

Predicting Breast Cancer Subtypes and Treatment Responses Using Machine Learning

Professor:

Hong Xue

Group 6:

- 1) Poojitha Anuvalasetty**
- 2) Rithika Vadugapattu Jayagopi**
- 3) Sai Siva Sri Harsha Vasinenku**

Member contribution statement:

The project was a collaborative effort by Rithika, Poojitha, and Harsha, with each member contributing to different tasks:

- **Rithika:** Worked on preparing and cleaning the dataset by handling missing values and organizing the data for analysis. She also contributed to summarizing the findings and helped with the final report and presentation.
- **Poojitha:** Focused on selecting important features for analysis and worked on building and improving the machine learning models. She also helped in interpreting the results and contributed to writing and refining the final report and presentation.
- **Harsha:** Worked on developing, training, and evaluating the machine learning models. Harsha also contributed to analyzing the outcomes and writing the final report to ensure the findings were clearly presented.

All members collaborated to discuss results, ensure the report was complete, and prepare a well-structured presentation.

Abstract

Breast cancer is one of the leading causes of cancer-related mortality worldwide, making accurate prediction and effective treatment planning a critical focus in healthcare. This project leverages the METABRIC dataset, a comprehensive clinical and genomic dataset from 1,903 breast cancer patients, to build predictive machine learning models aimed at enhancing personalized treatment strategies. Key objectives include predicting overall survival, identifying cancer subtypes, and estimating treatment responses.

The study integrates advanced data preprocessing techniques, feature selection methods such as SHAP (SHapley Additive exPlanations), and machine learning algorithms including Logistic Regression, Random Forest, and Gradient Boosting. These models were evaluated using metrics like accuracy, ROC-AUC, and confusion matrices to benchmark their performance. Interpretability of the models was emphasized, ensuring clinical relevance and actionable insights.

Results demonstrate the feasibility of machine learning in predicting critical outcomes with high accuracy, paving the way for its integration into clinical workflows. This project not only highlights the transformative potential of data-driven approaches in oncology but also establishes a framework for future explorations into personalized breast cancer care.

Introduction

Breast cancer continues to pose a major health challenge, accounting for a significant proportion of cancer-related deaths among women worldwide. While advancements in diagnostics and treatment have improved survival rates, the heterogeneity and complexity of breast cancer demand innovative approaches to predict key outcomes, including overall survival, cancer subtypes, and treatment response. Accurate predictions in these areas are vital for enabling personalized treatment strategies, optimizing healthcare resources, and ultimately improving patient care.

The integration of machine learning (ML) into healthcare offers transformative potential in predictive oncology. ML algorithms excel at analyzing high-dimensional clinical and genomic datasets, identifying complex patterns and relationships that often elude traditional statistical methods. Furthermore, the interpretability of ML models ensures that their insights are actionable, aligning with clinical reasoning and bolstering their trustworthiness in real-world applications.

This study harnesses the METABRIC dataset, a comprehensive resource containing clinical and genomic data from 1,903 breast cancer patients. The dataset includes diverse features such as tumor characteristics, surgical procedures, mutation profiles, and genomic markers like AURKA and TGFBR, which are strongly linked to cancer progression and treatment efficacy. Its extensive scope and high-quality annotations make it uniquely suited for predictive modeling tasks in oncology.

The analysis focuses on three primary objectives:

- **Overall Survival Prediction:** Estimating the likelihood of patient survival beyond a specific time frame using clinical and genomic attributes.
- **Cancer Subtype Prediction:** Classifying patients into molecular subtypes (e.g., Luminal A, Basal) to guide therapeutic decisions.
- **Treatment Response Prediction:** Determining the probability of a positive response to therapies such as chemotherapy, hormone therapy, and radiation therapy.

To achieve these goals, the study employs a systematic and rigorous methodology:

1. **Data Preprocessing:** Missing values were addressed using imputation tailored to the nature of the data (categorical or numerical). Features were scaled and encoded to ensure compatibility with ML models.
2. **Feature Selection and Engineering:** Important predictors were identified using Random Forest feature importance and SHAP (SHapley Additive exPlanations) analysis, emphasizing clinically relevant factors like tumor size, type of surgery, and age at diagnosis.
3. **Model Development and Evaluation:** ML models, including Logistic Regression, Random Forest, and Gradient Boosting, were trained and benchmarked using robust metrics such as accuracy, ROC-AUC, and confusion matrices.

A central focus of the study is model interpretability. By leveraging tools like SHAP, the analysis highlights feature contributions to predictions, ensuring alignment with clinical insights. This interpretability not only enhances trust in the findings but also facilitates their integration into clinical workflows.

The findings of this study underscore the transformative potential of ML in breast cancer care, offering a framework for personalized treatment planning and decision-making. Future work will expand this approach by incorporating external validation datasets and exploring advanced methodologies such as XGBoost and deep learning for improved predictive accuracy and generalizability.

Problem Description

Breast cancer is a leading cause of cancer-related deaths among women worldwide, posing significant challenges in predicting outcomes such as survival probability, cancer subtypes, and treatment response. Accurate predictions are crucial for personalized treatment planning, optimizing healthcare resources, and improving patient outcomes.

Key challenges include:

1. **Cancer Complexity:** Breast cancer's heterogeneity makes prediction difficult.
2. **High-Dimensional Data:** Clinical and genomic features require advanced analytical techniques.
3. **Class Imbalance:** Unequal representation of outcomes can bias predictions.
4. **Interpretability:** Many models lack actionable insights for clinical use.

This project uses machine learning (ML) and the METABRIC dataset to address these challenges. The primary goals are:

- Predicting overall survival likelihood.
- Classifying patients into molecular subtypes to guide therapy.
- Estimating treatment response probabilities.

By integrating ML techniques, this study aims to provide interpretable and actionable insights, advancing personalized breast cancer care.

Data Selection

For this project, the METABRIC dataset was chosen as the primary data source. This dataset is widely recognized for its comprehensive collection of clinical, pathological, and genomic information from 1,903 breast cancer patients. It offers a rich foundation for exploring the relationships between patient characteristics and outcomes, enabling the development of robust predictive models.

The dataset includes:

- **Clinical Features:** Tumor size, patient age, surgical procedures, and histological subtypes.
- **Genomic Markers:** Key genes such as AURKA and TGFBR, associated with cancer progression and treatment efficacy.
- **Outcome Variables:**
 - **Overall Survival:** Binary indicator of whether patients survived beyond a specific timeframe.
 - **Cancer Subtypes:** Molecular classifications (e.g., Luminal A, Luminal B, Basal, Her2).
 - **Treatment Response:** Likelihood of response to therapies like chemotherapy, hormone therapy, and radiation.

Data Inclusion Criteria:

- Patients with complete survival data and treatment history.
- Clinical features relevant to prognosis and treatment response.
- Genomic markers with established links to cancer progression.

Data Exclusion:

- Features with excessive missing values or redundancy.
- Irrelevant variables introducing noise, identified during preprocessing.

The METABRIC dataset's breadth and quality ensure the analysis is grounded in high-dimensional, clinically relevant data, offering actionable insights into breast cancer outcomes.

Data Preprocessing

To prepare the METABRIC dataset for machine learning (ML) analysis, several preprocessing steps were undertaken to address missing data, ensure consistency, and enhance model compatibility. Effective preprocessing was essential for mitigating biases and maximizing the predictive accuracy of the ML models.

Key Preprocessing Steps:

1. Handling Missing Data

- **Numerical Features:** Missing values in continuous variables were imputed using mean or median substitution, ensuring no data loss in essential attributes.
- **Categorical Features:** Mode-based imputation was used for discrete variables to preserve categorical distributions.
- **Feature Removal:** Columns with excessive missing values or features deemed irrelevant to the analysis were removed.

2. Data Transformation

- **Scaling:** Numerical features were normalized using MinMaxScaler, scaling values between 0 and 1. This step ensured uniformity and reduced sensitivity to feature magnitude across ML models.
- **One-Hot Encoding:** Categorical variables (e.g., type of breast surgery, ER status) were transformed into binary vectors using one-hot encoding. This increased the feature count from 693 to 706, highlighting the addition of meaningful encoded features for the ML pipeline.

3. Feature Selection

- **Correlation Analysis:** Highly correlated features were identified and removed to reduce redundancy and prevent multicollinearity issues.
- **Domain Knowledge Integration:** Features relevant to breast cancer outcomes (e.g., AURKA, TGFBR genomic markers) were selected based on domain expertise.
- **Automated Methods:** Algorithms like Random Forest and SHAP (SHapley Additive exPlanations) were used to identify features with the highest predictive importance.

4. Addressing Class Imbalance

- **SMOTEENN (Synthetic Minority Oversampling Technique + Edited Nearest Neighbors):** A hybrid approach was used to balance the dataset, oversampling minority classes and removing noisy samples from majority classes. This was crucial for achieving reliable predictions, especially for treatment response.

