MA981 DISSERTATION

# Tracking the Factors leading to Breast Cancer Staging and Survival: An Exploratory Analysis

**Sourav Ghosh Hansda**

Supervisor: **Lausen, Berthold**

December 22, 2023

Colchester

# Contents

**Abstract**

The heterogeneity of breast cancer underscores the urgency for accurate cancer staging and overall survival prediction. In this backdrop, the study has explored the clinical and genomic aspects of breast cancer patients based on the landmark TCGA-GDC database. It observes that the age group around 60 years is most vulnerable to breast cancer. The disease is found more prevalent among the whites. In a rigorous framework of logistic regression, the study brings out the lymph nodes examined count as the only clinical variable, which is statistically significant in the final model. Applying Kaplan-Meier Estimator and Log Rank Test, it finds that the groupings of patients in terms of vital status: alive and dead, tumor staging, and her2 status (ie, aggressive tumor development) are meaningful analytical categories for survival analysis. Given the limited utility of clinical data for survival analysis, it considers gene expression data and tackles the computational challenge by applying differential expression analysis and selecting 60 genes based on p-value out of about 60,000 genes without compromising the performance efficiency. It finds that the optimal number of clusters out of 60 genes may be two in consonance with the number of cancer staging being two: invasive and non-invasive, used in the study as a dependent variable. In view of the limitations of logistic regression in explaining cancer staging, the study carries out supervised learning, namely, random forest and neural network on the dataset. It finds that the random forest model works well with the dataset for explaining both cancer staging and overall survival. The results are found robust against resampling and random sampling of genes in place of differential expression analysis. Given the difficulty in interpreting machine learning-based results, the study suggests that future research may be directed at unraveling the black boxes like machine learning models with a view to improving their utility in the medical domain.

# Introduction

Breast cancer, the most prevalent malignancy among women worldwide, poses a significant public health concern. Accounting for around 25% of all female cancer diagnoses and 15% of all cancer-related deaths, breast cancer exhibits a heterogeneous nature, characterized by diverse clinical presentations, treatment responses, and prognoses. This heterogeneity underscores the urgency for accurate staging and overall survival prediction. Accurate prognostication enables tailored treatment strategies, optimizing patient outcomes and improving survival rates. The Cancer Genome Atlas (TCGA) project, a landmark initiative of the National Cancer Institute (NCI) and the National Human Genome Research Institute (NHGRI), provides a comprehensive resource for studying cancer at the molecular level. The TCGA database comprises clinical data, including patient demographics, tumor characteristics, and treatment information, as well as genomic data, such as gene expression, DNA copy number, and DNA methylation profiles. This wealth of data offers invaluable insights into the underlying biology of breast cancer and its response to treatment. The TCGA database is hosted at the Genomic Data Commons Data Portal (GDC) and henceforth is referred as TCGA-GDC database. This study aims to leverage the TCGA-GDC database to identify significant predictors of breast cancer staging and overall survival. By employing a battery of techniques in supervised and unsupervised learning, we will explore the complex interplay between clinical and genomic factors that contribute to breast cancer heterogeneity. The rest of the study is organized as follows. Chapter II delineates the conceptual underpinnings of cancer staging and overall survival (OS). Chapter III presents a brief survey of literature,

bringing out the learnings and guideposts relevant for the exercise at hand. Chapter IV elucidates the stylized facts and figures on the predictors and response variables based on the TCGA-GDC database. Chapter V discusses different dimensions of the battery of statistical methodologies employed to decode and decipher the big data set. Chapter VI applies the statistical techniques on the dataset and presents and interprets the results. It compares and, contrasts the results in terms of a comparative study of the methodologies employed on the dataset. The scope for refinement of the study in view of its limitations is explored in Chapter VII towards providing a direction for future research. The concluding observations from the study are provided as a summing up in Chapter VIII.

# Conceptual Underpinnings of Cancer Staging and Overall Survival

Cancer staging, which is a crucial prognostic factor, is the classification of (breast) cancer depending on the size and spread in respect of tumor, node, and metastasis (TNM). When such classification is carried out after (before) surgery of tumor, it's known as pathologic (clinical) staging. Pathologic staging is often considered more useful as it uses additional information contained inside the tumor. If, however, surgery is preceded by other treatments like chemo/radiation therapy, pathologic staging could be misleading. TCGA-GDC database provides data on the pathologic stage of breast cancer patients in line with the American Joint Committee on Cancer under the head: `ajcc_pathologic_stage`. According to Rocky Mountain Cancer Centers (2023) and DePolo (2023), breast cancer can be categorized into five major stages[3][6]. While stages I through IV may be considered invasive, stage 0 may be treated as non-invasive or in situ. However, the head: `ajcc_pathologic_stage` in the TCGA database has I to IV stages. As cancer is small and confined to one area in stage I, we take it as the non-invasive stage. As cancer becomes larger and/or extends to other parts of body, it moves to stages II, III and IV, which are taken as invasive stages in the current study. Depending on the invasive (as also its sub-categories) or non-invasive stage of the cancer, best treatments are designed for the patient. Overall survival (OS), a la Gobbini et al (2018), is defined as the duration of time from the date of initial diagnosis to the date of death from any cause, or the date of last follow-up for censored patients[10].

Often, it's presented in per cent of total number of cancer patients for a given number of years/days of survival. As the number of years/days of survival is raised, OS rate goes down. OS rates can vary depending on the cancer stage. Patients diagnosed with stage I or II breast cancers are expected to have higher overall survival rates than people diagnosed with stage III or IV breast cancers.

# A Brief Survey of Literature

Breast cancers with virtually identical TNM characteristics may exhibit highly contrasting behaviors due to divergent molecular profiles. Therefore, predictive models need to incorporate histopathological and molecular predictors to explain better the heterogeneity of the complex group of diseases, called cancer (Giacomelli et al, 2023)[9]. Cancer research (on breast) using survival analysis may be classified into two types: binary classification and risk regression. As part of binary classification approach, short-survival group is distinguished from long-survival group using a duration of survival, say 5 years. In risk regression approach, each patient is assigned a risk score. Using various classification methods like gradient boosting, random forest, artificial neural network, and support vector machine, Zhao et al (2018) attempted to predict 5-year breast cancer survival based on gene expression data and clinical/pathological factors[23]. As different classification methods produced similar accuracy and area under the curve (AUC), they highlighted the role of quality of the data itself. Similarly, Gevaert et al (2006) applied Bayesian networks and obtained higher AUC with the use of integrated dataset than the standalone clinical or gene expression dataset[8]. Using multiple kernel learning, Sun et al (2018) integrated genomic data and pathological imaging data, resulting in marginally higher AUC than the standalone genomic data[22]. In risk regression framework, Cox-proportional hazards model has often been used for survival analysis of breast cancer, wherein predictors are multiplicatively related to hazards like the event of death. With the introduction of machine learning like neural network, linear relationship in the Cox-proportional hazards model is replaced by a non-linear relationship

(Huang et al, 2019)[11]. Bismeijer et al (2018) used functional sparse-factor analysis to integrate different types of molecular data and identify dominant biological processes active in breast and lung tumors[2]. In view of missing data in the TCGA-GDC dataset, Sherafatian (2018) attempted to identify minimal set of biomarkers to classify cancer status and different subtypes of cancer[21]. Rendleman et al (2019) applied imputation techniques for missing data and also dimensionality reduction techniques to bring down computational cost while ensuring good performance[19].

# Preliminary Data Analysis – A Clinical Fact Sheet

To start with, we consider clinical data from the TCGA-GDC database and focus on the female category as breast cancer is overwhelmingly a disease faced by females. The age group of around 60 years is found to have the maximum number of patients who are diagnosed with breast cancer, followed by patients of about 50 years (Figure 4.1).



Figure 4.1: Age at Diagnosis

Breast cancer patients are predominantly white, followed by Black/African American and Asian in that order (Figure 4.2).



Figure 4.2: Ethnicity

The histological type of cancer provides details on the microscopic appearance of tumor. Based on this criterion, infiltrating carcinoma is found to have the highest frequency and is followed by mucinous carcinoma (Figure 4.3).



Figure 4.3: Histology Types

As per the pathologic tumor stage, the majority of patients is in stage II, for whom surgery alone may not be curative (Figure 4.4). The next stage in terms of number of patients turns out III, wherein the tumor size and/or its spread is even larger.



Figure 4.4: Stages

The dominant therapy in use against breast cancer has been chemotherapy, followed by hormone therapy (Table 4.1 and Figure 4.5). While breast cancer patients have generally tried various types of therapies, patents in stage II account for the maximum frequency under each therapy.



Figure 4.5: Stages in different Therapies

Table 4.1: Stages in different Therapies

| Stage | Ancillary | Chemotherapy | Hormone therapy | Immunotherapy | Targeted Molecular therapy | Vaccine |
|---|---|---|---|---|---|---|
| Stage I | 1 | 176 | 100 | 6 | 3 | 0 |
| Stage II | 12 | 920 | 390 | 19 | 18 | 0 |
| Stage III | 9 | 467 | 153 | 12 | 9 | 1 |
| Stage IV | 1 | 12 | 12 | 0 | 1 | 0 |

Source: TCGA database.

# Variables and Methodology

Terminologies used in the clinical dataset, data description and percentage missing values for 107 entries are presented in Appendix. Throughout the study, cancer staging and overall survival as explained in Chapter II are taken as response variables. Clearly, cancer staging can have binary outcome: invasive or non-invasive whereas overall survival is a continuous variable. After carefully examining the suitability of all clinical variables based on criteria such as clinical importance, absence of multi-collinearity and percentage missing values, we have selected the following nine as predictor variables: `age_at_diagnosis` , `menopause_status`, `margin_status` , `lymph_nodes_examined_count` , `er_status_by_ihc` , `pr_status_by_ihc` , `her2_status_by_ihc`, `history_other_malignancy` and `lymph_nodes_examined`. At a disaggregated level, the set of nine predictors b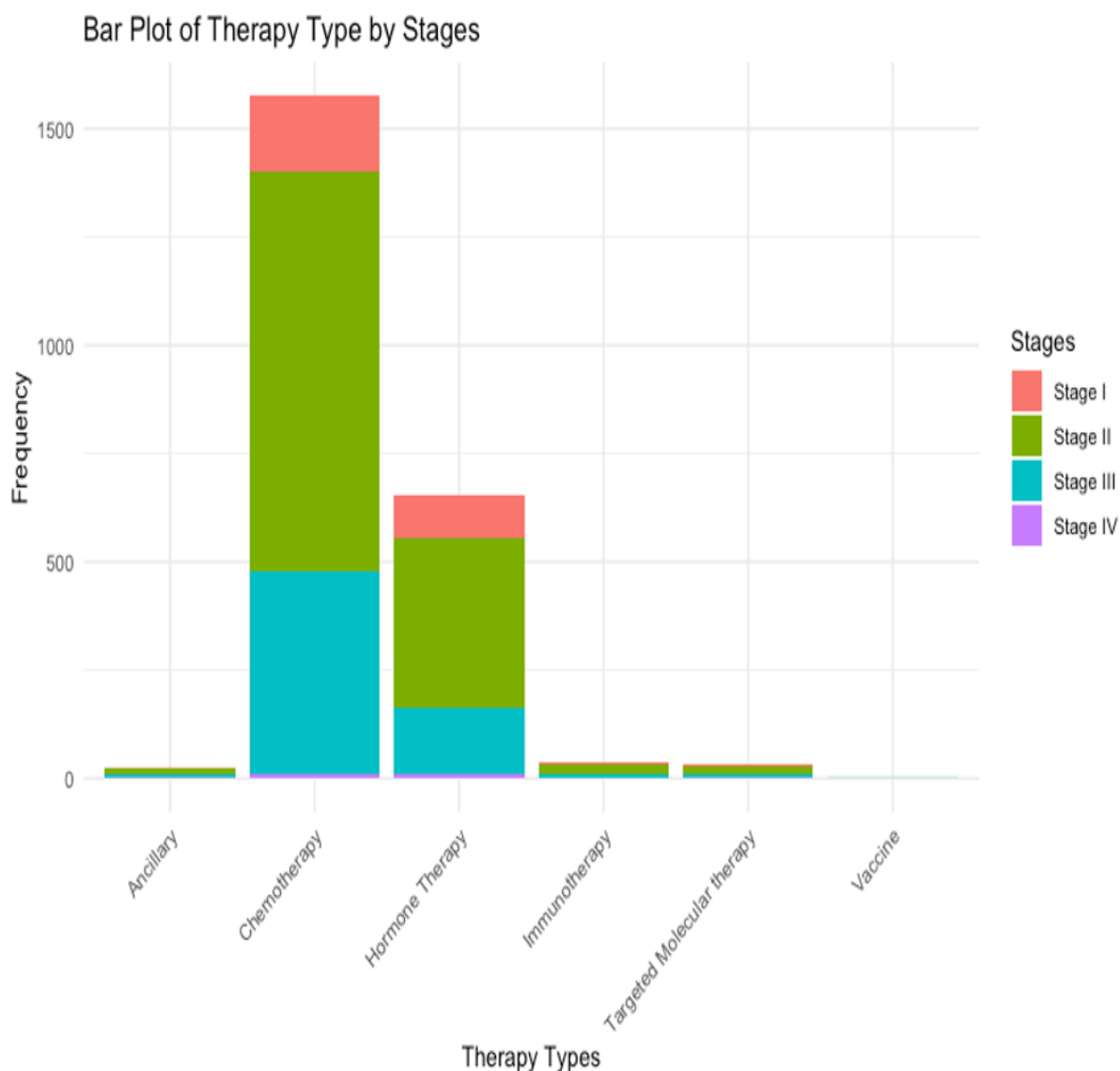oils down to a granular set of 24 clinical variables. Summary statistics of the select set of clinical predictors are presented in Appendix. To start with, we have attempted to fit a logistic regression on the clinical dataset, using cancer staging as the dependent categorical variable. Logistic regression estimates the probability of an event occurring based on a given set of independent variables. The coefficient in this model is commonly estimated through maximum likelihood estimation. After carrying out the logistic regression based on 24 independent variables at the disaggregated level, we perform a backward variable selection and drop variables, one at a time, if such removal helps improve the model's Akaike Information Criterion (AIC). For survival analysis, Kaplan-Meier Estimator, which is a non-parametric statistic, has been used to estimate the survival function in terms

