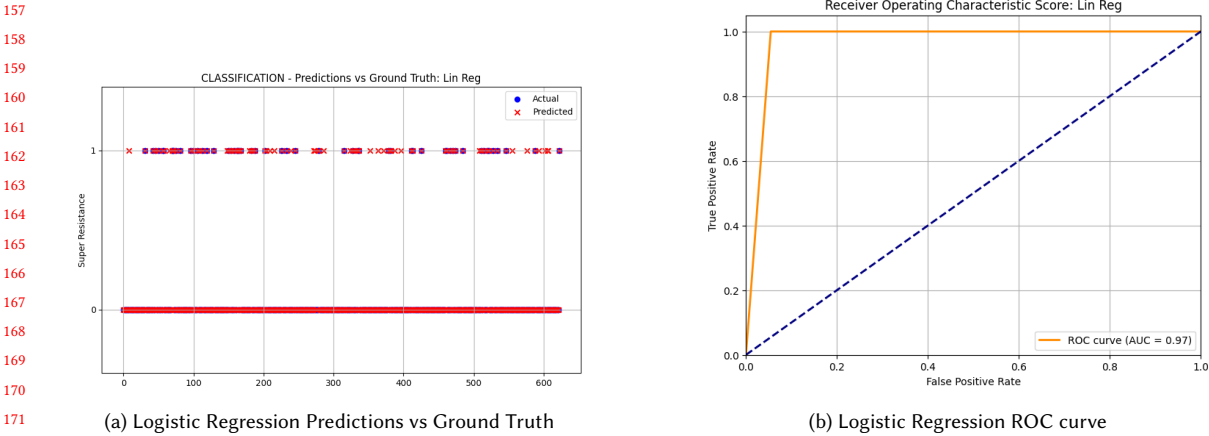


Fig. 3. Logistic Regression

Support Vector Machine (SVM)

We chose to include support vector machines as one of the models for our classification task. SVM is a powerful algorithm for binary classification problems. It works by finding an optimal hyperplane that separates the classes with the largest margin, making it effective in cases where the data is not linearly separable. SVM is well-suited to this problem as it can handle non-linear decision boundaries using various kernels. For our implementation, we used the default polynomial kernel and set the maximum number of iterations to 10 to ensure computational efficiency while maintaining acceptable performance. Despite the iteration limit, SVM performed reasonably ok, achieving an accuracy of approximately 87%. Our extended metrics for our SVM model are as follows:

Metric	Score
Accuracy	0.8732
Precision	0.9249
Recall	0.8732
F1 Score	0.8911

Table 4. Performance Metrics for Support Vector Machine

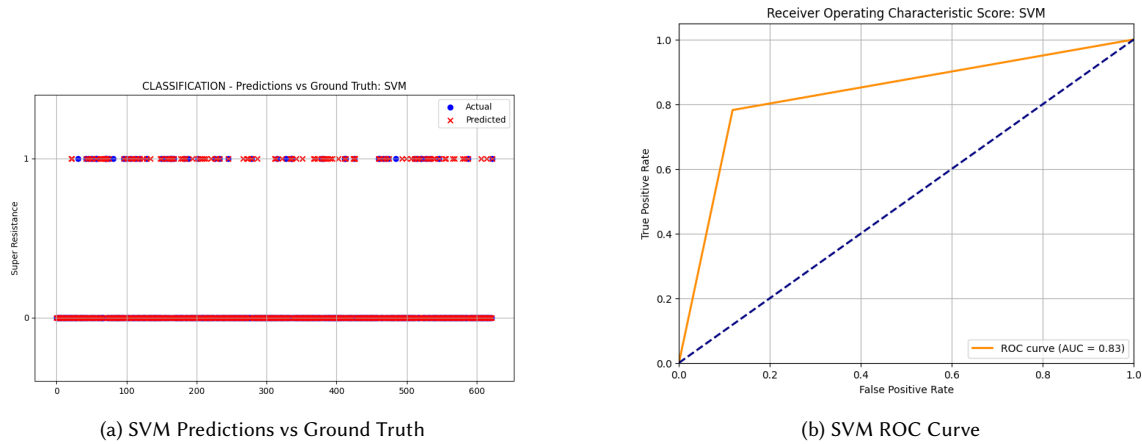


Fig. 4. Support Vector Machine

Neural Network

We included a custom Neural Network as one of the models for our classification task. Neural networks are highly flexible and capable of capturing complex patterns in data, making them suitable for problems where non-linear relationships between features and the target label exist. Given the relatively simple structure of our dataset, we designed a straightforward architecture to avoid overfitting and any other unnecessary complexities.