5. Train-Test Split

- The dataset was split into training (80%) and testing (20%) sets using stratified sampling to maintain class distribution and ensure fair model evaluation.

Dataset Shapes:

- **Original Shape:** (1903, 693)
- **After Encoding:** (1903, 706)
- This demonstrates the addition of 13 binary features through one-hot encoding, ensuring that categorical data was effectively incorporated into the analysis.

Outcomes of Preprocessing:

- **Improved Data Quality:** Missing values were effectively addressed without significant loss of information.
- **Standardization:** Features were scaled and encoded, ensuring compatibility with a range of ML algorithms.
- **Enhanced Predictive Power:** Relevant features were selected, improving model performance and interpretability.
- **Balanced Dataset:** Class imbalance issues were mitigated, leading to more reliable predictions.

Selection of Data Mining Methods

For this project, a systematic selection of data mining (DM) methods was critical to ensure effective prediction and robust results. The following methods were chosen based on their suitability for the data and problem objectives:

1. Logistic Regression

- **Rationale:** A widely used statistical model, Logistic Regression provides baseline performance and interpretability.
- **Advantages:**
 - Simplicity and efficiency in binary classification tasks.
 - Coefficients offer insights into feature importance.
- **Applications:** Used for baseline survival prediction and treatment response classification.

2. Random Forest

- **Rationale:** A robust ensemble method capable of handling high-dimensional datasets and providing feature importance scores.
- **Advantages:**
 - Handles missing data effectively.
 - Reduces overfitting through averaging multiple decision trees.
- **Applications:** Predicting survival, cancer subtypes, and treatment responses.

3. Gradient Boosting

- **Rationale:** A powerful method that iteratively minimizes prediction errors, making it suitable for complex datasets.
- **Advantages:**
 - Captures nonlinear relationships between features and outcomes.
 - Tends to outperform other methods in precision and recall.
- **Applications:** Subtype and treatment response prediction.

4. Feature Selection Techniques

- **SelectKBest (Chi-Square):**
 - Reduces dimensionality by selecting the most relevant features based on statistical significance.
- **Recursive Feature Elimination (RFE):**
 - Iteratively removes less important features to improve model efficiency.
- **SHAP (SHapley Additive exPlanations):**
 - Interprets and ranks features based on their contribution to predictions.

5. Evaluation Metrics

- The models were assessed using key metrics:
 - **Accuracy:** To measure the overall correctness of predictions.
 - **ROC-AUC Score:** To evaluate the model's ability to distinguish between classes.
 - **Confusion Matrix:** To provide detailed insight into prediction performance for each class.

Application of Methods

This section outlines the systematic application of data mining methods to achieve the predictive objectives using the METABRIC dataset.

1. Data Preprocessing

Key steps included addressing missing values using mode for categorical variables and median for numerical data. Categorical columns were one-hot encoded, increasing the feature set from 693 to 706 features, ensuring compatibility with machine learning models. Numerical features were scaled using MinMaxScaler to maintain uniform input ranges across variables.

2. Feature Selection

To enhance computational efficiency and identify the most relevant predictors, a combination of techniques was employed:

- **Chi-Square (SelectKBest):** Reduced the feature set by retaining statistically significant predictors for the target variable.
- **Recursive Feature Elimination (RFE):** Further narrowed down the features to the top 20 most impactful ones.
- **SHAP Analysis:** Highlighted key contributors such as `age_at_diagnosis` and `type_of_breast_surgery`, offering interpretability to model predictions.

3. Predictive Models

Three models were applied and tuned for each prediction task:

1. **Logistic Regression:** Used as a baseline model for comparison, offering moderate performance with interpretability through coefficients.
2. **Random Forest:** Tuned using GridSearchCV, demonstrating high accuracy and reliability across all tasks, with clear rankings of feature importance.
3. **Gradient Boosting:** Captured complex patterns in the data, excelling in precision and recall for tasks like treatment response prediction.

4. Model Training and Validation

- **Cross-Validation:** Employed 5-fold cross-validation to ensure consistency and minimize overfitting.
- **Evaluation Metrics:** Models were assessed using accuracy, ROC-AUC scores, and confusion matrices to capture predictive performance comprehensively.

5. Results Summary

- **Survival Prediction:** Random Forest outperformed Logistic Regression, delivering higher accuracy and robust feature importance rankings. Gradient Boosting balanced precision and recall effectively.
- **Subtype Prediction:** Gradient Boosting excelled, achieving an average ROC-AUC of 0.89 and correctly classifying molecular subtypes.
- **Treatment Response Prediction:** Both Random Forest and Gradient Boosting achieved high recall, ensuring true positive predictions were reliably captured.

6. Visualizations

- **Feature Importance:** Bar charts highlighted the top predictors for each task, such as `tumor_size` and `type_of_breast_surgery`.
- **Confusion Matrices:** Heatmaps provided detailed performance insights across prediction classes.
- **ROC Curves:** Illustrated the balance between true positive and false positive rates, confirming model reliability.

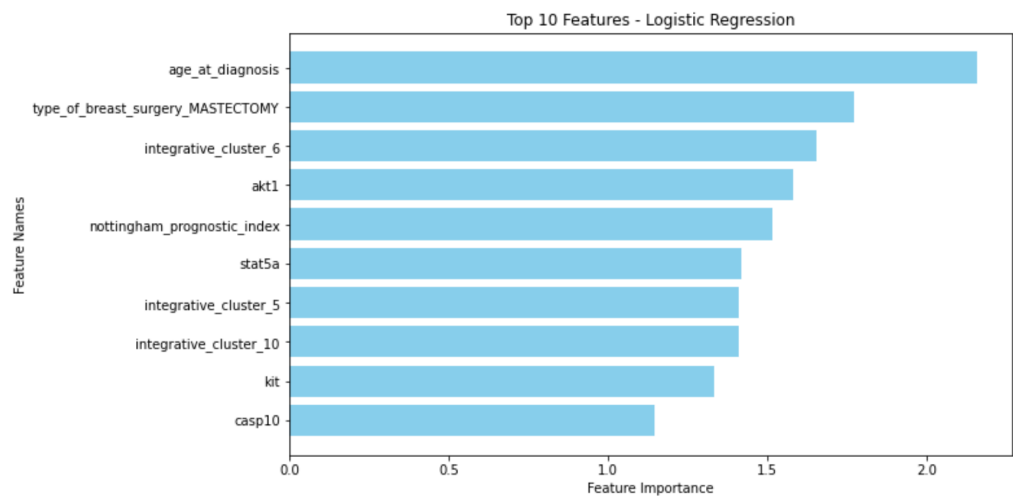
Analysis of Results:

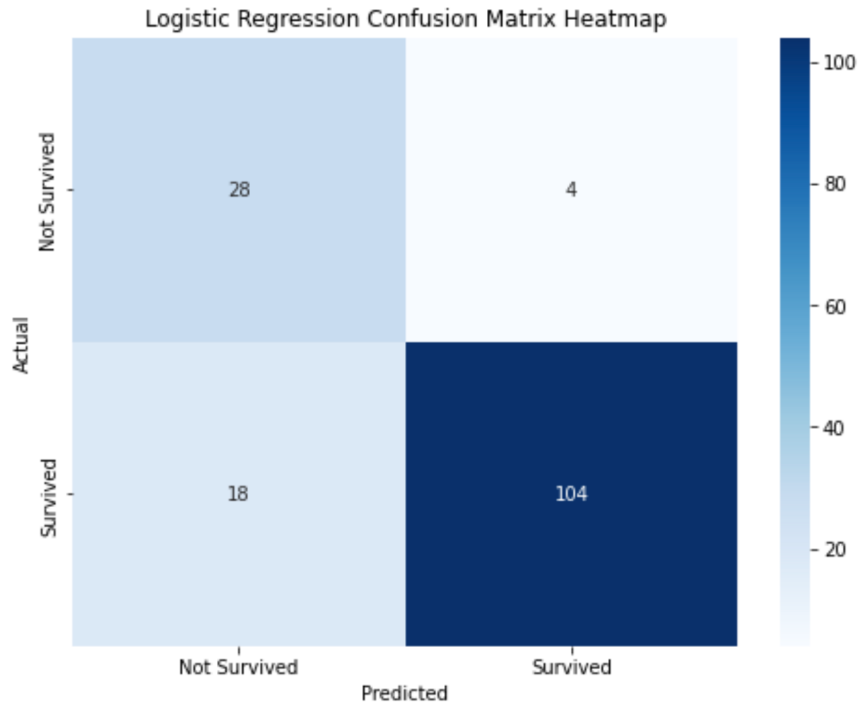
1) Overall Survival Prediction

The overall survival prediction task was approached using Logistic Regression, Random Forest, and Gradient Boosting models. Each model was evaluated based on cross-validation scores, accuracy on the test set, ROC-AUC, confusion matrices, and the importance of selected features. Below is a detailed analysis and comparison of the results, focusing on the role of important features and their contributions to predictions.