of percentage of patients surviving for a certain amount of time after diagnosis/treatment. The estimator is versatile enough to accommodate the event of drop-out of patients, data truncation or missing values. We have also used log rank test, which compares the survival curves of two or more groups. The null hypothesis for a log rank test, which uses a chi-square statistic, is that the groups have the same survival time. These statistics are useful when the variables are categorical or take small number of values. With the induction of gene expression data combined with clinical data, the data dimension increases manifold, given more than 1000 breast cancer patients and over 60,000 genes per patient. In view of the computational challenges in terms of cost/time involved in handling the big data set, we decide to explore and select a set of manageable number of genes, albeit without compromising the quality of analysis. Towards this direction, we first carry out normalization of the data and then perform differential expression (DE) analysis, following Maié and Manolov (2000)[14]. This brings out the genes that show differences in expression level between conditions or are associated with given predictors or responses. This is carried out with the help of a computational pipeline, called `limma_pipeline` in R, which returns the top (set at 60 in current study) differentially expressed genes sorted by p-value. Clustering partitions a dataset into clusters or groups such that data points within a cluster or group are more similar to each other than to those in other clusters. For clustering, we apply K-Means clustering, which partitions data into K clusters based on mean of data points. Deciding the optimal number of clusters, however, remains a challenge. A popular method for optimal cluster count is the Elbow method, which runs K-means clustering for a range of values of K and plots the within-cluster sum of squared distances (WCSS) between each point and the centroid of that cluster. This is also known as scree plot in clustering and Elbow of the plot is the point where the rate of decrease in WCSS slows down, implying that addition of more clusters does not improve the model's fit significantly. The optimal number of clusters can also be based on Silhouette score, which measures in a range of [-1 to +1] how similar an object is to its own cluster compared to other clusters. Scores close to 1 (-1) indicate high (low) similarity, implying well defined clusters. For K-Means clustering, Silhouette score is calculated for different values of K and we choose that K that produces the maximum Silhouette score. We now proceed with the set of 60 select genes and the two response variables to run a random

forest model on the dataset for more than 1000 breast cancer patients. Random forest is a supervised non-linear machine learning algorithm built through an ensemble of decision trees. It builds and combines multiple decision trees, trained on different parts of the same training set, with the goal of overcoming over-fitting problem of individual decision tree as was encountered in logistic regression. We also attempt to fit a neural network model on the 60 select genes and the two response variables. Neural network allows complex nonlinear relationships between the predictors (bottom layer) and the response variable (top layer). There would be intermediate/hidden layers in between the way human brain functions. The technique being extremely versatile has the ability to fit any functional relationship in the data. Finally, we have compared and contrasted the model performance of random forest and neural network, using the criteria like accuracy and precision.

# Results and Interpretation

## 6.1   Logistic Regression

The results of logistic regression when we consider the most disaggregated model are presented in Appendix. It may be noted that `lymph_nodes_examined_count`, and the disaggregated categories of `her2_status_by_ihc`: positive, negative and equivocal, are the only four clinical variables, which are statistically significant. In medical domain, `lymph_nodes_examined_count` is taken as indicative of possible presence of a malignant process and its localization or generalization. It is found statistically significant at 0.001 level of significance. Further, a positive `her2_status_by_ihc` implies that breast cancer tests positive for a protein, namely, human epidermal growth factor receptor 2 (HER2), which promotes the growth of cancer cells. Thus, a positive her2 status poses greater risk of relapse and metastasis than a negative her2 status. If her2 status is not clear, it is called equivocal. Though her2 status variables turn out statistically significant, the level of significance remains far higher than that for lymph nodes examined count. As only four variables are found statistically significant out of 24 independent variables, possible overfitting of the model is dealt with by applying a backward variable selection process based on AIC (Appendix). AIC makes a balance between the aspects of model complexity and goodness of fit and, helps retain variable/s that contribute more to model performance. Thus, lymph nodes examined count has turned out to be the single independent variable, which is statistically sig-

nificant at a very low level of significance of 0.001 and hence, sufficient to predict the binary cancer staging. The log odds of binary cancer staging is estimated to increase by 0.16941 with one unit increase in lymph nodes examined count. As AIC decreases with backward selection to 616.2 from 645.32 in the original disaggregated model, the single-variable reduced model appears to throw up a better fit to the data than the initial model. As all the clinical variables barring one are not statistically significant in the reduced (final) model, logistic regression framework based on clinical variables seems hardly satisfactory for explaining cancer staging.

## 6.2  Kaplan-Meier Estimator & Log Rank Test



Figure 6.1: Vital Status

While the binary vital status: alive and dead appears to be an analytically meaningful category, we check it rigorously using the technique of survival analysis. For the alive,

the survival curve remains horizontal all through, implying that the survival probability is unchanged at around 1 irrespective of the survival time, measured in number of days (Figure 6.1). In contrast, the survival curve for the dead has been downward sloping and lies below that for the alive group all through. The apparently higher survival probability of the alive than of the dead is checked for confirmation through the log rank test, which finds that the difference in means between the two curves is statistically significant, p-value being closer to zero. The lower panel of Figure 6.1, which provides the number of patients at risk shows that the number of patients declines gradually for both the groups and in tandem at least up to 4000 days.



Figure 6.2: Stages

Note: 'Filtered' qualifier indicates computation after dropping of missing values.

Patients with tumor in stage IV is found to have the lowest survival percentage all through whereas patients with tumor in stage I have the highest survival percentage at least up to 4000 days (Figure 6.2). In any case, there seems to be significant difference

in survival percentage across groups of patients with tumor in different stages. This is also buttressed by very low p-value. The relative number of patients at risk across stages also remain similar over time.



Figure 6.3: ER Status

Note: 'Filtered' qualifier indicates computation after dropping of missing values.

Clinically, patients with positive ER status are expected to respond better to hormone therapies like tamoxifen that block estrogen signaling. Such patients may have higher survival probability up to 4000 days (Figure 6.3). However, the utility of ER status category could not be confirmed by log rank test as we obtain high p-value.

Figure 6.4: her2 Status

Note: 'Filtered' qualifier indicates computation after dropping of missing values.

Grouping of patients as per her2 status – positive or negative – appears to matter, leading to lower probability of survival for her2 positive group (Figure 6.4). In such group, tumor tends to be more aggressive. This is also supported by log rank test with small p-value.

## 6.3   Differential Expression Analysis

Ahead of the DE analysis, we drop genes that have small counts of less than 10. Then, we apply trimmed mean of M-values (TMM) method to normalize the data and convert it to produce a similar variance as arrays, based on voom method, which estimates the mean-variance relationship in the data (Figure 6.5). Then, we carry out the differential analysis by comparing primary solid tumor samples with solid tissue normal as the reference group. In the process, we fit a series of linear models, one to each of the probes

and rank the most differentially expressed genes based on p-values. Accordingly, we pick up the top-60 genes for further analysis like clustering and machine learning. The first two principal components of the expression matrix of the voom-transformed data are presented in Figure 6.6, which shows that the two sample groups, namely, primary solid tumor and solid tissue normal have clearly distinguishable RNA expression profiles.



Figure 6.5: Mean Variance Trend

Figure 6.6: PCA

## 6.4    Clustering

The scree plot of clusters based on 60 identified genes shows that the elbow point is attained when the number of clusters, that is, K takes a value equal to 2 or 3, where adding more clusters does not significantly reduce the sum of squared distances (Figure 6.7).

Figure 6.7: Elbow Method

This is also confirmed by Silhouette score, which reaches its maximum of 0.989 when the number of clusters is set at 2 (Table 6.1). As the number of clusters increases from 2 to 3 and so on, the Silhouette score decreases slowly, signifying a gradual moderation in the degree of similarity within a group.

Table 6.1: Cluster and Silhouette Score

| Number of clusters | Score |
|:---:|:---:|
| 2 | 0.989 |
| 3 | 0.987 |
| 4 | 0.969 |
| 5 | 0.942 |
| 6 | 0.934 |
| 7 | 0.906 |
| 8 | 0.896 |

We may, therefore, fix the optimal number of clusters at 2 for the set of 60 genes. This is also consistent with the binary cancer staging considered in the study: invasive and non-invasive.

## 6.5   Random Forest

In the random forest output for overall survival as the dependent variable, R-squared score calculated for the training data (using 75% of the dataset) turns out low at 0.0715, indicating that the model might not be a good fit. However, a low Test RMSE for Best Model at 0.1300 suggests a reasonable predictive performance: predictions deviate by approximately 0.1300 units from actual values. Test RMSE increases marginally to 0.1309 if 60% data are used for training and 40% for testing, implying stability in results on resampling. When cancer staging is the dependent variable, the test accuracy stands at a high of 0.8324, far exceeding the 'rule of thumb' for a minimum accuracy of 0.70 for model acceptance. ROC curve, which indicates the trade-off between true positive rate and false positive rate at each threshold setting remains throughout above the diagonal line. However, AUC of 0.6288 suggests that the model's ability to discriminate between invasive and non-invasive cases is somewhat better than random chance (AUC=0.5) (Figure 6.8). Overall, the model stands out in predicting cancer staging. The test accuracy increases slightly to 0.8339 if 60% data are used for training and 40% for testing, implying stability in results on resampling.

Figure 6.8: ROC

## 6.6  Neural Network

When cancer staging is used as the target variable, we obtain a very low test accuracy of the model at 0.1676, using 75% and 25% of data for training and testing, respectively. Threshold optimization for classification also does not improve the model's performance. We do not recommend the model in its current form for predicting cancer staging. Test accuracy of neural network model with overall survival as the target variable turns out high at 0.7528. If we set the training at 60% of the dataset, test accuracy of the model decreases marginally to 0.1661 for cancer staging while test accuracy declines fractionally to 0.7515 for overall survival, indicating relative stability in results upon resampling.

## 6.7 Alternative Methodology for Gene Selection & Results

We check the robustness of results based on the DE analysis by employing an alternative methodology as under. We first check the presence of all the genes with respect to patient by doing a groupby. As the gene names are found equally spread out across patients, we sort the genes by case count in descending order, shuffle the data, reset the index and then take the first 60 genes to bring in randomness into the selection of genes. The rest of the exercise – clustering, fitting a random forest model and a neural network model, and comparing the model performance – is repeated with the new-found set of 60 genes.

## 6.8 Clustering

The scree plot of clusters based on 60 identified genes is presented in Figure 6.9, wherein the elbow point is attained at K=2 or 3. If we add more clusters, it may not significantly reduce the sum of squared distances. This is in a way confirmed by Silhouette score, which attains its maximum at K=2. We may, therefore, fix the optimal number of clusters at 2 for the set of 60 genes based on random sampling as has been the case based on the DE analysis.



Figure 6.9: Elbow Method

## 6.9   Random Forest

When we take overall survival as the dependent variable, R-squared score calculated for the training data (using 75% of the dataset) turns out low at 0.2220, even though it remains higher than that based on the DE analysis at 0.0715. Also, Test RMSE for Best Model at 0.0997 turns out lower than that of 0.1300 under the DE analysis. This suggests a reasonable predictive performance of the model wherein genes are picked up based on random sampling. Test RMSE increases marginally to 0.1097 if 60% data are used for training and 40% for testing, implying stability in results on resampling. When we consider cancer staging as the dependent variable, the test accuracy stands at a high of 0.8523, which is marginally higher than that at 0.8324 under the DE analysis. ROC curve based on the random sampling approach also remains throughout above the diagonal line. Moreover, AUC of 0.7236 turns out better than that at 0.6288 under the DE analysis (Figure 6.10). Overall, the model stands out in predicting cancer staging. The test accuracy decreases slightly to 0.8507 if 60% data are used for training and 40% for testing, implying stability in results on resampling.
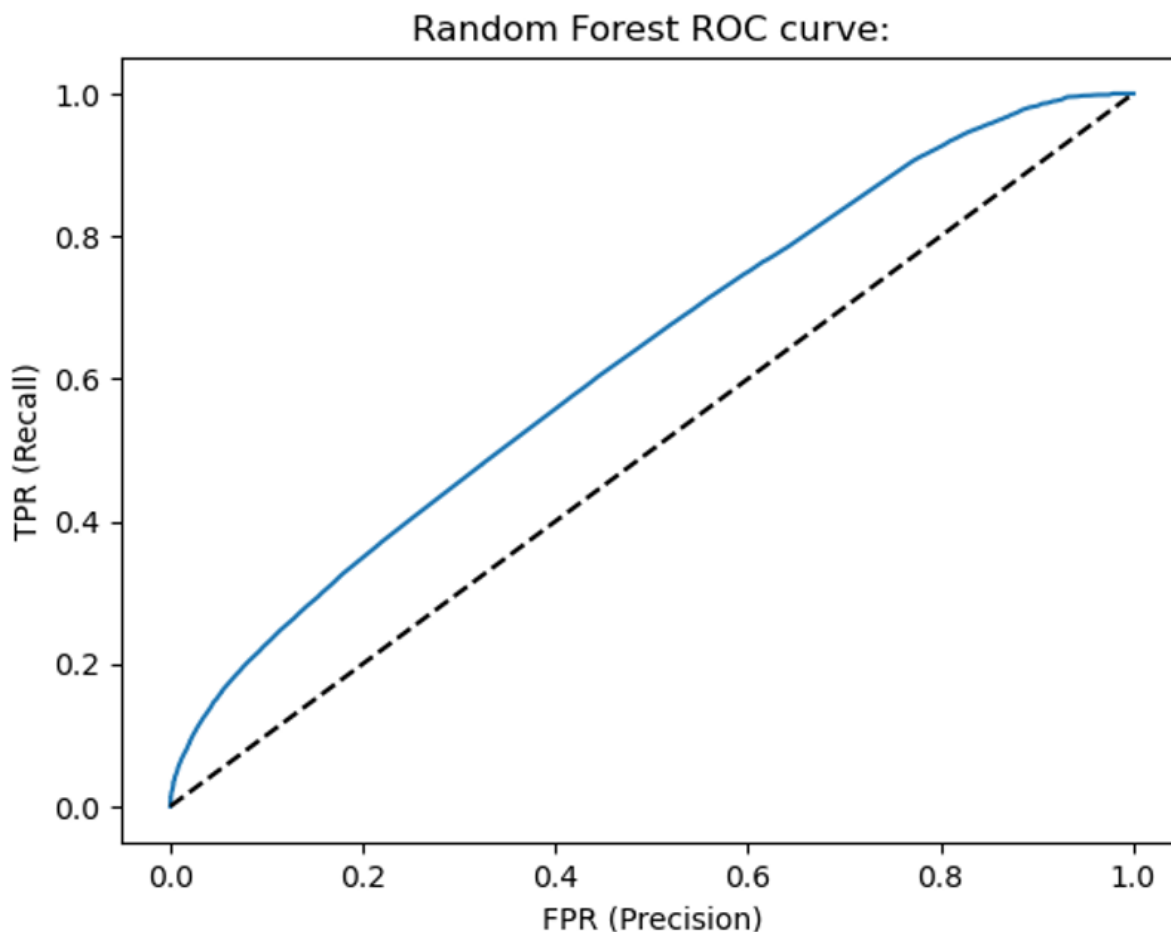


