

Internal_Lab_Assignment_Clean_Try_1

```
library(ggplot2)
library(dplyr)

## 
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
## 
##     filter, lag

## The following objects are masked from 'package:base':
## 
##     intersect, setdiff, setequal, union

library(mice)

## 
## Attaching package: 'mice'

## The following object is masked from 'package:stats':
## 
##     filter

## The following objects are masked from 'package:base':
## 
##     cbind, rbind

library(tidyr)
library(Rtsne)

## Warning: package 'Rtsne' was built under R version 4.4.3

library(ggplot2)
library(factoextra)

## Warning: package 'factoextra' was built under R version 4.4.3

## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```

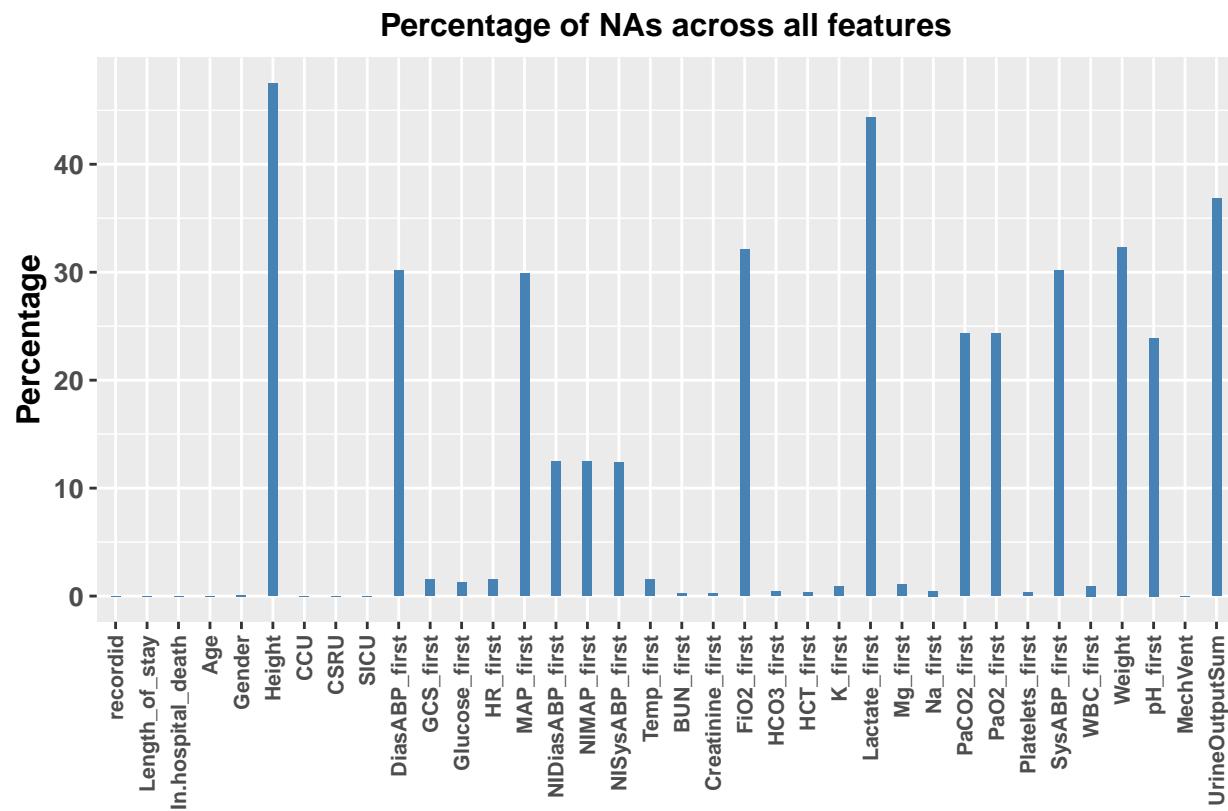
library(cluster)

file = 'Data/ICU_filtered.csv'
dfICU = read.csv(file, header = TRUE, stringsAsFactors = TRUE)

pMissDF = setNames(stack(sapply(dfICU, function(x){(sum(is.na(x))/length(x))*100}))[2:1], c('Feature', 'Value'))
p = ggplot(data = pMissDF, aes(x = Feature, y = Value)) +
  geom_bar(stat = 'identity', fill = 'steelblue', width = 0.3) +
  theme(text = element_text(size = 14, face = 'bold'),
        axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5)) +
  xlab('') + ylab('Percentage') +
  ggtitle('Percentage of NAs across all features') +
  theme(plot.title = element_text(size = 12, hjust = 0.5),
        axis.text = element_text(size = 8),
        axis.text.x = element_text(size = 8),
        axis.text.y = element_text(size = 10),
        axis.title = element_text(size = 12, face = "bold"))

p

```



```
dfICU = dfICU %>% select(-c(pMissDF[pMissDF['Value'] > 20, 'Feature']))
```

```

dfICU[dfICU['CCU'] == 1, 'ICU'] = 1
dfICU[dfICU['CSRU'] == 1, 'ICU'] = 2
dfICU[dfICU['SICU'] == 1, 'ICU'] = 3

```

```
dfICU[(dfICU['CCU'] == 0) & (dfICU['CSR'] == 0) & (dfICU['SICU'] == 0), 'ICU'] = 4
dfICU = dfICU %>% select(-c(CCU, CSR, SICU, recordid, In.hospital_death, Length_of_stay))
```

```
str(dfICU)
```

```
## 'data.frame':    7886 obs. of  20 variables:
##   $ Age          : int  54 76 44 68 88 64 68 78 64 74 ...
##   $ Gender       : int  0 1 0 1 0 1 0 0 0 1 ...
##   $ GCS_first    : int  15 3 7 15 15 7 15 15 15 10 ...
##   $ Glucose_first: int  205 105 141 129 113 264 94 132 113 106 ...
##   $ HR_first     : num  73 88 100 79 93 78 73 111 127 67 ...
##   $ NIDiasABP_first: int  65 38 84 63 41 89 88 51 71 57 ...
##   $ NIMAP_first  : num  92.3 49.3 100.3 86.7 75.3 ...
##   $ NISysABP_first: int  147 72 133 134 144 129 187 100 138 134 ...
##   $ Temp_first   : num  35.1 35.2 37.8 36.3 37.8 35.8 36.3 38 37.3 34.8 ...
##   $ BUN_first    : int  13 16 8 23 45 15 32 81 21 19 ...
##   $ Creatinine_first: num  0.8 0.8 0.4 0.9 1 1.4 3.4 0.9 0.7 1.1 ...
##   $ HCO3_first   : int  26 21 24 28 18 19 25 18 21 23 ...
##   $ HCT_first    : num  33.7 24.7 28.5 41.3 22.6 41.6 31.9 32.6 28.3 31.5 ...
##   $ K_first      : num  4.4 4.3 3.3 4 6 5.1 3.7 4.2 3.9 4.6 ...
##   $ Mg_first     : num  1.5 3.1 1.9 2.1 1.5 1.7 1.9 2.2 1.6 1.8 ...
##   $ Na_first     : int  137 139 137 140 140 141 140 141 139 141 ...
##   $ Platelets_first: int  221 164 72 391 109 276 325 91 696 141 ...
##   $ WBC_first    : num  11.2 7.4 4.2 11.5 3.8 24 6.2 16.1 15.2 9 ...
##   $ MechVent    : int  0 1 1 0 0 1 0 1 0 1 ...
##   $ ICU          : num  3 2 4 4 4 1 4 4 4 2 ...
```

```
dfICU = complete(mice(dfICU, m = 40, maxit = 10, pred = quickpred(dfICU, minpuc = 0.25), seed = 500))
```

```
##
## iter imp variable
##  1  1  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  2  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  3  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  4  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  5  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  6  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  7  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  8  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  9  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  10  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  11  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  12  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  13  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  14  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  15  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  16  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  17  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  18  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  19  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  20  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
##  1  21  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first  To
```



