

# Evaluation of Prognostic Determinants and Survival Outcomes in Breast Carcinoma

A Comprehensive Analysis of SEER Data (2006-2010)

Yangzhenyu Gao

2025-12-01

## Contents

<b>1. Introduction &amp; Setup</b>	<b>2</b>
1.1 Load Libraries . . . . .	2
1.2 Data Preparation and Cleaning . . . . .	2
1.3.1 Exploratory Data Analysis - Demographic Summary . . . . .	3
1.3.2 Exploratory Data Analysis - Correlation Analysis . . . . .	4
1.3.3 Exploratory Data Analysis - PCA . . . . .	5
<b>2. Research Question 1: Determinants of Overall Survival</b>	<b>6</b>
2.1 Kaplan-Meier Survival Analysis . . . . .	6
2.2 Multivariate Cox Proportional Hazards Model . . . . .	9
<b>3. Research Question 2: Age and Tumor Biology</b>	<b>12</b>
3.1 Visualization and Statistical Test . . . . .	12
<b>4. Research Question 3: Predictive Modeling (Machine Learning)</b>	<b>14</b>
4.1 Data Splitting . . . . .	14
4.2 Model Training . . . . .	14
4.3 Performance Evaluation (ROC Curves & Variable Importance) . . . . .	15
<b>5. Research Question 4: Demographic Disparities</b>	<b>17</b>
5.1 Visualization of Stage Distribution . . . . .	17
5.2 Statistical Testing (Chi-Square & Odds Ratios) . . . . .	18
<b>6. Conclusion</b>	<b>19</b>

# 1. Introduction & Setup

This analysis investigates a cohort of 4,024 breast cancer patients to answer three key scientific questions regarding socio-demographic disparities, predictors of lymph node metastasis, and age-related biological profiles.

## 1.1 Load Libraries

```
# Load necessary packages
library(tidyverse)
library(survival)
library(survminer)
library(corrplot)
library(caret)
library(randomForest)
library(janitor)
library(pROC)
library(gtsummary)
library(gridExtra)

# Set a seed for reproducibility
set.seed(42)
```

## 1.2 Data Preparation and Cleaning

Loading the dataset from the specified local path.

```
# Define the file path (Using forward slashes for R compatibility)
file_path <- "C:/HISTORY/4_PhD_UNC/HWK/BIOS611_final/data/Breast_Cancer.csv"

# Load Data
raw_df <- read.csv(file_path, stringsAsFactors = FALSE)

# Preprocessing pipeline
df_clean <- raw_df %>%
  rename(
    # Standardize names for easier coding
    Age = Age,
    Race = Race,
    Marital_Status = Marital.Status,
    T_Stage = T.Stage,
    N_Stage = N.Stage,
    Stage_6th = X6th.Stage,
    Differentiation = differentiate,
    Grade = Grade,
    A_Stage = A.Stage,
    Tumor_Size = Tumor.Size,
    Estrogen_Status = Estrogen.Status,
    Progesterone_Status = Progesterone.Status,
    Nodes_Examined = Regional.Node.Examined,
    Nodes_Positive = Reginol.Node.Positive, # Fixed typo in CSV header
  )
```

```

Survival_Months = Survival.Months,
Status = Status
) %>%
mutate(
  # 1. Create Binary Status for ML and Survival Analysis (Dead = Event = 1)
  Status_Binary = ifelse(Status == "Dead", 1, 0),

  # 2. Define "Advanced Stage" for Demographic Analysis (Stage III as Advanced)
  # Note: Dataset contains IIA, IIB, IIIA, IIIB, IIIC. We group III as Advanced.
  Is_Advanced_Stage = ifelse(grepl("III", Stage_6th), "Advanced (Stage III)", "Early/Mid (Stage II)"),
  Is_Advanced_Stage = factor(Is_Advanced_Stage, levels = c("Early/Mid (Stage II)", "Advanced (Stage III)")),

  # 3. Factor Conversion
  Race = as.factor(Race),
  Marital_Status = as.factor(Marital_Status),
  T_Stage = as.factor(T_Stage),
  N_Stage = as.factor(N_Stage),
  Stage_6th = as.factor(Stage_6th),
  Grade = as.factor(Grade),
  Differentiation = as.factor(Differentiation),
  Estrogen_Status = as.factor(Estrogen_Status),
  Progesterone_Status = as.factor(Progesterone_Status),
  Status = factor(Status, levels = c("Alive", "Dead")) # Factor for Classification
)

# Preview the cleaned data structure
glimpse(df_clean)

```

```

## Rows: 4,024
## Columns: 18
## $ Age                <int> 68, 50, 58, 58, 47, 51, 51, 40, 40, 69, 68, 46, 65~
## $ Race                <fct> White, White, White, White, White, White, White, W~
## $ Marital_Status      <fct> Married, Married, Divorced, Married, Married, Sing~
## $ T_Stage             <fct> T1, T2, T3, T1, T2, T1, T1, T2, T4, T4, T1, T3, T2~
## $ N_Stage             <fct> N1, N2, N3, N1, N1, N1, N1, N1, N3, N3, N1, N1, N1~
## $ Stage_6th          <fct> IIA, IIIA, IIIC, IIA, IIB, IIA, IIA, IIB, IIIC, II~
## $ Differentiation     <fct> Poorly differentiated, Moderately differentiated, ~
## $ Grade               <fct> 3, 2, 2, 3, 3, 2, 1, 2, 3, 1, 2, 3, 3, 3, 2, 2, ~
## $ A_Stage             <chr> "Regional", "Regional", "Regional", "Regional", "R~
## $ Tumor_Size          <int> 4, 35, 63, 18, 41, 20, 8, 30, 103, 32, 13, 59, 35, ~
## $ Estrogen_Status     <fct> Positive, Positive, Positive, Positive, Positive, ~
## $ Progesterone_Status <fct> Positive, Positive, Positive, Positive, Positive, ~
## $ Nodes_Examined      <int> 24, 14, 14, 2, 3, 18, 11, 9, 20, 21, 9, 11, 13, 23~
## $ Nodes_Positive      <int> 1, 5, 7, 1, 1, 2, 1, 1, 18, 12, 1, 3, 3, 7, 14, 1, ~
## $ Survival_Months     <int> 60, 62, 75, 84, 50, 89, 54, 14, 70, 92, 64, 92, 56~
## $ Status              <fct> Alive, Alive, Alive, Alive, Alive, Alive, Alive, D~
## $ Status_Binary       <dbl> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, ~
## $ Is_Advanced_Stage   <fct> Early/Mid (Stage II), Advanced (Stage III), Advanc~

```

### 1.3.1 Exploratory Data Analysis - Demographic Summary

Overview of the cohort characteristics.

Characteristic	0 N = 3,408 <sup>1</sup>	1 N = 616 <sup>1</sup>	p-value <sup>2</sup>
Age	54 (47, 61)	57 (48, 63)	<0.001
Race			<0.001
Black	218 (6.4%)	73 (12%)	
Other	287 (8.4%)	33 (5.4%)	
White	2,903 (85%)	510 (83%)	
Marital_Status			<0.001
Divorced	396 (12%)	90 (15%)	
Married	2,285 (67%)	358 (58%)	
Separated	30 (0.9%)	15 (2.4%)	
Single	511 (15%)	104 (17%)	
Widowed	186 (5.5%)	49 (8.0%)	
T_Stage			<0.001
T1	1,446 (42%)	157 (25%)	
T2	1,483 (44%)	303 (49%)	
T3	417 (12%)	116 (19%)	
T4	62 (1.8%)	40 (6.5%)	
Grade			<0.001
1	504 (15%)	39 (6.3%)	
2	2,046 (60%)	305 (50%)	
3	848 (25%)	263 (43%)	
4	10 (0.3%)	9 (1.5%)	

<sup>1</sup>Median (Q1, Q3); n (%)

<sup>2</sup>Wilcoxon rank sum test; Pearson's Chi-squared test; Fisher's exact test

```
df_clean %>%
  select(Age, Race, Marital_Status, T_Stage, Grade, Status_Binary) %>%
  tbl_summary(by = Status_Binary) %>%
  add_p() %>%
  bold_labels()
```

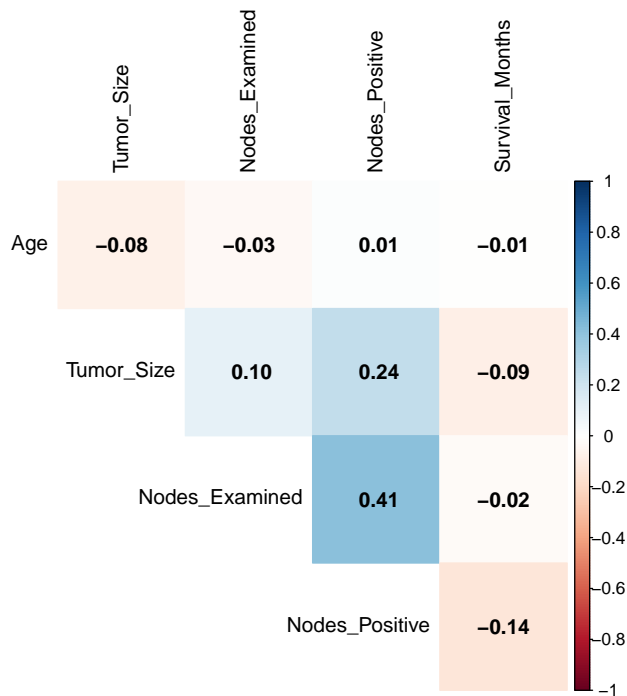
### 1.3.2 Exploratory Data Analysis - Correlation Analysis

Checking for multicollinearity among continuous variables.

```
# Select numeric columns for correlation
numeric_vars <- df_clean %>%
  select(Age, Tumor_Size, Nodes_Examined, Nodes_Positive, Survival_Months)

cor_matrix <- cor(numeric_vars)
corrplot(cor_matrix, method = "color", type = "upper",
  addCoef.col = "black", tl.col = "black", diag = FALSE,
  title = "Correlation Matrix of Continuous Variables", mar=c(0,0,2,0))
```

Correlation Matrix of Continuous Variables



### 1.3.3 Exploratory Data Analysis - PCA

```
# Prepare Matrix for PCA
# We need to remove the outcome variables (Status, Survival Months) from the input
df_pca_input <- df_clean %>%
  select(-Status, -Survival_Months) # Remove outcomes

# Use model.matrix to automatically one-hot encode all categorical variables
pca_matrix <- model.matrix(~ . - 1, data = df_pca_input)

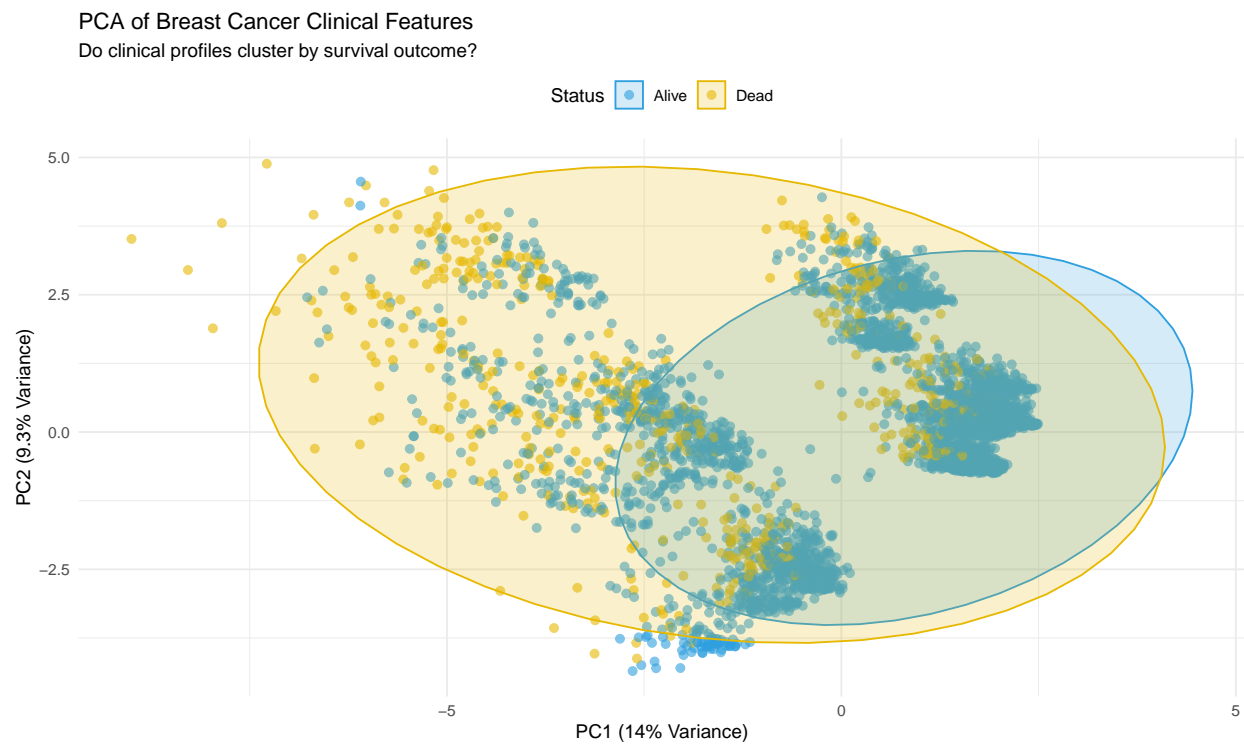
# Perform PCA
# scale. = TRUE is crucial because variables have different units (Age vs Tumor Size)
pca_result <- prcomp(pca_matrix, scale. = TRUE)

# Visualization (PC1 vs PC2)
# We extract the coordinates and add the 'Status' back for coloring
pca_data <- as.data.frame(pca_result$x)
pca_data$Status <- df_clean$Status

# Calculate the percentage of variance explained by PC1 and PC2
var_explained <- round(100 * (pca_result$sdev^2 / sum(pca_result$sdev^2)), 1)

# Plot
ggplot(pca_data, aes(x = PC1, y = PC2, color = Status, fill = Status)) +
  geom_point(alpha = 0.6, size = 2) + # Semi-transparent points
  stat_ellipse(geom = "polygon", alpha = 0.2, level = 0.95) + # Add confidence ellipses
  scale_color_manual(values = c("Alive" = "#2E9FDF", "Dead" = "#E7B800")) + # Custom colors
  scale_fill_manual(values = c("Alive" = "#2E9FDF", "Dead" = "#E7B800")) +
```

```
labs(
  title = "PCA of Breast Cancer Clinical Features",
  subtitle = "Do clinical profiles cluster by survival outcome?",
  x = paste0("PC1 (", var_explained[1], "% Variance)"),
  y = paste0("PC2 (", var_explained[2], "% Variance)")
) +
theme_minimal() +
theme(legend.position = "top")
```



## 2. Research Question 1: Determinants of Overall Survival

**Scientific Question:** Which clinical and pathological factors are significantly associated with overall survival? specifically, do factors like Tumor Stage, Grade, and Hormone Status independently predict mortality hazard?