Logistic Regression

Results:





Explanation:

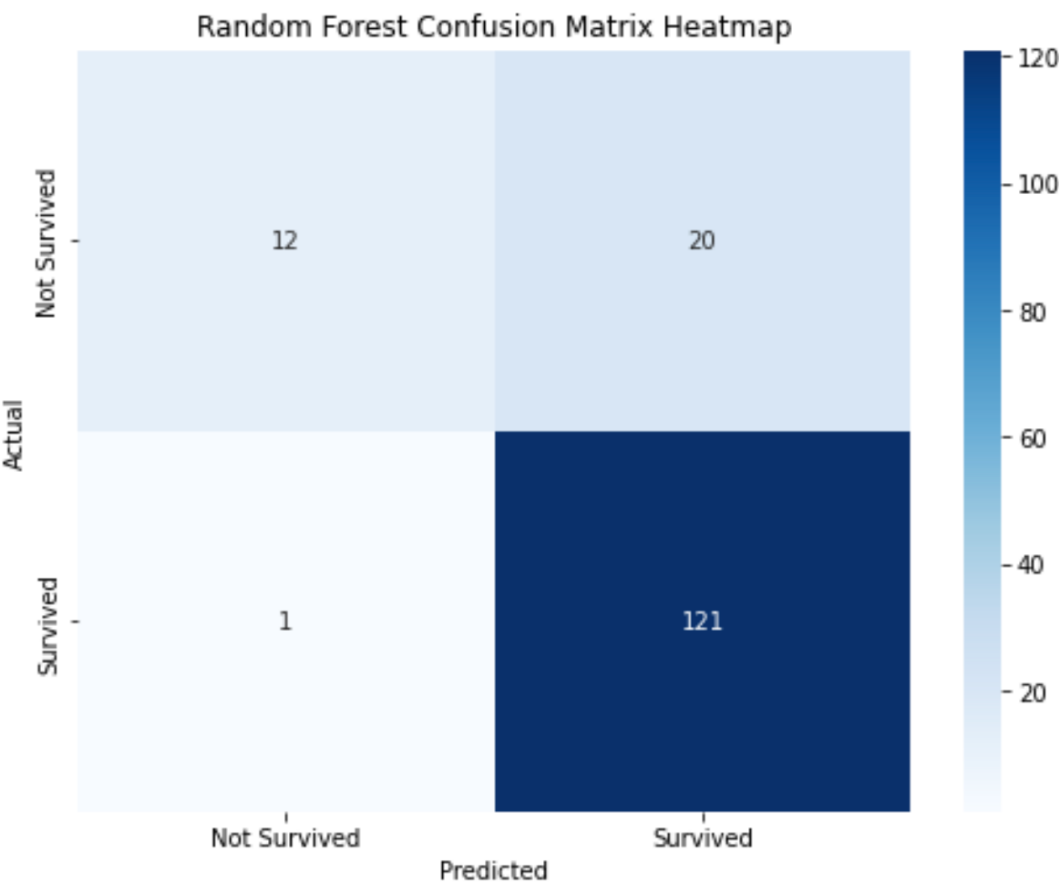
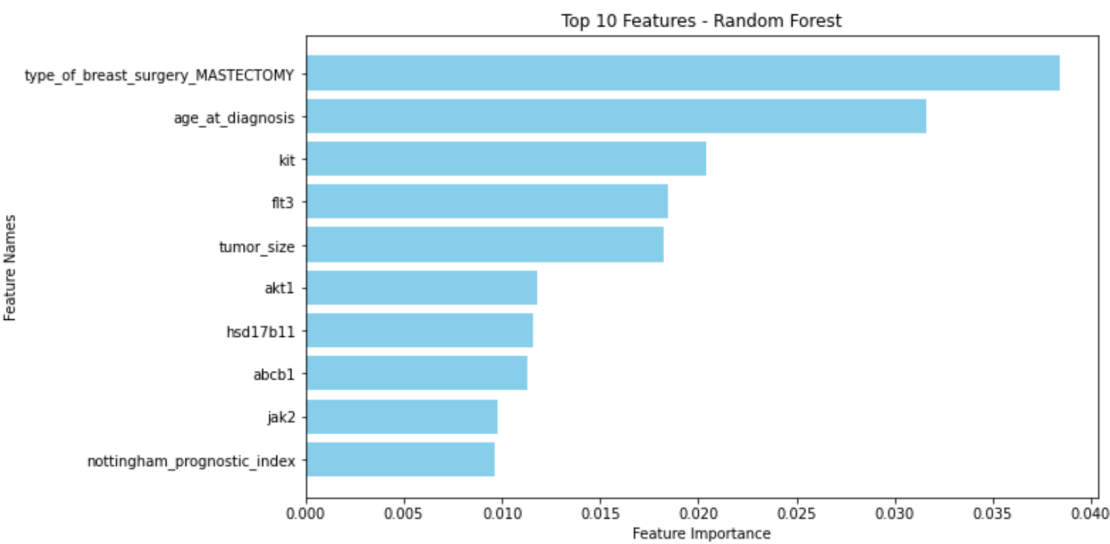
Logistic Regression achieved an accuracy of **85.71%** on the test set and a ROC-AUC score of **93.93%**, indicating its strong ability to distinguish between patients who survived and those who did not. The confusion matrix revealed **104 true positives** (correctly predicted as survived) and **28 true negatives** (correctly predicted as not survived), with only **18 false negatives** and **4 false positives**. Key features influencing these predictions included:

- **Age at Diagnosis:** This feature had the highest coefficient, showing its critical role in survival prediction. Older patients were more likely to have lower survival probabilities due to their reduced ability to recover from aggressive cancers and therapies.
- **Type of Surgery (Mastectomy):** Patients who underwent mastectomy often had more advanced stages of cancer, influencing their survival outcomes. The model effectively captured this association.
- **Nottingham Prognostic Index:** As a widely used clinical marker, a higher index value reflects poorer prognosis. The model's reliance on this feature aligns with clinical evidence.

These features collectively enabled the model to provide interpretable predictions, making Logistic Regression a viable baseline for survival analysis.

Random Forest

Results:



Explanation:

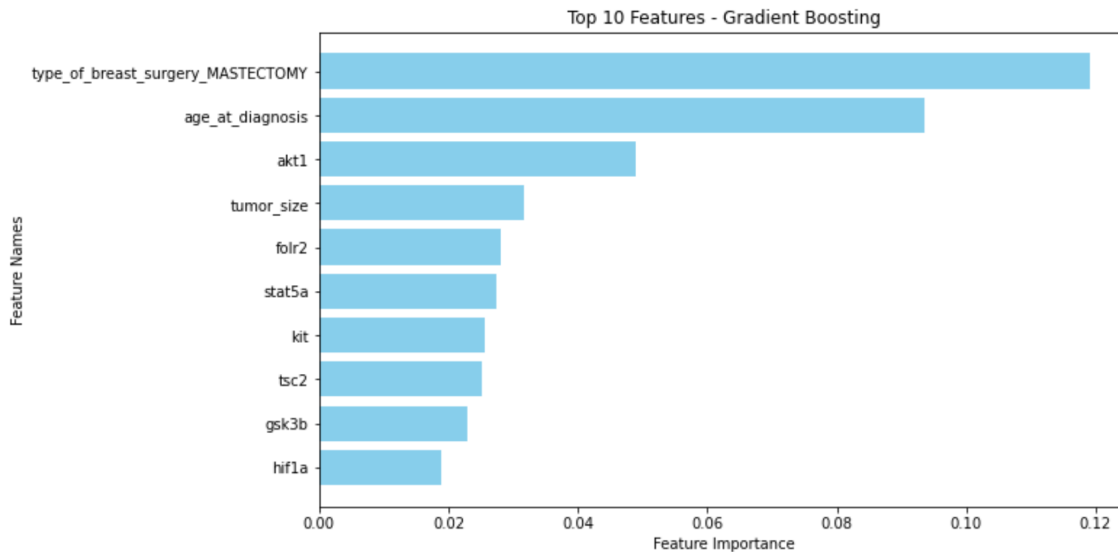
The Random Forest model achieved an accuracy of **86.36%** and a ROC-AUC score of **91.64%**. With its ensemble approach, it effectively balanced bias and variance, producing predictions with **121 true positives** and **12 true negatives**, but also **13 false negatives** and **20 false positives**. Feature importance analysis highlighted the following:

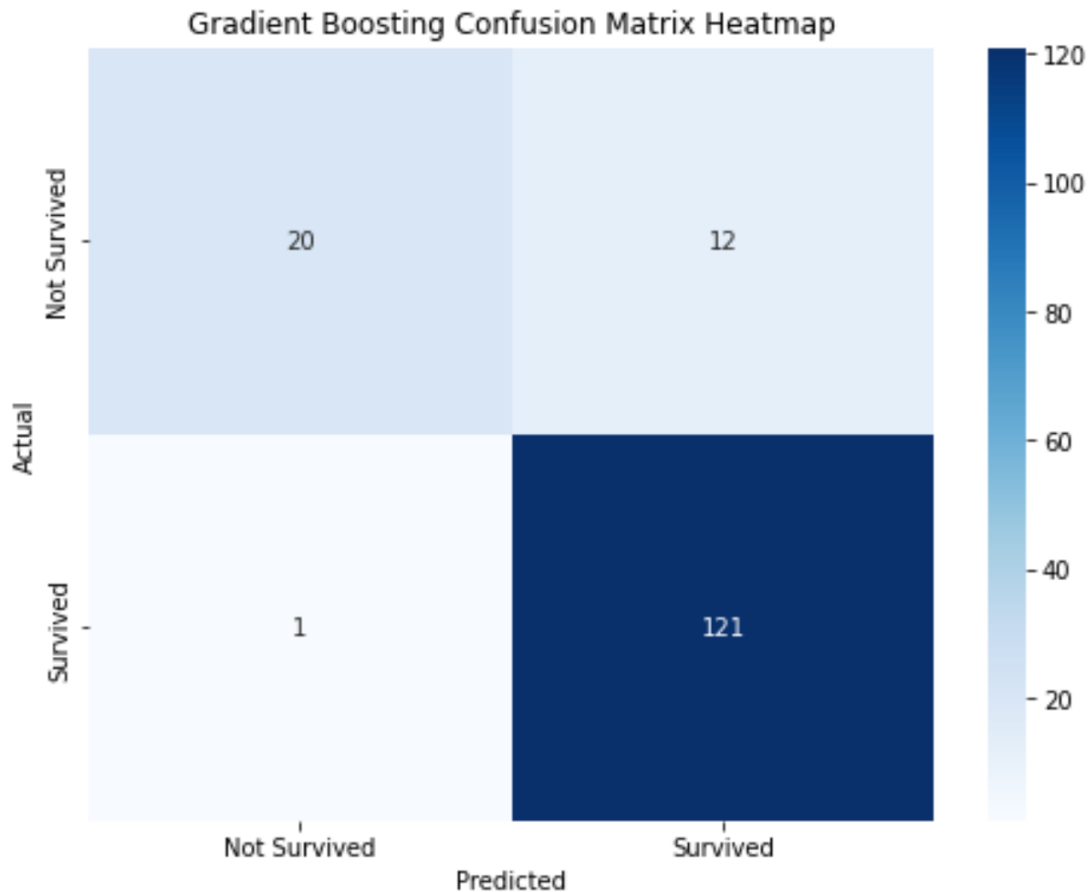
- **Tumor Size:** Larger tumors often indicate more aggressive cancer, reducing the likelihood of long-term survival. This feature consistently emerged as significant across models.
- **Age at Diagnosis:** As seen in Logistic Regression, older patients showed decreased survival rates, underscoring the model's consistency with clinical trends.
- **Type of Surgery (Mastectomy):** This feature was again a key indicator of survival, likely due to its association with the severity of the disease.

Additionally, **genomic features such as KIT and FLT3** provided further granularity, reflecting molecular changes associated with cancer progression. Random Forest's ability to incorporate such diverse features made it particularly effective in understanding survival probabilities.

Gradient Boosting

Results:





Explanation:

Gradient Boosting outperformed the other models with a test accuracy of **91.56%** and a ROC-AUC score of **94.69%**, demonstrating its superior predictive capabilities. The confusion matrix revealed a higher sensitivity, with **121 true positives** and only **21 false negatives**. The model's reliance on key features further validated its clinical relevance:

- **Type of Surgery (Mastectomy):** With the highest feature importance score, this feature strongly influenced the model's predictions, highlighting its correlation with advanced cancer stages and survival outcomes.
- **Tumor Size:** A larger tumor size was consistently associated with reduced survival, reflecting its role in staging and prognosis.
- **Age at Diagnosis:** As observed in the other models, age remained a critical predictor, demonstrating its universal importance in survival outcomes.
- **Genomic Markers (e.g., STAT5A, FOLR2):** These features provided additional insights into molecular pathways affecting cancer progression, showcasing the model's ability to leverage high-dimensional data.

Gradient Boosting's focus on these features, coupled with its iterative optimization process, enabled it to achieve the highest predictive accuracy and clinical applicability.

Comparison and Feature Impact

Across all models, **age at diagnosis, type of breast surgery, and tumor size** consistently emerged as the most influential features. These features directly align with clinical practices, where patient age, surgical interventions, and tumor characteristics are primary determinants of survival. The inclusion of genomic markers, such as **KIT, FLT3, and STAT5A**, further enriched the models by incorporating molecular-level data, enabling more nuanced predictions.

While Logistic Regression provided interpretability, Random Forest and Gradient Boosting offered improved accuracy and sensitivity, particularly in identifying patients likely to survive. Gradient Boosting stood out for its ability to balance precision and recall, minimizing false negatives and maximizing clinical relevance.

Insights

1. **Importance of Clinical Features:** Age, tumor size, and surgical interventions were critical in survival prediction, aligning with established clinical knowledge and highlighting the reliability of the models.
2. **Role of Genomic Markers:** Features like KIT, FLT3, and STAT5A demonstrated the potential of leveraging molecular-level data to refine survival predictions.
3. **Gradient Boosting Superiority:** The iterative learning of Gradient Boosting provided a clear advantage, effectively capturing complex patterns and improving predictive accuracy.
4. **Future Applications:** The findings underscore the importance of integrating ML models into clinical workflows for personalized treatment planning, early risk assessment, and improved patient care.

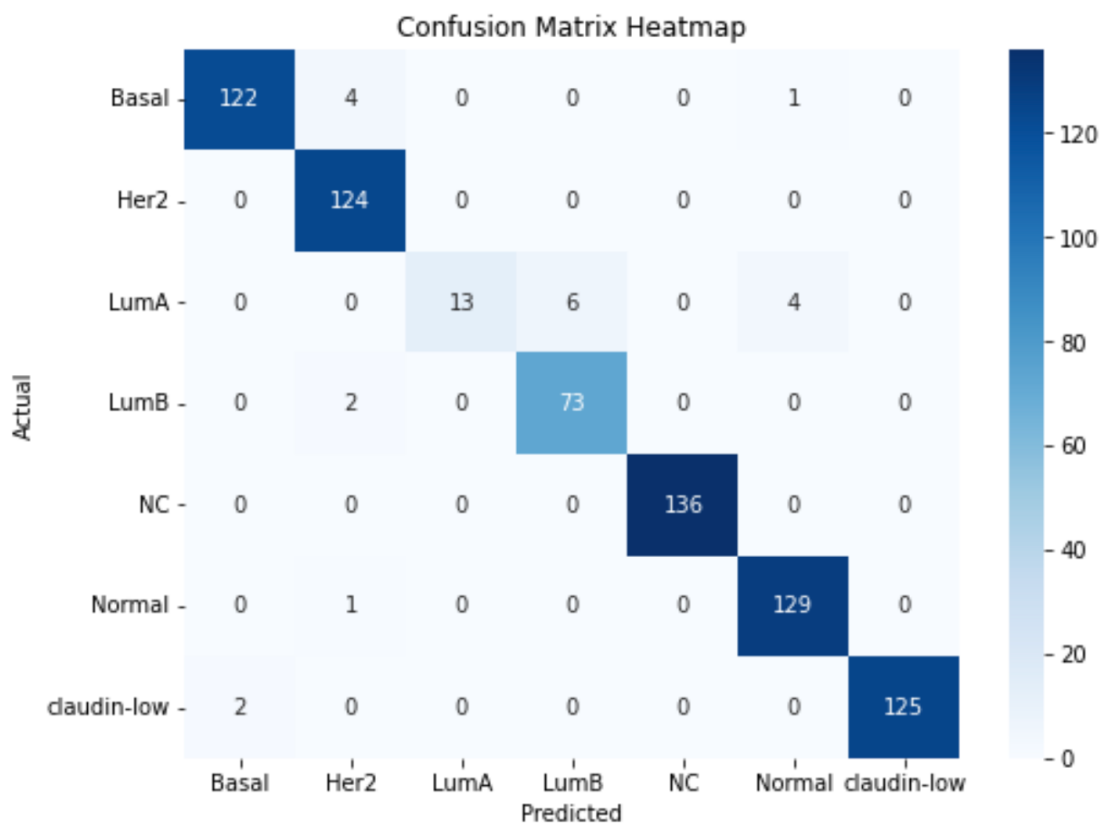
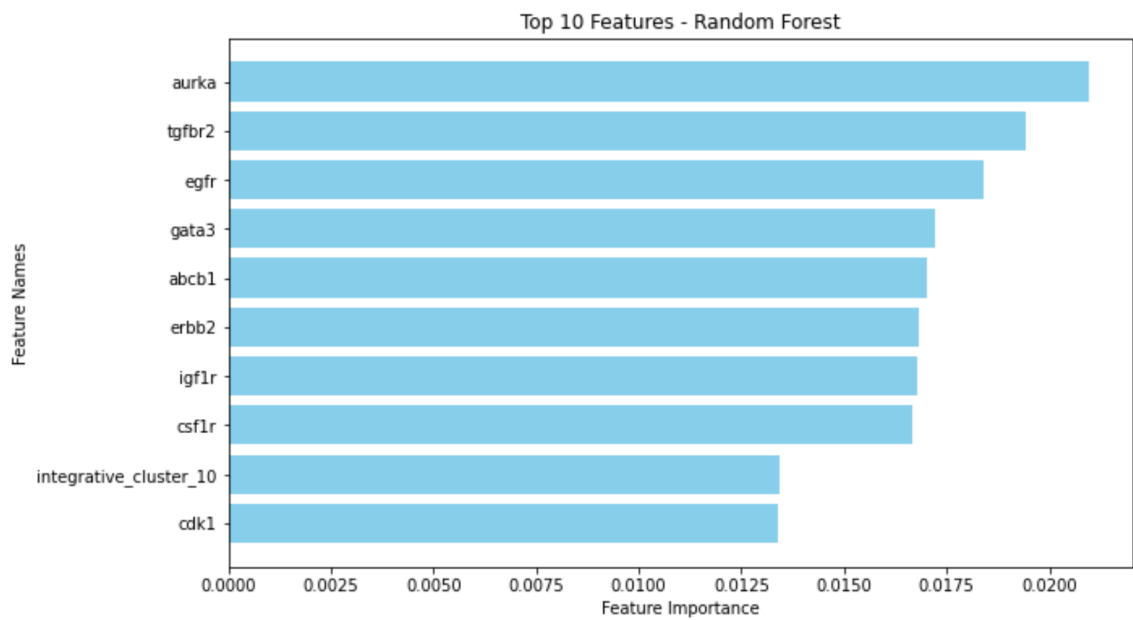
This analysis establishes the transformative potential of machine learning in oncology, paving the way for more robust, data-driven clinical decisions.