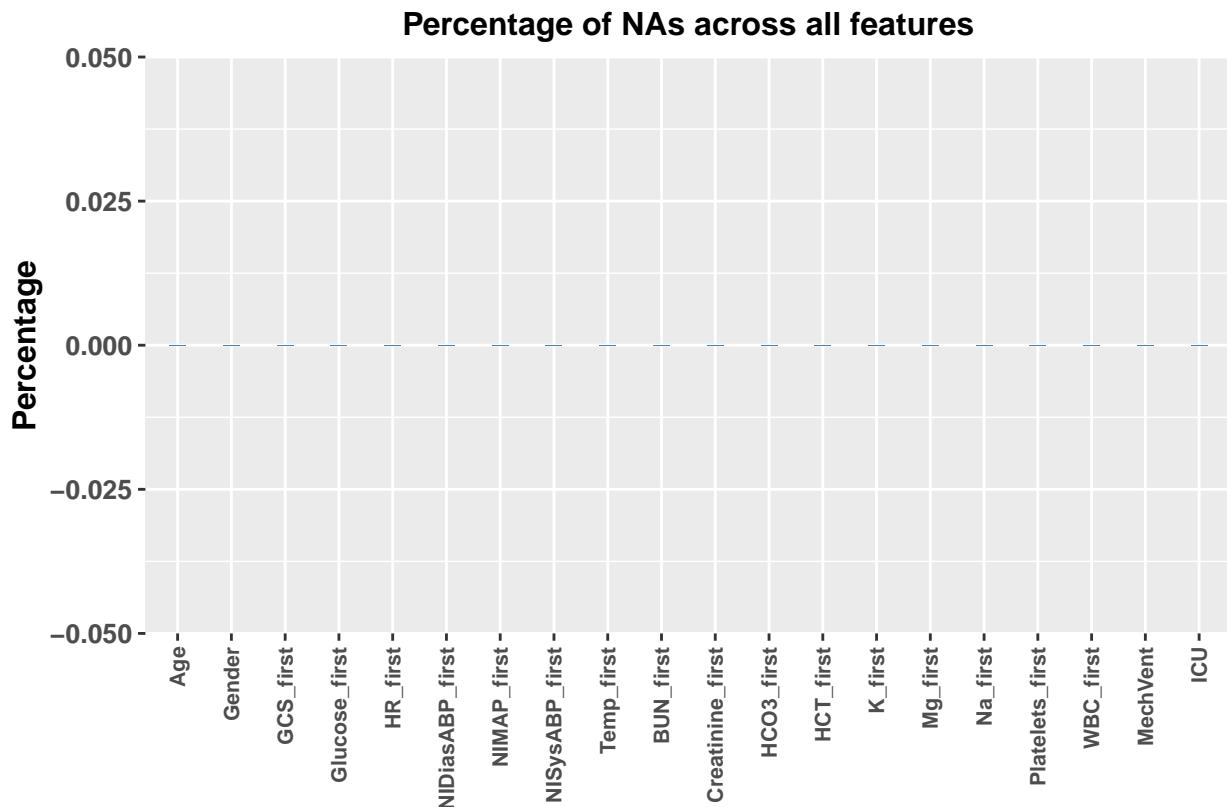
```

##    10   40  Gender  GCS_first  Glucose_first  HR_first  NIDiasABP_first  NIMAP_first  NISysABP_first

pMissDF = setNames(stack(sapply(dfICU, function(x){(sum(is.na(x))/length(x))*100}))[2:1], c('Feature', 'Value'))
p = ggplot(data = pMissDF, aes(x = Feature, y = Value)) +
  geom_bar(stat = 'identity', fill = 'steelblue', width = 0.3) +
  theme(text = element_text(size = 14, face = 'bold'),
        axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5)) +
  xlab('') + ylab('Percentage') +
  ggtitle('Percentage of NAs across all features') +
  theme(plot.title = element_text(size = 12, hjust = 0.5),
        axis.text = element_text(size = 8),
        axis.text.x = element_text(size = 8),
        axis.text.y = element_text(size = 10),
        axis.title = element_text(size = 12, face = "bold"))

p

```



```

features = colnames(dfICU)
categorical_features = c('Gender', 'GCS_first', 'MechVent', 'ICU')
continuous_features = features[c(!(colnames(dfICU) %in% categorical_features))]

dfICU_continuous = dfICU %>% select(continuous_features)

```

```

## Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
## i Please use 'all_of()' or 'any_of()' instead.

```

```

##  # Was:
##  data %>% select(continuous_features)
##
##  # Now:
##  data %>% select(all_of(continuous_features))
##
## See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

```

```
dfICU_categorical = dfICU %>% select(categorical_features)
```

```

## Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
## i Please use `all_of()` or `any_of()` instead.
##  # Was:
##  data %>% select(categorical_features)
##
##  # Now:
##  data %>% select(all_of(categorical_features))
##
## See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

```

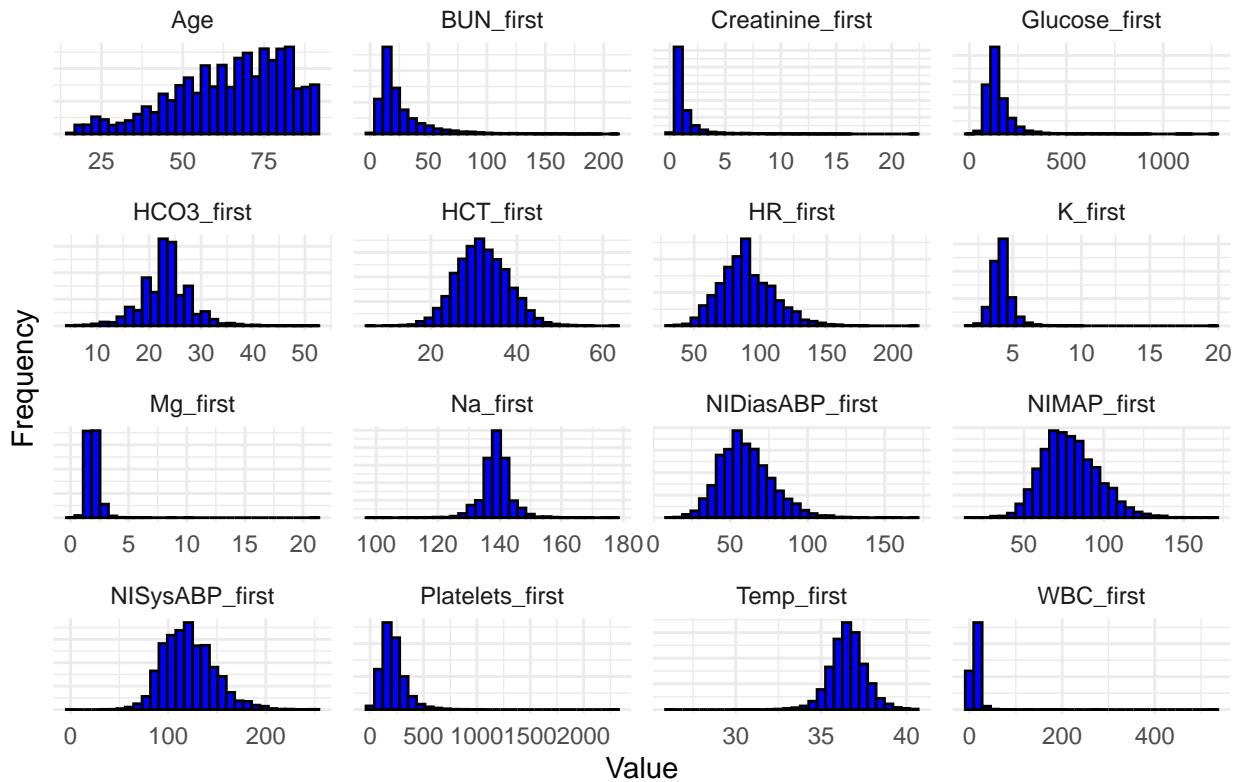
```

dfICU_long = dfICU_continuous %>% pivot_longer(cols = everything(), names_to = "variable", values_to = )

ggplot(dfICU_long, aes(x = value)) +
  geom_histogram(bins = 30, fill = "blue", color = "black") +
  facet_wrap(~ variable, scales = "free") +
  labs(title = "Histograms of Continuous Features",
       x = "Value",
       y = "Frequency") +
  theme_minimal() +
  theme(axis.text.y = element_blank(),
        axis.ticks.y = element_blank())