Our custom architecture consists of two hidden layers with 8 and 4 neurons, respectively, both using the RELU activation function. For the output layer, we used a single neuron with a sigmoid activation function, which is the norm for binary classification tasks. We compiled the model with the Adam optimizer and binary cross-entropy loss, ensuring efficient training with a learning rate of 0.1. Adam will automatically adjust the learning rate as it sees fit.

Despite the simplicity of the architecture, the neural network performed fine, achieving an accuracy of approximately 90%. We trained the model for 10 epochs, to enable easier comparisons with our other models which have had their maximum iterations set to 10. Our neural network did not perform as well as logistic regression despite having a more complicated structure and presumably higher computational cost. It is likely better to use logistic regression on our dataset compared to our neural network with our current choice of hyperparameters. I expect that if we were to increase the number of hidden units and epochs, we would see much better performance, though at a higher cost. Our extended metrics for this neural network model are as follows:

Metric	Score
Accuracy	0.9053
Precision	0.9543
Recall	0.9053
F1 Score	0.9192

Table 5. Performance Metrics for Neural Network

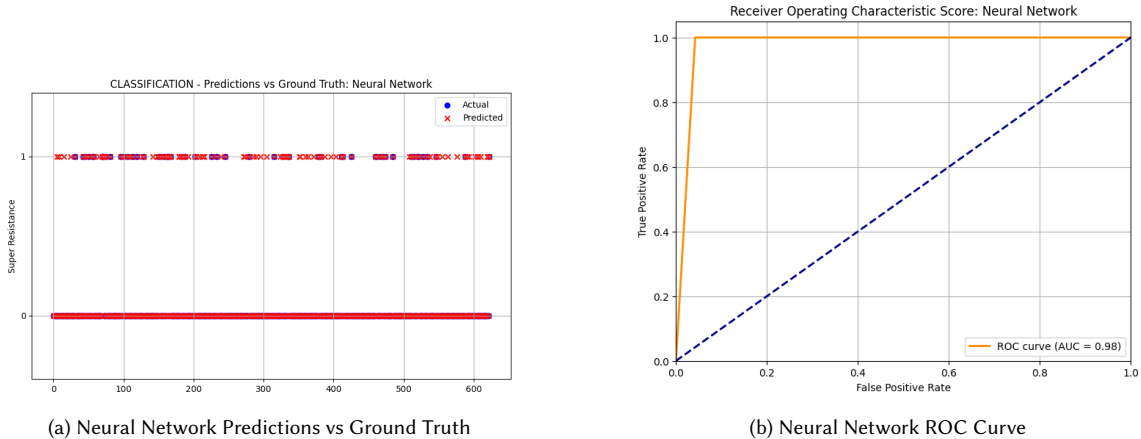


Fig. 5. Neural Network

XGBoost

We included Xgboost as one of the models for our classification task. Xgboost (Extreme Gradient Boosting) is a highly efficient and powerful implementation of gradient boosting, designed to optimize both speed and performance. It is particularly well-suited for tabular data such as ours and has become one of the go-to competition algorithms for many classification problems due to its ability to handle feature interactions, missing data, and non-linear relationships.

For our implementation, we used an Xgboost model with a logistic objective for binary classification. We set the number of estimators to 100, with a maximum depth of 3 and a learning rate of 0.01 to control for overfitting and ensure smooth learning.