Figure 6.10: ROC

## 6.10   Neural Network

When we take cancer staging as the target variable, we obtain a very low test accuracy at 0.1477, which stands out marginally lower than that at 0.1676 under the DE analysis, using 75% and 25% of data for training and testing, respectively. Threshold optimization for classification also does not improve the model's performance. We do not recommend the model in its current form for predicting cancer staging. Test accuracy of neural network model with overall survival as the target variable turns out higher at 0.7817 than at 0.7528 under the DE analysis. If we set the training at 60% of the dataset, test accuracy of the neural network model increases marginally to 0.1493 for cancer staging while test accuracy declines fractionally to 0.7735 for overall survival, indicating relative stability in results upon resampling. Random forest model seems to fit the given dataset well both for explaining cancer staging and overall survival regardless of the selection process of genes – either based on the DE analysis or the random sampling of genes. Both the selection processes of genes – one based on the DE analysis and the other based on random sampling – tend to produce very similar results based on random forest or neural network even though composition of the set of 60 genes is different depending on the selection process. Further, there is relative stability in results with respect to the size of training and testing datasets irrespective of the choice of model – random forest or neural network.

# 7

# Limitations of the Study & Scope for Future Research

Gene expression data has attracted lots of attention in biomarker research. Various statistical and machine learning models are used to identify the biomarkers related to breast cancer. The models used in our study have displayed a degree of robustness upon resampling by varying the size of training and testing datasets. Results have also remained invariant to the two gene selection processes used in the study. However, the models can display low stability of feature selection in view of high dimensionality of gene data and low sample size of patients (Kim et al, 2020)[13]. Indeed, neural network model did not perform well in the study for cancer staging and need to be taken up for further investigation and improvements. Machine learning models also pose challenges in terms of interpretation and are often dubbed as black boxes. Further, notwithstanding high prediction accuracy, they may or may not be clinically relevant. The models have also assumed away inter-dependence of genes, if any. Some of these issues could be taken up for exploration as part of future research.

# A Summing Up

Breast cancer accounts for around 25% of all female cancer diagnoses and 15% of all cancer-related deaths. It exhibits a heterogeneous nature, characterized by diverse clinical presentations, treatment responses, and prognoses. This heterogeneity underscores the urgency for accurate cancer staging and overall survival prediction. In this backdrop, the study has explored the clinical and genomic aspects of breast cancer patients based on the landmark TCGA-GDC database. It observes that the age group around 60 years is most vulnerable to breast cancer. The disease, which is known as the emperor of maladies (Mukherjee, 2010) is found more prevalent among the whites[16]. Infiltrating carcinoma is the histological type, often found amongst breast cancer patients. The majority of patients in the study belong to the invasive stage II of breast cancer, for whom surgery alone may not be curative. Most of the patients have gone for chemotherapy as part of treatment. In a rigorous framework of logistic regression, the study initially attempts to link cancer staging to a number of relevant clinical variables. It brings out the lymph nodes examined count as the only clinical variable, which is statistically significant in the final equation. Based on Kaplan-Meier Estimator and Log Rank Test, groupings of patients in terms of vital status: alive and dead, tumor staging, and her2 status (ie, aggressive tumor development) are found to be meaningful analytical categories for survival analysis while the grouping based on ER status (ie, sensitivity to hormones) is not so. Given the limited utility of clinical data for survival analysis, we consider the gene expression data. With the induction of gene data, the dimensionality of dataset increases. We handle the computational challenge by applying differential

expression analysis and selecting 60 genes based on p-value out of about 60,000 genes without compromising the performance efficiency. Then we carry out an unsupervised learning by clustering the 60 genes and find that the optimal number of clusters may be two. This is also consistent with the number of cancer staging being two: invasive and non-invasive, used in the study as a dependent variable. In view of the limitations of logistic regression in explaining cancer staging, we employ supervised learning, namely, random forest and neural network on the dataset. We find that the random forest model works out well with the dataset for explaining both cancer staging and overall survival. The results of the models remain robust against resampling, passing the litmus test of stability for assessment of results. They also remain largely invariant when we use random sampling for selection of genes in place of differential expression analysis. We are, however, aware that machine learning techniques often present black boxes which are not as interpretable as regression models rooted in explicit mathematical formulations. In this context, we suggest that future research can focus on unraveling the black boxes with a view to improving their utility in the medical domain.

# Definitions and Data Description

1. **Estrogen Receptor Status (ER status)**

   - ER status refers to whether the breast cancer cells have receptors for the hormone estrogen. This is typically measured by immunohistochemistry (IHC) on tumor samples.

   - If more than 1% of tumor cells stain positive, the cancer is deemed ER positive. This occurs in about 70% of breast cancers. ER+ cancers tend to respond better to hormone therapies like tamoxifen that block estrogen signaling. (1)

   - Example: A pathology report states the patient's breast cancer is ER positive, with 95% of tumor cells staining positive for ER by IHC. This indicates an ER+ tumor that will likely respond well to hormonal therapy.

2. **HER2 Status**

   - HER2 is a growth-promoting protein. Tumors with increased levels of HER2 (HER2 positive) tend to be more aggressive. (2)

   - HER2 status is measured by IHC and/or fluorescent in situ hybridization (FISH). FISH counts the number of HER2 gene copies within cancer cells.

   - Example: A patient's tumor is scored 3+ on IHC staining, indicating strong, uniform staining in >10% of tumor cells. This confirms HER2 positivity.

Alternatively, FISH showing HER2 gene amplification (HER2:Chr17 ratio > 2) would confirm a HER2+ result.

3. **Lymph Node Involvement**

   - Lymph nodes near the breast are often the first place breast cancer spreads. Assessing if cancer is present helps determine prognosis and need for additional treatments. (3)

   - Example: A pathology report indicates 2 of 15 lymph nodes examined contained metastatic breast carcinoma. This suggests the cancer has started spreading to the lymph nodes. More aggressive adjuvant chemotherapy may be warranted.

4. **Treatment Implications**

   - Early stage cancers are treated with surgery and radiation. Additional treatments like chemotherapy or hormonal therapy are considered for node positive or higher risk cancers to reduce recurrence risk. (4)

   - Example: A large tumor with close surgical margins and cancer in 3 lymph nodes has a higher risk of recurrence. Adjuvant chemotherapy followed by 5 years of hormone therapy is typically recommended.

5. **`form_completion_date`**

   - This field typically indicates the date when a particular form or document related to cancer research or medical records was completed. It's essential for tracking the timing of data collection.

6. **`prospective_collection`**

   - This field is used to indicate whether the data was collected prospectively, meaning it was gathered in real-time as events occurred, or retrospectively, meaning it was collected from historical records.

   - Example: Prospective collection may involve ongoing clinical trials where data is collected as patients receive treatment.

7. **`retrospective_collection`**

- Similar to prospective collection, this field indicates whether the data was collected retrospectively, often from existing medical records or databases.

- Example: Retrospective collection might involve studying the medical records of cancer patients over the past decade to analyze treatment outcomes.

8. **menopause_status**

   - This field refers to a patient's menopause status, whether they are pre-menopausal, postmenopausal, or perimenopausal.

   - Example: Understanding menopause status is crucial in breast cancer research, as hormone receptor status can influence treatment decisions.

9. **history_other_malignancy**

   - Indicates whether the patient has a history of other malignancies or cancer types apart from the current diagnosis.

   - Example: A patient with a history of lung cancer who is now diagnosed with breast cancer would have a history of other malignancy.

10. **history_neoadjuvant_treatment**

    - This field indicates whether the patient received neoadjuvant treatment (treatment before surgery).

    - Example: A breast cancer patient receiving chemotherapy before surgery to shrink the tumor.

11. **radiation_treatment_adjuvant**

    - Specifies whether adjuvant radiation therapy was administered after primary treatment.

    - Example: Post-surgery radiation therapy for breast cancer patients to reduce the risk of recurrence.

12. **pharmaceutical_tx_adjuvant**

    - Indicates the use of adjuvant pharmaceutical treatments, such as hormone therapy or targeted therapy.

- Example: Adjuvant hormone therapy for hormone receptor-positive breast cancer.

13. **`histologic_diagnosis_other`**

    - If the standard histologic diagnosis doesn't apply, this field captures other specific histologic diagnoses.

    - Example: In rare cases, a cancer type might have unique histological features not covered by standard classifications.

14. **`initial_pathologic_dx_year`**

    - The year in which the initial pathologic diagnosis was made.

    - Example: A patient diagnosed with breast cancer in 2020 would have an `initial_pathologic_dx_year` of 2020.

15. **`method_initial_path_dx`**

    - Describes the method used for the initial pathologic diagnosis, which can include methods like biopsy or imaging.

    - Example: A breast cancer diagnosis may involve a core needle biopsy for pathologic diagnosis.

16. **`method_initial_path_dx_other`**

    - If the method used for the initial pathologic diagnosis is not covered by standard categories, this field allows for specifying other methods.

    - Example: A specialized imaging technique used for diagnosing rare tumors.

17. **`surgical_procedure_first`**

    - Indicates the type of surgical procedure performed as the first treatment.

    - Example: Lumpectomy or mastectomy for breast cancer.

18. **`first_surgical_procedure_other`**

    - If the first surgical procedure is not covered by standard categories, this field allows for specifying other procedures.

- Example: Complex reconstructive surgery following tumor removal.

19. **margin_status**

    - Specifies whether surgical margins were clear or involved by tumor cells.

    - Example: Negative margin status indicates clear margins, while positive margin status means tumor cells were found at the edges of the surgical specimen.

20. **surgery_for_positive_margins**

    - Indicates whether additional surgery was performed due to positive margin status.

    - Example: Re-excision surgery following positive margin status in breast cancer surgery.

21. **surgery_for_positive_margins_other**

    - Allows for specifying other surgical procedures performed in case of positive margins if they don't fit standard categories.

    - Example: A specialized resection technique used for specific tumor types.

22. **margin_status_re-excision**

    - Describes the margin status after re-excision surgery.

    - Example: After re-excision, margin status may change from positive to negative.

23. **axillary_staging_method**

    - Refers to the method used for staging axillary lymph nodes in breast cancer or similar procedures in other cancers.

    - Example: Sentinel lymph node biopsy as an axillary staging method in breast cancer.

24. **axillary_staging_method_other**

- Allows for specifying other methods used for axillary staging if they are not covered by standard categories.

- Example: Rarely used axillary staging techniques.

25. **micromet_detection_by_ihc**

- Indicates whether micrometastases were detected using immunohistochemistry (IHC).

- Example: Detecting micrometastases in lymph nodes through IHC staining.

26. **lymph_nodes_examined**

- The total number of lymph nodes examined during surgery.

- Example: If 20 lymph nodes were removed and examined, this field would contain the value 20.

27. **lymph_nodes_examined_count**

- Specifies the count of lymph nodes examined, which can be crucial for staging.

- Example: If only 10 out of 20 lymph nodes examined were found to be positive for cancer, this field would contain the value 10.

28. **lymph_nodes_examined_he_count**

- This field might represent the count of lymph nodes examined using hematoxylin and eosin (HE) staining.

- Example: Count of lymph nodes examined using HE staining to determine metastasis.

29. **lymph_nodes_examined_ihc_count**

- Indicates the count of lymph nodes examined using immunohistochemistry (IHC) staining.

- Example: The number of lymph nodes examined with IHC to detect micrometastases.

30. **ajcc_staging_edition**

   - Refers to the edition of the American Joint Committee on Cancer (AJCC) staging system used for cancer staging.

   - Example: AJCC 8th edition for breast cancer staging.

31. **ajcc_tumor_pathologic_pt**

   - Represents the pathologic T stage in the AJCC staging system, indicating the primary tumor size.

   - Example: pT2 for a breast tumor larger than 2 cm but not larger than 5 cm.

32. **ajcc_nodes_pathologic_pn**

   - Indicates the pathologic N stage in the AJCC staging system, describing the extent of regional lymph node involvement.

   - Example: pN1a for breast cancer with metastasis to movable ipsilateral axillary lymph nodes.

33. **ajcc_metastasis_pathologic_pm**

   - Refers to the pathologic M stage in the AJCC staging system, indicating the presence or absence of distant metastases.

   - Example: pM0 for no distant metastases in breast cancer.

34. **ajcc_pathologic_tumor_stage**

   - Represents the overall pathologic stage of cancer according to the AJCC staging system.

   - Example: Stage IIA for breast cancer with specific pT and pN stages.

35. **metastasis_site**

   - Specifies the site(s) of metastasis in cases where cancer has spread to distant locations.

   - Example: Liver, lung, bone, or brain metastasis in various cancer types.

36. **metastasis_site_other**

- Allows for specifying other sites of metastasis not covered by standard categories.

- Example: Rare sites of metastasis that may require special consideration.

37. **er_status_by_ihc**

    - Indicates estrogen receptor (ER) status determined by immunohistochemistry (IHC) staining.

    - Example: ER-positive or ER-negative status in breast cancer.

38. **er_status_ihc_Percent_Positive**

    - Represents the percentage of ER-positive cells determined by IHC.

    - Example: If 80% of tumor cells are ER-positive, this field would contain the value 80.

39. **er_positivity_scale_used**

    - Specifies the scale used to measure ER positivity.

    - Example: The Allred scoring system or H-score.

40. **er_ihc_score**

    - Represents the ER IHC score based on staining intensity and percentage of positive cells.

    - Example: An ER IHC score of 7 on the Allred scale.

41. **er_positivity_scale_other**

    - Allows for specifying other ER positivity scales if they are used.

    - Example: A novel ER positivity scoring system developed for research purposes.

42. **er_positivity_method**

    - Describes the method used to determine ER positivity.

    - Example: Immunohistochemistry (IHC) or reverse transcription-polymerase chain reaction (RT-PCR).

43. **pr_status_by_ihc**

   - Indicates progesterone receptor (PR) status determined by immunohisto-chemistry (IHC).

   - Example: PR-positive or PR-negative status in breast cancer.

44. **pr_status_ihc_percent_positive**

   - Represents the percentage of PR-positive cells determined by IHC.

   - Example: If 60% of tumor cells are PR-positive, this field would contain the value 60.

45. **pr_positivity_scale_used**

   - Specifies the scale used to measure PR positivity.

   - Example: The Allred scoring system or H-score.

46. **pr_positivity_ihc_intensity_score**

   - Represents the PR IHC intensity score, which is part of the PR positivity assessment.

   - Example: An PR IHC intensity score of 3 indicating strong staining.

47. **pr_positivity_scale_other**

   - Allows for specifying other PR positivity scales if used.

   - Example: A custom PR positivity scoring system developed for specific research.

48. **pr_positivity_define_method**

   - Describes the method used to define PR positivity.

   - Example: Immunohistochemistry (IHC) or gene expression profiling.

49. **her2_status_by_ihc**

   - Indicates human epidermal growth factor receptor 2 (HER2) status deter-mined by immunohistochemistry (IHC).

- Example: HER2-positive or HER2-negative status in breast cancer.

50. **her2_ihc_percent_positive**

    - Represents the percentage of cells with HER2-positive staining determined by IHC.

    - Example: If 90% of tumor cells show HER2-positive staining, this field would contain the value 90.

51. **her2_ihc_score**

    - Represents the HER2 IHC score based on staining intensity and percentage of positive cells.

    - Example: An HER2 IHC score of 3+ indicating strong staining.

52. **her2_positivity_scale_other**

    - Allows for specifying other HER2 positivity scales if used.

    - Example: A modified HER2 positivity scoring system for research purposes.

53. **her2_positivity_method_text**

    - Describes the method used to assess HER2 positivity.

    - Example: Immunohistochemistry (IHC), fluorescence in situ hybridization (FISH), or gene expression profiling.

54. **her2_fish_status**

    - Indicates HER2 status determined by fluorescence in situ hybridization (FISH).

    - Example: HER2 amplified or not amplified status in breast cancer.

55. **her2_copy_number**

    - Represents the HER2 copy number determined by FISH.

    - Example: If the HER2 copy number is 5, this field would contain the value 5.

56. **cent17_copy_number**

- Represents the cent17 copy number determined by FISH.

- Example: If the cent17 copy number is 2, this field would contain the value 2.

57. **her2_and_cent17_cells_count**

- Indicates the count of cells in which both HER2 and cent17 were counted during FISH analysis.

- Example: If 100 cells were counted, this field would contain the value 100.

58. **her2_cent17_ratio**

- Represents the ratio of HER2 copy number to cent17 copy number in FISH analysis.

- Example: If the HER2/cent17 ratio is 2.5, this field would contain the value 2.5.

59. **her2_and_cent17_scale_other**

- Allows for specifying other scales or methods used in HER2 and cent17 copy number analysis.

- Example: A novel scale developed for specific research purposes.

60. **her2_fish_method**

- Describes the method used for HER2 fluorescence in situ hybridization (FISH) analysis.

- Example: Standard FISH or alternative FISH techniques.

61. **new_tumor_event_dx_indicator**

- Indicates whether a new tumor event or recurrence was diagnosed.

- Example: Used to mark whether a patient has experienced a cancer recurrence.

62. **nte_er_status**

- ER status of the new tumor event, if applicable.

- Example: To determine if the recurrent tumor retains its ER-positive or ER-negative status.

63. **nte_er_status_ihc__positive**

    - Indicates whether the ER status of the recurrent tumor is positive.

    - Example: To track changes in ER status between the primary tumor and recurrence.

64. **nte_er_ihc_intensity_score**

    - Represents the intensity score for ER staining in the recurrent tumor.

    - Example: To assess changes in ER staining intensity in recurrent cancer.

65. **nte_er_positivity_other_scale**

    - Allows for specifying alternative scales or methods used to assess ER positivity in the recurrent tumor.

    - Example: A modified scoring system for ER positivity in recurrent cancer research.

66. **nte_er_positivity_define_method**

    - Describes the method used to define ER positivity in the recurrent tumor.

    - Example: Immunohistochemistry (IHC) or gene expression profiling in recurrent cancer analysis.

67. **nte_pr_status_by_ihc**

    - PR status of the new tumor event, if applicable.

    - Example: To determine if the recurrent tumor retains its PR-positive or PR-negative status.

68. **nte_pr_status_ihc__positive**

    - Indicates whether the PR status of the recurrent tumor is positive.

    - Example: To track changes in PR status between the primary tumor and recurrence.

69. **nte_pr_ihc__score**

   - Represents the intensity score for PR staining in the recurrent tumor.

   - Example: To assess changes in PR staining intensity in recurrent cancer.

70. **nte_pr_positivity_other_scale**

   - Allows for specifying alternative scales or methods used to assess PR positivity in the recurrent tumor.

   - Example: A modified scoring system for PR positivity in recurrent cancer research.

71. **nte_pr_positivity_define_method**

   - Describes the method used to define PR positivity in the recurrent tumor.

   - Example: Immunohistochemistry (IHC) or gene expression profiling in recurrent cancer analysis.

72. **nte_her2_status**

   - HER2 status of the new tumor event, if applicable.

   - Example: To determine if the recurrent tumor retains its HER2-positive or HER2-negative status.

73. **nte_her2_status_ihc__positive**

   - Indicates whether the HER2 status of the recurrent tumor is positive.

   - Example: To track changes in HER2 status between the primary tumor and recurrence.

74. **nte_her2_positivity_ihc_score**

   - Represents the HER2 IHC score in the recurrent tumor.

   - Example: To assess changes in HER2 staining intensity in recurrent cancer.

75. **nte_her2_positivity_other_scale**

   - Allows for specifying alternative scales or methods used to assess HER2 positivity in the recurrent tumor.

- Example: A modified scoring system for HER2 positivity in recurrent cancer research.

76. **nte_her2_positivity_method**

- Describes the method used to assess HER2 positivity in the recurrent tumor.

- Example: Immunohistochemistry (IHC), fluorescence in situ hybridization (FISH), or gene expression profiling in recurrent cancer analysis.

77. **nte_her2_fish_status**

- Indicates HER2 status of the new tumor event determined by fluorescence in situ hybridization (FISH).

- Example: To determine if the recurrent tumor retains its HER2 amplified or not amplified status.

78. **nte_her2_signal_number**

- Represents the HER2 signal number in FISH analysis of the recurrent tumor.

- Example: To assess changes in HER2 copy number in recurrent cancer.

79. **nte_cent_17_signal_number**

- Represents the cent17 signal number in FISH analysis of the recurrent tumor.

- Example: To assess changes in cent17 copy number in recurrent cancer.

80. **her2_cent17_counted_cells_count**

- Indicates the count of cells in which both HER2 and cent17 were counted during FISH analysis of the recurrent tumor.

- Example: If 150 cells were counted, this field would contain the value 150.

81. **nte_cent17_her2_ratio**

- Represents the ratio of HER2 copy number to cent17 copy number in FISH analysis of the recurrent tumor.

- Example: To assess changes in the HER2/cent17 ratio in recurrent cancer.

82. **nte_cent17_her2_other_scale**

    - Allows for specifying alternative scales or methods used in HER2 and cent17 copy number analysis of the recurrent tumor.

    - Example: A novel scale developed for specific research purposes.

83. **nte_her2_fish_define_method**

    - Describes the method used for HER2 fluorescence in situ hybridization (FISH) analysis in the recurrent tumor.

    - Example: Standard FISH or alternative FISH techniques for recurrent cancer research.

84. **anatomic_neoplasm_subdivision**

    - Specifies the anatomical subdivision or location of the neoplasm (tumor).

    - Example: In breast cancer, this could specify whether the tumor is in the upper, lower, inner, or outer quadrant of the breast.

85. **clinical_M**

    - Represents the clinical M stage, indicating the presence or absence of distant metastases based on clinical evaluation.

    - Example: cM0 for no clinical evidence of distant metastasis.

86. **clinical_N**

    - Represents the clinical N stage, describing the extent of regional lymph node involvement based on clinical evaluation.

    - Example: cN1 for clinical evidence of metastasis to movable ipsilateral axillary lymph nodes.

87. **clinical_T**

    - Represents the clinical T stage, indicating the size and extent of the primary tumor based on clinical evaluation.

    - Example: cT2 for a clinically determined tumor size larger than 2 cm but not larger than 5 cm.

88. **clinical_stage**

   - Specifies the clinical stage of cancer based on the combined assessment of clinical T, N, and M stages.

   - Example: Stage IIIB for advanced breast cancer with specific cT, cN, and cM stages.

89. **days_to_initial_pathologic_diagnosis**

   - Indicates the number of days from a relevant reference point to the initial pathologic diagnosis.

   - Example: The number of days from the date of biopsy to the date of pathologic diagnosis.

90. **extranodal_involvement**

   - Specifies whether cancer has spread to extranodal sites outside the lymph nodes.

   - Example: Metastasis to extranodal sites like liver, lung, or bone.

91. **histological_type**

   - Describes the specific histological type of cancer, providing details about the tumor's microscopic appearance.

   - Example: Ductal carcinoma in situ (DCIS) for a type of breast cancer.

92. **icd_10**

   - Refers to the International Classification of Diseases, 10th Revision, code used for coding the diagnosis of cancer.

   - Example: C50.0 for malignant neoplasm of the nipple and areola of the female breast.

93. **icd_o_3_histology**

   - Specifies the histology code from the International Classification of Diseases for Oncology, 3rd Edition.

- Example: 8500/3 for invasive ductal carcinoma of the breast.

94. **icd_o_3_site**

    - Specifies the site code from the International Classification of Diseases for Oncology, 3rd Edition.

    - Example: C50.0 for the site of the nipple and areola of the female breast.

95. **informed_consent_verified**

    - Indicates whether informed consent for treatment or participation in research has been verified.

    - Example: Verifying that patients have provided informed consent before participating in a clinical trial.

96. **metastatic_tumor_indicator**

    - Specifies whether the tumor is metastatic, meaning it has spread to distant sites.

    - Example: Used to distinguish between primary and metastatic tumors.

97. **site_of_primary_tumor_other**

    - Allows for specifying the primary tumor site if it doesn't fit standard categories.

    - Example: Rare primary tumor sites not covered by standard classifications.

98. **stage_other**

    - Provides additional details or descriptors related to cancer staging if needed.

    - Example: Additional information about the extent of disease not covered by standard staging criteria.

99. **tissue_source_site**

    - Indicates the source or location where the tumor tissue was obtained for analysis.

    - Example: A specific hospital or research center where the tumor biopsy or specimen was collected.

# Codes

## B.1 Data analysis of clinical baseline data

[4][1]

```r
setwd("/Users/souravghoshhansda/Library/CloudStorage/
    OneDrive-UniversityofEssex/Dissertation")
library("TCGAbiolinks")
library("glmnet")
library("factoextra")
library("FactoMineR")
library("caret")
library("SummarizedExperiment")
library("ggplot2")
library("RColorBrewer")
library("gProfileR")
library("GenomeInfoDbData")
library("keras")
library("tensorflow")
library("dplyr")
library("DT")
library("MultiAssayExperiment")
```

```
17 library("maftools")
18 library("ComplexHeatmap")
19 library("ggplot2")
20 library("psych")
21 library("corrplot")
22 library("survival")
23 library("glmnet")
24 library("survminer")
25 library("writexl")
26 # get a list of projects
27 TCGAbiolinks:::getProjectSummary("TCGA-BRCA")
28 query1<- GDCquery(
29   project = "TCGA-BRCA",
30   data.category = "Clinical",
31   data.type = "Clinical Supplement",
32   data.format = "BCR Biotab"
33 )
34 output_query1<-getResults(query1)
35 GDCdownload(query1)
36 clinical.BCRtab.all <- GDCprepare(query1,
      summarizedExperiment = TRUE)
37 clinical.BCRtab.all <- na.omit(clinical.BCRtab.all)
38 names(clinical.BCRtab.all)
39 #rm(patient_stage)
40 #patient
41 patient <- rbind(clinical.BCRtab.all$clinical_patient_brca
      )
42 patient<-patient[-2, ]
43 patient<-patient[-1, ]
44 #checking missing values
45 values_to_check <- c("[Not Available]", "[Not Applicable]"
      , "[Not Evaluated]","[Unknown]","Indeterminate","
      Equivocal","[Discrepancy]","not amplified")
46 occurrences <- sapply(values_to_check, function(value) sum
```

```
     (patient == value, na.rm = TRUE))
47 # Display the results
48 result_table <- data.frame(Value = values_to_check,
      Occurrences = occurrences)
49 print(result_table)
50 values_to_check <- c("[Not Available]", "[Not Applicable]"
      , "[Not Evaluated]","[Unknown]","Indeterminate","
      Equivocal","[Discrepancy]","not amplified")
51 calculate_percentage <- function(column) {
52   total_rows <- length(column)
53   count <- sum(column %in% values_to_check)
54   percentage <- round((count / total_rows) * 100, 0)
55   return(percentage)
56 }
57 percentage_results <- sapply(patient, calculate_percentage
      )
58 result_table <- data.frame(Column = names(percentage_
      results), Percentage = as.integer(percentage_results))
59 print(result_table)
60 #readr::write_csv(result_table, file = "/Users/
      souravghoshhansda/Library/CloudStorage/OneDrive-
      UniversityofEssex/Dissertation/complete/missing_values_
      percentage.csv")
61 patient <- na.omit(patient)
62 dim(patient)
63 patient <- patient[patient$gender != "MALE", ]
64 patient <- patient %>%
65   mutate(
66     grouped_stage = case_when(
67       ajcc_pathologic_tumor_stage %in% c('Stage I', 'Stage
            IA', 'Stage IB') ~ 'Stage I',
68       ajcc_pathologic_tumor_stage %in% c('Stage II','Stage
            IIA', 'Stage IIB') ~ 'Stage II',
69       ajcc_pathologic_tumor_stage %in% c('Stage III','
```

```
            Stage IIIA', 'Stage IIIB', 'Stage IIIC') ~ 'Stage
             III',
70        ajcc_pathologic_tumor_stage == 'Stage IV' ~ 'Stage
             IV',
71        TRUE ~ ajcc_pathologic_tumor_stage
72      )
73    )
74 any(colnames(patient) %in% c("vital_status","last_contact_
     days_to","death_days_to"))
75 which(colnames(patient) %in% c("vital_status","last_
     contact_days_to","death_days_to"))
76 patient[,c(14,15,16)]
77 patient$deceased<-ifelse(patient$vital_status=="Alive",
     FALSE,TRUE)
78 table(patient$deceased)
79 table(patient$vital_status)
80 patient$overall_survival<-ifelse(patient$vital_status=="
     Alive",patient$last_contact_days_to,  patient$death_
     days_to)
81 sum(is.na(patient$overall_survival))
82 class(patient$overall_survival)
83 max(patient$overall_survival)
84 table(patient$ajcc_pathologic_tumor_stage)
85 patient$cancer_staging <- ifelse(
86    patient$ajcc_pathologic_tumor_stage %in% c("Stage II",
          "Stage IIA", "Stage IIB", "Stage III", "Stage IIIA
          ", "Stage IIIB", "Stage IIIC", "Stage IV"),
87    "Invasive",
88    ifelse(patient$ajcc_pathologic_tumor_stage %in% c("
          Stage I", "Stage IA", "Stage IB"),
89         "Non-Invasive",
90         NA)
91    )
92 table(patient$cancer_staging)
```

```r
93  patient$binary_cancer_staging <- ifelse(patient$cancer_
        staging == "Invasive", 1, 0)
94  table(patient$binary_cancer_staging)
95  #readr::write_csv(patient, file = "/Users/
        souravghoshhansda/Library/CloudStorage/OneDrive-
        UniversityofEssex/Dissertation/complete/patient_female.
        csv")
96  table(patient$grouped_stage)
97  #grouped_stage_patient_female
98  ggplot(data = patient[!(patient$grouped_stage %in% c("[Not
        Available]", "[Discrepancy]","Stage X")), ], aes(x =
        grouped_stage)) +
99    geom_bar(fill = "blue") +
100   labs(title = "Distribution of ajcc pathological tumor
        stages", x = "Stages", y = "Frequency") +
101   theme_minimal() +
102   theme(axis.text.x = element_text(angle = 45, hjust = 1))
103 #Race_patient_female
104 table(patient$race)
105 ggplot(data = patient[!(patient$race %in% c("[Not
       Available]", "[Not Evaluated]")), ],
106      aes(x = race)) +
107   geom_bar(fill = "lightblue") +
108   labs(title = "Bar Plot of Ethnicity", x = "Ethnicity", y
        = "Frequency") +
109   theme_minimal() +
110   theme(axis.text.x = element_text(angle = 45, hjust = 1))
111 #histology_grouped_patient_female
112 table(patient$histological_type)
113 patient <- patient %>%
114   mutate(
115     histology_grouped = case_when(
116       histological_type %in% c('Infiltrating Carcinoma NOS
            ', 'Infiltrating Ductal Carcinoma', 'Infiltrating
```

```r
             Lobular Carcinoma') ~ 'Infiltrating Carcinoma',
117       TRUE ~ histological_type
118     )
119   )
120 table(patient$histology_grouped)
121 ggplot(data = patient[!(patient$histology_grouped %in% c("
      [Not Available]", "Mixed Histology (please specify)","
      Other, specify")), ]
122 , aes(x = histology_grouped)) +
123   geom_bar(fill = "lightblue") +
124   labs(title = "Bar Plot of Histology", x = "Histology
        Type", y = "Frequency") +
125   theme_minimal() +
126   theme(axis.text.x = element_text(angle = 45, hjust = 1))
127 #age_at_diagnosis_patient_female
128 class(patient$age_at_diagnosis)
129 patient$age_at_diagnosis <- as.numeric(patient$age_at_
      diagnosis)
130 ggplot(patient, aes(x = age_at_diagnosis)) +
131   geom_histogram(binwidth = 5, fill = "blue", color = "
        black") +
132   labs(title = "Age at Diagnosis Histogram", x = "Age", y
        = "Frequency")
133 female_selected <- c("bcr_patient_uuid", "bcr_patient_
      barcode", "grouped_stage")
134 patient_female_filtered <- patient[, female_selected]
135 #drug
136 drug<-rbind(clinical.BCRtab.all$clinical_drug_brca)
137 drug<-drug[-2, ]
138 drug<-drug[-1, ]
139 drug <- na.omit(drug)
140 names(drug)
141 drug_selected<-c("bcr_patient_uuid", "bcr_patient_barcode"
      , "pharmaceutical_therapy_type")
```

```r
142 drug_filtered <- drug[, drug_selected]
143 #combination_patient_drug
144 patient_drug <- inner_join(patient_female_filtered, drug_
        filtered, by = c("bcr_patient_uuid", "bcr_patient_
        barcode"))
145 table(patient_drug$grouped_stage)
146 names(patient_drug)
147 dim(patient_drug)
148 table(patient_drug$pharmaceutical_therapy_type)
149 condition <- !(patient_drug$pharmaceutical_therapy_type %
        in% c("[Not Available]", "[Discrepancy]","Other,
        specify in notes","Chemotherapy|Hormone Therapy")) &
150   !(patient_drug$grouped_stage %in% c("[Not Available]", "
        [Discrepancy]","Stage X"))
151 filtered_patient_drug <- patient_drug[condition, ]
152 my_table<-table(filtered_patient_drug$grouped_stage,
        filtered_patient_drug$pharmaceutical_therapy_type)
153 #csv_file_path <-"/Users/souravghoshhansda/Library/
        CloudStorage/OneDrive-UniversityofEssex/Dissertation/
        complete/therapy_stage.csv"
154 #write.table(my_table, file = csv_file_path, sep = ",",
        col.names = NA, quote = FALSE)
155 cross_table<-as.data.frame(table(filtered_patient_drug$
        grouped_stage, filtered_patient_drug$pharmaceutical_
        therapy_type))
156 names(cross_table) <- c("Stages", "Pharmaceutical_Therapy_
        Types", "Frequency")
157 ggplot(cross_table, aes(x = Pharmaceutical_Therapy_Types,
        y = Frequency, fill = Stages)) +
158   geom_bar(stat = "identity", position = "stack") +
159   labs(title = "Bar Plot of Therapy Type by Stages",
         x = "Therapy Types", y = "Frequency") +
161   theme_minimal() +
162   theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

```r
163  #survival analysis
164  max(patient$overall_survival)
165  sum(is.na(patient$overall_survival))
166  patient$overall_survival <- as.numeric(patient$overall_
         survival)
167  survival_formula<-Surv(patient$overall_survival, patient$
         deceased)
168  #vital_status_survival
169  fit1 = survfit(survival_formula ~ vital_status, data=
         patient)
170  print(fit1)
171  pval1= surv_pvalue(fit1, data=patient)$pval
172  print(pval1)
173  ggsurvplot(fit1, data=patient, pval=T, risk.table=T, risk.
         table.col="strata", xlab="Days")
174  #tumor_stage_survival
175  patient$grouped_stage_filtered <- factor(patient$grouped_
         stage, levels = c("Stage I", "Stage II", "Stage III", "
         Stage IV"))
176  fit2 = survfit(survival_formula ~ grouped_stage_filtered,
         data=patient)
177  print(fit2)
178  ggsurvplot(fit2, data=patient, pval=T, risk.table=T, risk.
         table.col="strata",xlab="Days")
179  pval2= surv_pvalue(fit2, data=patient)$pval
180  print(pval2)
181  #er_status_by_ihc_survival
182  patient$er_status_filtered <- factor(patient$er_status_by_
         ihc, exclude = c("[Not Evaluated]", "Indeterminate"))
183  fit4 = survfit(survival_formula~ er_status_filtered, data=
         patient)
184  print(fit4)
185  ggsurvplot(fit4, data=patient, pval=T, risk.table=T, risk.
         table.col="strata",xlab="Days")
```

```r
186 pval4= surv_pvalue(fit4, data=patient)$pval
187 print(pval4)
188 #her2_status_by_ihc_survival
189 patient$her2_status_filtered <- factor(patient$her2_status
        _by_ihc, exclude = c("[Not Available]", "[Not Evaluated
        ]", "Equivocal", "Indeterminate"))
190 fit5 = survfit(survival_formula ~ her2_status_filtered,
        data=patient)
191 print(fit5)
192 ggsurvplot(fit5, data=patient, pval=T, risk.table=T, risk.
        table.col="strata",xlab="Days")
193 pval5= surv_pvalue(fit5, data=patient)$pval
194 print(pval5)
195 #summary statistics
196 predictors_response<-patient[, c("age_at_diagnosis", "
        menopause_status","margin_status","lymph_nodes_examined
        _count","er_status_by_ihc","pr_status_by_ihc","history_
        other_malignancy","lymph_nodes_examined","her2_status_
        by_ihc")]
197 describe(predictors_response)
198 #write.csv(describe(predictors_response), file = "/Users/
        souravghoshhansda/Library/CloudStorage/OneDrive-
        UniversityofEssex/Dissertation/complete/summary_
        statistics.csv")
199 #logistic regression
200 patient$lymph_nodes_examined_count<-as.numeric(patient$
        lymph_nodes_examined_count)
201 set.seed(123)
202 trainingRowIndex <- sample(1:nrow(patient), 0.8 * nrow(
        patient))
203 trainingData <- patient[trainingRowIndex, ]
204 testData <- patient[-trainingRowIndex, ] %>% na.omit()  #
        Remove NAs for test data
205 # Fit the initial logistic regression model
```

```r
206 model1 <- glm(binary_cancer_staging ~ age_at_diagnosis +
207                 menopause_status +
208                 margin_status +
209                 lymph_nodes_examined_count +
210                 er_status_by_ihc +
211                 pr_status_by_ihc+
212                 history_other_malignancy+
213                 lymph_nodes_examined +
214                 her2_status_by_ihc, data = trainingData,
                    family = "binomial")
215 summary(model1)
216 # backward variable selection
217 backward_model <- step(model1, direction = "backward",
       trace = 0)
218 summary(backward_model)
219 #output_text <- capture.output(summary(backward_model))
220 #csv_file_path <- "/Users/souravghoshhansda/Library/
       CloudStorage/OneDrive-UniversityofEssex/Dissertation/
       complete/backward_model_summary_text.txt"
221 #writeLines(output_text, con = csv_file_path)
222 % Code sources:
223 % - ChatGPT: https://www.chatbot.com/
224 %TCGABOILINKS:https://bioconductor.org/packages/devel/bioc
       /vignettes/TCGAbiolinks/inst/doc/download_prepare.html
```

## B.2 Differential Expression Analysis on Gene Data

[14]

```r
setwd("/Users/souravghoshhansda/Library/CloudStorage/
    OneDrive-UniversityofEssex/Dissertation/test56")
#Load required Libraries
library("TCGAbiolinks")
library("limma")
library("edgeR")
library("glmnet")
library("factoextra")
library("FactoMineR")
library("caret")
library("SummarizedExperiment")
library("gplots")
library("survival")
library("survminer")
library("RColorBrewer")
library("gProfileR")
library("gprofiler2")
library("genefilter")
library(neuralnet)
library("dplyr")
library("tidyr")
library("tibble")
library("data.table")
query_TCGA = GDCquery(
  project = "TCGA-BRCA",
  data.category = "Transcriptome Profiling", # parameter
      enforced by GDCquery
  experimental.strategy = "RNA-Seq",
  workflow.type = "STAR - Counts",
  sample.type = c("Primary Tumor","Solid Tissue Normal"),
  data.type = "Gene Expression Quantification",
```

```r
30    access = "open")
31 brca_res = getResults(query_TCGA) # make results as table
32 # head(brca_res) # data of the first 6 patients.
33 colnames(brca_res) # columns present in the table
34 head(brca_res$sample_type)
35 summary(factor(brca_res$sample_type))
36 GDCdownload(query = query_TCGA)
37 tcga_data = GDCprepare(query_TCGA)
38 dim(tcga_data)
39 colnames(colData(tcga_data))
40 table(tcga_data@colData$vital_status)
41 table(tcga_data@colData$gender)
42 table(tcga_data@colData$race)
43 table(tcga_data$definition)
44 dim(assay(tcga_data))
45 head(assay(tcga_data)[,1:10])
46 head(rowData(tcga_data))
47 saveRDS(object = tcga_data, file = "tcga_data.RDS",
       compress = FALSE)
48 # To load the previous data
49 tcga_data = readRDS(file = "tcga_data.RDS")
50 clinical_data = colData(tcga_data)
51 any(colnames(clinical_data) %in% c("vital_status","days_to
       _last_follow_up","days_to_death"))
52 which(colnames(clinical_data) %in% c("vital_status","days_
       to_last_follow_up","days_to_death"))
53 clinical_data[,c(28,50,63)]
54 clinical_data$deceased<-ifelse(clinical_data$vital_status
       =="Alive",FALSE,TRUE)
55 table(clinical_data$deceased)
56 table(clinical_data$vital_status)
57 clinical_data$overall_survival<ifelse(clinical_data$vital_
       status=="Alive",clinical_data$days_to_last_follow_up,
58 clinical_data$days_to_death)
```