```

Histograms of Continuous Features



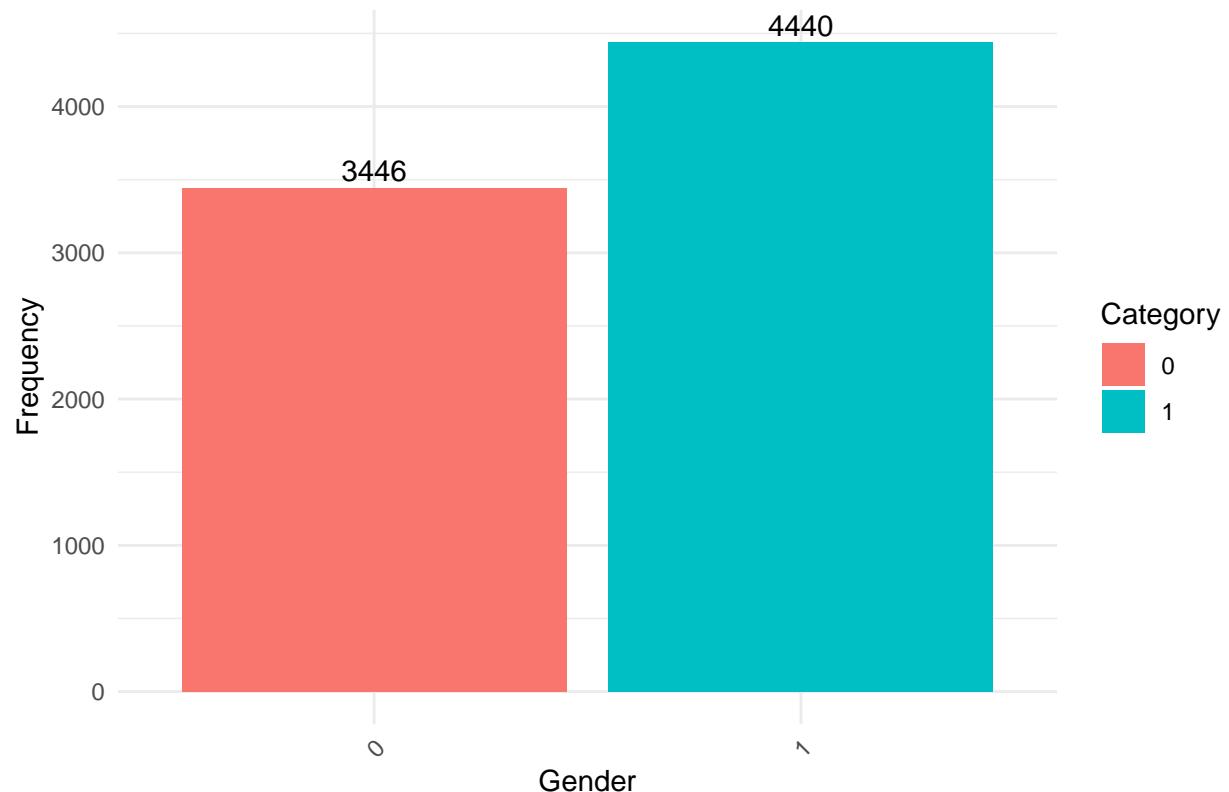
```
categorical_cols <- names(dfICU_categorical)

for (col in categorical_cols) {
  freq_table <- table(dfICU_categorical[[col]])
  freq_df <- as.data.frame(freq_table)
  colnames(freq_df) <- c("Category", "Frequency")

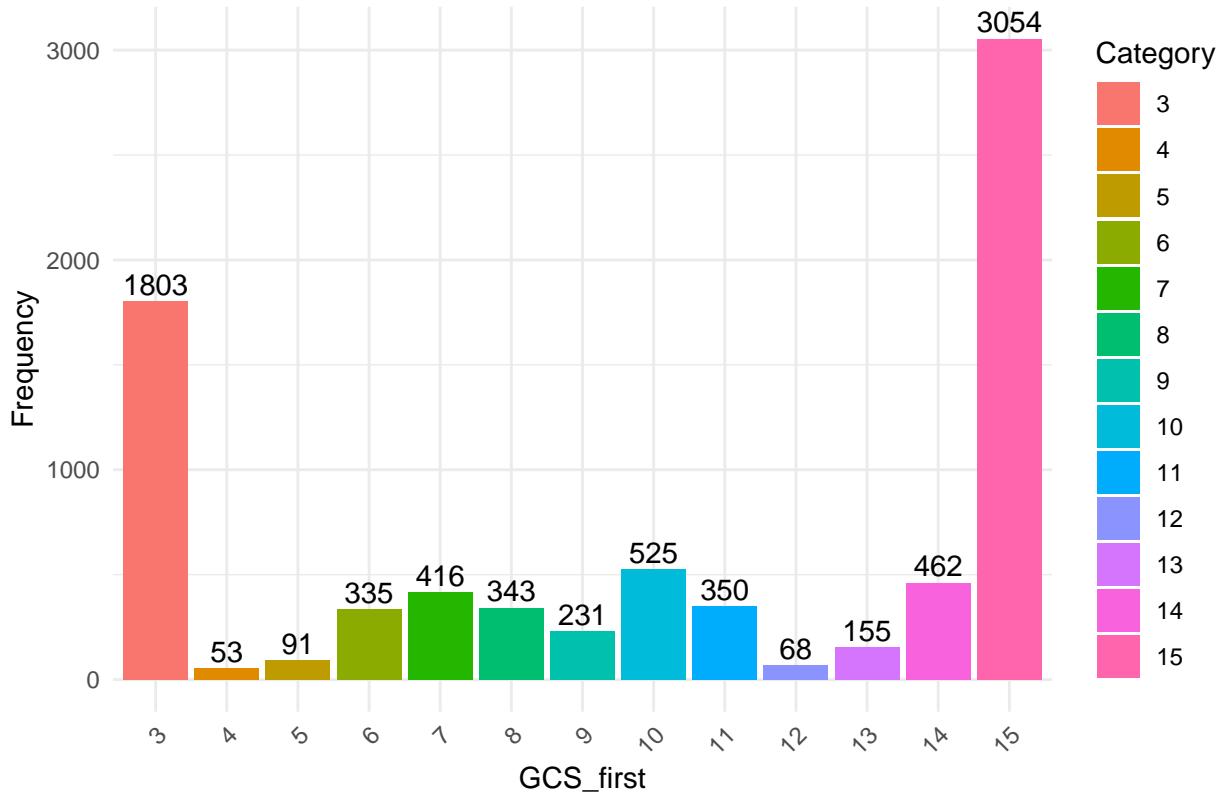
  p <- ggplot(freq_df, aes(x = Category, y = Frequency, fill = Category)) +
    geom_bar(stat = "identity") +
    geom_text(aes(label = Frequency), vjust = -0.3) +
    labs(title = paste("Distribution of", col),
        x = col,
        y = "Frequency") +
    theme_minimal() +
    theme(axis.text.x = element_text(angle = 45, hjust = 1))

  print(p)
}
```