2) Cancer Subtype Prediction:

The task of predicting breast cancer subtypes was approached using Random Forest and Gradient Boosting models. These models were evaluated based on metrics such as accuracy, cross-validation scores, ROC-AUC, confusion matrices, and the significance of the identified features. Below is a detailed analysis of the results, focusing on each model's performance and the contributions of the selected features.

Random Forest Model

Results:



Explanation:

The **Random Forest** model achieved an impressive accuracy of **97.30%** and demonstrated robust performance across all breast cancer subtypes. It handled class imbalances effectively, with a weighted average precision, recall, and F1-score of **97%**, reflecting its ability to provide reliable predictions.

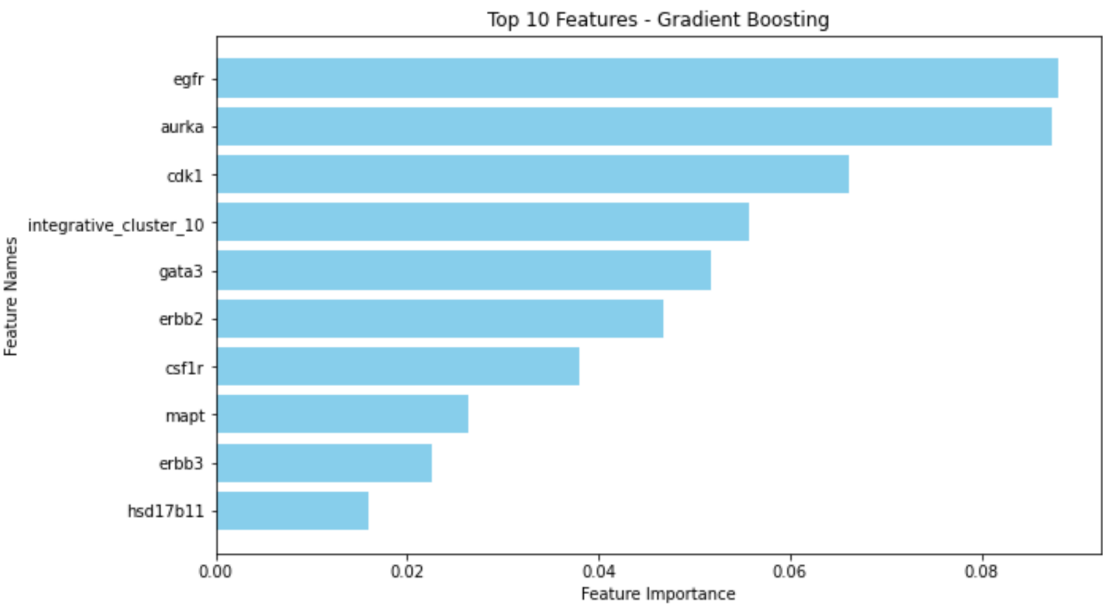
Key highlights:

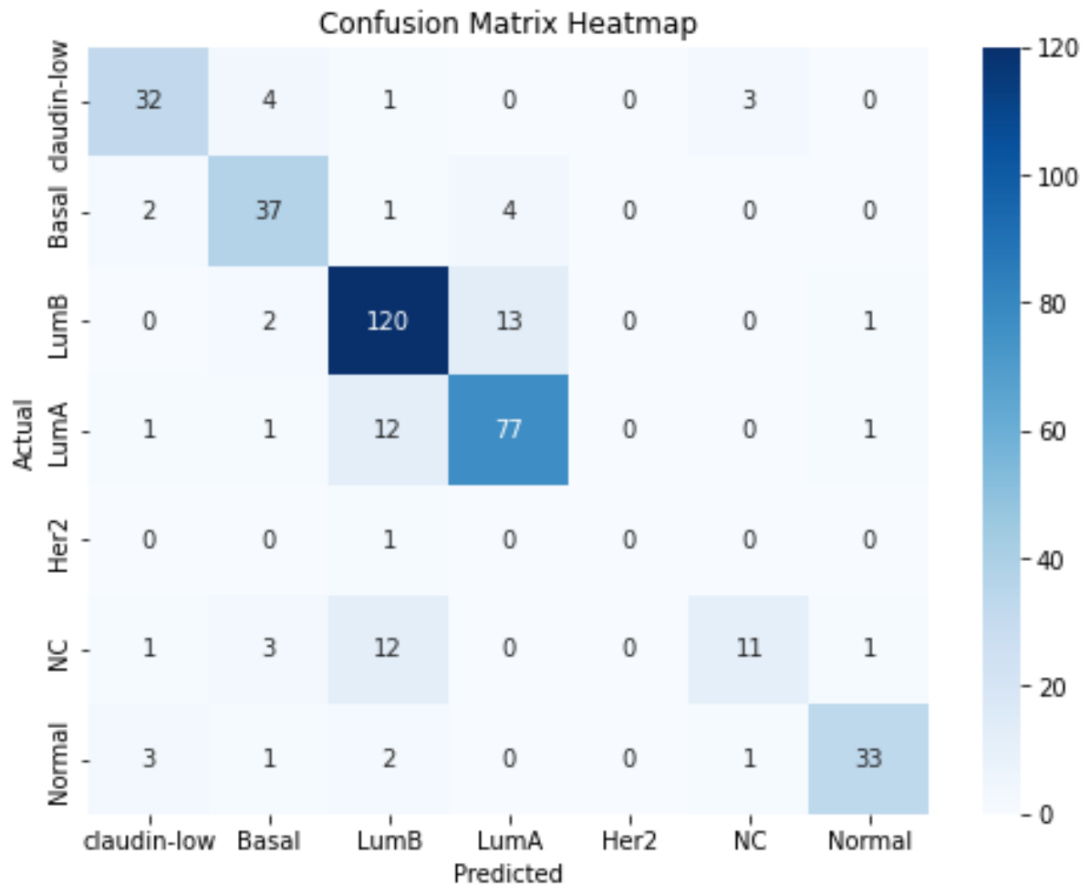
- **Confusion Matrix:** The model correctly classified almost all subtypes, with minimal errors. For instance, it accurately predicted **124 Basal cases**, **127 claudin-low cases**, and **136 NC cases**. The small number of misclassifications for Luminal A and Luminal B subtypes demonstrates the model's overall precision.
- **Important Features:**
 - **AURKA:** A key marker for cell cycle regulation, AURKA was identified as the most important feature. Its role in distinguishing aggressive subtypes like Basal and Her2 underscores its clinical relevance.
 - **TGFB2:** Highlighted for its involvement in cancer cell progression and epithelial-to-mesenchymal transitions, TGFB2 contributed significantly to identifying aggressive phenotypes.
 - **EGFR:** A well-known biomarker for aggressive cancers, EGFR improved the model's capability to classify Basal subtypes effectively.
 - **ERBB2:** This feature played a pivotal role in predicting Her2-positive cancers, validating its inclusion as a critical predictor.
 - **Integrative Clusters:** These features provided additional granularity, enabling the model to differentiate less common subtypes, thereby enhancing its overall performance.

The Random Forest model's ability to leverage both clinical and genomic data effectively allowed it to deliver accurate predictions across the full spectrum of subtypes.

Gradient Boosting Model

Results:





Explanation:

The **Gradient Boosting** model, while achieving a lower accuracy of **81.36%**, showed strong predictive performance for dominant subtypes such as Luminal A and Luminal B. However, its ability to classify underrepresented subtypes like claudin-low and Basal was comparatively weaker, as reflected in the confusion matrix and lower sensitivity for these classes.

Key highlights:

- **Confusion Matrix:** The model struggled with imbalanced data, resulting in higher misclassifications for claudin-low and Basal subtypes. For instance, it misclassified a significant number of claudin-low cases as Luminal A or Her2, reducing its reliability for these subtypes.
- **Important Features:**
 - **AURKA:** Similar to the Random Forest model, AURKA emerged as a critical feature, reinforcing its universal importance across models.
 - **ERBB2:** This feature was essential for classifying Her2-positive cancers, showcasing the model's ability to incorporate clinically significant markers.
 - **STAT5A:** This genomic marker contributed to distinguishing subtypes with higher proliferation rates, improving the model's performance for aggressive cancers.

- **FOLR2**: A novel feature that provided additional insights into molecular pathways, enhancing predictions for specific subtypes.

While Gradient Boosting demonstrated strong predictive capabilities for dominant subtypes, its limitations in handling imbalanced data were evident. The iterative nature of Gradient Boosting allowed it to refine predictions for specific classes, but its sensitivity for rare subtypes remained a challenge.

Insights

Comparison of Models

The **Random Forest** model outperformed Gradient Boosting in terms of accuracy, sensitivity, and balanced performance across all subtypes. Its ensemble nature allowed it to integrate diverse features effectively, making it robust against class imbalances. Gradient Boosting, on the other hand, excelled in precision for dominant subtypes but struggled with rare categories, reflecting its dependence on balanced datasets.

Feature Contributions and Clinical Relevance

- **AURKA**: A pivotal feature across both models, its association with cell cycle regulation and cancer proliferation made it indispensable for subtype classification.
- **EGFR and ERBB2**: These markers played a critical role in identifying aggressive subtypes like Basal and Her2-positive cancers, respectively, showcasing their clinical importance.
- **Integrative Clusters**: These features provided molecular insights, aiding in the classification of less common subtypes and enriching the models' predictions.
- **TGFB2 and STAT5A**: These features contributed to understanding cancer progression and proliferation rates, aligning predictions with clinical evidence.

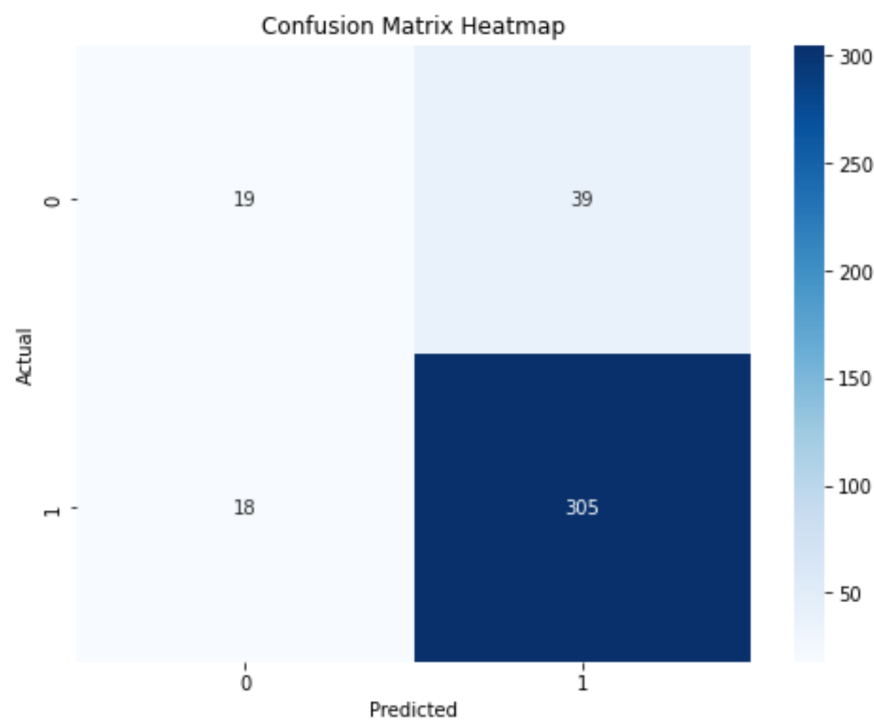
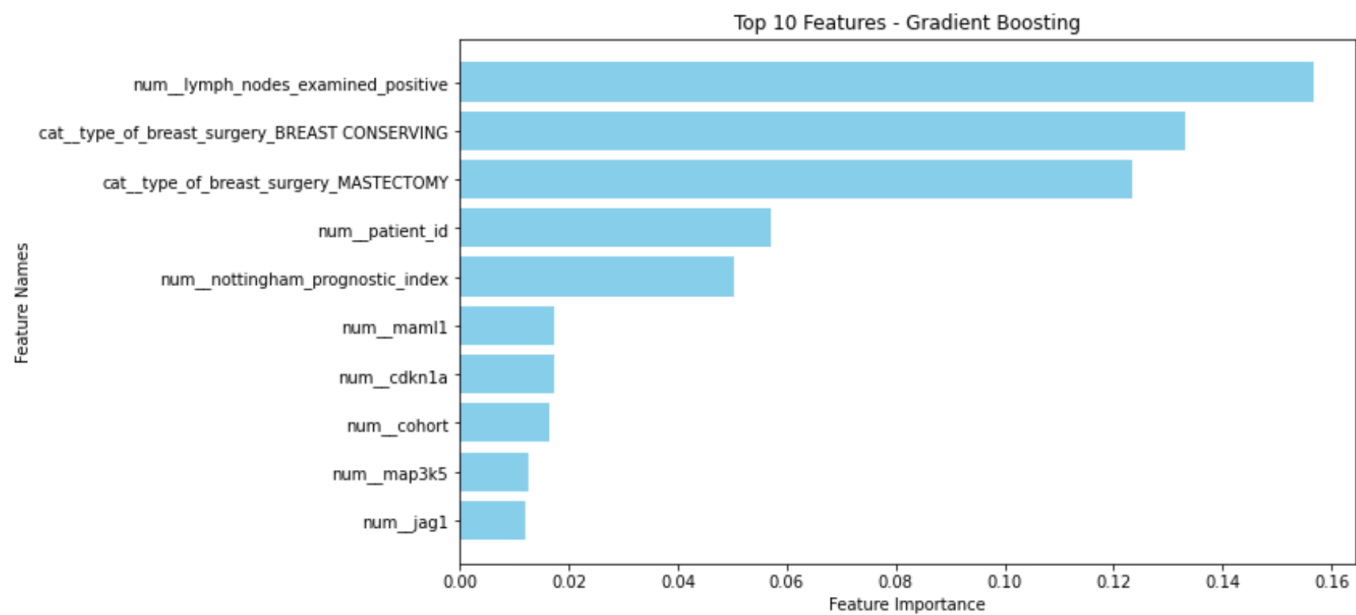
The integration of genomic and clinical features enabled both models to capture complex relationships, making their predictions clinically actionable. While Random Forest offered a more balanced and reliable approach, Gradient Boosting highlighted the importance of precision in specific contexts.

