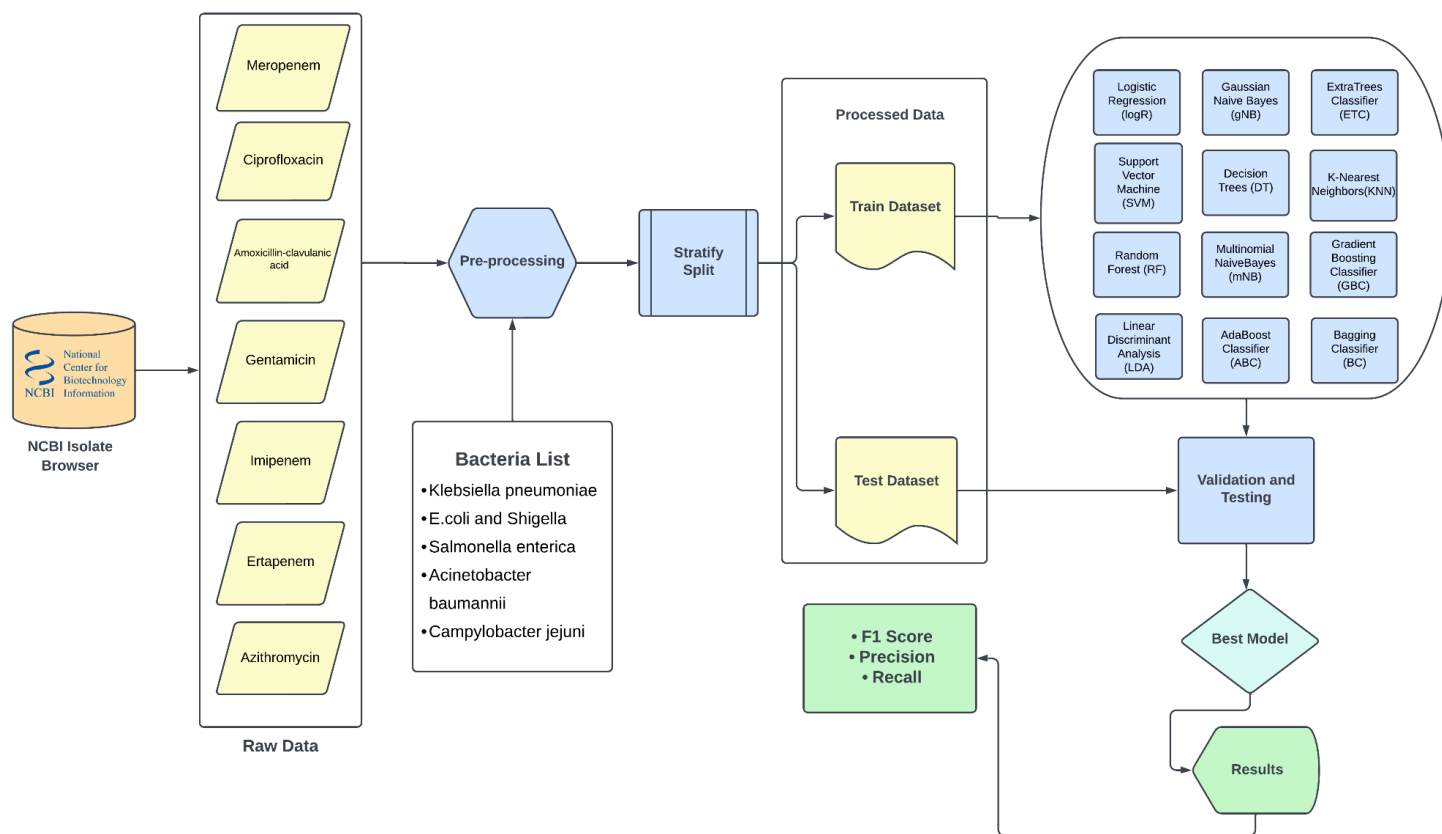


# Antimicrobial Resistance Prediction Using Statistical ML Models



View Diagram here: [Click Here](#)

## Datasets:

Bacterial strains, AMR genotypes and AST phenotypes

Website: [Isolates Browser - Pathogen Detection - NCBI](#)

Raw Datasets:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/tree/master/raw\\_data\\_from\\_NCBI](https://github.com/RAGUL1902/AMR_Prediction_Project/tree/master/raw_data_from_NCBI)

Processed Datasets:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/tree/master/datasets](https://github.com/RAGUL1902/AMR_Prediction_Project/tree/master/datasets)

Information about the datasets:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/blob/master/datasets\\_info.csv](https://github.com/RAGUL1902/AMR_Prediction_Project/blob/master/datasets_info.csv)

The NCBI Isolate Browser is a web-based tool developed by the National Center for Biotechnology Information (NCBI) that allows users to explore and analyze microbial genomes and their associated metadata.

Bacterial genus and species were selected in the 'organism group', with filters checked for 'has AMR genotypes' and 'has AST phenotypes'. Then the isolates were downloaded for different antibiotics separately.

The Collected Antibiotics are as follows:

- Meropenem
- Ciprofloxacin
- Amoxicillin-clavulanic acid
- Gentamicin
- Imipenem
- Ertapenem
- Azithromycin

Once this data was collected, the following list of bacteria isolates were segregated from each antibiotic dataset.

- *Klebsiella pneumoniae*
- *E.coli* and *Shigella*
- *Salmonella enterica*
- *Acinetobacter baumannii*
- *Campylobacter jejuni*

The combination of these bacteria and antibiotics were separated and made into different datasets. Based on the AMR genotypes and AST phenotypes, binary matrices of genotypes (0 for absence and 1 for presence of an AMR gene) and relevant antibiotics' phenotypes (0 for susceptibility and 1 for resistance to an antibiotic) were created.

## **Prediction - Machine Learning Models**

Code for statistical models:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/blob/master/Statistical\\_models\\_on\\_all\\_datasets.ipynb](https://github.com/RAGUL1902/AMR_Prediction_Project/blob/master/Statistical_models_on_all_datasets.ipynb)

Machine learning in Python with Scikit-learn (<https://scikit-learn.org/stable/>) was used to evaluate the performance of 12 machine learning algorithms, namely,

- Logistic Regression (logR),
- Gaussian Naive Bayes (gNB),
- Support Vector Machine (SVM),
- Decision Trees (DT),
- Random Forest (RF),
- K-Nearest Neighbors(KNN),
- Linear Discriminant Analysis (LDA),
- Multinomial NaiveBayes (mNB),
- AdaBoost Classifier (ABC),
- Gradient Boosting Classifier (GBC),
- ExtraTrees Classifier (ETC),
- Bagging Classifier (BC).

We used the bacterial strains' AMR genotype and AST phenotype data for training and testing the machine learning algorithms. The performance metrics include precision, recall, F1-score and Confusion Matrix.

## Results

Results - Tabular format:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/blob/master/Best\\_statistical\\_model\\_on\\_all\\_datasets.csv](https://github.com/RAGUL1902/AMR_Prediction_Project/blob/master/Best_statistical_model_on_all_datasets.csv)

Results - Plots:

[https://github.com/RAGUL1902/AMR\\_Prediction\\_Project/blob/master/plot\\_results.ipynb](https://github.com/RAGUL1902/AMR_Prediction_Project/blob/master/plot_results.ipynb)

Assessment of different machine learning algorithms on the entire dataset of bacterial strains with known genotypes and phenotypes showed there is no single optimal machine learning model for predicting resistance phenotype across all bacterial species considered for the prediction. Added the Results as bar plots:

