

GPH 2338 Project

Yi Yang

2023-03-13

Data Upload

```
library(haven)
library(psych)
library(caret)
library(tidyverse)
library(ggplot2)
library(psych)
library(pander)
library(corrplot)
library(pander)
library(readr)
library(r02pro)
library(plyr)
library(tree)
library(gbm)
library(caret)
library(leaps)
library(readr)
```

Data Preparation

```
df<-readr::read_tsv("brca_metabric_clinical_data.tsv")
```

```
## Rows: 2509 Columns: 39
## -- Column specification -----
## Delimiter: "\t"
## chr (27): Study ID, Patient ID, Sample ID, Type of Breast Surgery, Cancer Ty...
## dbl (12): Age at Diagnosis, Cohort, Neoplasm Histologic Grade, Lymph nodes e...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
colnames(df)
```

```
## [1] "Study ID"           "Patient ID"
## [3] "Sample ID"          "Age at Diagnosis"
## [5] "Type of Breast Surgery" "Cancer Type"
## [7] "Cancer Type Detailed" "Cellularity"
## [9] "Chemotherapy"       "Pam50 + Claudin-low subtype"
```

```
## [11] "Cohort" "ER status measured by IHC"
## [13] "ER Status" "Neoplasm Histologic Grade"
## [15] "HER2 status measured by SNP6" "HER2 Status"
## [17] "Tumor Other Histologic Subtype" "Hormone Therapy"
## [19] "Inferred Menopausal State" "Integrative Cluster"
## [21] "Primary Tumor Laterality" "Lymph nodes examined positive"
## [23] "Mutation Count" "Nottingham prognostic index"
## [25] "Oncotree Code" "Overall Survival (Months)"
## [27] "Overall Survival Status" "PR Status"
## [29] "Radio Therapy" "Relapse Free Status (Months)"
## [31] "Relapse Free Status" "Number of Samples Per Patient"
## [33] "Sample Type" "Sex"
## [35] "3-Gene classifier subtype" "TMB (nonsynonymous)"
## [37] "Tumor Size" "Tumor Stage"
## [39] "Patient's Vital Status"
```

```
colSums(is.na(df))
```

```
## Study ID Patient ID
## 0 0
## Sample ID Age at Diagnosis
## 0 11
## Type of Breast Surgery Cancer Type
## 554 0
## Cancer Type Detailed Cellularity
## 0 592
## Chemotherapy Pam50 + Claudin-low subtype
## 529 529
## Cohort ER status measured by IHC
## 11 83
## ER Status Neoplasm Histologic Grade
## 40 121
## HER2 status measured by SNP6 HER2 Status
## 529 529
## Tumor Other Histologic Subtype Hormone Therapy
## 135 529
## Inferred Menopausal State Integrative Cluster
## 529 529
## Primary Tumor Laterality Lymph nodes examined positive
## 639 266
## Mutation Count Nottingham prognostic index
## 151 222
## Oncotree Code Overall Survival (Months)
## 0 528
## Overall Survival Status PR Status
## 528 529
## Radio Therapy Relapse Free Status (Months)
## 529 121
## Relapse Free Status Number of Samples Per Patient
## 21 0
## Sample Type Sex
## 0 0
## 3-Gene classifier subtype TMB (nonsynonymous)
## 745 0
```

```
##           Tumor Size           Tumor Stage
##           149           721
## Patient's Vital Status
##           529
```

```
print(df)
```

```
## # A tibble: 2,509 x 39
##   'Study ID'   Patie~1 Sampl~2 Age a~3 Type ~4 Cance~5 Cance~6 Cellu~7 Chemo~8
##   <chr>       <chr>   <chr>   <dbl> <chr>   <chr>   <chr>   <chr>   <chr>
## 1 brca_metabric MB-0000 MB-0000  75.6 MASTEC~ Breast~ Breast~ <NA>   NO
## 2 brca_metabric MB-0002 MB-0002  43.2 BREAST~ Breast~ Breast~ High   NO
## 3 brca_metabric MB-0005 MB-0005  48.9 MASTEC~ Breast~ Breast~ High   YES
## 4 brca_metabric MB-0006 MB-0006  47.7 MASTEC~ Breast~ Breast~ Modera~ YES
## 5 brca_metabric MB-0008 MB-0008  77.0 MASTEC~ Breast~ Breast~ High   YES
## 6 brca_metabric MB-0010 MB-0010  78.8 MASTEC~ Breast~ Breast~ Modera~ NO
## 7 brca_metabric MB-0014 MB-0014  56.4 BREAST~ Breast~ Breast~ Modera~ YES
## 8 brca_metabric MB-0020 MB-0020  70   MASTEC~ Breast~ Breast~ High   YES
## 9 brca_metabric MB-0022 MB-0022  89.1 BREAST~ Breast~ Breast~ Modera~ NO
## 10 brca_metabric MB-0025 MB-0025  76.2 <NA>   Breast~ Breast~ <NA>   <NA>
## # ... with 2,499 more rows, 30 more variables:
## #   'Pam50 + Claudin-low subtype' <chr>, Cohort <dbl>,
## #   'ER status measured by IHC' <chr>, 'ER Status' <chr>,
## #   'Neoplasm Histologic Grade' <dbl>, 'HER2 status measured by SNP6' <chr>,
## #   'HER2 Status' <chr>, 'Tumor Other Histologic Subtype' <chr>,
## #   'Hormone Therapy' <chr>, 'Inferred Menopausal State' <chr>,
## #   'Integrative Cluster' <chr>, 'Primary Tumor Laterality' <chr>, ...
```

```
df1 <- df %>% na.omit()
colSums(is.na(df1))
```

```
##           Study ID           Patient ID
##           0           0
##           Sample ID           Age at Diagnosis
##           0           0
##           Type of Breast Surgery           Cancer Type
##           0           0
##           Cancer Type Detailed           Cellularity
##           0           0
##           Chemotherapy Pam50 + Claudin-low subtype
##           0           0
##           Cohort ER status measured by IHC
##           0           0
##           ER Status Neoplasm Histologic Grade
##           0           0
## HER2 status measured by SNP6 HER2 Status
##           0           0
## Tumor Other Histologic Subtype Hormone Therapy
##           0           0
## Inferred Menopausal State Integrative Cluster
##           0           0
## Primary Tumor Laterality Lymph nodes examined positive
##           0           0
```

```
##           Mutation Count      Nottingham prognostic index
##           0                0
##           Oncotree Code      Overall Survival (Months)
##           0                0
##           Overall Survival Status      PR Status
##           0                0
##           Radio Therapy      Relapse Free Status (Months)
##           0                0
##           Relapse Free Status      Number of Samples Per Patient
##           0                0
##           Sample Type      Sex
##           0                0
##           3-Gene classifier subtype      TMB (nonsynonymous)
##           0                0
##           Tumor Size      Tumor Stage
##           0                0
##           Patient's Vital Status
##           0
```

```
df2<-df1[,-c(1, 2, 3, 6)] # remove StudyID, Patient ID, Sample ID, Cancer Type
print(df2)
```

```
## # A tibble: 1,092 x 35
##   Age at Diagn~1 Type ~2 Cance~3 Cellu~4 Chemo~5 Pam50~6 Cohort ER st~7 ER St~8
##   <dbl> <chr> <chr> <chr> <chr> <chr> <dbl> <chr> <chr>
## 1 43.2 BREAST~ Breast~ High NO LumA 1 Positive Positi~
## 2 77.0 MASTEC~ Breast~ High YES LumB 1 Positive Positi~
## 3 78.8 MASTEC~ Breast~ Modera~ NO LumB 1 Positive Positi~
## 4 86.4 BREAST~ Breast~ Modera~ NO LumB 1 Positive Positi~
## 5 84.2 MASTEC~ Breast~ High NO Her2 1 Negati~ Positi~
## 6 85.5 MASTEC~ Breast~ Modera~ NO LumA 1 Positive Positi~
## 7 45.4 BREAST~ Breast~ High YES LumB 1 Positive Positi~
## 8 61.5 BREAST~ Breast~ High NO LumB 1 Positive Positi~
## 9 68.7 MASTEC~ Breast~ Low YES Basal 1 Negati~ Negati~
## 10 46.9 MASTEC~ Breast~ Modera~ NO Normal 1 Positive Positi~
## # ... with 1,082 more rows, 26 more variables:
## # 'Neoplasm Histologic Grade' <dbl>, 'HER2 status measured by SNP6' <chr>,
## # 'HER2 Status' <chr>, 'Tumor Other Histologic Subtype' <chr>,
## # 'Hormone Therapy' <chr>, 'Inferred Menopausal State' <chr>,
## # 'Integrative Cluster' <chr>, 'Primary Tumor Laterality' <chr>,
## # 'Lymph nodes examined positive' <dbl>, 'Mutation Count' <dbl>,
## # 'Nottingham prognostic index' <dbl>, 'Oncotree Code' <chr>, ...
```

```
dim(df2)
```

```
## [1] 1092 35
```

```
is_categorical <- sapply(df2, is.character)
is_categorical
```

```
##           Age at Diagnosis      Type of Breast Surgery
##           FALSE                TRUE
```

```

##          Cancer Type Detailed          Cellularity
##                TRUE                TRUE
##          Chemotherapy    Pam50 + Claudin-low subtype
##                TRUE                TRUE
##                Cohort    ER status measured by IHC
##                FALSE                TRUE
##                ER Status    Neoplasm Histologic Grade
##                TRUE                FALSE
##    HER2 status measured by SNP6    HER2 Status
##                TRUE                TRUE
## Tumor Other Histologic Subtype    Hormone Therapy
##                TRUE                TRUE
##    Inferred Menopausal State    Integrative Cluster
##                TRUE                TRUE
##    Primary Tumor Laterality    Lymph nodes examined positive
##                TRUE                FALSE
##                Mutation Count    Nottingham prognostic index
##                FALSE                FALSE
##                Oncotree Code    Overall Survival (Months)
##                TRUE                FALSE
##    Overall Survival Status    PR Status
##                TRUE                TRUE
##                Radio Therapy    Relapse Free Status (Months)
##                TRUE                FALSE
##                Relapse Free Status    Number of Samples Per Patient
##                TRUE                FALSE
##                Sample Type    Sex
##                TRUE                TRUE
##    3-Gene classifier subtype    TMB (nonsynonymous)
##                TRUE                FALSE
##                Tumor Size    Tumor Stage
##                FALSE                FALSE
##    Patient's Vital Status
##                TRUE

```

Create Dummy Variable for Categorical Variable

```

df2$Type_of_Breast_Surgery = ifelse(df2$`Type of Breast Surgery` == "BREAST CONSERVING",1,0)
df2$new_Chemotherapy = ifelse(df2$Chemotherapy == "NO",0,1)
for (i in 1:nrow(df2)) {
  if (df2[i,]$Cellularity == "High") {
    df2$new_Cellularity[i] <- 3
  }
  else if (df2[i,]$Cellularity == "Moderate") {
    df2$new_Cellularity[i] <- 2
  }
  else {
    df2$new_Cellularity[i] <- 1
  }
}

df2$Pam50_Claudin_low_subtype_Luma = ifelse(df2$`Pam50 + Claudin-low subtype` == "LumA",1,0)
df2$Pam50_Claudin_low_subtype_LumB = ifelse(df2$`Pam50 + Claudin-low subtype` == "LumB",1,0)
df2$Pam50_Claudin_low_subtype_Her2 = ifelse(df2$`Pam50 + Claudin-low subtype` == "Her2",1,0)

```

```

df2$Pam50_Claudin_low_subtype_Basal = ifelse(df2$`Pam50 + Claudin-low subtype` == "Basal",1,0)
df2$Pam50_Claudin_low_subtype_Normal = ifelse(df2$`Pam50 + Claudin-low subtype` == "Normal",1,0)
df2$Pam50_Claudin_low_subtype_claudin = ifelse(df2$`Pam50 + Claudin-low subtype` == "claudin-low",1,0)
df2$Pam50_Claudin_low_subtype_NC = ifelse(df2$`Pam50 + Claudin-low subtype` == "NC",1,0)

df2$ER_status_measured_by_IHC = ifelse(df2$`ER status measured by IHC`=="Positive",1,0)
df2$ER_Status_Positive = ifelse(df2$`ER Status`=="Positive",1,0)

for (i in 1:nrow(df2)) {
  if (df2[i,]$`HER2 status measured by SNP6` == "NEUTRAL") {
    df2$HER2_status_measured_by_SNP6[i] <- 4
  }
  else if (df2[i,]$`HER2 status measured by SNP6` == "GAIN") {
    df2$HER2_status_measured_by_SNP6[i] <- 3
  }
  else if (df2[i,]$`HER2 status measured by SNP6` == "LOSS"){
    df2$HER2_status_measured_by_SNP6[i] <- 2
  }
  else if (df2[i,]$`HER2 status measured by SNP6` == "UNDEF"){
    df2$HER2_status_measured_by_SNP6[i] <- 1
  }
}
df2$HER2_Status_Positive = ifelse(df2$`HER2 Status`=="Positive",1,0)
df2$Tumor_Other_Histologic_Subtype_Ductal = ifelse(df2$`Tumor Other Histologic Subtype` == "Ductal/NST",1,0)
df2$Tumor_Other_Histologic_Subtype_Tubular = ifelse(df2$`Tumor Other Histologic Subtype` == "Tubular/ c",1,0)
df2$Tumor_Other_Histologic_Subtype_Medullary = ifelse(df2$`Tumor Other Histologic Subtype` == "Medullary",1,0)

#df2$Oncotree_Code_IDC = ifelse(df2$`Oncotree Code` == "IDC",1,0)
#df2$Oncotree_Code_MLDC = ifelse(df2$`Oncotree Code` == "MLDC",1,0)
#df2$Oncotree_Code_ILC = ifelse(df2$`Oncotree Code` == "ILC",1,0)
#df2$Oncotree_Code_IMMC = ifelse(df2$`Oncotree Code` == "IMMC",1,0)
#df2$Oncotree_Code_BREAST = ifelse(df2$`Oncotree Code` == "BREAST",1,0)
df2$Hormone_Therapy = ifelse(df2$`Hormone Therapy`=="NO",0,1)
df2$Inferred_Menopausal_State = ifelse(df2$`Inferred Menopausal State`=="Pre",0,1)
df2$Primary_Tumor_Laterality = ifelse(df2$`Primary Tumor Laterality`=="Right",0,1)
df2$Overall_Survival_Status = ifelse(df2$`Overall Survival Status`=="1:DECEASED",1,0)
df2$PR_Status = ifelse(df2$`PR Status`=="Positive",1,0)
df2$Radio_Therapy = ifelse(df2$`Radio Therapy`=="NO",0,1)
df2$Relapse_Free_Status = ifelse(df2$`Relapse Free Status`=="0:Not Recurred",0,1)

df2$Gene_classifier_subtype_ERH = ifelse(df2$`3-Gene classifier subtype` == "ER+/HER2- High Prolif",1,0)
df2$Gene_classifier_subtype_ERL = ifelse(df2$`3-Gene classifier subtype` == "ER+/HER2- Low Prolif",1,0)
df2$Gene_classifier_subtype_ERM = ifelse(df2$`3-Gene classifier subtype` == "ER-/HER2-",1,0)
df2$Gene_classifier_subtype_ERP = ifelse(df2$`3-Gene classifier subtype` == "HER2+",1,0)

df2$Patients_Vital_Status_Living = ifelse(df2$`Patient's Vital Status` == "Living",1,0)
df2$Patients_Vital_Status_Died = ifelse(df2$`Patient's Vital Status` == "Died of Disease",1,0)
df2$Patients_Vital_Status_Do = ifelse(df2$`Patient's Vital Status` == "Died of Other Causes",1,0)

for (i in 1:nrow(df2)) {

```

```

if (df2[i,]$`Integrative Cluster` == "4ER+") {
  df2$Integrative_Cluster[i] <- 4.5
}
else if (df2[i,]$`Integrative Cluster` == "9") {
  df2$Integrative_Cluster[i] <- 9
}
else if (df2[i,]$`Integrative Cluster` == "7"){
  df2$Integrative_Cluster[i] <- 7
}
else if (df2[i,]$`Integrative Cluster` == "3"){
  df2$Integrative_Cluster[i] <- 3
}
else if (df2[i,]$`Integrative Cluster` == "10"){
  df2$Integrative_Cluster[i] <- 10
}
else if (df2[i,]$`Integrative Cluster` == "8"){
  df2$Integrative_Cluster[i] <- 8
}
else if (df2[i,]$`Integrative Cluster` == "6"){
  df2$Integrative_Cluster[i] <- 6
}
else if (df2[i,]$`Integrative Cluster` == "1"){
  df2$Integrative_Cluster[i] <- 1
}
else if (df2[i,]$`Integrative Cluster` == "2"){
  df2$Integrative_Cluster[i] <- 2
}
else if (df2[i,]$`Integrative Cluster` == "5"){
  df2$Integrative_Cluster[i] <- 5
}
else if (df2[i,]$`Integrative Cluster` == "4ER-"){
  df2$Integrative_Cluster[i] <- 3.5
}
}
names(df2)[names(df2) == "Cancer Type Detailed"] <- "Cancer_Type_Detailed"
df2<- df2[df2$Cancer_Type_Detailed == "Breast Invasive Ductal Carcinoma",]

df2

```

```

## # A tibble: 859 x 67
##   Age at Diagn~1 Type ~2 Cance~3 Cellu~4 Chemo~5 Pam50~6 Cohort ER st~7 ER St~8
##           <dbl> <chr>  <chr>  <chr>  <chr>  <chr>  <dbl> <chr>  <chr>
## 1         43.2 BREAST~ Breast~ High   NO     LumA      1 Positive Positi~
## 2         78.8 MASTEC~ Breast~ Modera~ NO     LumB      1 Positive Positi~
## 3         86.4 BREAST~ Breast~ Modera~ NO     LumB      1 Positive Positi~
## 4         85.5 MASTEC~ Breast~ Modera~ NO     LumA      1 Positive Positi~
## 5         45.4 BREAST~ Breast~ High   YES    LumB      1 Positive Positi~
## 6         61.5 BREAST~ Breast~ High   NO     LumB      1 Positive Positi~
## 7         68.7 MASTEC~ Breast~ Low    YES    Basal     1 Negati~ Negati~
## 8         49.9 MASTEC~ Breast~ Modera~ YES    LumA      1 Positive Positi~
## 9         54.2 MASTEC~ Breast~ High   NO     LumA      1 Positive Positi~
## 10        48.6 MASTEC~ Breast~ Low    NO     LumA      1 Positive Positi~

```

```
## # ... with 849 more rows, 58 more variables: 'Neoplasm Histologic Grade' <dbl>,
## # 'HER2 status measured by SNP6' <chr>, 'HER2 Status' <chr>,
## # 'Tumor Other Histologic Subtype' <chr>, 'Hormone Therapy' <chr>,
## # 'Inferred Menopausal State' <chr>, 'Integrative Cluster' <chr>,
## # 'Primary Tumor Laterality' <chr>, 'Lymph nodes examined positive' <dbl>,
## # 'Mutation Count' <dbl>, 'Nottingham prognostic index' <dbl>,
## # 'Oncotree Code' <chr>, 'Overall Survival (Months)' <dbl>, ...
```

Create Dummy Variable for Numeric Variable

```
names(df2)[names(df2) == "Age at Diagnosis"] <- "Age"
names(df2)[names(df2) == "Neoplasm Histologic Grade"] <- "Neo_Grade"
names(df2)[names(df2) == "Lymph nodes examined positive"] <- "Lymph"
names(df2)[names(df2) == "Mutation Count"] <- "Mutation"
names(df2)[names(df2) == "Nottingham prognostic index"] <- "Nottingham"
names(df2)[names(df2) == "Overall Survival (Months)"] <- "Overall_Month"
names(df2)[names(df2) == "Relapse Free Status (Months)"] <- "Relapse_Month"
names(df2)[names(df2) == "TMB (nonsynonymous)"] <- "TMB"
names(df2)[names(df2) == "Tumor Size"] <- "Tumor_Size"
names(df2)[names(df2) == "Tumor Stage"] <- "Tumor_Stage"
```

```
df3 <- df2[,-c(2,3,4:6,8,9,11:17,21,23:25,27:32,35)]
df3
```

```
## # A tibble: 859 x 42
##   Age Cohort Neo_Grade Lymph Mutation Notti~1 Overa~2 Relap~3 Tumor~4 Tumor~5
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 43.2      1      3      0      2      4.02  84.6   83.5     10      1
## 2 78.8      1      3      0      4      4.06   7.8    2.89     31      4
## 3 86.4      1      3      1      4      5.03  36.6   36.1     16      2
## 4 85.5      1      2      0      1      3.04 132.   123.     22      4
## 5 45.4      1      3      0      5      4.05 141.   139.     23      2
## 6 61.5      1      2      1      3      4.03 157.   155.     16      2
## 7 68.7      1      3      0      1      4.08   8.07   7.83     39      2
## 8 49.9      1      1      5      4      4.14  85.3   84.2     70      3
## 9 54.2      1      1      0      4      2.05 127.   125.     27      2
## 10 48.6      1      2      0      3      3.06 13.4   13.2     30      2
## # ... with 849 more rows, 32 more variables: Type_of_Breast_Surgery <dbl>,
## # new_Chemotherapy <dbl>, new_Cellularity <dbl>,
## # Pam50_Claudin_low_subtype_Luma <dbl>, Pam50_Claudin_low_subtype_LumB <dbl>,
## # Pam50_Claudin_low_subtype_Her2 <dbl>,
## # Pam50_Claudin_low_subtype_Basal <dbl>,
## # Pam50_Claudin_low_subtype_Normal <dbl>,
## # Pam50_Claudin_low_subtype_claudin <dbl>, ...
```

```
pander(summary(df3),caption='Descriptive Statistics of The Data')
```

Table 1: Descriptive Statistics of The Data (continued below)

| Age | Cohort | Neo_Grade | Lymph |
|-------------|-------------|-------------|--------------|
| Min. :26.36 | Min. :1.000 | Min. :1.000 | Min. : 0.000 |

| Age | Cohort | Neo_Grade | Lymph |
|---------------|---------------|---------------|----------------|
| 1st Qu.:49.99 | 1st Qu.:1.000 | 1st Qu.:2.000 | 1st Qu.: 0.000 |
| Median :60.62 | Median :2.000 | Median :3.000 | Median : 0.000 |
| Mean :60.05 | Mean :2.191 | Mean :2.517 | Mean : 1.916 |
| 3rd Qu.:69.75 | 3rd Qu.:3.000 | 3rd Qu.:3.000 | 3rd Qu.: 2.000 |
| Max. :96.29 | Max. :5.000 | Max. :3.000 | Max. :41.000 |

Table 2: Table continues below

| Mutation | Nottingham | Overall_Month | Relapse_Month |
|----------------|---------------|----------------|----------------|
| Min. : 1.000 | Min. :2.018 | Min. : 0.10 | Min. : 0.10 |
| 1st Qu.: 3.000 | 1st Qu.:3.080 | 1st Qu.: 58.05 | 1st Qu.: 40.10 |
| Median : 5.000 | Median :4.050 | Median :115.30 | Median : 98.42 |
| Mean : 5.423 | Mean :4.216 | Mean :124.05 | Mean :109.35 |
| 3rd Qu.: 7.000 | 3rd Qu.:5.050 | 3rd Qu.:186.32 | 3rd Qu.:172.12 |
| Max. :46.000 | Max. :6.360 | Max. :337.03 | Max. :296.91 |

Table 3: Table continues below

| Tumor_Size | Tumor_Stage | Type_of_Breast_Surgery | new_Chemotherapy |
|---------------|---------------|------------------------|------------------|
| Min. : 1.0 | Min. :1.000 | Min. :0.0000 | Min. :0.0000 |
| 1st Qu.: 17.0 | 1st Qu.:1.000 | 1st Qu.:0.0000 | 1st Qu.:0.0000 |
| Median : 22.0 | Median :2.000 | Median :0.0000 | Median :0.0000 |
| Mean : 25.7 | Mean :1.767 | Mean :0.4319 | Mean :0.2538 |
| 3rd Qu.: 30.0 | 3rd Qu.:2.000 | 3rd Qu.:1.0000 | 3rd Qu.:1.0000 |
| Max. :180.0 | Max. :4.000 | Max. :1.0000 | Max. :1.0000 |

Table 4: Table continues below

| new_Cellularity | Pam50_Claudin_low_subtype_Luma |
|-----------------|--------------------------------|
| Min. :1.000 | Min. :0.0000 |
| 1st Qu.:2.000 | 1st Qu.:0.0000 |
| Median :3.000 | Median :0.0000 |
| Mean :2.421 | Mean :0.3667 |
| 3rd Qu.:3.000 | 3rd Qu.:1.0000 |
| Max. :3.000 | Max. :1.0000 |

Table 5: Table continues below

| Pam50_Claudin_low_subtype_LumB | Pam50_Claudin_low_subtype_Her2 |
|--------------------------------|--------------------------------|
| Min. :0.0000 | Min. :0.0000 |
| 1st Qu.:0.0000 | 1st Qu.:0.0000 |
| Median :0.0000 | Median :0.0000 |
| Mean :0.2608 | Mean :0.1036 |
| 3rd Qu.:1.0000 | 3rd Qu.:0.0000 |
| Max. :1.0000 | Max. :1.0000 |

Table 6: Table continues below

| Pam50_Claudin_low_subtype_Basal | Pam50_Claudin_low_subtype_Normal |
|---------------------------------|----------------------------------|
| Min. :0.0000 | Min. :0.00000 |
| 1st Qu.:0.0000 | 1st Qu.:0.00000 |
| Median :0.0000 | Median :0.00000 |
| Mean :0.1141 | Mean :0.05122 |
| 3rd Qu.:0.0000 | 3rd Qu.:0.00000 |
| Max. :1.0000 | Max. :1.00000 |

Table 7: Table continues below

| Pam50_Claudin_low_subtype_claudin | Pam50_Claudin_low_subtype_NC |
|-----------------------------------|------------------------------|
| Min. :0.0000 | Min. :0.000000 |
| 1st Qu.:0.0000 | 1st Qu.:0.000000 |
| Median :0.0000 | Median :0.000000 |
| Mean :0.1013 | Mean :0.002328 |
| 3rd Qu.:0.0000 | 3rd Qu.:0.000000 |
| Max. :1.0000 | Max. :1.000000 |

Table 8: Table continues below

| ER_status_measured_by_IHC | ER_Status_Positive | HER2_status_measured_by_SNP6 |
|---------------------------|--------------------|------------------------------|
| Min. :0.0000 | Min. :0.0000 | Min. :1.000 |
| 1st Qu.:0.0000 | 1st Qu.:0.0000 | 1st Qu.:3.000 |
| Median :1.0000 | Median :1.0000 | Median :4.000 |
| Mean :0.7404 | Mean :0.7404 | Mean :3.631 |
| 3rd Qu.:1.0000 | 3rd Qu.:1.0000 | 3rd Qu.:4.000 |
| Max. :1.0000 | Max. :1.0000 | Max. :4.000 |

Table 9: Table continues below

| HER2_Status_Positive | Tumor_Other_Histologic_Subtype_Ductal |
|----------------------|---------------------------------------|
| Min. :0.0000 | Min. :0.0000 |
| 1st Qu.:0.0000 | 1st Qu.:1.0000 |
| Median :0.0000 | Median :1.0000 |
| Mean :0.1444 | Mean :0.9674 |
| 3rd Qu.:0.0000 | 3rd Qu.:1.0000 |
| Max. :1.0000 | Max. :1.0000 |

Table 10: Table continues below

| Tumor_Other_Histologic_Subtype_Tubular |
|--|
| Min. :0.00000 |
| 1st Qu.:0.00000 |
| Median :0.00000 |
| Mean :0.01746 |

| Tumor_Other_Histologic_Subtype_Tubular |
|--|
| 3rd Qu.:0.00000 |
| Max. :1.00000 |