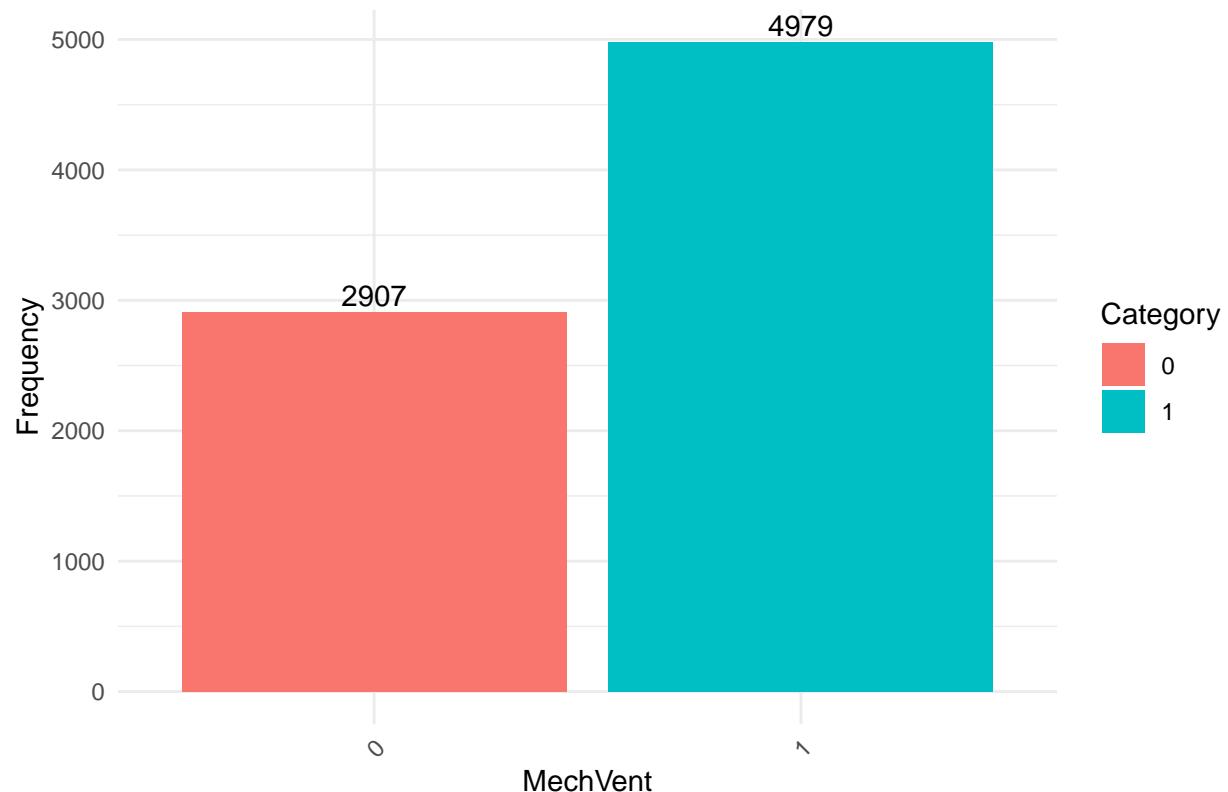
Distribution of Gender



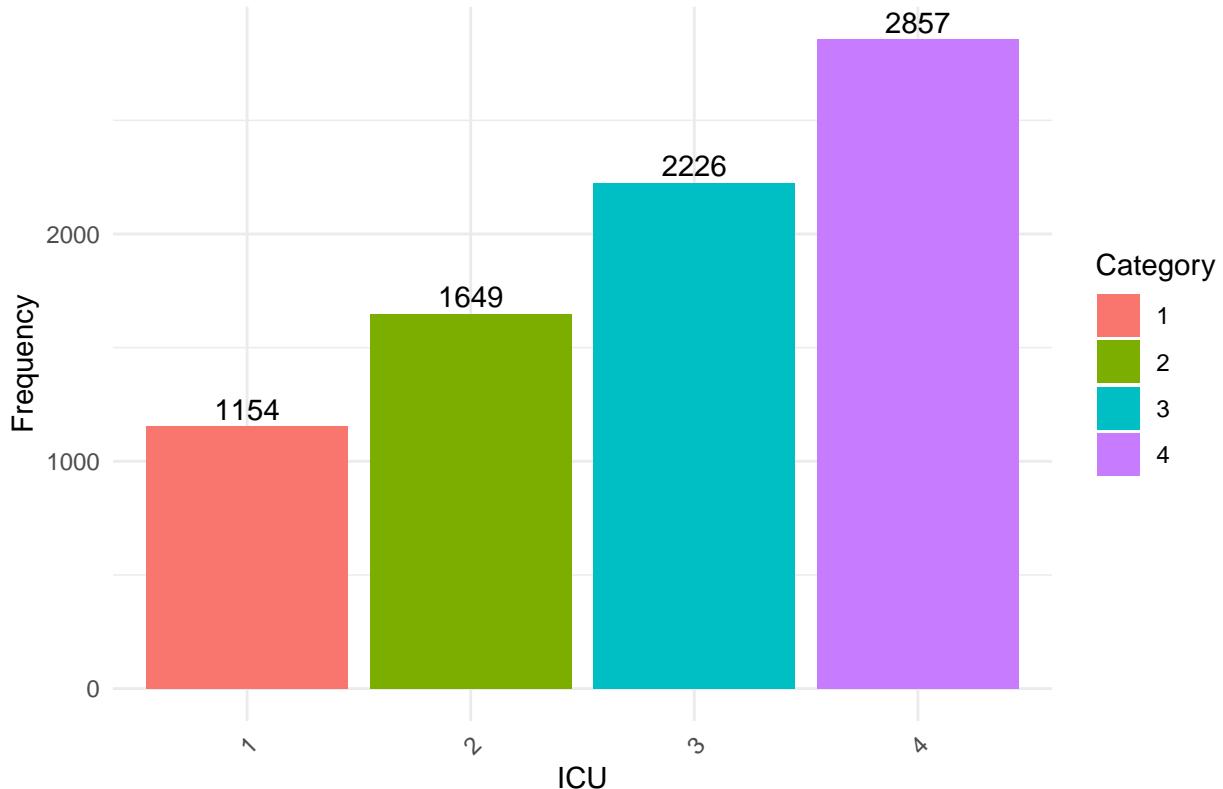
Distribution of GCS_first



Distribution of MechVent



Distribution of ICU



```

# Define categorical features explicitly
categorical_features <- c('Gender', 'GCS_first', 'MechVent', 'ICU')

# Convert categorical columns to factors
dfICU[, categorical_features] <- lapply(dfICU[, categorical_features], as.factor)

# Compute Gower distance across all rows and all features
gower_dist <- daisy(dfICU, metric = "gower")

# Convert to a full distance matrix (optional, but useful for inspection or clustering)
gower_matrix <- as.matrix(gower_dist)

# Compute average distance to all other rows (excluding self)
gower_avg_dist <- apply(as.matrix(gower_dist), 1, function(x) mean(x[x != 0]))

# Add to dfICU
dfICU_continuous$GowerAvgDist <- gower_avg_dist

str(dfICU_continuous)

## 'data.frame':    7886 obs. of  17 variables:
## $ Age           : int  54 76 44 68 88 64 68 78 64 74 ...
## $ Glucose_first : int  205 105 141 129 113 264 94 132 113 106 ...
## $ HR_first      : num  73 88 100 79 93 78 73 111 127 67 ...
## $ NIDiasABP_first : int  65 38 84 63 41 89 88 51 71 57 ...