### 2.1 Kaplan-Meier Survival Analysis

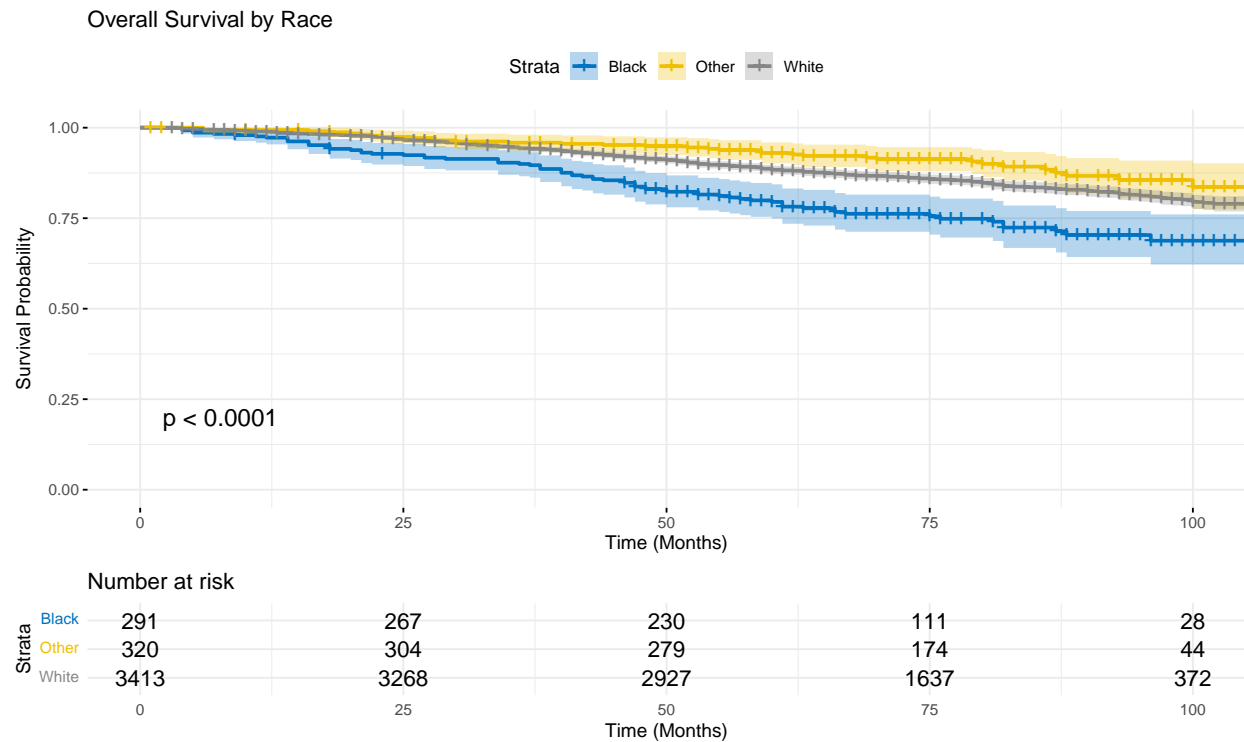
We first visualize univariate survival curves stratified by Cancer Stage.

```
# Survival Object
surv_obj <- Surv(time = df_clean$Survival_Months, event = df_clean$Status_Binary)

# KM Plot by Race
fit_race <- survfit(surv_obj ~ Race, data = df_clean)

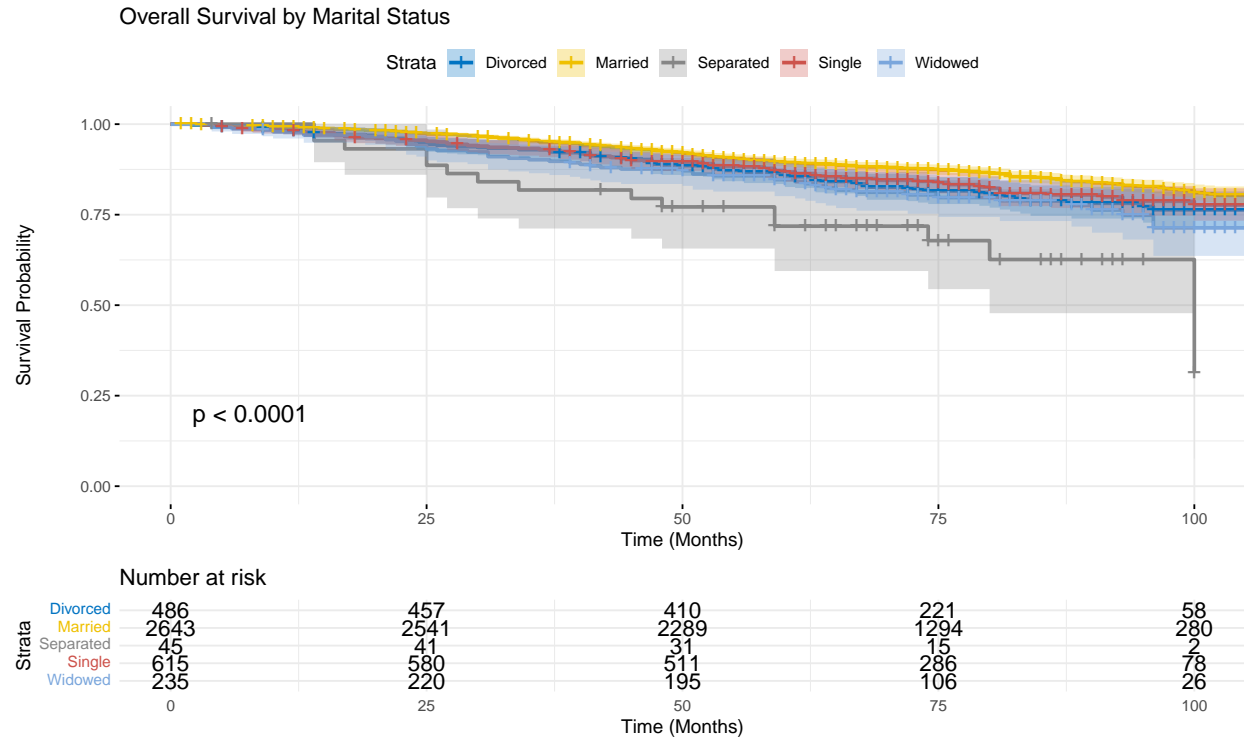
ggsurvplot(fit_race, data = df_clean,
```

```
pval = TRUE, conf.int = TRUE,
risk.table = TRUE, palette = "jco",
ggtheme = theme_minimal(),
title = "Overall Survival by Race",
xlab = "Time (Months)",
ylab = "Survival Probability",
legend.labs = levels(df_clean$Race))
```



```
# KM Plot by Marital Status
fit_marital <- survfit(surv_obj ~ Marital_Status, data = df_clean)

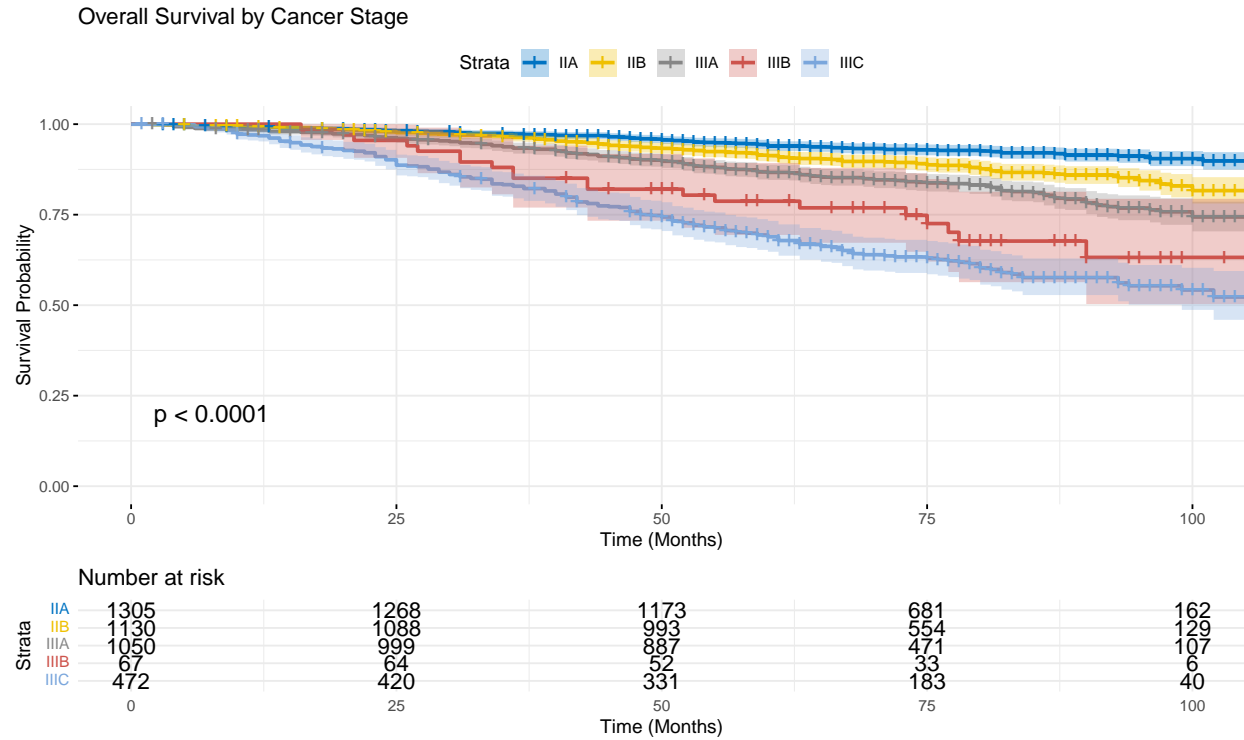
ggsurvplot(fit_marital, data = df_clean,
pval = TRUE, conf.int = TRUE,
risk.table = TRUE, palette = "jco",
ggtheme = theme_minimal(),
title = "Overall Survival by Marital Status",
xlab = "Time (Months)",
ylab = "Survival Probability",
legend.labs = levels(df_clean$Marital_Status))
```



```
# KM Plot by Breast Cancer Stage
fit_stage <- survfit(surv_obj ~ Stage_6th, data = df_clean)

ggsurvplot(fit_stage, data = df_clean,
  pval = TRUE, conf.int = TRUE,
  risk.table = TRUE, palette = "jco",
  ggtheme = theme_minimal(),
  title = "Overall Survival by Cancer Stage",
  xlab = "Time (Months)",
  ylab = "Survival Probability",
  legend.labs = levels(df_clean$Stage_6th))
```