Table 11: Table continues below

| Tumor_Other_Histologic_Subtype_Medullary | Hormone_Therapy |
|--|-----------------|
| Min. :0.00000 | Min. :0.0000 |
| 1st Qu.:0.00000 | 1st Qu.:0.0000 |
| Median :0.00000 | Median :1.0000 |
| Mean :0.01513 | Mean :0.6042 |
| 3rd Qu.:0.00000 | 3rd Qu.:1.0000 |
| Max. :1.00000 | Max. :1.0000 |

Table 12: Table continues below

| Inferred_Menopausal_State | Primary_Tumor_Laterality | Overall_Survival_Status |
|---------------------------|--------------------------|-------------------------|
| Min. :0.0000 | Min. :0.000 | Min. :0.0000 |
| 1st Qu.:0.5000 | 1st Qu.:0.000 | 1st Qu.:0.0000 |
| Median :1.0000 | Median :1.000 | Median :1.0000 |
| Mean :0.7497 | Mean :0.525 | Mean :0.5565 |
| 3rd Qu.:1.0000 | 3rd Qu.:1.000 | 3rd Qu.:1.0000 |
| Max. :1.0000 | Max. :1.000 | Max. :1.0000 |

Table 13: Table continues below

| PR_Status | Radio_Therapy | Relapse_Free_Status |
|----------------|----------------|---------------------|
| Min. :0.0000 | Min. :0.0000 | Min. :0.0000 |
| 1st Qu.:0.0000 | 1st Qu.:0.0000 | 1st Qu.:0.0000 |
| Median :0.0000 | Median :1.0000 | Median :0.0000 |
| Mean :0.4994 | Mean :0.6799 | Mean :0.4214 |
| 3rd Qu.:1.0000 | 3rd Qu.:1.0000 | 3rd Qu.:1.0000 |
| Max. :1.0000 | Max. :1.0000 | Max. :1.0000 |

Table 14: Table continues below

| Gene_classifier_subtype_ERH | Gene_classifier_subtype_ERL |
|-----------------------------|-----------------------------|
| Min. :0.000 | Min. :0.0000 |
| 1st Qu.:0.000 | 1st Qu.:0.0000 |
| Median :0.000 | Median :0.0000 |
| Mean :0.362 | Mean :0.3132 |
| 3rd Qu.:1.000 | 3rd Qu.:1.0000 |
| Max. :1.000 | Max. :1.0000 |

Table 15: Table continues below

| Gene_classifier_subtype_ERM | Gene_classifier_subtype_ERP |
|-----------------------------|-----------------------------|
| Min. :0.0000 | Min. :0.0000 |
| 1st Qu.:0.0000 | 1st Qu.:0.0000 |
| Median :0.0000 | Median :0.0000 |
| Mean :0.1921 | Mean :0.1327 |
| 3rd Qu.:0.0000 | 3rd Qu.:0.0000 |
| Max. :1.0000 | Max. :1.0000 |

Table 16: Table continues below

| Patients_Vital_Status_Living | Patients_Vital_Status_Died |
|------------------------------|----------------------------|
| Min. :0.0000 | Min. :0.0000 |
| 1st Qu.:0.0000 | 1st Qu.:0.0000 |
| Median :0.0000 | Median :0.0000 |
| Mean :0.4435 | Mean :0.3481 |
| 3rd Qu.:1.0000 | 3rd Qu.:1.0000 |
| Max. :1.0000 | Max. :1.0000 |

| Patients_Vital_Status_Do | Integrative_Cluster |
|--------------------------|---------------------|
| Min. :0.0000 | Min. : 1.000 |
| 1st Qu.:0.0000 | 1st Qu.: 3.500 |
| Median :0.0000 | Median : 5.000 |
| Mean :0.2084 | Mean : 5.861 |
| 3rd Qu.:0.0000 | 3rd Qu.: 8.000 |
| Max. :1.0000 | Max. :10.000 |

```
pander(head(df3),caption='Head of data selection')
```

Table 18: Head of data selection (continued below)

| Age | Cohort | Neo_Grade | Lymph | Mutation | Nottingham | Overall_Month |
|-------|--------|-----------|-------|----------|------------|---------------|
| 43.19 | 1 | 3 | 0 | 2 | 4.02 | 84.63 |
| 78.77 | 1 | 3 | 0 | 4 | 4.062 | 7.8 |
| 86.41 | 1 | 3 | 1 | 4 | 5.032 | 36.57 |
| 85.49 | 1 | 2 | 0 | 1 | 3.044 | 132 |
| 45.43 | 1 | 3 | 0 | 5 | 4.046 | 140.9 |
| 61.49 | 1 | 2 | 1 | 3 | 4.032 | 157.4 |

Table 19: Table continues below

| Relapse_Month | Tumor_Size | Tumor_Stage | Type_of_Breast_Surgery |
|---------------|------------|-------------|------------------------|
| 83.52 | 10 | 1 | 1 |
| 2.89 | 31 | 4 | 0 |
| 36.09 | 16 | 2 | 1 |

| Relapse_Month | Tumor_Size | Tumor_Stage | Type_of_Breast_Surgery |
|---------------|------------|-------------|------------------------|
| 123.3 | 22 | 4 | 0 |
| 139 | 23 | 2 | 1 |
| 155.4 | 16 | 2 | 1 |

Table 20: Table continues below

| new_Chemotherapy | new_Cellularity | Pam50_Claudin_low_subtype_Luma |
|------------------|-----------------|--------------------------------|
| 0 | 3 | 1 |
| 0 | 2 | 0 |
| 0 | 2 | 0 |
| 0 | 2 | 1 |
| 1 | 3 | 0 |
| 0 | 3 | 0 |

Table 21: Table continues below

| Pam50_Claudin_low_subtype_LumB | Pam50_Claudin_low_subtype_Her2 |
|--------------------------------|--------------------------------|
| 0 | 0 |
| 1 | 0 |
| 1 | 0 |
| 0 | 0 |
| 1 | 0 |
| 1 | 0 |

Table 22: Table continues below

| Pam50_Claudin_low_subtype_Basal | Pam50_Claudin_low_subtype_Normal |
|---------------------------------|----------------------------------|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |

Table 23: Table continues below

| Pam50_Claudin_low_subtype_claudin | Pam50_Claudin_low_subtype_NC |
|-----------------------------------|------------------------------|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |

Table 24: Table continues below

| ER_status_measured_by_IHC | ER_Status_Positive | HER2_status_measured_by_SNP6 |
|---------------------------|--------------------|------------------------------|
| 1 | 1 | 4 |
| 1 | 1 | 4 |
| 1 | 1 | 3 |
| 1 | 1 | 4 |
| 1 | 1 | 4 |
| 1 | 1 | 4 |

Table 25: Table continues below

| HER2_Status_Positive | Tumor_Other_Histologic_Subtype_Ductal |
|----------------------|---------------------------------------|
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |

Table 26: Table continues below

| Tumor_Other_Histologic_Subtype_Tubular |
|--|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

Table 27: Table continues below

| Tumor_Other_Histologic_Subtype_Medullary | Hormone_Therapy |
|--|-----------------|
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |

Table 28: Table continues below

| Inferred_Menopausal_State | Primary_Tumor_Laterality | Overall_Survival_Status |
|---------------------------|--------------------------|-------------------------|
| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

| Inferred_Menopausal_State | Primary_Tumor_Laterality | Overall_Survival_Status |
|---------------------------|--------------------------|-------------------------|
| 0 | 0 | 0 |
| 1 | 1 | 0 |

Table 29: Table continues below

| PR_Status | Radio_Therapy | Relapse_Free_Status | Gene_classifier_subtype_ERH |
|-----------|---------------|---------------------|-----------------------------|
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 |
| 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 0 | 1 |

Table 30: Table continues below

| Gene_classifier_subtype_ERL | Gene_classifier_subtype_ERM |
|-----------------------------|-----------------------------|
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 1 | 0 |
| 0 | 0 |
| 0 | 0 |

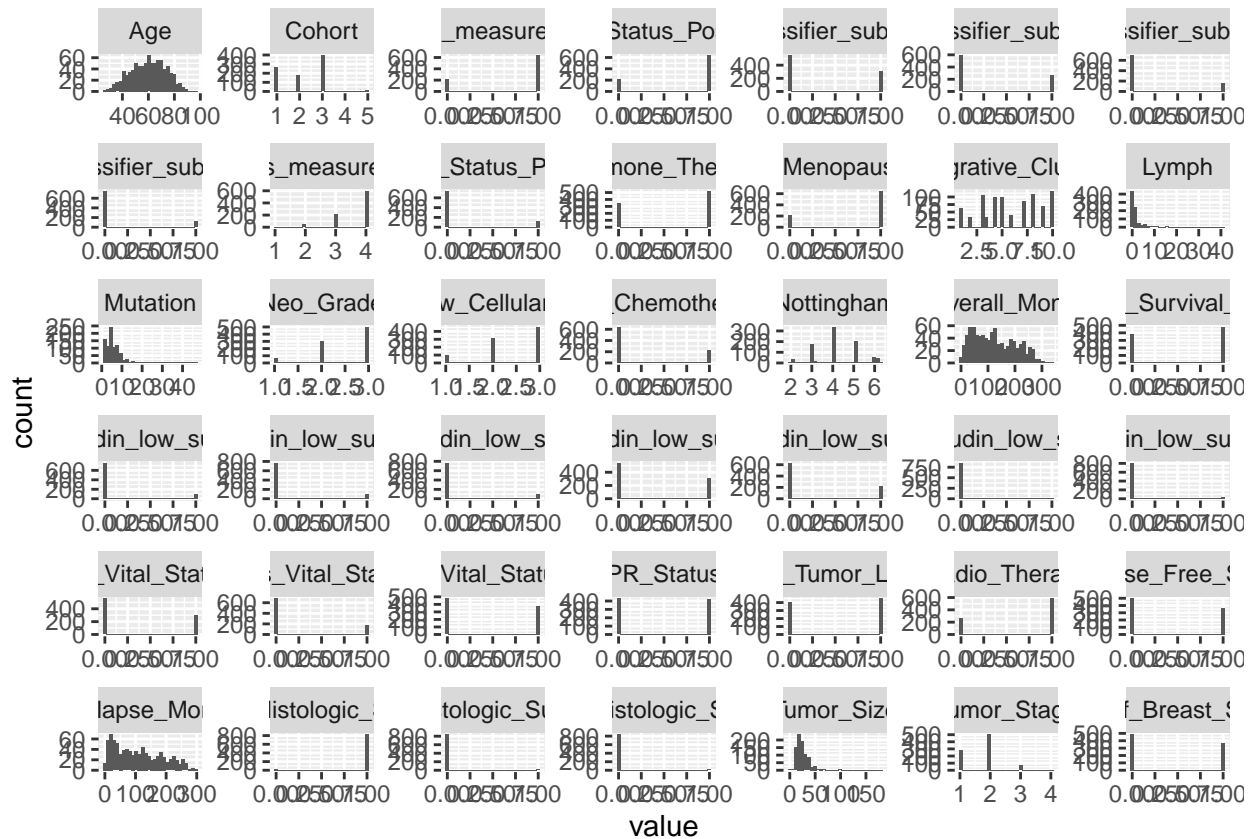
Table 31: Table continues below

| Gene_classifier_subtype_ERP | Patients_Vital_Status_Living |
|-----------------------------|------------------------------|
| 0 | 1 |
| 0 | 0 |
| 0 | 0 |
| 0 | 0 |
| 0 | 1 |
| 0 | 1 |

| Patients_Vital_Status_Died | Patients_Vital_Status_Do | Integrative_Cluster |
|----------------------------|--------------------------|---------------------|
| 0 | 0 | 4.5 |
| 1 | 0 | 7 |
| 0 | 1 | 9 |
| 1 | 0 | 3 |
| 0 | 0 | 10 |
| 0 | 0 | 7 |

```
df3 %>%
  keep(is.numeric) %>%
  gather() %>%
  ggplot(aes(value)) +
  facet_wrap(~ key, scales = "free") +
  geom_histogram()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
tr_ind <- 1:(nrow(df3) * 0.7)
df3_tr <- df3[tr_ind, ]
nrow(df3_tr)
```

```
## [1] 601
```

```
df3_te <- df3[-tr_ind, ]
nrow(df3_te)
```

```
## [1] 258
```

variable selection

```
set.seed(0)
fit_BIC <- regsubsets(Overall_Survival_Status ~ ., data = df3_tr)
```

```
## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax, force.in =
## force.in, : 4 linear dependencies found
```

```
## Reordering variables and trying again:
```



```
summary_BIC <- summary(fit_BIC)
```

```
## Warning in log(vr): NaNs produced
```

```
summary_BIC
```

```
## Subset selection object
## Call: regsubsets.formula(Overall_Survival_Status ~ ., data = df3_tr)
## 41 Variables (and intercept)
##
```