```

```

## $ NIMAP_first      : num  92.3 49.3 100.3 86.7 75.3 ...
## $ NISysABP_first   : int  147 72 133 134 144 129 187 100 138 134 ...
## $ Temp_first       : num  35.1 35.2 37.8 36.3 37.8 35.8 36.3 38 37.3 34.8 ...
## $ BUN_first        : int  13 16 8 23 45 15 32 81 21 19 ...
## $ Creatinine_first : num  0.8 0.8 0.4 0.9 1 1.4 3.4 0.9 0.7 1.1 ...
## $ HCO3_first       : int  26 21 24 28 18 19 25 18 21 23 ...
## $ HCT_first        : num  33.7 24.7 28.5 41.3 22.6 41.6 31.9 32.6 28.3 31.5 ...
## $ K_first          : num  4.4 4.3 3.3 4 6 5.1 3.7 4.2 3.9 4.6 ...
## $ Mg_first         : num  1.5 3.1 1.9 2.1 1.5 1.7 1.9 2.2 1.6 1.8 ...
## $ Na_first         : int  137 139 137 140 140 141 140 141 139 141 ...
## $ Platelets_first  : int  221 164 72 391 109 276 325 91 696 141 ...
## $ WBC_first        : num  11.2 7.4 4.2 11.5 3.8 24 6.2 16.1 15.2 9 ...
## $ GowerAvgDist    : num  0.19 0.192 0.198 0.178 0.2 ...

# Original column list
cols_to_transform <- c("Age", "BUN_first", "categorical_avg_distance",
                      "categorical_dist_to_ref", "Creatinine_first",
                      "Glucose_first", "K_first", "Mg_first",
                      "Platelets_first", "WBC_first")

# Filter out non-existent columns
cols_to_transform_valid <- intersect(cols_to_transform, colnames(dfICU_continuous))

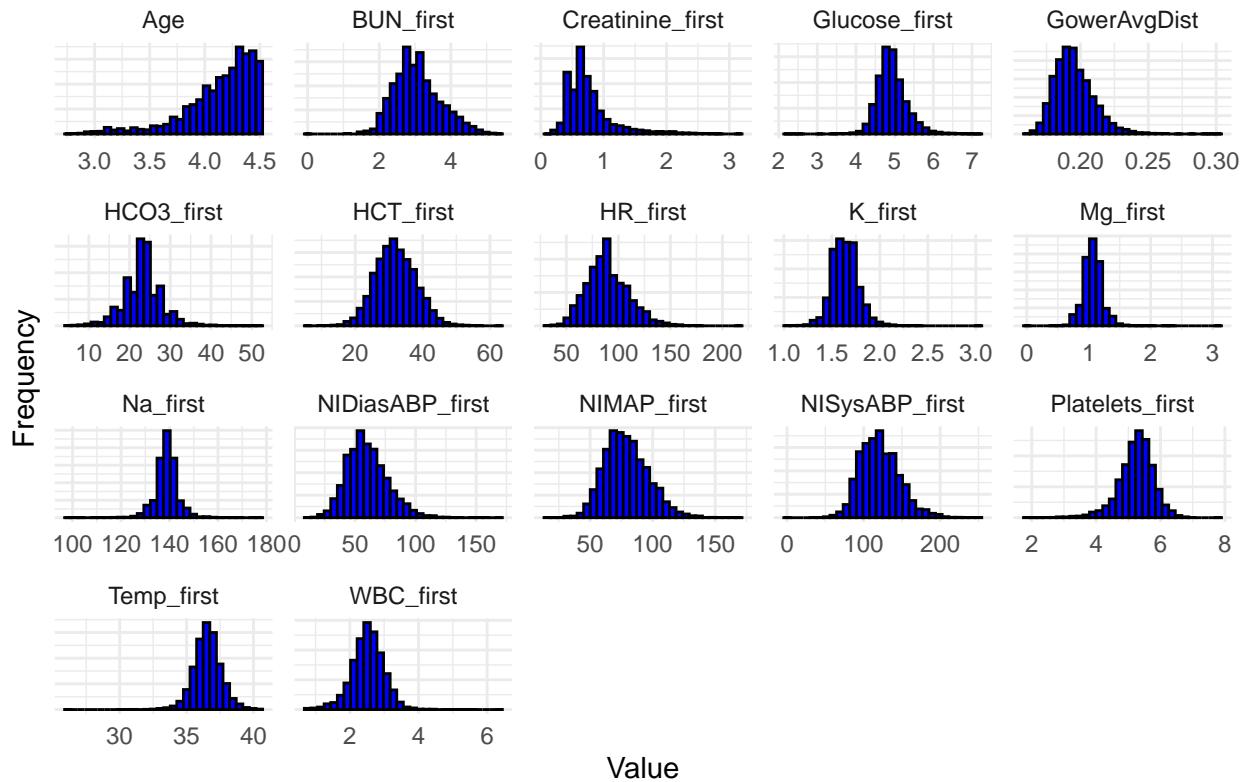
# Apply log1p transformation only on valid columns
dfICU_transformed <- dfICU_continuous %>%
  mutate(across(all_of(cols_to_transform_valid), ~ log1p(.)))

# Pivot to long format
dfICU_long_transformed <- dfICU_transformed %>%
  pivot_longer(cols = everything(), names_to = "variable", values_to = "value")

# Plot histograms
ggplot(dfICU_long_transformed, aes(x = value)) +
  geom_histogram(bins = 30, fill = "blue", color = "black") +
  facet_wrap(~ variable, scales = "free") +
  labs(title = "Histograms of Continuous Features (Log-Transformed Where Applicable)",
       x = "Value",
       y = "Frequency") +
  theme_minimal() +
  theme(axis.text.y = element_blank(),
        axis.ticks.y = element_blank())

```

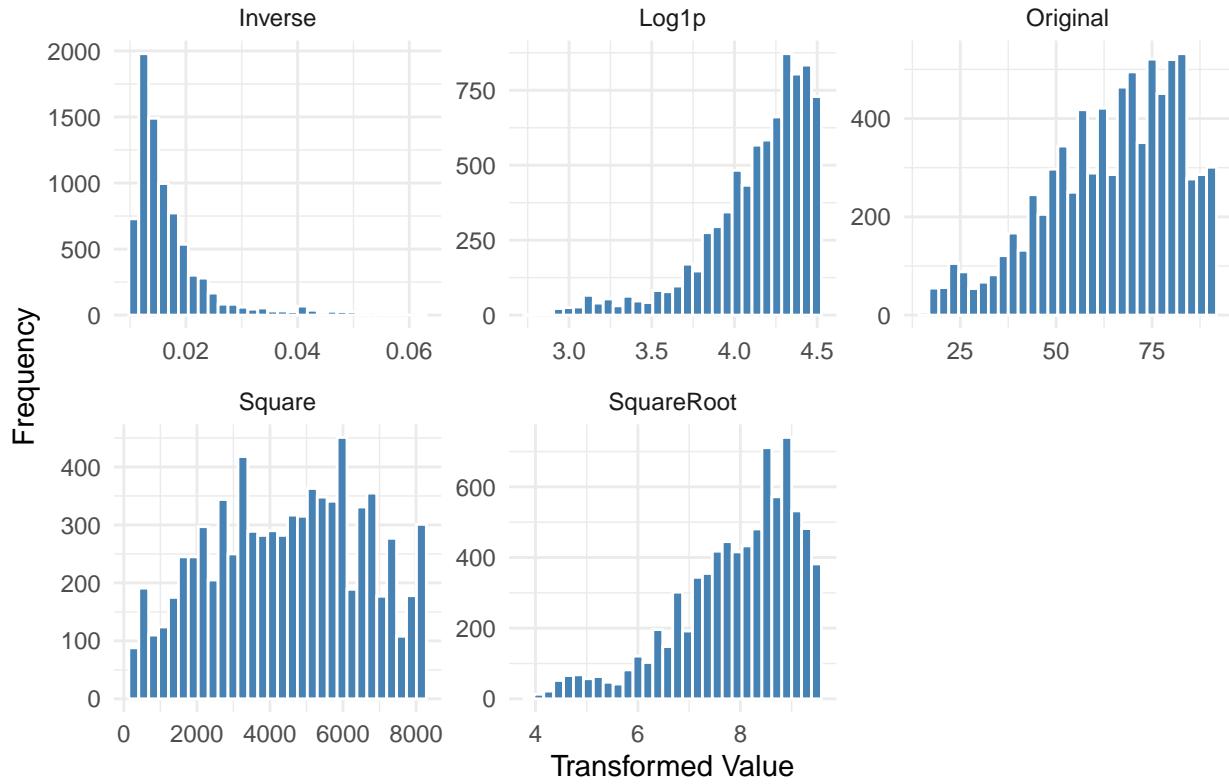
Histograms of Continuous Features (Log-Transformed Where Applicable)



```
# Create a new data frame with Age transformations
df_age_transformed <- dfICU_continuous %>%
  transmute(
    Original = Age,
    Log1p = log1p(Age),
    SquareRoot = sqrt(Age),
    Square = Age^2,
    Inverse = 1 / (Age + 1) # avoid divide-by-zero just in case
  ) %>%
  pivot_longer(cols = everything(), names_to = "Transformation", values_to = "Value")

# Plot histograms of all transformations of Age
ggplot(df_age_transformed, aes(x = Value)) +
  geom_histogram(bins = 30, fill = "steelblue", color = "white") +
  facet_wrap(~ Transformation, scales = "free") +
  labs(title = "Transformations of Age Feature",
       x = "Transformed Value",
       y = "Frequency") +
  theme_minimal()
```

Transformations of Age Feature



```

# Clean and scale the data
df_clean <- na.omit(dfICU_transformed)
df_scaled <- scale(df_clean)

# ----- 1. PCA Scores -----
pca_result <- prcomp(df_scaled, center = TRUE, scale. = TRUE)

# Get PCA scores (coordinates in the principal component space)
pca_scores <- as.data.frame(pca_result$x)

# View first few rows of PCA scores
head(pca_scores)

```

```

##          PC1         PC2         PC3         PC4         PC5         PC6         PC7
## 1 -1.240807 -0.2788929  0.6051838  0.5926341  1.9308423 -0.08773139  0.05782075
## 2  3.261596 -1.1195016  0.7937084  0.1085575 -0.1186429  0.58398941  0.99979868
## 3 -2.341497 -1.4264655  0.6447906 -2.5971002 -0.1864642  0.70820115 -0.53431537
## 4 -1.027961 -0.2581753  0.8354423  2.0913222 -0.3147096 -0.19566675  0.19433030
## 5  1.426883  0.7581013  1.0301603 -2.0691249 -0.3649144  0.61069540 -1.26806320
## 6 -1.796893  1.8587572 -0.8556086  1.3608710  1.9222207 -0.55099963  0.64285212
##          PC8         PC9         PC10        PC11        PC12        PC13        PC14
## 1  0.2566928  0.4149775  0.1488125 -0.9811031 -0.06078672 -0.4689355 -0.5330022
## 2  1.0126072 -0.2807118  1.0470405  1.0414343 -0.33877168  1.0775435  0.3807416
## 3  0.6380201  0.6432960  0.3262887  0.9225268  0.57006181  0.1948819  0.3244991
## 4 -0.1020917  0.2112627 -0.6989318  0.1562970 -0.39468687  0.4475666 -0.7322000
## 5 -0.1245768 -0.4358454 -0.7099052 -2.1083191 -0.32259366  0.7960597  1.9722530

```

```

## 6  0.8261914 -0.2427970 -1.0012560 -0.5614419  0.22396452 -0.2737824  0.3605981
##          PC15         PC16         PC17
## 1  0.4624728 -0.085079608 -0.04952401
## 2 -0.3241927 -0.132126000 -0.09156405
## 3 -0.5214371  0.003997787 -0.07661933
## 4  0.4431521  0.271877982 -0.07185698
## 5  0.7765094  0.934066046  0.01223954
## 6 -1.0623206 -0.480091549 -0.06682070

# Standard deviations of the principal components
std_devs <- pca_result$sdev

# Variance explained by each PC
var_explained <- std_devs^2
prop_var_explained <- var_explained / sum(var_explained)

# Create a data frame of variance explained
pca_variance_df <- data.frame(
  PC = paste0("PC", 1:length(prop_var_explained)),
  Variance_Explained = prop_var_explained,
  Cumulative_Variance = cumsum(prop_var_explained)
)

# View the variance explained
print(pca_variance_df)

##      PC Variance_Explained Cumulative_Variance
## 1  PC1        0.16257075    0.1625708
## 2  PC2        0.133007319   0.2955781
## 3  PC3        0.093050056   0.3886281
## 4  PC4        0.082768551   0.4713967
## 5  PC5        0.068800754   0.5401974
## 6  PC6        0.064807270   0.6050047
## 7  PC7        0.058375365   0.6633801
## 8  PC8        0.051632032   0.7150121
## 9  PC9        0.050391883   0.7654040
## 10 PC10       0.046679712   0.8120837
## 11 PC11       0.040262674   0.8523464
## 12 PC12       0.037811435   0.8901578
## 13 PC13       0.034738974   0.9248968
## 14 PC14       0.030174315   0.9550711
## 15 PC15       0.023677299   0.9787484
## 16 PC16       0.013501665   0.9922501
## 17 PC17       0.007749939   1.0000000

pca_mahalanobis_analysis <- function(df, n_components = 13) {
  # Ensure we're working with numeric data only
  numeric_cols <- sapply(df, is.numeric)
  data_numeric <- df[, numeric_cols]

  # Remove NA values if any
  if(any(is.na(data_numeric))) {
    warning("NA values detected. Using complete cases only.")
  }
}