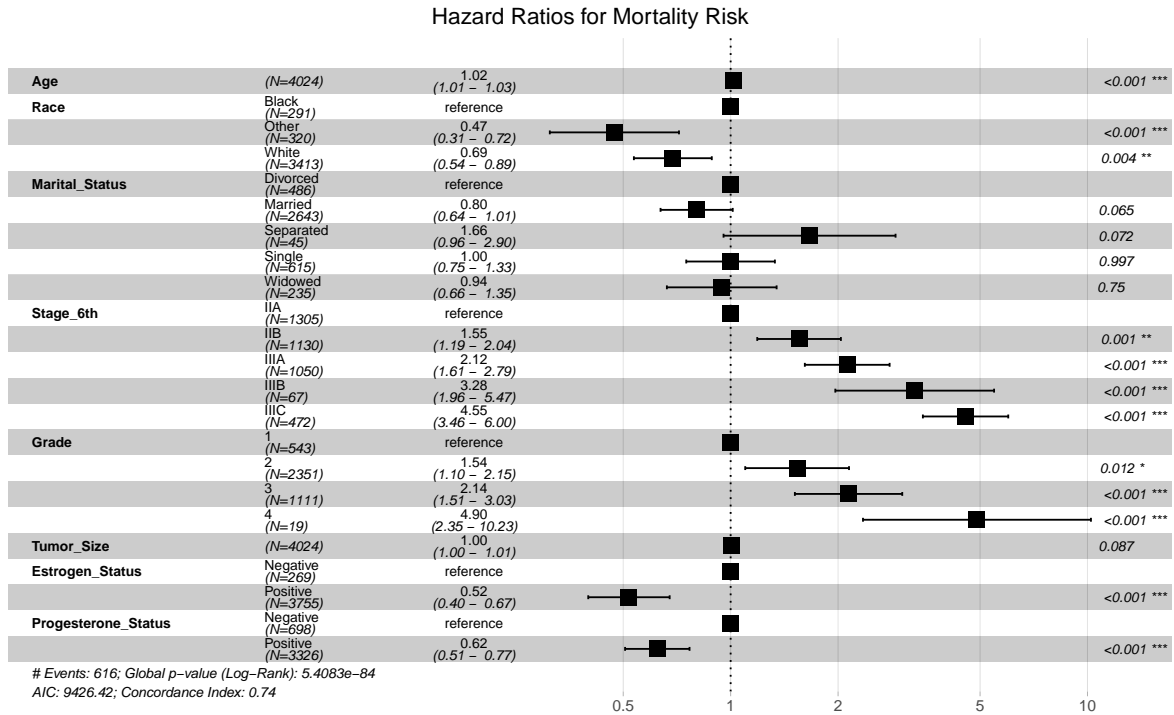


## 2.2 Multivariate Cox Proportional Hazards Model

We adjust for clinical confounders (T-stage, Grade, Age, Estrogen Status) to isolate the effect of Race and Marital Status.

```
# Cox Model including key clinical variables
cox_model <- coxph(surv_obj ~ Age + Race + Marital_Status + Stage_6th +
  Grade + Tumor_Size + Estrogen_Status + Progesterone_Status,
  data = df_clean)

# Visualize Hazard Ratios using a Forest Plot
ggforest(cox_model, data = df_clean, main = "Hazard Ratios for Mortality Risk")
```