| | Forced in | Forced out |
|---|-----------|------------|
| ## Age | FALSE | FALSE |
| ## Cohort | FALSE | FALSE |
| ## Neo_Grade | FALSE | FALSE |
| ## Lymph | FALSE | FALSE |
| ## Mutation | FALSE | FALSE |
| ## Nottingham | FALSE | FALSE |
| ## Overall_Month | FALSE | FALSE |
| ## Relapse_Month | FALSE | FALSE |
| ## Tumor_Size | FALSE | FALSE |
| ## Tumor_Stage | FALSE | FALSE |
| ## Type_of_Breast_Surgery | FALSE | FALSE |
| ## new_Chemotherapy | FALSE | FALSE |
| ## new_Cellularity | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_Luma | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_LumB | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_Her2 | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_Basal | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_Normal | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_claudin | FALSE | FALSE |
| ## ER_status_measured_by_IHC | FALSE | FALSE |
| ## ER_Status_Positive | FALSE | FALSE |
| ## HER2_status_measured_by_SNP6 | FALSE | FALSE |
| ## HER2_Status_Positive | FALSE | FALSE |
| ## Tumor_Other_Histologic_Subtype_Ductal | FALSE | FALSE |
| ## Tumor_Other_Histologic_Subtype_Tubular | FALSE | FALSE |
| ## Hormone_Therapy | FALSE | FALSE |
| ## Inferred_Menopausal_State | FALSE | FALSE |
| ## Primary_Tumor_Laterality | FALSE | FALSE |
| ## PR_Status | FALSE | FALSE |
| ## Radio_Therapy | FALSE | FALSE |
| ## Relapse_Free_Status | FALSE | FALSE |
| ## Gene_classifier_subtype_ERH | FALSE | FALSE |
| ## Gene_classifier_subtype_ERL | FALSE | FALSE |
| ## Gene_classifier_subtype_ERM | FALSE | FALSE |
| ## Patients_Vital_Status_Living | FALSE | FALSE |
| ## Patients_Vital_Status_Died | FALSE | FALSE |
| ## Integrative_Cluster | FALSE | FALSE |
| ## Pam50_Claudin_low_subtype_NC | FALSE | FALSE |
| ## Tumor_Other_Histologic_Subtype_Medullary | FALSE | FALSE |
| ## Gene_classifier_subtype_ERP | FALSE | FALSE |
| ## Patients_Vital_Status_Do | FALSE | FALSE |

```
## 1 subsets of each size up to 9
```

```

## Selection Algorithm: exhaustive
##      Age Cohort Neo_Grade Lymph Mutation Nottingham Overall_Month
## 1 ( 1 ) " " " " " " " " " " " "
## 2 ( 1 ) " " " " " " " " " " " "
## 3 ( 1 ) "*" " " " " " " " " " "
## 4 ( 1 ) "*" " " " " " " " " " "
## 5 ( 1 ) "*" "*" "*" " " " " " "
## 6 ( 1 ) "*" "*" "*" "*" " " " "
## 7 ( 1 ) "*" "*" "*" "*" " " " "
## 8 ( 1 ) "*" "*" "*" "*" "*" "*" "*"
## 9 ( 1 ) "*" "*" "*" "*" "*" "*" "*"
##      Relapse_Month Tumor_Size Tumor_Stage Type_of_Breast_Surgery
## 1 ( 1 ) " " " " " " " "
## 2 ( 1 ) " " " " " " " "
## 3 ( 1 ) " " " " " " " "
## 4 ( 1 ) " " " " " " " "
## 5 ( 1 ) " " " " " " " "
## 6 ( 1 ) " " " " " " " "
## 7 ( 1 ) " " " " " " " "
## 8 ( 1 ) " " " " " " " "
## 9 ( 1 ) " " "*" " " " " "
##      new_Chemotherapy new_Cellularity Pam50_Claudin_low_subtype_Luma
## 1 ( 1 ) " " " " " "
## 2 ( 1 ) " " " " " "
## 3 ( 1 ) " " " " " "
## 4 ( 1 ) " " " " " "
## 5 ( 1 ) " " " " " "
## 6 ( 1 ) " " " " " "
## 7 ( 1 ) " " " " " "
## 8 ( 1 ) " " " " " "
## 9 ( 1 ) " " " " " "
##      Pam50_Claudin_low_subtype_LumB Pam50_Claudin_low_subtype_Her2
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " "*"
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      Pam50_Claudin_low_subtype_Basal Pam50_Claudin_low_subtype_Normal
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) "*" " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      Pam50_Claudin_low_subtype_claudin Pam50_Claudin_low_subtype_NC
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "

```

```

## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      ER_status_measured_by_IHC ER_Status_Positive
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      HER2_status_measured_by_SNP6 HER2_Status_Positive
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      Tumor_Other_Histologic_Subtype_Ductal
## 1 ( 1 ) " "
## 2 ( 1 ) " "
## 3 ( 1 ) " "
## 4 ( 1 ) " "
## 5 ( 1 ) " "
## 6 ( 1 ) " "
## 7 ( 1 ) " "
## 8 ( 1 ) " "
## 9 ( 1 ) " "
##      Tumor_Other_Histologic_Subtype_Tubular
## 1 ( 1 ) " "
## 2 ( 1 ) " "
## 3 ( 1 ) " "
## 4 ( 1 ) " "
## 5 ( 1 ) " "
## 6 ( 1 ) " "
## 7 ( 1 ) " "
## 8 ( 1 ) " "
## 9 ( 1 ) " "
##      Tumor_Other_Histologic_Subtype_Medullary Hormone_Therapy
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "

```

```

## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      Inferred_Menopausal_State Primary_Tumor_Laterality PR_Status
## 1 ( 1 ) " " " " " "
## 2 ( 1 ) " " " " " "
## 3 ( 1 ) " " " " " "
## 4 ( 1 ) " " " " " "
## 5 ( 1 ) " " " " " "
## 6 ( 1 ) " " " " " "
## 7 ( 1 ) " " " " " "
## 8 ( 1 ) " " " " " "
## 9 ( 1 ) " " " " " "
##      Radio_Therapy Relapse_Free_Status Gene_classifier_subtype_ERH
## 1 ( 1 ) " " " " " "
## 2 ( 1 ) " " "*" " " "
## 3 ( 1 ) " " "*" " " "
## 4 ( 1 ) " " "*" " " "
## 5 ( 1 ) " " " " " "
## 6 ( 1 ) " " " " " "
## 7 ( 1 ) " " " " " "
## 8 ( 1 ) " " " " " "
## 9 ( 1 ) " " " " " "
##      Gene_classifier_subtype_ERL Gene_classifier_subtype_ERM
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " " "
##      Gene_classifier_subtype_ERP Patients_Vital_Status_Living
## 1 ( 1 ) " " "*"
## 2 ( 1 ) " " "*"
## 3 ( 1 ) " " "*"
## 4 ( 1 ) " " "*"
## 5 ( 1 ) " " " "
## 6 ( 1 ) " " " "
## 7 ( 1 ) " " " "
## 8 ( 1 ) " " " "
## 9 ( 1 ) " " "*"
##      Patients_Vital_Status_Died Patients_Vital_Status_Do
## 1 ( 1 ) " " " "
## 2 ( 1 ) " " " "
## 3 ( 1 ) " " " "
## 4 ( 1 ) " " " "
## 5 ( 1 ) "*" "*"
## 6 ( 1 ) "*" "*"
## 7 ( 1 ) "*" "*"
## 8 ( 1 ) "*" "*"
## 9 ( 1 ) " " " "
##      Integrative_Cluster