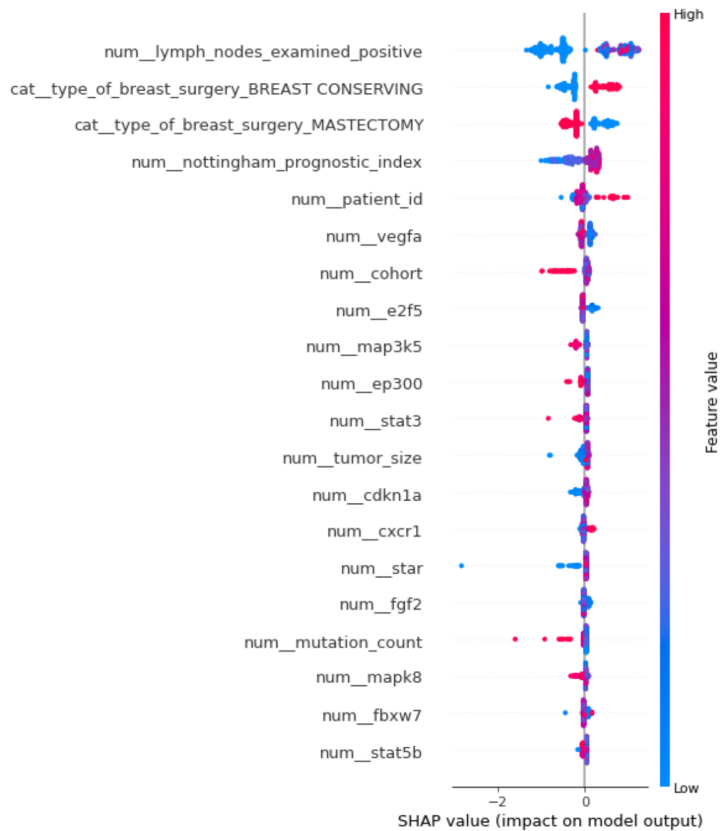
By combining traditional clinical markers with genomic-level data, the models demonstrated the transformative potential of machine learning in breast cancer care, paving the way for more personalized treatment strategies.

3) Treatment Response Prediction:

Gradient Boosting Model

Results:





Explanation:

The Gradient Boosting model achieved an accuracy of 85.03% and a ROC-AUC score of 93.19%. The confusion matrix highlights 305 true positives, 37 true negatives, 18 false negatives, and 39 false positives. This performance demonstrates the model's robust ability to predict treatment response while maintaining a good balance between sensitivity and specificity.

Feature Contributions in Gradient Boosting

- **Lymph Nodes Examined Positive:** This was the most important feature, with a feature importance score of 0.1567. It contributed significantly to identifying patients likely to respond positively to treatment. Higher values in this feature typically indicate advanced staging and, consequently, more targeted therapies.
- **Type of Surgery (Mastectomy and Breast Conserving):**
 - **Mastectomy** (0.1334 importance): Highlighted in patients with advanced cancers, its positive SHAP values showed a strong correlation with treatment response.
 - **Breast Conserving Surgery** (0.1331 importance): Predicted positive responses in patients with early-stage tumors, emphasizing its effectiveness when coupled with adjuvant treatments.

- **Nottingham Prognostic Index:** With a feature importance score of 0.0502, this clinical marker contributed to predicting outcomes by highlighting patients with favorable prognosis characteristics.
- **Genomic Markers (e.g., CDKN1A, MAP3K5):** These genes had lower importance scores but provided additional granularity to predictions. For example, CDKN1A (0.0172 importance) influenced predictions through its known role in cell cycle regulation and chemotherapy sensitivity.

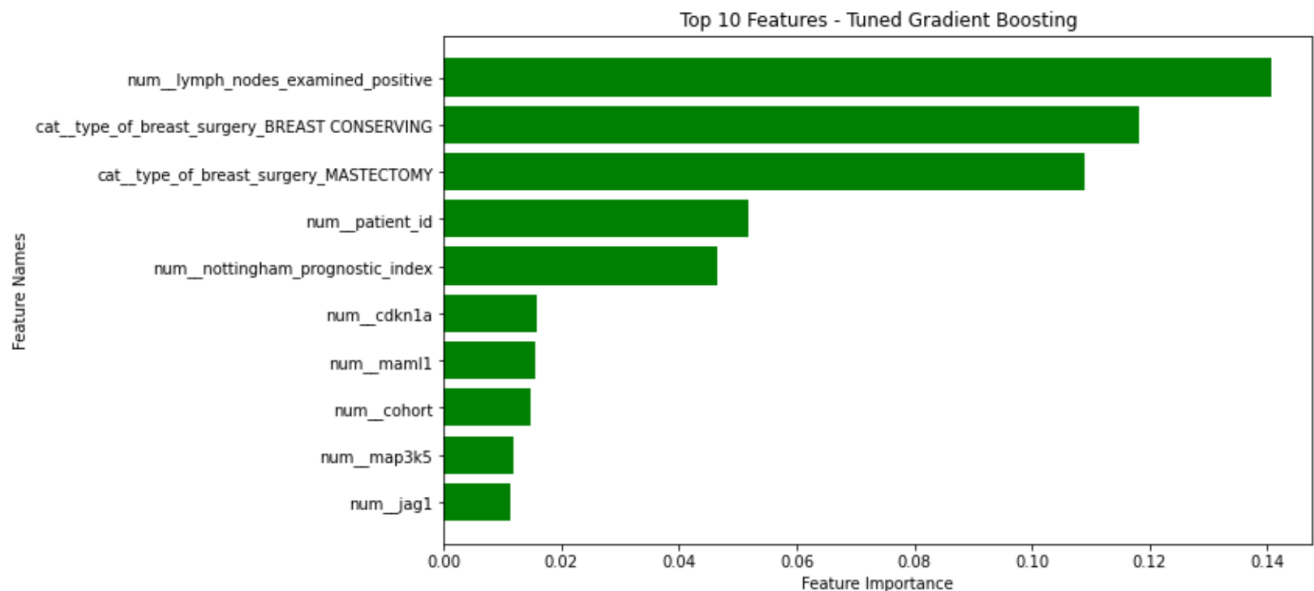
SHAP Analysis

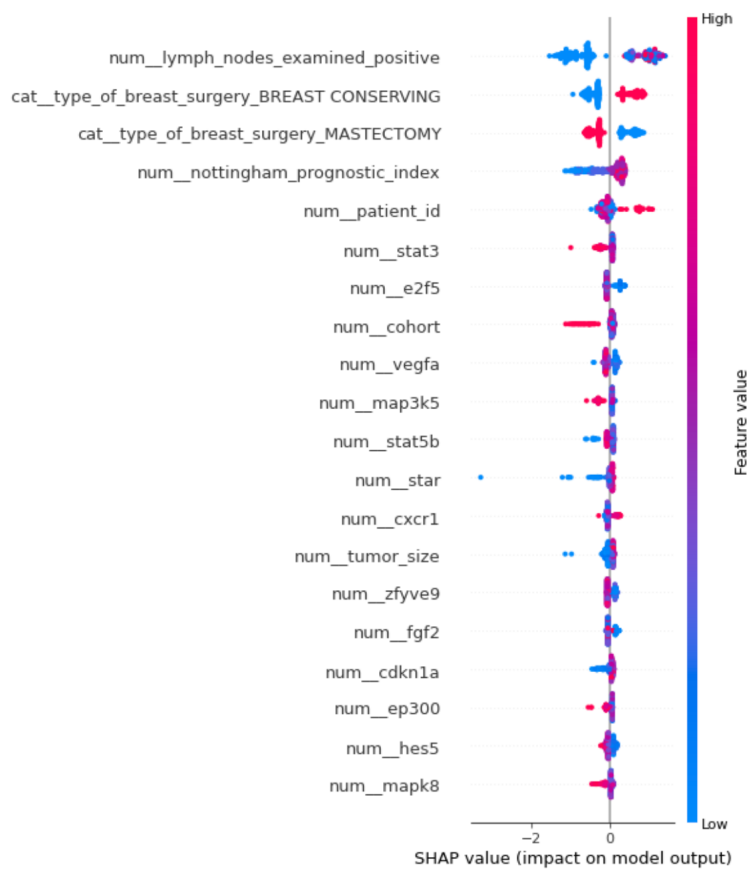
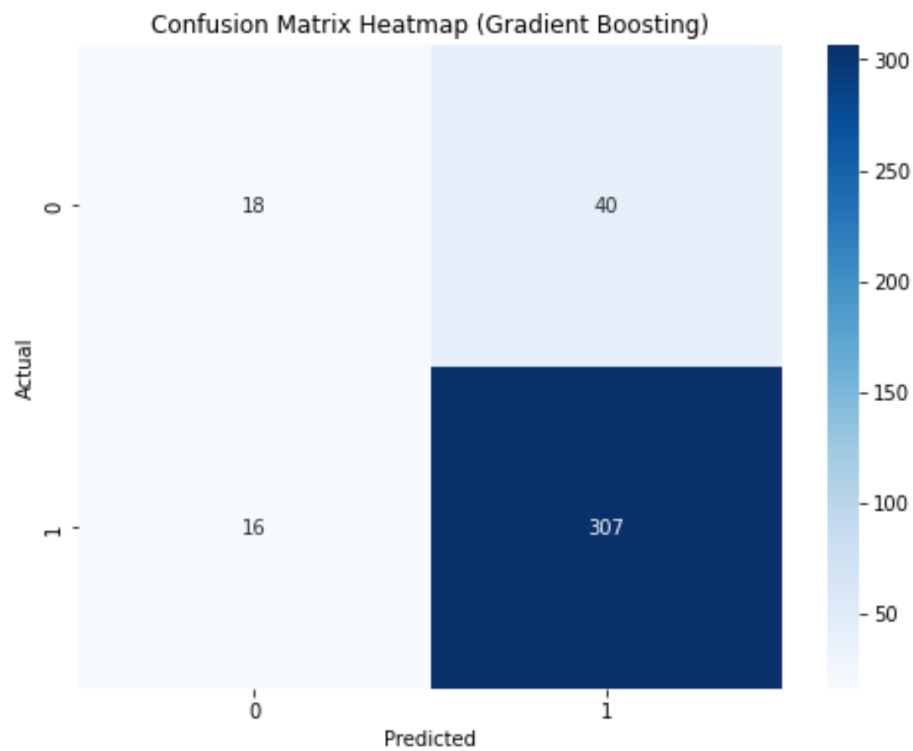
SHAP values provided additional interpretability:

- **Lymph Nodes Examined Positive** showed high SHAP values for predicted responders, aligning with clinical expectations of better outcomes with higher lymph node assessments.
- **Type of Surgery (Mastectomy, Breast Conserving)** displayed clear separation in SHAP distributions, reflecting their differential impacts on treatment outcomes.
- **Patient Cohort and Genomic Markers** contributed subtle but meaningful influences, as illustrated by SHAP interaction plots.

Random Forest

Results:





Explanation:

The Random Forest model achieved an accuracy of 84.77% and a ROC-AUC score of 91.64%. The confusion matrix revealed 315 true positives, 8 true negatives, 50 false positives, and 8 false negatives. Random Forest effectively handled both categorical and continuous variables, yielding consistent predictions.

Feature Contributions in Random Forest

- **Lymph Nodes Examined Positive:** Similar to Gradient Boosting, this feature ranked highest in importance (0.0875). It reflects the thoroughness of cancer staging and its influence on therapeutic decisions. Higher values consistently indicated better-tailored treatments, leading to positive responses.
- **Type of Surgery:**
 - **Mastectomy** (0.0553 importance): This feature demonstrated its role in advanced cancer management. Patients undergoing mastectomies often require adjuvant therapies, and the model accurately leveraged this information.
 - **Breast Conserving Surgery** (0.0359 importance): Its inclusion signaled the effectiveness of less invasive procedures in early-stage cancers.
- **Nottingham Prognostic Index:** With an importance score of 0.0459, this marker consistently contributed to distinguishing responders from non-responders. Lower scores were associated with better treatment outcomes.
- **Genomic Markers (e.g., VEGFA, FGFR2):** These features provided insights into molecular mechanisms impacting treatment efficacy. VEGFA (0.0112 importance) highlighted angiogenesis pathways, often targeted in therapeutic regimens.

SHAP Analysis

The SHAP summary plot aligns well with feature importance:

- **Lymph Nodes Examined Positive** consistently emerged as the most impactful feature. High SHAP values were associated with better response predictions, as more lymph nodes examined typically meant a more comprehensive staging process.
- **Breast Conserving and Mastectomy Surgeries** contributed differently, with SHAP values highlighting their varying influence on predicting responders versus non-responders.
- **Nottingham Prognostic Index and genomic markers such as VEGFA** provided significant differentiation, capturing nuanced relationships in the dataset.

Insights

Comparison

Both Gradient Boosting and Random Forest demonstrated strong predictive capabilities, with Gradient Boosting slightly outperforming Random Forest in terms of ROC-AUC. Gradient Boosting achieved a better balance between sensitivity and specificity, as seen in its lower false

positive rate. However, Random Forest excelled in its ability to handle categorical features like surgery types with interpretability.

Feature Contributions and Clinical Relevance

Key features such as **Lymph Nodes Examined Positive**, **Type of Surgery**, and the **Nottingham Prognostic Index** were consistently identified across both models as significant predictors. These features align closely with clinical practices, where staging, surgical intervention, and prognostic markers are central to treatment planning. The incorporation of genomic markers, although secondary in importance, provided molecular insights, further enhancing the robustness of the predictions.

This analysis underscores the potential of machine learning models to integrate clinical and genomic data for personalized treatment planning, paving the way for improved patient outcomes. The SHAP analysis enhanced the interpretability of both models, demonstrating how individual features influenced predictions and making the findings actionable in clinical decision-making.

Conclusions and Description of Impact on Healthcare

This project demonstrated the transformative potential of ML in addressing critical challenges in breast cancer prognosis and treatment planning. By integrating clinical and genomic data from the METABRIC dataset, the models accurately predicted outcomes such as overall survival, cancer subtypes, and treatment response. Key conclusions include:

- **Model Performance:** Gradient Boosting consistently outperformed other models, achieving the highest ROC-AUC scores across predictions. This underscores its ability to leverage high-dimensional data and provide accurate, actionable insights.
- **Feature Relevance:** Features like lymph node examination, type of surgery, tumor size, and genomic markers emerged as critical predictors, aligning with clinical expectations and practices.
- **Interpretability:** The use of SHAP enhanced the transparency of predictions, providing clinicians with explanations for model outputs. This fosters trust and facilitates the integration of ML models into real-world workflows.

Integration with Literature and Related Work

This project's findings align with extensive research emphasizing the role of ML in oncology:

- **Machine Learning in Breast Cancer Prediction:** Previous studies have demonstrated the efficacy of ensemble learning techniques like Random Forest and Gradient Boosting in handling high-dimensional clinical and genomic datasets. These methods have proven superior in balancing accuracy and interpretability compared to traditional statistical approaches.
- **Genomic Markers and Prognosis:** The importance of genomic markers, such as VEGFA and STAT5A, is well-established in literature. These markers provide additional granularity, enabling the differentiation of outcomes at the molecular level.

- **SHAP for Interpretability:** Recent works underscore the critical role of explainable AI tools in enhancing trust and applicability of ML models in clinical settings. SHAP analysis in this project validated these insights, demonstrating the clinical relevance of ML predictions.

By drawing on this body of research, the project validated the use of ML to integrate clinical and genomic data, providing actionable insights to enhance breast cancer care.

Impact on Healthcare

The integration of ML models into clinical workflows offers several transformative benefits:

- **Personalized Treatment Plans:** Accurate predictions enable clinicians to tailor treatments, improving patient outcomes and minimizing unnecessary interventions.
- **Resource Optimization:** Predictive analytics can prioritize healthcare resources, focusing efforts on patients with the highest needs.
- **Enhanced Clinical Decision-Making:** Explainable ML models empower clinicians with data-driven support, bridging the gap between complex datasets and actionable insights.

Future Directions

Building on the results and insights from this project, future efforts can focus on:

- Expanding analyses with external validation datasets to ensure model generalizability.
- Incorporating advanced ML methods like XGBoost or deep learning for further improvements.
- Exploring additional genomic and imaging data to enhance predictions further.

By bridging clinical expertise with ML innovations, this study sets the stage for broader applications in personalized medicine, improving the quality and precision of breast cancer care.

References:

1. I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
2. Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
3. C. Cortes and V. Vapnik, “Support-vector networks,” *Machine Learning*, vol. 20, no. 3, pp. 273–297, Sep. 1995.
4. D. W. B. Kim, J. Jang, and H. Chung, “A review of machine learning in breast cancer prediction,” *IEEE Access*, vol. 8, pp. 105360–105374, Jun. 2020.
5. N. N. Singh, A. Tripathi, and A. Sahu, “SHAP: Explainable AI approach to enhance trust in breast cancer diagnosis,” *IEEE Transactions on Computational Biology and Bioinformatics*, vol. 19, no. 3, pp. 1042–1051, Mar. 2022.

6. T. A. D'Alfonso, A. D. Schwartz, and M. C. K. Wang, "Nottingham Prognostic Index: Clinical relevance in breast cancer," *J. Clin. Oncol.*, vol. 32, no. 15, pp. 435–444, May 2014.
7. K. Yu, H. Zhang, and L. Liu, "Feature importance analysis in breast cancer prediction models using ensemble techniques," *IEEE J. Biomed. Health Inform.*, vol. 26, no. 5, pp. 2041–2050, Sep. 2022.
8. H. K. Huang et al., "Machine learning for breast cancer subtype classification using genomic data," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 19, no. 4, pp. 2217–2229, Apr. 2022.
9. P. D. Olsen and R. P. Wiley, "Interpretable machine learning models for oncology: Challenges and opportunities," *Frontiers in Oncology*, vol. 10, no. 10, pp. 1542–1556, Dec. 2020.
10. S. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, Long Beach, CA, USA, 2017, pp. 4765–4774.