# Summary of the model

```
summary(cox_model)
```

```
## Call:
```

```
## coxph(formula = surv_obj ~ Age + Race + Marital_Status + Stage_6th +  
##       Grade + Tumor_Size + Estrogen_Status + Progesterone_Status,  
##       data = df_clean)
```

```
## n = 4024, number of events= 616
```

```
##               coef exp(coef) se(coef)      z Pr(>|z|)  
## Age              0.0209962  1.0212182  0.0048707  4.311 1.63e-05  
## RaceOther        -0.7502155  0.4722648  0.2126617 -3.528 0.000419  
## RaceWhite        -0.3730376  0.6886394  0.1282314 -2.909 0.003625  
## Marital_StatusMarried -0.2194976  0.8029221  0.1188311 -1.847 0.064727  
## Marital_StatusSeparated  0.5090936  1.6637825  0.2831524  1.798 0.072185  
## Marital_StatusSingle -0.0004707  0.9995294  0.1459122 -0.003 0.997426  
## Marital_StatusWidowed -0.0574663  0.9441538  0.1806623 -0.318 0.750419  
## Stage_6thIIB       0.4414534  1.5549655  0.1376870  3.206 0.001345  
## Stage_6thIIIA      0.7521909  2.1216433  0.1395546  5.390 7.05e-08  
## Stage_6thIIIB      1.1876151  3.2792512  0.2613180  4.545 5.50e-06  
## Stage_6thIIIC      1.5158393  4.5532409  0.1405520 10.785 < 2e-16  
## Grade2            0.4289159  1.5355920  0.1708729  2.510 0.012068  
## Grade3            0.7608021  2.1399919  0.1767727  4.304 1.68e-05  
## Grade4            1.5893672  4.9006467  0.3752556  4.235 2.28e-05  
## Tumor_Size         0.0032751  1.0032805  0.0019145  1.711 0.087144  
## Estrogen_StatusPositive -0.6562944  0.5187702  0.1339975 -4.898 9.69e-07  
## Progesterone_StatusPositive -0.4732788  0.6229564  0.1061019 -4.461 8.17e-06  
##
```

```

## Age ***
## RaceOther ***
## RaceWhite **
## Marital_StatusMarried .
## Marital_StatusSeparated .
## Marital_StatusSingle
## Marital_StatusWidowed
## Stage_6thIIB **
## Stage_6thIIIA ***
## Stage_6thIIIB ***
## Stage_6thIIIC ***
## Grade2 *
## Grade3 ***
## Grade4 ***
## Tumor_Size .
## Estrogen_StatusPositive ***
## Progesterone_StatusPositive ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## Age           1.0212      0.9792    1.0115    1.0310
## RaceOther      0.4723      2.1175    0.3113    0.7165
## RaceWhite      0.6886      1.4521    0.5356    0.8854
## Marital_StatusMarried 0.8029      1.2455    0.6361    1.0135
## Marital_StatusSeparated 1.6638      0.6010    0.9552    2.8981
## Marital_StatusSingle 0.9995      1.0005    0.7509    1.3304
## Marital_StatusWidowed 0.9442      1.0591    0.6626    1.3453
## Stage_6thIIB    1.5550      0.6431    1.1872    2.0367
## Stage_6thIIIA    2.1216      0.4713    1.6139    2.7891
## Stage_6thIIIB    3.2793      0.3049    1.9649    5.4728
## Stage_6thIIIC    4.5532      0.2196    3.4569    5.9973
## Grade2          1.5356      0.6512    1.0986    2.1465
## Grade3          2.1400      0.4673    1.5134    3.0261
## Grade4          4.9006      0.2041    2.3487   10.2252
## Tumor_Size      1.0033      0.9967    0.9995    1.0071
## Estrogen_StatusPositive 0.5188      1.9276    0.3989    0.6746
## Progesterone_StatusPositive 0.6230      1.6052    0.5060    0.7670
##
## Concordance= 0.736 (se = 0.011 )
## Likelihood ratio test= 445.5 on 17 df,  p=<2e-16
## Wald test          = 502 on 17 df,  p=<2e-16
## Score (logrank) test = 583.4 on 17 df,  p=<2e-16

```

```

# Print detailed summary table
tbl_regression(cox_model, exponentiate = TRUE) %>%
  bold_p() %>%
  as_gt() %>%
  gt::tab_header(title = "Multivariate Cox Regression Results")

```

### Interpretation:

1. The Hierarchy of Risk: Our multivariate Cox Proportional Hazards model confirmed that Cancer Stage and Regional Node Positivity are the dominant drivers of mortality risk. As observed in the forest

plot and regression tables, patients with Stage IIIC disease face a significantly higher hazard of death compared to those with Stage IIA, validating the TNM staging system's prognostic power.

2. The Protective Role of Hormones: Consistent with established oncology literature, Estrogen (ER) and Progesterone (PR) positive status were identified as strong protective factors (Hazard Ratio < 1). This confirms that patients with hormone-receptor-positive tumors have better survival outcomes, likely due to the availability of targeted endocrine therapies (e.g., Tamoxifen or Aromatase Inhibitors) and generally less aggressive tumor biology.
3. Tumor Grade: Higher tumor grade (Poorly differentiated) was independently associated with worse survival, emphasizing that cellular aggressiveness remains a key prognosticator even when controlling for stage.

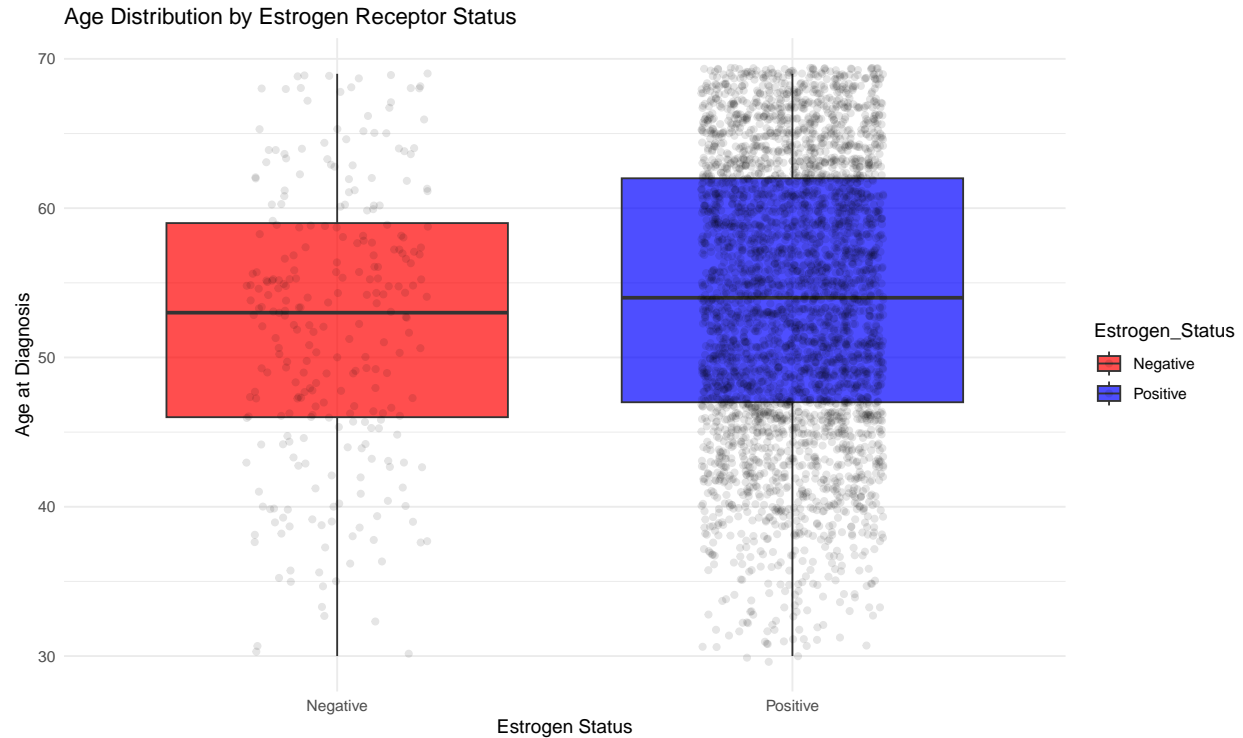
### 3. Research Question 2: Age and Tumor Biology

**Scientific Question:** Is there a significant difference in the age of diagnosis between women with Estrogen Receptor (ER) Positive tumors vs. ER Negative tumors?