```

```

    data_numeric <- na.omit(data_numeric)
}

# Store row indices
original_rows <- as.integer(rownames(data_numeric))

# Perform PCA
pca_result <- prcomp(data_numeric, scale. = TRUE)
pc_scores <- pca_result$x[, 1:n_components]

# Mahalanobis distance
mean_vec <- colMeans(pc_scores)
cov_mat <- cov(pc_scores)
mahalanobis_dist <- mahalanobis(pc_scores, mean_vec, cov_mat)

result_df <- data.frame(
  original_index = original_rows,
  mahalanobis_dist = mahalanobis_dist
)

return(list(
  result_df = result_df,
  pca_result = pca_result,
  pc_scores = pc_scores
))
}

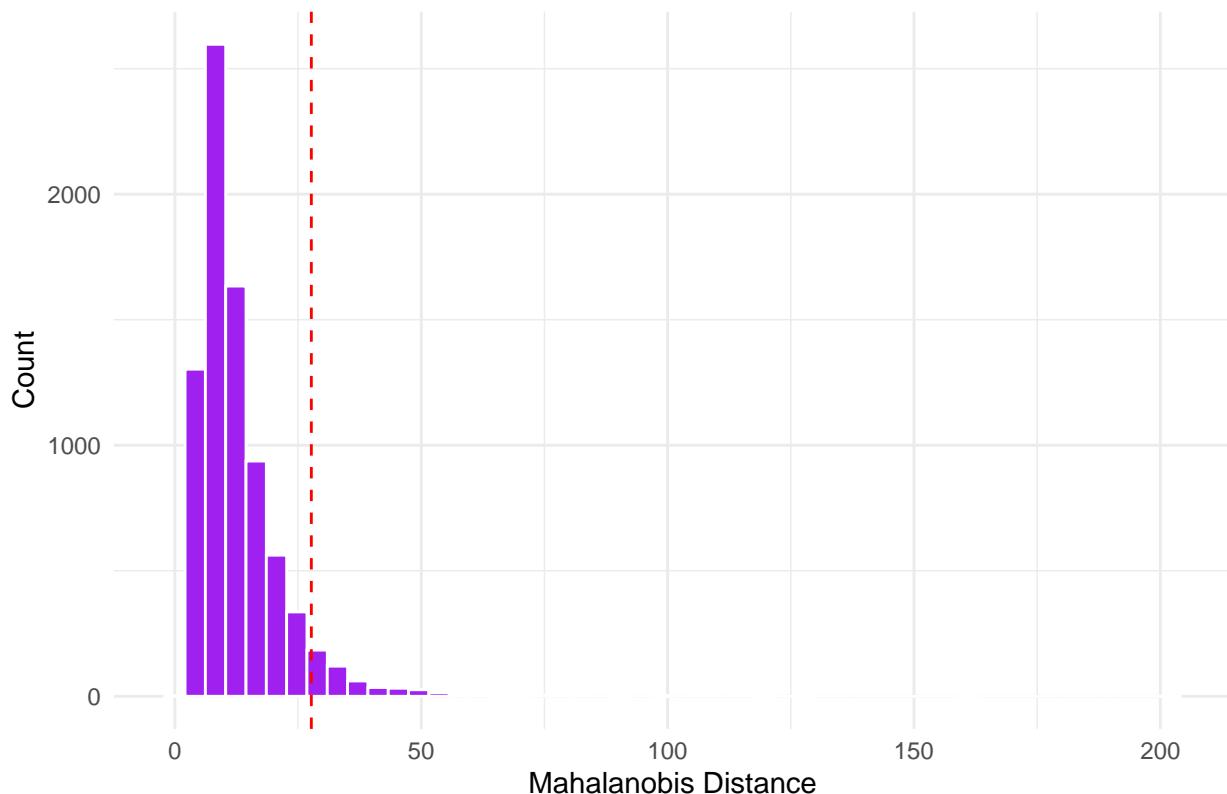
# Run on your transformed data
pca_result <- pca_mahalanobis_analysis(dfICU_transformed, n_components = 13)

# Add Mahalanobis distance as a new feature to dfICU_transformed
dfICU_transformed$MahalanobisDist <- NA
dfICU_transformed[pca_result$result_df$original_index, "MahalanobisDist"] <-
  pca_result$result_df$mahalanobis_dist

# Plot Mahalanobis distance
ggplot(dfICU_transformed, aes(x = MahalanobisDist)) +
  geom_histogram(bins = 50, fill = "purple", color = "white") +
  geom_vline(xintercept = qchisq(0.99, df = 13), col = "red", linetype = "dashed") +
  labs(title = "Mahalanobis Distance from PCA (13 Components)",
       x = "Mahalanobis Distance",
       y = "Count") +
  theme_minimal()

```

Mahalanobis Distance from PCA (13 Components)



```
# Filter anomalies with a chosen threshold (e.g. 95%)
threshold <- qchisq(0.99, df = 13)
df_anomalies <- dfICU_transformed %>% filter(MahalanobisDist > threshold)
```

```
# Count the number of outliers for the 95% threshold
outlier_count_99 <- nrow(df_anomalies)
cat("Number of outlier points (95% threshold):", outlier_count_99, "\n")
```

```
## Number of outlier points (95% threshold): 460
```

```
library(ggplot2)
library(gridExtra)
```

```
## Warning: package 'gridExtra' was built under R version 4.4.3
```

```
##
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
##
##     combine
```

```

# Extract PC scores
pc_df <- as.data.frame(pca_result$pc_scores)

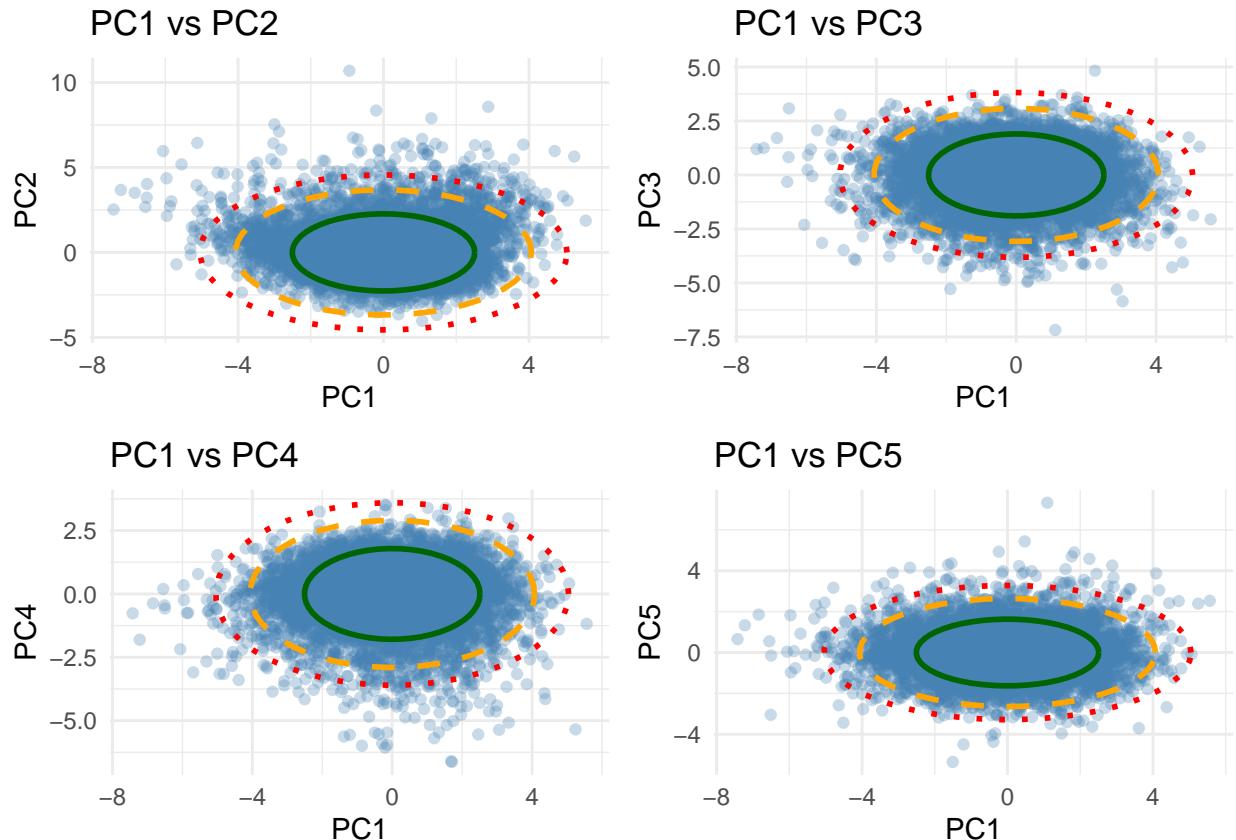
# List of plots PC1 vs PC2 to PC5
plots <- list()
for (i in 2:5) {
  p <- ggplot(pc_df, aes_string(x = "PC1", y = paste0("PC", i))) +
    geom_point(alpha = 0.3, color = "steelblue") +
    stat_ellipse(type = "norm", level = 0.68, color = "darkgreen", linetype = "solid", size = 1) +
    stat_ellipse(type = "norm", level = 0.95, color = "orange", linetype = "dashed", size = 1) +
    stat_ellipse(type = "norm", level = 0.99, color = "red", linetype = "dotted", size = 1) +
    labs(title = paste0("PC1 vs PC", i),
         x = "PC1",
         y = paste0("PC", i)) +
    theme_minimal()
  plots[[i - 1]] <- p
}

## Warning: `aes_string()` was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with `aes()`'.
## i See also `vignette("ggplot2-in-packages")` for more information.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

## Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use `linewidth` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

# Arrange all plots
gridExtra::grid.arrange(grobs = plots, ncol = 2)