```

```
## 1 ( 1 ) " "
## 2 ( 1 ) " "
## 3 ( 1 ) " "
## 4 ( 1 ) " "
## 5 ( 1 ) " "
## 6 ( 1 ) " "
## 7 ( 1 ) " "
## 8 ( 1 ) " "
## 9 ( 1 ) " "
```

```
min_BIC <- which.min(summary_BIC$bic)
min_BIC
```

```
## [1] 8
```

```
coef_BIC = coef(fit_BIC,min_BIC)
```

```
## Warning in log(vr): NaNs produced
```

```
coef_BIC
```

```
##                (Intercept)
##                -0.677629271
##                Age
##                0.010622924
##                Cohort
##                0.145974676
##                Neo_Grade
##                0.028670848
##                Lymph
##                0.018462621
##                Mutation
##                0.004245253
##                Nottingham
##                0.037713480
## Tumor_Other_Histologic_Subtype_Medullary
##                0.166695765
##                Gene_classifier_subtype_ERP
##                0.103582951
```

```
#Forward Stepwise Selection with Adjusted R squared
```

```
#fit_FORWARD <- regsubsets(Overall_Survival_Status ~ ., data = df3_tr, method = "forward", numax = 10)
#as.factor(df3_tr$Cancer_Type_Detailed)
#summary_FORWARD <- summary(fit_FORWARD)
#max_FORWARD <- which.max(summary_FORWARD$adjr2)
#max_FORWARD
#coef_FORWARD = coef(fit_FORWARD, max_FORWARD)
#coef_FORWARD
```

```
Backward Stepwise Selection with Cp
```

```
fit_BACKWARD <- regsubsets(Overall_Survival_Status ~ ., data = df3_tr, method = "backward", nvmax = ncol(df3_tr) - 1)
```

```
## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax, force.in =  
## force.in, : 4 linear dependencies found
```

```
## Reordering variables and trying again:
```

```
## Warning in rval$lopt[] <- rval$vorder[rval$lopt]: number of items to replace is  
## not a multiple of replacement length
```

```
summary_BACKWARD <- summary(fit_BACKWARD)  
mix_BACKWARD <- which.min(summary_BACKWARD$cp)  
mix_BACKWARD
```

```
## [1] 19
```

```
coef_BACKWARD = coef(fit_BACKWARD, mix_BACKWARD)  
coef_BACKWARD
```

```
##              (Intercept)  
##              0.101529777  
##              Age  
##              0.008725082  
##              Cohort  
##              0.205473176  
##              Lymph  
##              0.014394673  
##              Overall_Month  
##              -0.001669157  
##              Relapse_Month  
##              -0.001376028  
##              Tumor_Stage  
##              -0.037640439  
## Pam50_Claudin_low_subtype_Luma  
##              -0.118693753  
## Pam50_Claudin_low_subtype_LumB  
##              -0.069504986  
## Pam50_Claudin_low_subtype_Basal  
##              -0.065317856  
## Pam50_Claudin_low_subtype_claudin  
##              -0.117549258  
##              HER2_Status_Positive  
##              0.022130151  
## Tumor_Other_Histologic_Subtype_Tubular  
##              0.033264067  
##              Hormone_Therapy  
##              -0.057071936  
##              PR_Status  
##              0.021133168  
## Gene_classifier_subtype_ERH  
##              0.090267981
```

```
##           Gene_classifier_subtype_ERL
##                               -0.003118641
##           Gene_classifier_subtype_ERM
##                               0.006531345
##           Pam50_Claudin_low_subtype_NC
##                               0.412900231
## Tumor_Other_Histologic_Subtype_Medullary
##                               0.146141437
```

```
set.seed(0)
predict_BIC = glm(Overall_Survival_Status~Age + Cohort + Neo_Grade + Lymph + Mutation + Nottingham + Tumor_Other_Histologic_Subtype_Medullary, data=df3_te)
pred_BIC = round(predict(predict_BIC,df3_te,type = "response"))
error_BIC = mean((df3_te$Overall_Survival_Status - pred_BIC)^2)
error_BIC
```

```
## [1] 0.3333333
```

```
#predict_FORWARD = glm(Overall_Survival_Status ~ Cohort + Neo_Grade + Overall_Month + Relapse_Month + Nottingham + Tumor_Other_Histologic_Subtype_Medullary, data=df3_te)
#pred_FORWARD = round(predict(predict_FORWARD,df3_te,type = "response"))
#pred_FORWARD
#error_FORWARD = mean((df3_te$Overall_Survival_Status - pred_FORWARD)^2)
#error_FORWARD
```

```
predict_BACKWARD = glm(Overall_Survival_Status ~ Age+ Cohort + Lymph + Overall_Month + Relapse_Month + Nottingham + Tumor_Other_Histologic_Subtype_Medullary, data=df3_tr)
pred_BACKWARD = round(predict(predict_BACKWARD,df3_tr,type = "response"))
error_BACKWARD = mean((df3_te$Overall_Survival_Status - pred_BACKWARD)^2)
error_BACKWARD
```

```
## [1] 0.5291181
```

```
which.min(data.frame(error_BIC, error_BACKWARD)) # All the errors listed
```

```
## error_BIC
##           1
```

```
formula = Overall_Survival_Status~Age + Cohort + Neo_Grade + Lymph + Mutation + Nottingham + Tumor_Other_Histologic_Subtype_Medullary
```

Model 1

```
set.seed(0)
library(randomForest)
rf_model <- randomForest(formula, data = df3_tr,ntree = 10,importance = T)
predictions <- predict(rf_model, df3_te)
error_rf <- mean((df3_te$Overall_Survival_Status - predictions)^2)
error_rf
```

```
## [1] 0.217086
```

Model 2

```

set.seed(0)
knn_test_error <- vector()
new <- data.frame(k=numeric(), knn_test_error)

for(i in seq(from = 1, to = 100, by = 5)) {
  knnmodel <- knn3(formula, df3_tr, k = i)
  knn_test <- predict(knnmodel, newdata = df3_te)
  knn_test_error <- mean((knn_test - df3_te$Overall_Survival_Status)^2)
  new[i,] <- c(i, knn_test_error)
}
frame = data.frame(knn_test_error,1:24)
error_knn <- data.frame(knn_test_error,1:24)[which.min(frame$knn_test_error),]
error_knn

```

```

## knn_test_error X1.24
## 1 0.2786025 1

```

Model 3 Ridge

```

x_tr<-as.matrix(df3_tr[,c(1:23, 25:ncol(df3_tr))])
y_tr<-as.matrix(df3_tr[,24])
x_te<-as.matrix(df3_te[,c(1:23, 25:ncol(df3_te))])
y_te<-as.matrix(df3_te[,24])

```

```
library(glmnet)
```

```
## Warning: package 'glmnet' was built under R version 4.2.2
```

```
## Warning: package 'Matrix' was built under R version 4.2.2
```

```

set.seed(0)
ridge<-glmnet(x=x_tr,y=y_tr,alpha=0)
#plot(ridge,xvar='lambda')

ridge_cv<-cv.glmnet(x=x_tr,y=y_tr,type.measure='mse',nfold=10,alpha=0)

best_ridge<-coef(ridge_cv, s = ridge_cv$lambda.min)

result_ridge<-predict(ridge_cv,newx=x_te,interval='prediction')
error_ridge <- (err_ridge<-mean((y_te-result_ridge)^2))
error_ridge

```

```
## [1] 0.04608354
```

Model 4 Lasso

```

set.seed(0)
lasso<-glmnet(x=x_tr,y=y_tr,alpha=1)
#plot(lasso,xvar='lambda')

```



```

lasso_cv<-cv.glmnet(x=x_tr,y=y_tr,type.measure='mse',nfold=10,alpha=1,keep=T)
#lasso_cv$lambda.min

result_la<-predict(lasso_cv,newx=x_te,interval='prediction')
error_lasso <- (err_la<-mean((y_te-result_la)^2))
error_lasso

```

```
## [1] 0.04679185
```

Model Selection

```
which.min(data.frame(error_rf, error_knn, error_ridge, error_lasso)) # All the errors listed
```

```
## error_ridge
##           4
```

Model Present cross validation