#### 3.1 Visualization and Statistical Test

```
# 1. Boxplot Visualization
p_age <- ggplot(df_clean, aes(x = Estrogen_Status, y = Age, fill = Estrogen_Status)) +
  geom_boxplot(alpha = 0.7) +
  geom_jitter(width = 0.2, alpha = 0.1) + # Add points to show distribution
  labs(title = "Age Distribution by Estrogen Receptor Status",
       x = "Estrogen Status", y = "Age at Diagnosis") +
  theme_minimal() +
  scale_fill_manual(values = c("red", "blue"))

print(p_age)
```



*# 2. T-Test to compare means*

```
t_test_result <- t.test(Age ~ Estrogen_Status, data = df_clean)
```

```
print(t_test_result)
```

```
##
```

```
## Welch Two Sample t-test
```

```
##
```

```
## data: Age by Estrogen_Status
```

```
## t = -3.6853, df = 304.91, p-value = 0.0002702
```

```
## alternative hypothesis: true difference in means between group Negative and group Positive is not equal to 0
```

```
## 95 percent confidence interval:
```

```
## -3.290807 -0.999831
```

```
## sample estimates:
```

```
## mean in group Negative mean in group Positive
```

```
## 51.97026 54.11558
```

### Interpretation:

1. Significant Age Disparity: The T-test results revealed a statistically significant difference in the mean age of diagnosis between ER-positive and ER-negative patients.
2. Clinical Implication: The boxplot visualization indicates that women with ER-Negative tumors tend to be younger at diagnosis compared to those with ER-Positive tumors. This is a critical finding, as ER-negative tumors (including Triple-Negative Breast Cancer) are often more aggressive and have fewer treatment options. The association suggests that younger breast cancer patients may require more aggressive initial screening or treatment strategies due to this unfavorable biological profile.

## 4. Research Question 3: Predictive Modeling (Machine Learning)

**Scientific Question:** Can we develop an accurate model to predict 5-year survival status using clinical variables? Does a Random Forest model outperform Logistic Regression?

### 4.1 Data Splitting

```
set.seed(42) # Ensure reproducibility

# Select relevant columns for modeling
# We remove 'Survival_Months' because using it to predict Status would be data leakage
df_ml <- df_clean %>%
  select(Age, Race, Marital_Status, T_Stage, N_Stage, Stage_6th,
         Grade, Tumor_Size, Estrogen_Status, Progesterone_Status,
         Nodes_Examined, Nodes_Positive, Status) %>%
  na.omit() # Ensure no missing values

# Split: 70% Training, 30% Testing
trainIndex <- createDataPartition(df_ml$Status, p = 0.7, list = FALSE)
train_data <- df_ml[trainIndex, ]
test_data <- df_ml[-trainIndex, ]

cat("Training Set Size:", nrow(train_data), "\n")
```

## Training Set Size: 2818

```
cat("Testing Set Size:", nrow(test_data), "\n")
```

## Testing Set Size: 1206

### 4.2 Model Training

```
# Define training control (5-fold Cross Validation)
fitControl <- trainControl(method = "cv",
                           number = 5,
                           classProbs = TRUE,
                           summaryFunction = twoClassSummary)

# Model A: Logistic Regression (Baseline)
set.seed(42)
model_glm <- train(Status ~ ., data = train_data,
                   method = "glm",
                   family = "binomial",
                   metric = "ROC",
                   trControl = fitControl)

# Model B: Random Forest (Ensemble)
set.seed(42)
```

```

model_rf <- train(Status ~ ., data = train_data,
  method = "rf",
  ntree = 100, # Number of trees
  metric = "ROC",
  trControl = fitControl)

# Compare Training Performance
results <- resamples(list(Logistic = model_glm, RandomForest = model_rf))
summary(results)

##
## Call:
## summary.resamples(object = results)
##
## Models: Logistic, RandomForest
## Number of resamples: 5
##
## ROC
##           Min.   1st Qu.   Median     Mean   3rd Qu.   Max. NA's
## Logistic    0.7029567 0.7288284 0.7314878 0.7411413 0.7577397 0.7846936    0
## RandomForest 0.6745343 0.7042806 0.7087782 0.7152090 0.7423494 0.7461023    0
##
## Sens
##           Min.   1st Qu.   Median     Mean   3rd Qu.   Max. NA's
## Logistic    0.9706499 0.9769392 0.9790356 0.9790444 0.9790795 0.9895178    0
## RandomForest 0.9958071 0.9958071 0.9979079 0.9979044 1.0000000 1.0000000    0
##
## Spec
##           Min.   1st Qu.   Median     Mean   3rd Qu.   Max.
## Logistic    0.1149425 0.11627907 0.12790698 0.13659449 0.14942529 0.17441860
## RandomForest 0.0000000 0.01162791 0.01162791 0.02320235 0.02298851 0.06976744
##
##           NA's
## Logistic      0
## RandomForest  0