XGBoost delivered impressive results, achieving an accuracy of approximately 95.5%. This demonstrates the model’s ability to generalize well while maintaining strong predictive performance. Our extended metrics for this XGBoost model are as follows:

Metric	Score
Accuracy	0.9551
Precision	0.9702
Recall	0.9551
F1 Score	0.9590

Table 6. Performance Metrics for XGBoost

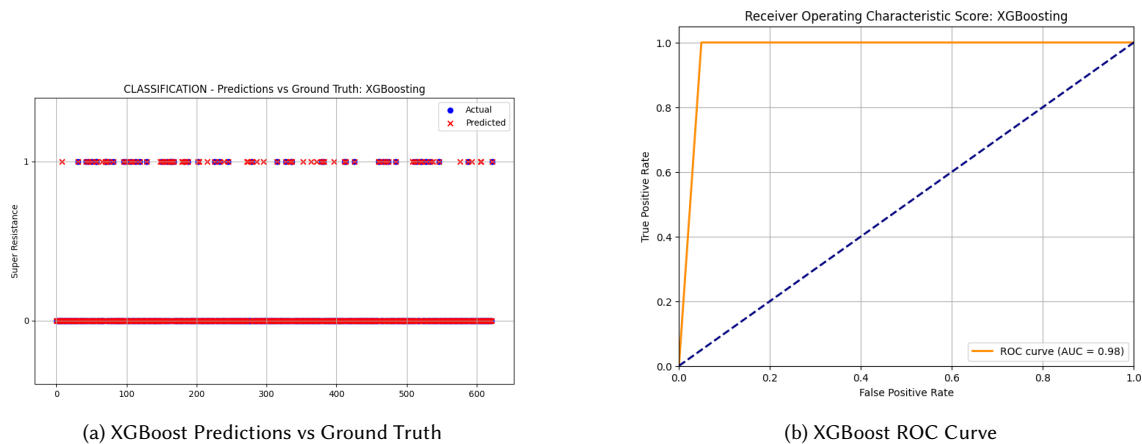


Fig. 6. XGBoost

2 Analysis

K-Nearest Neighbors, with 5 neighbors demonstrated the highest overall performance, achieving an accuracy of 96.47%, with a precision of 97.14% and an F1 score of 96.67%. This shows that KNN was effective in correctly predicting both true positives and negatives, offering a balanced performance across all metrics. The high precision shows that KNN minimizes false positives well.

Naive bayes, though simple and computationally efficient, had a lower accuracy of 91.97% compared to other models, with a precision of 95.45%. Despite this, Naive Bayes still performed relatively well, particularly with a strong F1 score of 92.99%, showing that it is capable of achieving good results even with simplified assumptions of feature independence. However, the lower recall shows that it may miss some true positives compared to the higher-performing models.

Logistic regression, with a max iteration limit of 10, provided decent performance with an accuracy of 95.02%, a precision of 96.82%, and an F1 score of 95.50%. Logistic regression maintains a good balance between efficiency and predictive performance, offering a reliable model that consistently performs well across various metrics. The slight drop compared to KNN in recall suggests that Logistic Regression may have a higher chance to miss some true positives.

Support vector machine, also capped at 10 iterations, had the lowest overall performance. It had an accuracy of 87.32% and an F1 score of 89.11%. While SVM is powerful for high-dimensional data, its performance here indicates that it may have struggled with the data's complexity or structure, especially with the limited number of iterations, resulting in a higher false positive rate compared to other models.

The neural network, with a simple architecture, performed reasonably well with an accuracy of 90.53%, precision of 95.43%, and an F1 score of 91.92%. Neural networks typically perform best at capturing non-linear relationships. While the model performed fine, it did not outperform the more straightforward models like KNN or Logistic Regression in this instance, possibly due to the limited number of epochs or the simplicity of the architecture.

Finally, XGBoost, known for its performance in tabular data, performed on par with KNN, achieving an accuracy of 95.51%, precision of 97.02%, and an F1 score of 95.90%. XGBoost's ability to handle complex feature interactions and provide strong generalization explains why it performed at the top, closely competing with KNN in terms of predictive power. It also maintained a good balance between precision and recall, making it a very good choice for this classification task despite the computational cost compared to KNN.

In conclusion, KNN and XGBoost emerged as the top performers, with slightly higher precision and accuracy than the other models, making them ideal for this dataset. Logistic Regression also provided solid performance, while Naive Bayes and Neural Network offered decent results but fell slightly short of the two leading models. SVM had the lowest performance, likely requires further tuning or more iterations to improve its results.

3 Citations for outside models

- 1 Neural Network - Prior Knowledge. Harvey Kwong has experience building simple neural network architectures in various ML frameworks including Keras.
- 2 XGBoost - https://xgboost.readthedocs.io/en/stable/get_started.html