```r
sum(is.na(clinical_data$overall_survival))
class(clinical_data$overall_survival)
max(clinical_data$overall_survival)
table(clinical_data$ajcc_pathologic_stage)
clinical_data$cancer_staging <- ifelse(
  clinical_data$ajcc_pathologic_stage %in% c("Stage II","
      Stage IIA","Stage IIB","Stage III", "Stage IIIA", "
      Stage IIIB", "Stage IIIC", "Stage IV"),
  "Invasive",
  ifelse(clinical_data$ajcc_pathologic_stage %in% c("Stage
      I", "Stage IA", "Stage IB"),
          "Non-Invasive",
          NA)
)
table(clinical_data$cancer_staging)
clinical_data$binary_cancer_staging <- ifelse(clinical_
    data$cancer_staging == "Invasive", 1, 0)
table(clinical_data$binary_cancer_staging)
names(clinical_data)
group = factor(clinical_data$definition)
group = relevel(group, ref="Solid Tissue Normal")
design = model.matrix(~group)
head(design)
dge = DGEList( # creating a DGEList object
  counts=assay(tcga_data),
  samples=colData(tcga_data),
  genes=as.data.frame(rowData(tcga_data)))
# filtering
keep = filterByExpr(dge,design) # defining which genes to
    keep
dge = dge[keep,,keep.lib.sizes=FALSE] # filtering the dge
    object
rm(keep) #  use rm() to remove objects from memory if you
    don't need them anymore
```

```r
86 dge = calcNormFactors(dge,method="TMM")
87 v = voom(dge,design,plot=TRUE)
88 fit = lmFit(v, design)
89 fit = eBayes(fit)
90 topGenes = topTable(fit, coef=1, sort.by="p")
91 print(topGenes)
92 limma_pipeline = function(tcga_data, condition_variable,
     reference_group = NULL) {
93 design_factor = colData(tcga_data)[, condition_variable,
     drop = TRUE]
94   group = factor(design_factor)
95   if (!is.null(reference_group)) {
96     group = relevel(group, ref = reference_group)
97   }
98   design = model.matrix(~ group)
99   dge = DGEList(counts = assay(tcga_data),
100                samples = colData(tcga_data),
101                genes = as.data.frame(rowData(tcga_data)))
102   # Filtering
103   keep = filterByExpr(dge, design)
104   dge = dge[keep, , keep.lib.sizes = FALSE]
105   rm(keep)
106   # Normalization (TMM followed by voom)
107   dge = calcNormFactors(dge)
108   v = voom(dge, design = design, plot = TRUE)  # Pass the
       design matrix directly
109   # Fit model to data given design
110   fit = lmFit(v, design)
111   fit = eBayes(fit)
112   # Show top genes
113   topGenes = topTable(fit, coef = ncol(design), number =
       60, sort.by = "p")
114   return(
115     list(
```

```r
      voomObj = v,    # normalized data
      fit = fit,      # fitted model and statistics
      topGenes = topGenes  # the 22000 most differentially
          expressed genes
    )
  )
}
limma_res = limma_pipeline(
  tcga_data = tcga_data,
  condition_variable = "definition",
  reference_group = "Solid Tissue Normal"
)
saveRDS(object = limma_res,
        file = "limma_res.RDS",
        compress = FALSE)
limma_res = readRDS(file = "limma_res.RDS")
plot_PCA = function(voomObj, condition_variable){
  group = factor(voomObj$targets[, condition_variable])
  pca = prcomp(t(voomObj$E))
  # Take PC1 and PC2 for the plot
  plot(pca$x[,1:2],col=group, pch=19)
  # include a legend for points
  legend("bottomright", inset=.01, levels(group), pch=19,
      col=1:length(levels(group)))
  return(pca)
}
res_pca = plot_PCA(limma_res$voomObj, "definition")
head(limma_res$topGenes)
geneid <- limma_res$topGenes$gene_id
genename <- limma_res$voomObj$genes[geneid, "gene_name"]
genecounts <- assay(tcga_data)[geneid, ]
length(geneid)
length(genename)
GenesData <- data.frame(
```

```r
148   gene_id = geneid,
149   gene_name = genename,
150   counts = as.vector(genecounts),
151   case_id = colnames(genecounts)
152   )
153 #readr::write_csv(GenesData, file = "/Users/
       souravghoshhansda/Library/CloudStorage/OneDrive-
       UniversityofEssex/Dissertation/complete/GenesData.csv")
154 clinical_data <- as.data.frame(clinical_data)
155 combined_data_final <- GenesData %>%
156   merge(clinical_data, by.x = 'case_id', by.y = 'barcode')
157 names(combined_data_final)
158 #readr::write_csv(combined_data_final, file = "/Users/
       souravghoshhansda/Library/CloudStorage/OneDrive-
       UniversityofEssex/Dissertation/complete/combined_data_
       final.csv")
159 combined_data_final_selected <- c("overall_survival", "
       cancer_staging", "submitter_id", "binary_cancer_staging
       ","gene_name","gene_id","counts")
160 selected_combined_data_final <- combined_data_final[,
       combined_data_final_selected]
161 sum(is.na(selected_combined_data_final))
162 max(selected_combined_data_final$overall_survival)
163 selected_combined_data_final<-na.omit(selected_combined_
       data_final)
164 selected_combined_data_final$overall_survival<-selected_
       combined_data_final$overall_survival/8605
165 #readr::write_csv(selected_combined_data_final, file = "/
       Users/souravghoshhansda/Library/CloudStorage/OneDrive-
       UniversityofEssex/Dissertation/complete/selected_
       combined_data_final_60.csv"
166 gene_expression_data<-as.data.frame(assay(tcga_data))
167 gene_metadata <- as.data.frame(rowData(tcga_data))
168 combined_data <- gene_expression_data %>%
```

```r
169    as.data.frame() %>%
170    rownames_to_column(var = 'gene_id') %>%
171    gather(key = 'case_id', value = 'counts', -gene_id) %>%
172    left_join(., gene_metadata, by = "gene_id")
173 combined_data$case_id <- gsub('-01.*', '', combined_data$
       case_id)
174 combined_data_selected <- c("case_id", "counts", "gene_
       name")
175 selected_combined_data <- combined_data[, combined_data_
       selected]
176 clinical_data_selected <- c("overall_survival", "cancer_
       staging", "submitter_id", "binary_cancer_staging")
177 selected_clinical_data <- clinical_data[, clinical_data_
       selected]
178 combined_data_final_60000 <- selected_combined_data %>%
179    merge(selected_clinical_data, by.x = 'case_id', by.y = '
         submitter_id')
180 sum(is.na(combined_data_final_60000$overall_survival))
181 sum(is.na(combined_data_final_60000))
182 max(combined_data_final_60000$overall_survival)
183 combined_data_final_60000<-na.omit(combined_data_final_
       60000)
184 combined_data_final_60000$overall_survival<-combined_data_
       final_60000$overall_survival/8605
185 head(combined_data_final_60000)
186 #readr::write_csv(combined_data_final_60000, file = "/
       Users/souravghoshhansda/Library/CloudStorage/OneDrive-
       UniversityofEssex/Dissertation/complete/combined_data_
       final_60000.csv")
187 gene_expression_data<-as.data.frame(assay(tcga_data))
188 gene_metadata <- as.data.frame(rowData(tcga_data))
189 combined_data <- gene_expression_data %>%
190    as.data.frame() %>%
191    rownames_to_column(var = 'gene_id') %>%
```

```
192  gather(key = 'case_id', value = 'counts', -gene_id) %>%
193    left_join(., gene_metadata, by = "gene_id")
194 combined_data$case_id <- gsub('-01.*', '', combined_data$
         case_id)
195 combined_data_selected <- c("case_id", "counts", "gene_
         name")
196 selected_combined_data <- combined_data[, combined_data_
         selected]
197 clinical_data_selected <- c("overall_survival", "cancer_
         staging", "submitter_id", "binary_cancer_staging")
198 selected_clinical_data <- clinical_data[, clinical_data_
         selected]
199 combined_data_final_60000 <- selected_combined_data %>%
200    merge(selected_clinical_data, by.x = 'case_id', by.y = '
         submitter_id')
201 sum(is.na(combined_data_final_60000$overall_survival))
202 sum(is.na(combined_data_final_60000))
203 max(combined_data_final_60000$overall_survival)
204 combined_data_final_60000<-na.omit(combined_data_final_
         60000)
205 combined_data_final_60000$overall_survival<-combined_data_
         final_60000$overall_survival/8605
206 head(combined_data_final_60000)
207 #readr::write_csv(combined_data_final_60000, file = "/
         Users/souravghoshhansda/Library/CloudStorage/OneDrive-
         UniversityofEssex/Dissertation/complete/combined_data_
         final_60000.csv")
208 % Code sources:
209 % - ChatGPT: https://www.chatbot.com/
210 % - :https://www.costalab.org/wp-content/uploads/2020/11/R
         _class_D3.html#3_RNASeq_Normalization
```

[14]

## B.3   Gene Selection by Random Sampling

[4]

```python
import pandas as pd
import numpy as np
import seaborn as sns
import tensorflow as tf
from tensorflow import keras
from matplotlib
import pyplot as plt
%matplotlib inline
import warningswarnings.filterwarnings("ignore")
df= pd.read_csv("/content/combined_data_final_60000.csv")
df.head(2)
df.shape
gene_counts = df.groupby('gene_name')['case_id'].nunique()
    .reset_index()gene_counts.columns = ['Gene_Name', 'Case
    _id_Count']

gene_counts = gene_counts.sort_values(by='Case_id_Count',
    ascending=False)

gene_counts

 df_shuffled = gene_counts.sample(frac=1, random_state=42)

 only_60_genes = df_shuffled.iloc[:60]

 gene_names_list = only_60_genes['Gene_Name'].tolist()

 gene_names_list

 filtered_df = df[df['gene_name'].isin(gene_names_list)]
```

```
29
30  filtered_df
31
32  unique_gene_names = filtered_df['gene_name'].nunique()
33  print("Number of unique gene names:", unique_gene_names)
34
35  filtered_df.to_csv('/content/filtered_df.csv', index=
        False)
36
37  import os
38  print(os.getcwd())
39
40  gene_counts = filtered_df.groupby('gene_name')['case_id'
        ].nunique().reset_index()
41  gene_counts.columns = ['Gene_Name', 'Case_id_Count']
42  gene_counts = gene_counts.sort_values(by='Case_id_Count',
        ascending=False)
43
44  gene_counts.shape
45
46  df2 = filtered_df[['counts','gene_name']]
47  df2.head(2)
48  % Code sources:
49  % - ChatGPT: https://www.chatbot.com/
```

## B.4 Supervised learning/predictive modelling:Neural Networks

[4][12][15][18]

```
1  import pandas as pd
2  import numpy as np
3  import seaborn as sns
```

```python
4  import tensorflow as tf
5  from tensorflow import keras
6  from matplotlib import pyplot as plt
7  %matplotlib inline
8
9  df = pd.read_csv("/Users/souravghoshhansda/Library/
       CloudStorage/OneDrive-UniversityofEssex/Dissertation/
       complete/selected_combined_data_final_60.csv")
10
11 df.head(2)
12
13 sns.countplot(x = df.cancer_staging)
14
15 from imblearn.over_sampling import SMOTE #generating
       synthetic samples of the minority class
16 #from imblearn.under_sampling import RandomUnderSampler
17 from imblearn.pipeline import Pipeline
18 from collections import Counter
19
20 df.columns
21
22 df.info()
23
24 df.isnull().sum()
25
26 X = df[['counts']]
27 y1 = pd.get_dummies(df[['cancer_staging']], drop_first =
       True)
28
29 y1.value_counts()
30
31 from keras.models import Sequential
32 from keras.layers import Dense
33
```

```python
34 classifier = Sequential()
35 classifier.add(Dense(10, activation = 'relu', input_dim =
      30))
36 classifier.add(Dense(10, activation = 'relu'))
37 classifier.add(Dense(1,  activation = 'sigmoid'))
38 classifier.compile(optimizer = 'adam', loss = 'binary_
      crossentropy', metrics = ['accuracy'])
39 classifier.summary()
40
41 from sklearn.model_selection import train_test_split
42 X_train, X_test, y_train, y_test = train_test_split(X, y1,
       test_size = 0.25, random_state = 0)
43
44 # Create the ANN model
45 model = Sequential()
46 model.add(Dense(10, input_dim=X.shape[1], activation='relu
      '))  # Input layer with 10 neurons
47 model.add(Dense(8, activation='relu'))  # Hidden layer
      with 8 neurons
48 model.add(Dense(y1.shape[1], activation='softmax'))  #
      Output layer with number of classes as neurons
49
50 # Compile the model
51 model.compile(loss='categorical_crossentropy', optimizer='
      adam', metrics=['accuracy'])
52
53 # Train the model
54 model.fit(X_train, y_train, epochs=1, batch_size=5,
      verbose=1, validation_split=0.1)
55
56 # Evaluate the model on test data
57 loss, accuracy = model.evaluate(X_test, y_test)
58 print(f'Test Accuracy: {accuracy:.4f}')
59
```

```python
60 from sklearn import metrics
61 from sklearn.metrics import precision_score,recall_score
62
63 y_pred = model.predict(X_test)
64
65 precision_score(y_test, y_pred)
66
67 recall_score(y_test, y_pred)
68
69 def make_confusion_matrix(y_actual,y_predict,title):
70     fig, ax = plt.subplots(1, 1)
71
72     cm = confusion_matrix(y_actual, y_predict, labels
          =[0,1])
73     disp = ConfusionMatrixDisplay(confusion_matrix=cm,
74                             display_labels=["No","Yes"
                                    ])
75     disp.plot(cmap='Greens',colorbar=True,ax=ax)
76     ax.set_title(title)
77     plt.tick_params(axis=u'both', which=u'both',length=0)
78     #plt.grid(b=None,axis='both',which='both',visible=
          False)
79     plt.show()
80
81 from sklearn.metrics import confusion_matrix,
      classification_report
82 from sklearn.metrics import accuracy_score,precision_score
      ,recall_score,f1_score
83 from statsmodels.stats.outliers_influence import variance_
      inflation_factor
84 from sklearn import metrics
85 #AUC ROC curve
86 from sklearn.metrics import roc_auc_score
87 from sklearn.metrics import roc_curve
```

```python
88  from sklearn.metrics import precision_recall_curve

89

90  from sklearn.metrics import confusion_matrix,
        ConfusionMatrixDisplay

91

92  def get_metrics_score(model,X_train_df,X_test_df,y_train_
        pass,y_test_pass,statsklearn,threshold=0.5,flag=True,
        roc=False):