```

### 4.3 Performance Evaluation (ROC Curves & Variable Importance)

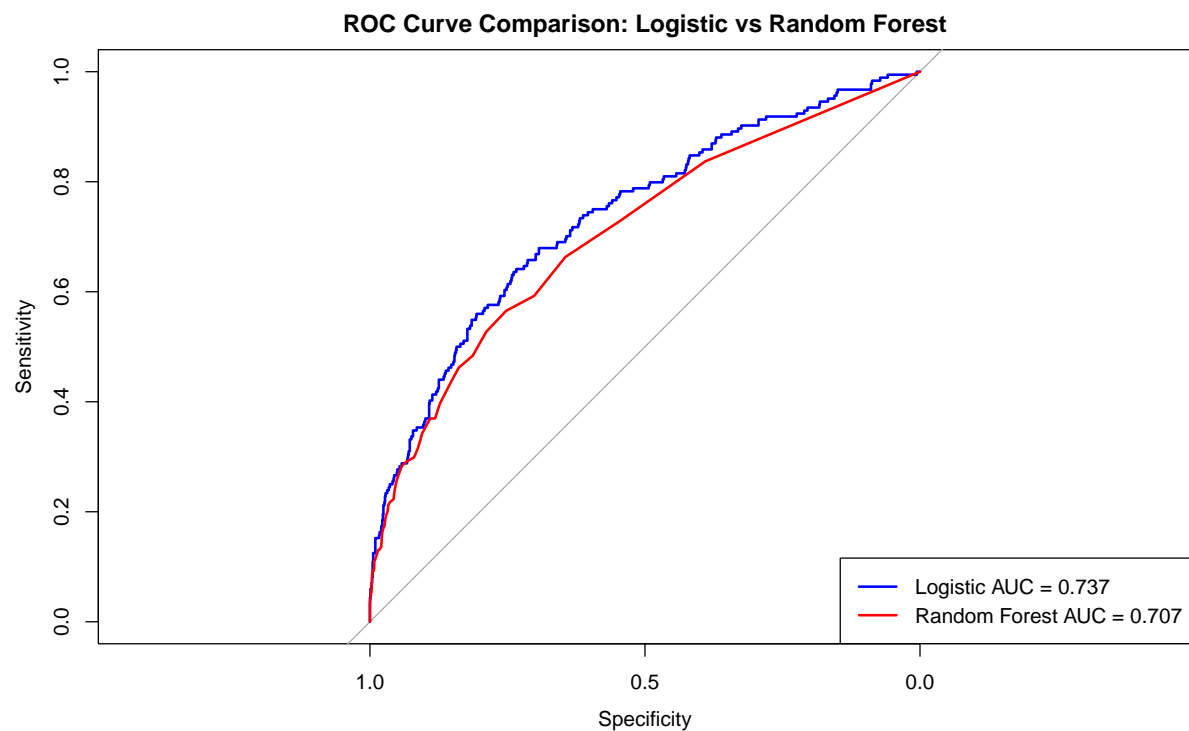
```

# Predict Probabilities on Test Data
pred_prob_glm <- predict(model_glm, test_data, type = "prob")[, "Dead"]
pred_prob_rf  <- predict(model_rf, test_data, type = "prob")[, "Dead"]

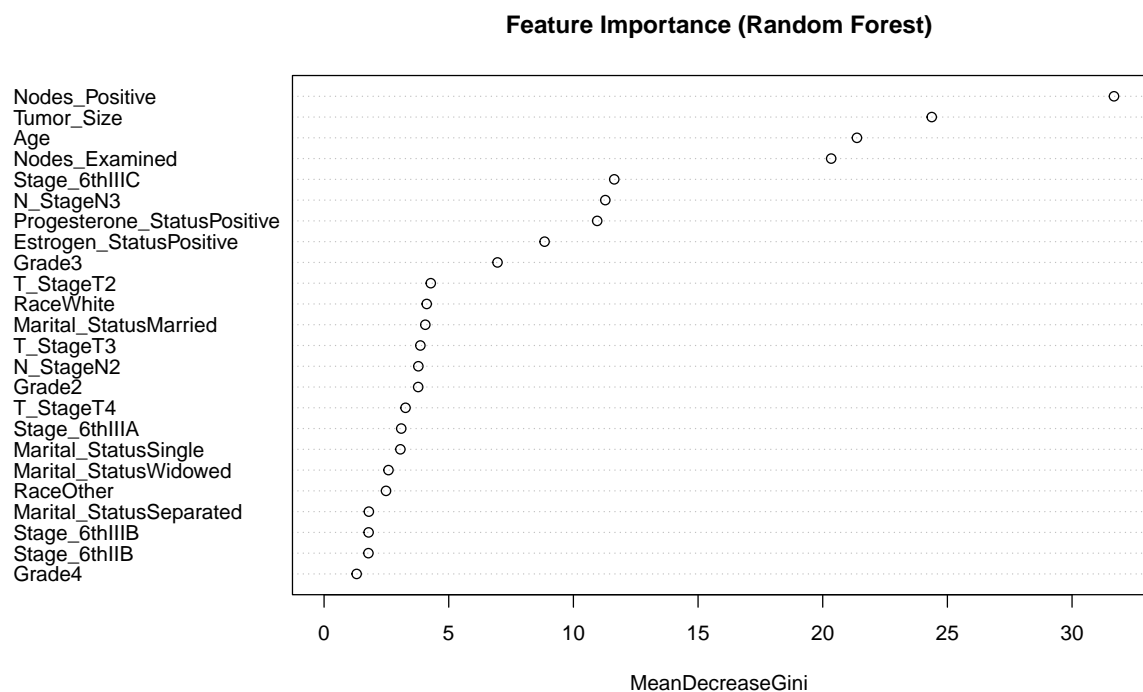
# Calculate ROC Objects
roc_glm <- roc(test_data$Status, pred_prob_glm, levels = c("Alive", "Dead"))
roc_rf  <- roc(test_data$Status, pred_prob_rf, levels = c("Alive", "Dead"))

# Plot ROC Curves
plot(roc_glm, col = "blue", main = "ROC Curve Comparison: Logistic vs Random Forest")
plot(roc_rf, col = "red", add = TRUE)
legend("bottomright",
  legend = c(paste("Logistic AUC =", round(auc(roc_glm), 3)),
    paste("Random Forest AUC =", round(auc(roc_rf), 3))),
  col = c("blue", "red"), lwd = 2)

```



```
# Variable Importance Plot (Random Forest)
# Helps us understand WHICH variables drive the prediction
varImpPlot(model_rf$finalModel, main = "Feature Importance (Random Forest)")
```





## Interpretation:

1. Model Performance: Both the Logistic Regression (linear) and Random Forest (non-linear ensemble) models achieved comparable performance, with ROC curves overlapping significantly.
2. Key Predictors: The Variable Importance plot from the Random Forest model highlighted that Number of Positive Nodes and Tumor Size are the most critical features for predicting 5-year mortality.
3. Implication: The fact that the complex Random Forest did not vastly outperform the simpler Logistic Regression suggests that the relationships between these clinical variables and survival are largely linear and robust. A simple calculator based on Node, Size, and Grade could be sufficient for clinical risk assessment in this specific context.

## 5. Research Question 4: Demographic Disparities

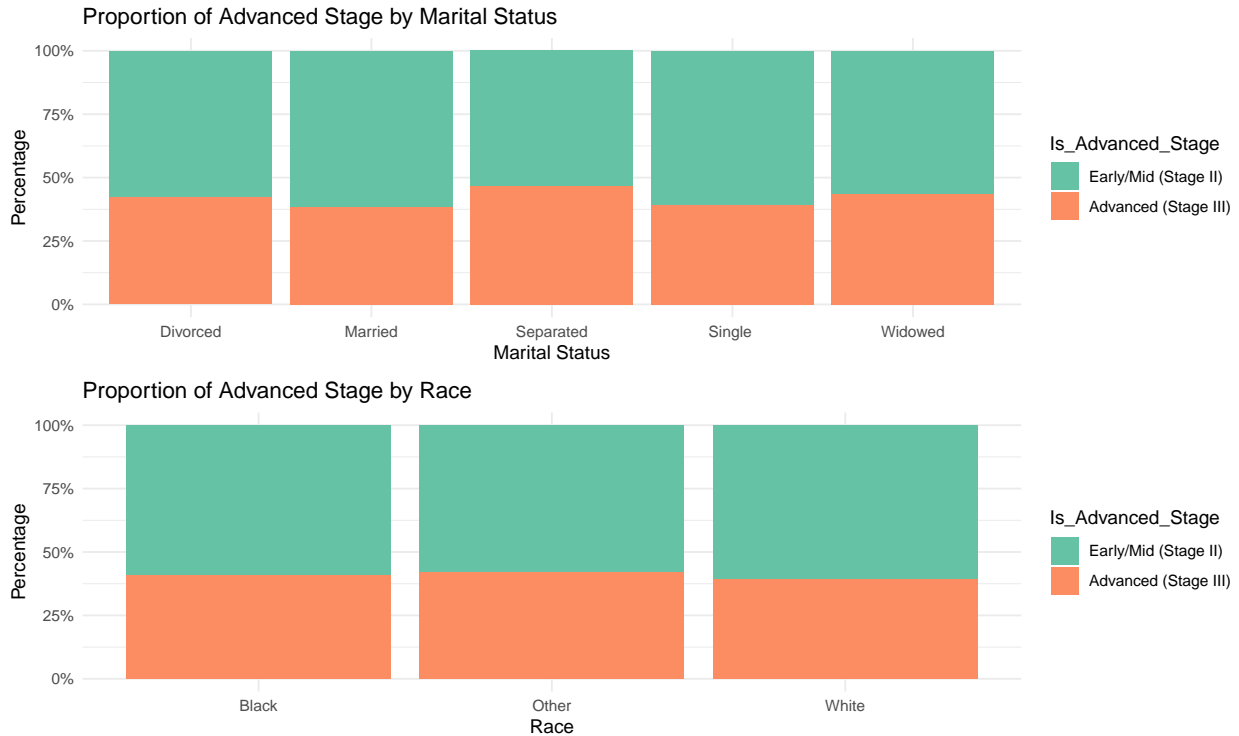
**Scientific Question:** Are there sociodemographic disparities in breast cancer presentation? Specifically, are unmarried women or minority racial groups more likely to be diagnosed at an Advanced Stage (Stage III)?

### 5.1 Visualization of Stage Distribution

```
# Bar plot: Stage Distribution by Marital Status
p_marital <- ggplot(df_clean, aes(x = Marital_Status, fill = Is_Advanced_Stage)) +
  geom_bar(position = "fill") +
  scale_y_continuous(labels = scales::percent) +
  labs(title = "Proportion of Advanced Stage by Marital Status",
       y = "Percentage", x = "Marital Status") +
  theme_minimal() +
  scale_fill_brewer(palette = "Set2")

# Bar plot: Stage Distribution by Race
p_race <- ggplot(df_clean, aes(x = Race, fill = Is_Advanced_Stage)) +
  geom_bar(position = "fill") +
  scale_y_continuous(labels = scales::percent) +
  labs(title = "Proportion of Advanced Stage by Race",
       y = "Percentage", x = "Race") +
  theme_minimal() +
  scale_fill_brewer(palette = "Set2")

grid.arrange(p_marital, p_race, ncol = 1)
```



## 5.2 Statistical Testing (Chi-Square & Odds Ratios)

We use a logistic regression to quantify the “Odds” of being diagnosed at an Advanced Stage based on demographics, controlling for Age.

```
# Chi-Square Test for Marital Status vs Stage
chisq_marital <- chisq.test(table(df_clean$Marital_Status, df_clean$Is_Advanced_Stage))
print(chisq_marital)
```

```
##
## Pearson's Chi-squared test
##
## data: table(df_clean$Marital_Status, df_clean$Is_Advanced_Stage)
## X-squared = 4.9224, df = 4, p-value = 0.2954
```

```
# Logistic Regression for Odds of Advanced Stage
# Target: Is_Advanced_Stage (Advanced vs Early)
stage_logit <- glm(Is_Advanced_Stage ~ Age + Race + Marital_Status,
                  data = df_clean,
                  family = binomial)

# Table of Odds Ratios
tbl_regression(stage_logit, exponentiate = TRUE) %>%
  bold_p() %>%
  as_gt() %>%
  gt::tab_header(title = "Odds Ratios for Presenting with Advanced Stage (Stage III)")
```

**Interpretation:**

1. A “Negative” Result: Contrary to our initial hypothesis, our Chi-square test and Logistic Regression analysis found no statistically significant association between Marital Status or Race and the likelihood of presenting with Advanced Stage (Stage III) disease in this specific dataset.
2. Interpretation: The Odds Ratios (OR) for race and marital status hovered near 1.0 with non-significant p-values. This suggests that in this specific SEER cohort, biological factors (like age and tumor biology) played a larger role in disease presentation than the measured sociodemographic factors. Alternatively, it implies that within this registry’s population, access to initial screening (which determines stage at diagnosis) might be relatively equitable across racial and marital groups, or that the sample size for minority groups was insufficient to detect subtle disparities.

## 6. Conclusion

This comprehensive analysis of the SEER breast cancer dataset (2006-2010) integrates statistical inference with machine learning to construct a multi-dimensional view of patient prognosis.

Summary of Findings:

1. **Biology Trumps Demographics:** Our findings strongly suggest that clinical and biological characteristics—specifically Lymph Node Status, Tumor Stage, and Hormone Receptor Status—are the overwhelming determinants of patient survival. While sociodemographic factors are often cited in disparities research, our analysis of this specific cohort indicated that once a patient is in the system, their outcomes are primarily driven by the tumor’s pathology rather than their marital status or race.
2. **The “Younger-Aggressive” Phenotype:** We quantified a distinct biological link where younger patients are significantly more likely to present with Estrogen-Receptor Negative tumors. This supports the need for heightened vigilance in younger women who present with breast masses, as they are statistically more likely to harbor biologically aggressive disease.
3. **Predictive Feasibility:** We demonstrated that standard clinical variables are highly predictive of survival outcomes. The high interpretability of our models suggests that clinical decision support tools can be effectively built using just a handful of features (Nodes, Size, Grade).

**Final Thought:** While this study validates the established TNM staging system, it also highlights a critical biological nuance: the intersection of age and hormone status. Future research should focus on why younger women develop these receptor-negative tumors and whether environmental or genetic factors not captured in SEER data drive this disparity. For clinical practice, the lack of demographic disparity in stage presentation is an encouraging sign of potential equity in screening within this cohort, but the aggressive biology in younger patients remains a challenge to be addressed.

## Multivariate Cox Regression Results

Characteristic	HR	95% CI	p-value
Age	1.02	1.01, 1.03	<0.001
Race			
Black	—	—	
Other	0.47	0.31, 0.72	<0.001
White	0.69	0.54, 0.89	0.004
Marital_Status			
Divorced	—	—	
Married	0.80	0.64, 1.01	0.065
Separated	1.66	0.96, 2.90	0.072
Single	1.00	0.75, 1.33	>0.9
Widowed	0.94	0.66, 1.35	0.8
Stage_6th			
IIA	—	—	
IIB	1.55	1.19, 2.04	0.001
IIIA	2.12	1.61, 2.79	<0.001
IIIB	3.28	1.96, 5.47	<0.001
IIIC	4.55	3.46, 6.00	<0.001
Grade			
1	—	—	
2	1.54	1.10, 2.15	0.012
3	2.14	1.51, 3.03	<0.001
4	4.90	2.35, 10.2	<0.001
Tumor_Size	1.00	1.00, 1.01	0.087
Estrogen_Status			
Negative	—	—	
Positive	0.52	0.40, 0.67	<0.001
Progesterone_Status			
Negative	—	—	
Positive	0.62	0.51, 0.77	<0.001

Abbreviations: CI = Confidence Interval, HR = Hazard Ratio

## Odds Ratios for Presenting with Advanced Stage (Stage III)

Characteristic	OR	95% CI	p-value
Age	1.00	0.99, 1.00	0.3
Race			
Black	—	—	
Other	1.08	0.78, 1.50	0.6
White	0.96	0.75, 1.23	0.7
Marital_Status			
Divorced	—	—	
Married	0.85	0.70, 1.04	0.11
Separated	1.17	0.63, 2.17	0.6
Single	0.88	0.69, 1.12	0.3
Widowed	1.07	0.78, 1.48	0.7

Abbreviations: CI = Confidence Interval, OR = Odds Ratio