93

94      # defining an empty list to store train and test
            results

95

96      score_list=[]

97

98      if statsklearn==0:
99          pred_train = model.predict(X_train_df)
100         pred_test = model.predict(X_test_df)
101     else:
102         pred_train = (model.predict(X_train_df)>threshold)
103         pred_test = (model.predict(X_test_df)>threshold)

104

105

106     pred_train = np.round(pred_train)
107     pred_test = np.round(pred_test)

108

109     train_acc = accuracy_score(y_train_pass,pred_train)
110     test_acc = accuracy_score(y_test_pass,pred_test)

111

112     train_recall = recall_score(y_train_pass,pred_train)
113     test_recall = recall_score(y_test_pass,pred_test)

114

115     train_precision = precision_score(y_train_pass,pred_
            train)
116     test_precision = precision_score(y_test_pass,pred_test
```

```python
        )

117
118     train_f1 = f1_score(y_train_pass,pred_train)
119     test_f1 = f1_score(y_test_pass,pred_test)
120     #till here - all metrics are assigned variables
            separately

121

122

123     score_list.extend((train_acc,test_acc,train_recall,
            test_recall,train_precision,test_precision,train_f1
            ,test_f1))

124

125     if flag == True:
126         print("\x1b[0;30;47m \033[1mMODEL PERFORMANCE\x1b
                [0m")

127

128         print("\x1b[0;30;47m \033[1mAccuracy   : Train:\
                x1b[0m",
129                 round(accuracy_score(y_train_pass,pred_train
                    ),3),
130                 "\x1b[0;30;47m \033[1mTest:\x1b[0m ",
131                 round(accuracy_score(y_test_pass,pred_test)
                    ,3))

132

133         print("\x1b[0;30;47m \033[1mRecall     : Train:\
                x1b[0m"
134                 ,round(recall_score(y_train_pass,pred_train)
                    ,3),
135                 "\x1b[0;30;47m \033[1mTest:\x1b[0m" ,
136                 round(recall_score(y_test_pass,pred_test),3)
                    )

137

138         print("\x1b[0;30;47m \033[1mPrecision  : Train:\
                x1b[0m",
```

```python
139            round(precision_score(y_train_pass,pred_
                   train),3),
140            "\x1b[0;30;47m \033[1mTest:\x1b[0m ",
141            round(precision_score(y_test_pass,pred_test)
                   ,3))
142
143        print("\x1b[0;30;47m \033[1mF1         : Train:\
               x1b[0m",
144            round(f1_score(y_train_pass,pred_train),3),
145            "\x1b[0;30;47m \033[1mTest:\x1b[0m",
146            round(f1_score(y_test_pass,pred_test),3))
147
148        make_confusion_matrix(y_train_pass,pred_train,"
               Confusion Matrix for Train")
149        make_confusion_matrix(y_test_pass,pred_test,"
               Confusion Matrix for Test")
150
151    if roc == True:
152
153        print("\x1b[0;30;47m \033[1mROC-AUC Score  :Train
               :\x1b[0m: ",
154            round(roc_auc_score(y_train_pass,pred_train)
                   ,3),
155            "\x1b[0;30;47m \033[1mTest:\x1b[0m: ",
156            round(roc_auc_score(y_test_pass,pred_test)
                   ,3))
157
158    return score_list # returning the list with train and
           test scores
159
160 # # defining empty lists to add train and test results
161 acc_train = []
162 acc_test = []
163 recall_train = []
```

```python
164  recall_test = []
165  precision_train = []
166  precision_test = []
167  f1_train = []
168  f1_test = []
169
170  def add_score_model(score):
171       #'''Add scores to list so that we can compare all
              models score together'''
172      acc_train.append(score[0])
173      acc_test.append(score[1])
174      recall_train.append(score[2])
175      recall_test.append(score[3])
176      precision_train.append(score[4])
177      precision_test.append(score[5])
178      f1_train.append(score[6])
179      f1_test.append(score[7])
180
181  statmodel=1
182  scores_Sklearn = get_metrics_score(model,X_train, X_test,
         y_train, y_test,statmodel)
183
184  optimal_threshold_curve = 0.3
185
186  scores_opt_curve = get_metrics_score(model,X_train, X_test
         , y_train, y_test,statmodel,threshold=optimal_threshold
         _curve,roc=True)
187  add_score_model(scores_opt_curve)
188
189  from tensorflow.keras.models import Sequential
190  from tensorflow.keras.layers import Embedding, SimpleRNN,
         Dense
191  from tensorflow.keras.preprocessing import sequence
192
```

```python
193 y = df[['overall_survival']]
194
195 from sklearn.model_selection import train_test_split
196 X_train1, X_test1, y_train1, y_test1 = train_test_split(X,
        y, test_size = 0.25, random_state = 0)
197
198 # Create the ANN model
199 model = Sequential()
200 model.add(Dense(10, input_dim=X.shape[1], activation='relu
        '))  # Input layer with 10 neurons
201 model.add(Dense(8, activation='relu'))  # Hidden layer
        with 8 neurons
202 model.add(Dense(y.shape[1], activation='softmax'))  #
        Output layer with number of classes as neurons
203
204 # Compile the model
205 model.compile(loss='categorical_crossentropy', optimizer='
        adam', metrics=['mse'])
206
207 # Train the model
208 model.fit(X_train1, y_train1, epochs=1, batch_size=5,
        verbose=1, validation_split=0.1)
209
210 # Evaluate the model on test data
211 loss, accuracy = model.evaluate(X_test1, y_test1)
212 print(f'Test Accuracy: {accuracy:.4f}')
213 % Code sources:
214 % - ChatGPT: https://www.chatbot.com/
215 % - Kanncaa1's deep learning tutorial: https://www.kaggle.
        com/code/kanncaa1/deep-learning-tutorial-for-beginners
216 % - Machine Learning Mastery tutorial: https://
        machinelearningmastery.com/tutorial-first-neural-
        network-python-keras/
217 % - Real Python tutorial: https://realpython.com/python-ai
```

```
    -neural-network/
```

## B.5 Supervised learning/predictive modelling:Random Forest

[4][7][20][17][5]

```python
1  import pandas as pd
2  import numpy as np
3  import seaborn as sns
4  import tensorflow as tf
5  from tensorflow import keras
6  from matplotlib import pyplot as plt
7  %matplotlib inline
8
9  df = pd.read_csv("/Users/souravghoshhansda/Library/
       CloudStorage/OneDrive-UniversityofEssex/Dissertation/
       complete/selected_combined_data_final_60.csv")
10
11 df.head(2)
12
13 df.columns
14
15 df.info()
16
17 df.isnull().sum()
18
19 X = df[['counts']]
20 y = df[['overall_survival']]
21
22 y
23
24 from sklearn.model_selection import train_test_split
```

```
25 X_train, X_test, y_train, y_test = train_test_split(X, y,
       test_size = 0.25, random_state = 0)

26

27 from sklearn.ensemble import RandomForestRegressor
28 RForest = RandomForestRegressor(n_estimators=300)
29 RForest.fit(X_train,y_train)
30 print(RForest.score(X_train,y_train))

31

32 from sklearn.ensemble import RandomForestRegressor
33 from sklearn.model_selection import GridSearchCV, train_
       test_split
34 from sklearn.metrics import mean_squared_error

35

36 rf_regressor = RandomForestRegressor()

37

38 # Define the parameters for grid search
39 parameters = {
40     'n_estimators': [100, 200, 300],  # Number of trees in
           the forest
41     'max_depth': [None, 5, 10, 15],  # Maximum depth of a
           tree
42     'min_samples_split': [2, 5, 10]  # Minimum number of
           samples required to split a node
43 }

44

45 # Perform Grid Search with cross-validation
46 grid_search = GridSearchCV(estimator=rf_regressor, param_
       grid=parameters,
47  scoring='neg_mean_squared_error', cv=5, n_jobs=-1)
48 grid_search.fit(X_train, y_train)

49

50 # Get the best parameters and best score
51 best_params = grid_search.best_params_
52 best_score = grid_search.best_score_
```

```python
53
54 print("Best Parameters:", best_params)
55 print("Best Negative Mean Squared Error:", best_score)
56
57 # Evaluate the model with the best parameters on the test
       set
58 best_model = grid_search.best_estimator_
59 y_pred = best_model.predict(X_test)
60 test_rmse = mean_squared_error(y_test, y_pred, squared=
       False)
61 print(f"Test RMSE with Best Model: {test_rmse:.4f}")
62
63 y1 = pd.get_dummies(df[['cancer_staging']], drop_first =
       True)
64
65 X = df[['counts']]
66
67 from sklearn.model_selection import train_test_split
68 X_train, X_test, y_train, y_test = train_test_split(X, y1,
       test_size = 0.25, random_state = 0)
69
70 from sklearn.ensemble import RandomForestClassifier
71 from sklearn.model_selection import GridSearchCV
72 from sklearn.metrics import accuracy_score
73
74 # Define the Random Forest classifier
75 rf_classifier = RandomForestClassifier()
76
77 # Define the parameters for grid search
78 param_grid = {
79     'n_estimators': [100, 200, 300],
80     'max_depth': [None, 5, 10, 15],
81     'min_samples_split': [2, 5, 10]
82 }
```

```python
83
84 # Perform Grid Search with cross-validation
85 grid_search = GridSearchCV(estimator=rf_classifier, param_
       grid=param_grid,
86                             scoring='accuracy', cv=5, n_
                                 jobs=-1)
87 grid_search.fit(X_train, y_train)
88
89 # Get the best parameters and best score
90 best_params = grid_search.best_params_
91 best_score = grid_search.best_score_
92
93 print("Best Parameters:", best_params)
94 print("Best Accuracy Score:", best_score)
95
96 # Evaluate the model with the best parameters on the test
       set
97 best_model = grid_search.best_estimator_
98 y_pred = best_model.predict(X_test)
99 test_accuracy = accuracy_score(y_test, y_pred)
100 print(f"Test Accuracy with Best Model: {test_accuracy:.4f}
       ")
101
102 from sklearn import metrics
103 from sklearn.metrics import precision_score,recall_score
104
105 precision_score(y_test, y_pred)
106
107 precision_score(y_test, y_pred)
108
109 recall_score(y_test, y_pred)
110
111 recall_score(y_test, y_pred)
112
```

```
113 y_train.value_counts()

114

115 y_train.value_counts()

116

117 from sklearn.metrics import  roc_auc_score,precision_
        recall_curve, auc, roc_curve

118

119 y_score = best_model.predict_proba(X_train)[:,1]

120 fpr, tpr, _ = roc_curve(y_train, y_score)

121

122 plt.title('Random Forest ROC curve:')

123 plt.xlabel('FPR (Precision)')

124 plt.ylabel('TPR (Recall)')

125

126 plt.plot(fpr,tpr)

127 plt.plot((0,1), ls='dashed',color='black')

128 plt.show()

129 print ('Area under curve (AUC): ', auc(fpr,tpr))

130

131 % Code sources:

132 % - GeeksforGeeks Random Forest Classifier tutorial: https
        ://www.geeksforgeeks.org/random-forest-classifier-using
        -scikit-learn/

133 % - Scikit-learn RandomForestClassifier documentation:
        https://scikit-learn.org/stable/modules/generated/
        sklearn.ensemble.RandomForestClassifier.html

134 % - Kaggle Tutorial by prashant111: https://www.kaggle.com
        /code/prashant111/random-forest-classifier-tutorial

135 % - DataCamp:https://www.datacamp.com/tutorial/random-
        forests-classifier-python
```

## B.6   Unsupervised learning/clustering of Gene expression:Elbow method and Silhouette Score

[4]

```python
import pandas as pd
import numpy as np
import seaborn as sns
import tensorflow as tf
from tensorflow import keras
from matplotlib import pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
df2= pd.read_csv("/content/selected_combined_data_final_
    60.csv")
df2.head()
df2.shape
df2 = df2[['counts','gene_name']]
df2.head(2)
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from scipy.cluster.hierarchy import linkage
from scipy.cluster.hierarchy import dendrogram
from scipy.cluster.hierarchy import cut_tree
class MultiColumnLabelEncoder:

    def __init__(self, columns=None):
        self.columns = columns # array of column names to
            encode


    def fit(self, X, y=None):
        self.encoders = {}
```

```python
28          columns = X.columns if self.columns is None else
                self.columns
29          for col in columns:
30              self.encoders[col] = LabelEncoder().fit(X[col
                    ])
31          return self
32
33
34      def transform(self, X):
35          output = X.copy()
36          columns = X.columns if self.columns is None else
                self.columns
37          for col in columns:
38              output[col] = self.encoders[col].transform(X[
                    col])
39          return output
40
41
42      def fit_transform(self, X, y=None):
43          return self.fit(X,y).transform(X)
44
45
46      def inverse_transform(self, X):
47          output = X.copy()
48          columns = X.columns if self.columns is None else
                self.columns
49          for col in columns:
50              output[col] = self.encoders[col].inverse_
                    transform(X[col])
51          return output
52          from sklearn.preprocessing import LabelEncoder
53          multi = MultiColumnLabelEncoder(columns=['gene_
                name'])
54
```

```python
55  df2 = multi.fit_transform(df2)
56  #from sklearn.preprocessing import LabelEncoder
57  #le = LabelEncoder()
58  #df2.gene_name = le.fit_transform(df2.gene_name)
59  df2
60  kmeans = KMeans(n_clusters=5, max_iter=50)
61  kmeans.fit(df2)
62  ssd = []
63  range_n_clusters = [2, 3, 4, 5, 6, 7, 8, 9
        ,10,20,30,40,50,59]
64  for num_clusters in range_n_clusters:
65      kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
66      kmeans.fit(df2)
67
68      ssd.append(kmeans.inertia_)
69
70  # plot the SSDs for each n_clusters
71  plt.plot(ssd)
72  # Silhouette analysis
73  range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
74
75  for num_clusters in range_n_clusters:
76
77      # intialise kmeans
78      kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
79      kmeans.fit(df2)
80
81      cluster_labels = kmeans.labels_
82
83      # silhouette score
84      silhouette_avg = silhouette_score(df2, cluster_labels)
85      print("For n_clusters={0}, the silhouette score is {1}
            ".format(num_clusters, silhouette_avg))
```

# **Outputs**

## C.1 Descriptive Statistics of clinical predictors

| Predictors | age_at_diagnosis | menopause_status | margin_status | lymph_nodes_examined_count | er_status_by_ihc | pr_status_by_ihc | history_other_malignancy | lymph_nodes_examined | her2_status_by_ihc |
|---|---|---|---|---|---|---|---|---|---|
| mean | 58.41751152 | 5.788940092 | 3.858986175 | 16.35576037 | 3.641474654 | 3.539170507 | 2.060829493 | 2.941935484 | 4.317050691 |
| sd | 13.22359052 | 1.354460746 | 0.795179715 | 13.48049967 | 0.710707278 | 0.729950495 | 0.242954701 | 1.414651175 | 1.363376588 |
| median | 58 | 6 | 4 | 14 | 4 | 4 | 2 | 4 | 5 |
| trimmed | 58.19102417 | 6.085155351 | 4 | 15.13463751 | 3.791714614 | 3.66858458 | 2 | 3.051783659 | 4.405063291 |
| mad | 13.3434 | 0 | 0 | 14.826 | 0 | 0 | 0 | 0 | 0 |
| min | 26 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| max | 90 | 7 | 5 | 43 | 4 | 4 | 3 | 4 | 6 |
| range | 64 | 6 | 4 | 42 | 3 | 3 | 2 | 3 | 5 |
| skew | 0.134456532 | -2.44490975 | -2.63591618 | 0.684645338 | -2.408460729 | -1.941511446 | 3.475559351 | -0.62405394 | -0.669572211 |
| kurtosis | -0.506111292 | 5.864472036 | 7.441841308 | -0.857287599 | 5.863199787 | 4.014942872 | 11.12733448 | -1.585187903 | -0.930901952 |
| se | 0.40145283 | 0.041119853 | 0.024140731 | 0.409252293 | 0.021576246 | 0.022160448 | 0.007375822 | 0.042947164 | 0.041390528 |

Figure C.1: Descriptive Statistics

## C.2   Logistic Regression

```
Call:

glm(formula = binary_cancer_staging ~ age_at_diagnosis+

menopause_status + margin_status + lymph_nodes_examined_count

+ er_status_by_ihc + pr_status_by_ihc + history_other_malignancy

+ lymph_nodes_examined + her2_status_by_ihc, family = "binomial",

data = trainingData)
```

Coefficients:

| | Estimate | Std. Error | z value | Pr(>\|z\|) |
|---|---|---|---|---|
| (Intercept) | -0.76 | 1.42 | -0.54 | 0.590 |
| age_at_diagnosis | -0.01 | 0.01 | -1.03 | 0.304 |
| menopause_status[Not Evaluated] | 15.59 | 1278.34 | 0.01 | 0.990 |
| menopause_status[Unknown] | -0.92 | 0.92 | -0.99 | 0.322 |
| menopause_statusIndeterminate | -0.78 | 0.87 | -0.90 | 0.367 |
| menopause_statusPeri | -0.51 | 0.69 | -0.74 | 0.458 |
| menopause_statusPost | -0.25 | 0.50 | -0.51 | 0.611 |
| menopause_statusPre | -0.40 | 0.55 | -0.74 | 0.459 |
| margin_status[Unknown] | 15.54 | 1417.95 | 0.01 | 0.991 |
| margin_statusClose | -0.35 | 0.78 | -0.45 | 0.652 |
| margin_statusNegative | 0.24 | 0.56 | 0.43 | 0.664 |
| margin_statusPositive | 0.34 | 0.70 | 0.48 | 0.628 |
| lymph_nodes_examined_count | 0.17 | 0.02 | 7.93 | 2.26e-15*** |
| er_status_by_ihcNegative | 14.70 | 2399.54 | 0.01 | 0.995 |
| er_status_by_ihcPositive | 14.70 | 2399.54 | 0.01 | 0.995 |
| pr_status_by_ihcIndeterminate | -0.19 | 2675.00 | 0.00 | 0.999 |
| pr_status_by_ihcNegative | -14.62 | 2399.54 | -0.01 | 0.995 |
| pr_status_by_ihcPositive | -14.93 | 2399.54 | -0.01 | 0.995 |
| history_other_malignancyYes | 0.18 | 0.45 | 0.40 | 0.687 |
| lymph_nodes_examinedYES | 0.10 | 0.25 | 0.39 | 0.699 |
| her2_status_by_ihc[Not Evaluated] | 1.41 | 0.92 | 1.54 | 0.125 |

```
her2_status_by_ihcEquivocal        2.01        0.90      2.24      0.025*

her2_status_by_ihcIndeterminate    1.99        1.40      1.42      0.156

her2_status_by_ihcNegative         1.72        0.87      1.98      0.047*

her2_status_by_ihcPositive         1.91        0.91      2.11      0.035*


---

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


(Dispersion parameter for binomial family taken to be 1)


Null deviance: 723.68  on 753  degrees of freedom
Residual deviance: 595.32  on 729  degrees of freedom
(114 observations deleted due to missingness)
AIC: 645.32


Number of Fisher Scoring iterations: 15
```

## C.3   Backward Variable Selection

```
Call:
glm(formula = binary_cancer_staging ~ lymph_nodes_examined_count,
     family = "binomial", data = trainingData)


Coefficients:
                             Estimate Std. Error z value Pr(>|z|)
(Intercept)                   0.20019    0.14993   1.335    0.182
lymph_nodes_examined_count    0.16941    0.02075   8.165 3.21e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


(Dispersion parameter for binomial family taken to be 1)


Null deviance: 723.68  on 753  degrees of freedom
Residual deviance: 612.20  on 752  degrees of freedom
  (114 observations deleted due to missingness)
AIC: 616.2


Number of Fisher Scoring iterations: 6
```

## C.4   Missing Values

Table C.1: Missing Values

| Predictors | Percentage |
| --- | --- |
| birth_days_to | 1 |
| gender | 0 |
| menopause_status | 8 |
| race | 9 |
| ethnicity | 16 |
| history_other_malignancy | 0 |
| history_neoadjuvant_treatment | 0 |
| tumor_status | 11 |
| vital_status | 0 |
| last_contact_days_to | 9 |
| death_days_to | 91 |
| radiation_treatment_adjuvant | 85 |
| pharmaceutical_tx_adjuvant | 85 |
| histologic_diagnosis_other | 93 |
| initial_pathologic_dx_year | 0 |
| age_at_diagnosis | 0 |
| method_initial_path_dx | 9 |
| method_initial_path_dx_other | 94 |
| surgical_procedure_first | 5 |

Table C.2: Missing Values

| Predictors | Percentage |
|---|---|
| first_surgical_procedure_other | 75 |
| margin_status | 6 |
| surgery_for_positive_margins | 95 |
| surgery_for_positive_margins_other | 99 |
| margin_status_reexcision | 95 |
| axillary_staging_method | 20 |
| axillary_staging_method_other | 98 |
| micromet_detection_by_ihc | 33 |
| lymph_nodes_examined | 34 |
| lymph_nodes_examined_count | 11 |
| lymph_nodes_examined_he_count | 15 |
| lymph_nodes_examined_ihc_count | 69 |
| ajcc_staging_edition | 13 |
| ajcc_tumor_pathologic_pt | 0 |
| ajcc_nodes_pathologic_pn | 0 |
| ajcc_metastasis_pathologic_pm | 0 |
| ajcc_pathologic_tumor_stage | 1 |
| metastasis_site | 98 |
| metastasis_site_other | 99 |
| er_status_by_ihc | 5 |
| er_status_ihc_Percent_Positive | 57 |
| er_positivity_scale_used | 87 |
| er_ihc_score | 80 |

Table C.3: Missing Values

| Predictors | Percentage |
|---|---|
| er_positivity_scale_other | 78 |
| er_positivity_method | 80 |
| pr_status_by_ihc | 5 |
| pr_status_ihc_percent_positive | 61 |
| pr_positivity_scale_used | 87 |
| pr_positivity_ihc_intensity_score | 81 |
| pr_positivity_scale_other | 80 |
| pr_positivity_define_method | 81 |
| her2_status_by_ihc | 34 |
| her2_ihc_percent_positive | 81 |
| her2_ihc_score | 43 |
| her2_positivity_scale_other | 99 |
| her2_positivity_method_text | 92 |
| her2_fish_status | 63 |
| her2_copy_number | 90 |
| cent17_copy_number | 90 |
| her2_and_cent17_cells_count | 91 |
| her2_cent17_ratio | 79 |
| her2_and_cent17_scale_other | 100 |
| her2_fish_method | 95 |
| new_tumor_event_dx_indicator | 82 |
| nte_er_status | 99 |
| nte_er_status_ihc__positive | 99 |

Table C.4: Missing Values

| Predictors | Percentage |
|:---:|:---:|
| nte_er_ihc_intensity_score | 100 |
| nte_er_positivity_other_scale | 100 |
| nte_er_positivity_define_method | 100 |
| nte_pr_status_by_ihc | 99 |
| nte_pr_status_ihc__positive | 100 |
| nte_pr_ihc_intensity_score | 100 |
| nte_pr_positivity_other_scale | 100 |
| nte_pr_positivity_define_method | 100 |
| nte_her2_status | 99 |
| nte_her2_status_ihc__positive | 100 |
| nte_her2_positivity_ihc_score | 100 |
| nte_her2_positivity_other_scale | 100 |
| nte_her2_positivity_method | 100 |
| nte_her2_fish_status | 100 |
| nte_her2_signal_number | 100 |
| nte_cent_17_signal_number | 100 |
| her2_cent17_counted_cells_count | 100 |
| nte_cent_17_her2_ratio | 100 |
| nte_cent17_her2_other_scale | 100 |
| nte_her2_fish_define_method | 100 |
| anatomic_neoplasm_subdivision | 0 |
| clinical_M | 100 |

Table C.5: Missing Values

| Predictors | Percentage |
|:---:|:---:|
| clinical_N | 100 |
| clinical_T | 100 |
| clinical_stage | 100 |
| days_to_initial_pathologic_diagnosis | 0 |
| days_to_patient_progression_free | 100 |
| days_to_tumor_progression | 100 |
| disease_code | 100 |
| extranodal_involvement | 100 |
| histological_type | 0 |
| icd_10 | 0 |
| icd_o_3_histology | 0 |
| icd_o_3_site | 0 |
| informed_consent_verified | 0 |
| metastatic_tumor_indicator | 64 |
| patient_id | 0 |
| project_code | 100 |
| site_of_primary_tumor_other | 100 |
| stage_other | 100 |
| tissue_source_site | 0 |
| tumor_tissue_site | 0 |

# Bibliography

[1] TCGAbiolinks package on bioconductor. Accessed on December 18, 2023.

[2] T. Bismeijer, S. Canisius, and L.F.A. Wessels. Molecular characterization of breast and lung tumors by integration of multiple data types with functional sparse-factor analysis. *PLoS Computational Biology*, 14:1–28, 2018.

[3] Rocky Mountain Cancer Centers. What's the difference between invasive and non-invasive breast cancers?, 2023. Accessed on December 8, 2023.

[4] ChatGPT. Chatgpt. https://www.chatbot.com/.

[5] DataCamp. DataCamp Tutorial, Year. Accessed on Date.

[6] Jamie DePolo. Breast cancer stages, 2023. Accessed on December 8, 2023.

[7] GeeksforGeeks. GeeksforGeeks Random Forest Classifier Tutorial, Year. Accessed on Date.

[8] O. Gevaert, F.D. Smet, D. Timmerman, Y. Moreau, and B.D. Moor. Predicting the prognosis of breast cancer by integrating clinical and microarray data with bayesian networks. *Bioinformatics*, 22(14):e184–90, 2006.

[9] L. Giacomelli, R. Sacco, S. Papa, and B.I. Carr. Understanding the drawbacks of the current tumor staging systems: How to improve? *Cancers*, 15(4):1242, 2023.

[10] E. Gobbini et al. Time trends of overall survival among metastatic breast cancer patients in the real-life esme cohort. *European Journal of Cancer*, 96, June 2018.

[11] Z. Huang, X. Zhan, S. Xiang, et al. Salmon: survival analysis learning with multi-omics neural networks on breast cancer. *Frontiers in Genetics*, 10, 2019.

[12] Kanncaa1. Kanncaa1's Deep Learning Tutorial, Year. Accessed on Date.

[13] A.A. Kim, S.R. Zaim, and V. Subbian. Assessing reproducibility and veracity across machine learning techniques in biomedicine: A case study using tcga data. *International Journal of Medical Informatics*, 141, September 2020.

[14] Tiago Maié and Martin Manolov. Analysis of cancer genome atlas in r. Accessed on December 19, 2023.

[15] Machine Learning Mastery. Machine Learning Mastery Tutorial, Year. Accessed on Date.

[16] Siddhartha Mukherjee. *The Emperor of All Maladies: A Biography of Cancer*. Scribner, 2010.

[17] Prashant111. Kaggle Tutorial by prashant111, Year. Accessed on Date.

[18] Real Python. Real Python Tutorial, Year. Accessed on Date.

[19] M.C. Rendleman, J.M. Buatti, T.A. Braun, et al. Machine learning with the tcga-hnsc dataset: Improving usability by addressing inconsistency, sparsity, and high-dimensionality. *BMC Bioinformatics*, 20:1–9, 2019.

[20] Scikit-Learn. Scikit-Learn RandomForestClassifier Documentation, Year. Accessed on Date.

[21] M. Sherafatian. Tree-based machine learning algorithms identified a minimal set of mirna biomarkers for breast cancer diagnosis and molecular subtyping. *Gene*, 677:111–118, 2018.

[22] D. Sun, A. Li, B. Tang, and M. Wang. Integrating genomic data and pathological images to effectively predict breast cancer clinical outcome. *Computer Methods and Programs in Biomedicine*, 161:45–53, 2018.

[23] M. Zhao, Y. Tang, H. Kim, and K. Hasegawa. Machine learning with k-means dimensional reduction for predicting survival outcomes in patients with breast cancer. *Cancer Informatics*, 17:1176935118810215, 2018.