```



```

library(ggplot2)

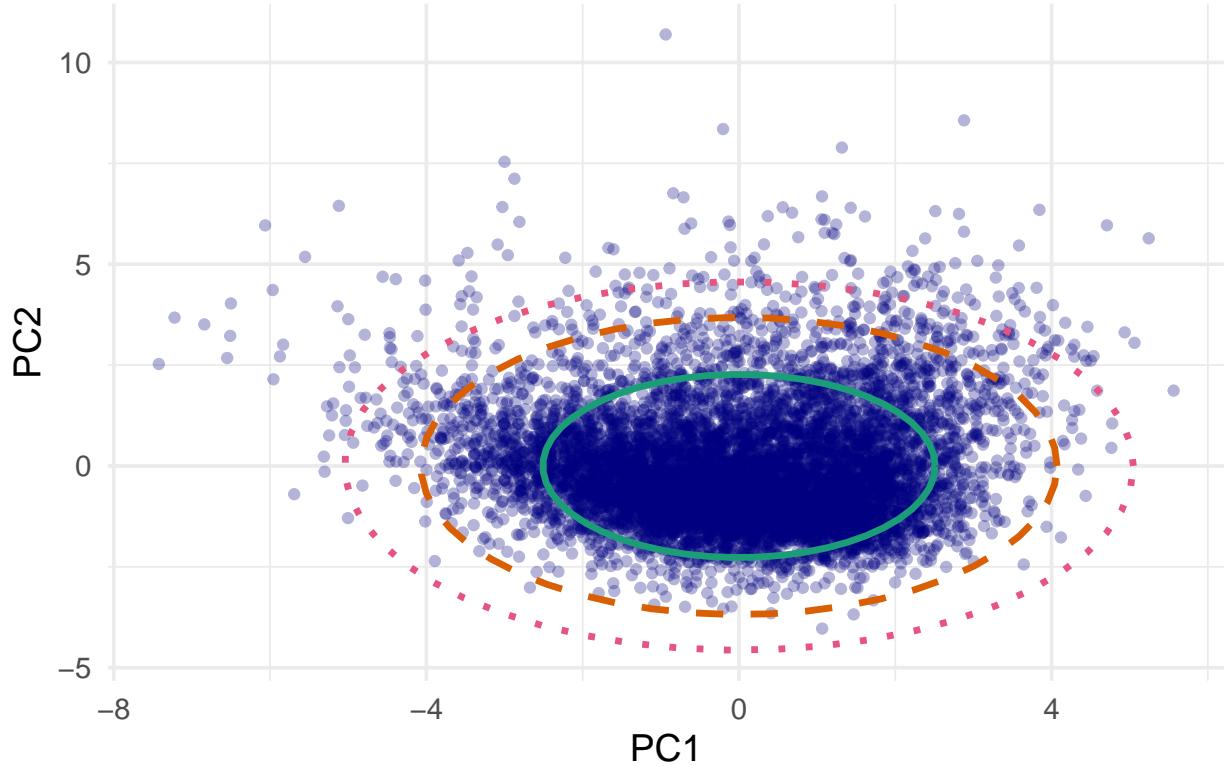
# Convert PCA scores to a data frame
pc_df <- as.data.frame(pca_result$pc_scores)

# Helper function to plot PC1 vs PCx with ellipses
plot_pca_ellipse <- function(pc_df, x = "PC1", y = "PC2") {
  ggplot(pc_df, aes_string(x = x, y = y)) +
    geom_point(alpha = 0.3, color = "navy") +
    stat_ellipse(type = "norm", level = 0.68, color = "#1b9e77", size = 1.2) + # Green
    stat_ellipse(type = "norm", level = 0.95, color = "#d95f02", size = 1.2, linetype = "dashed") + # Yellow
    stat_ellipse(type = "norm", level = 0.99, color = "#E75480", size = 1.2, linetype = "dotted") + # Red
    labs(title = paste(x, "vs", y, "with Confidence Ellipses"),
         x = x,
         y = y) +
    theme_minimal(base_size = 14)
}

# Draw each plot separately
plot_pca_ellipse(pc_df, "PC1", "PC2")

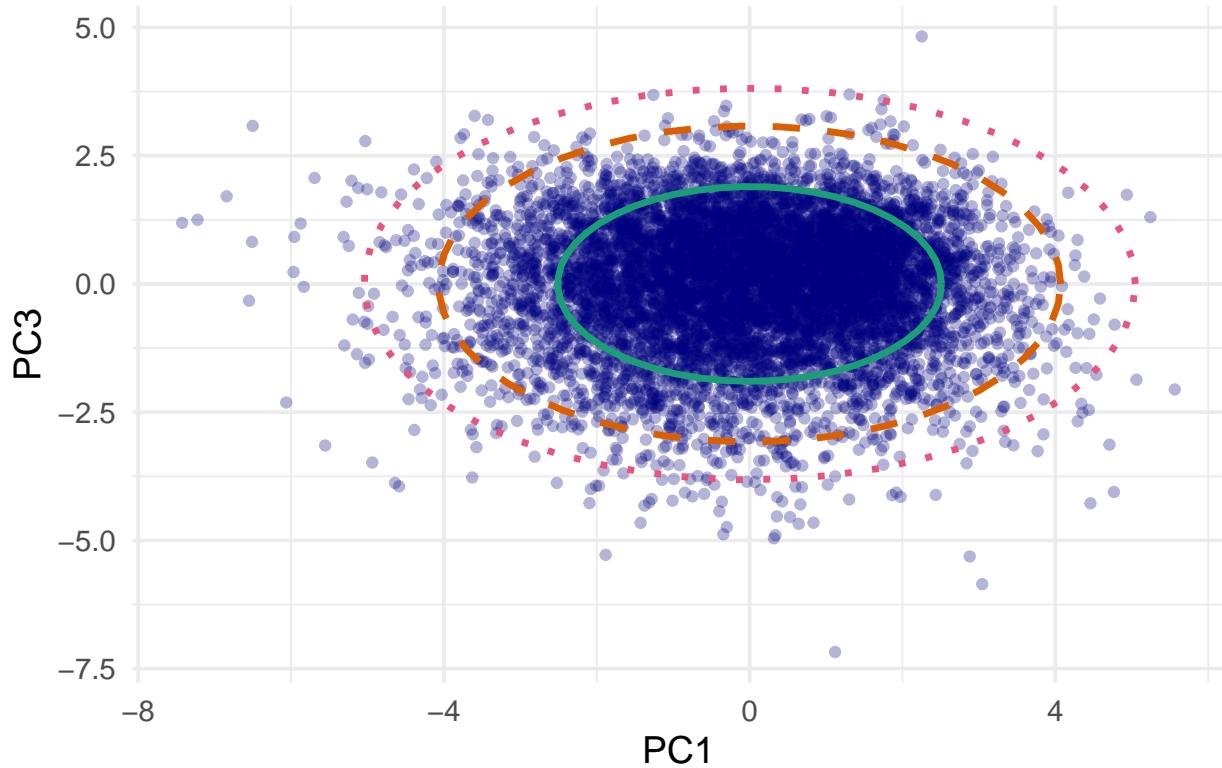
```

PC1 vs PC2 with Confidence Ellipses



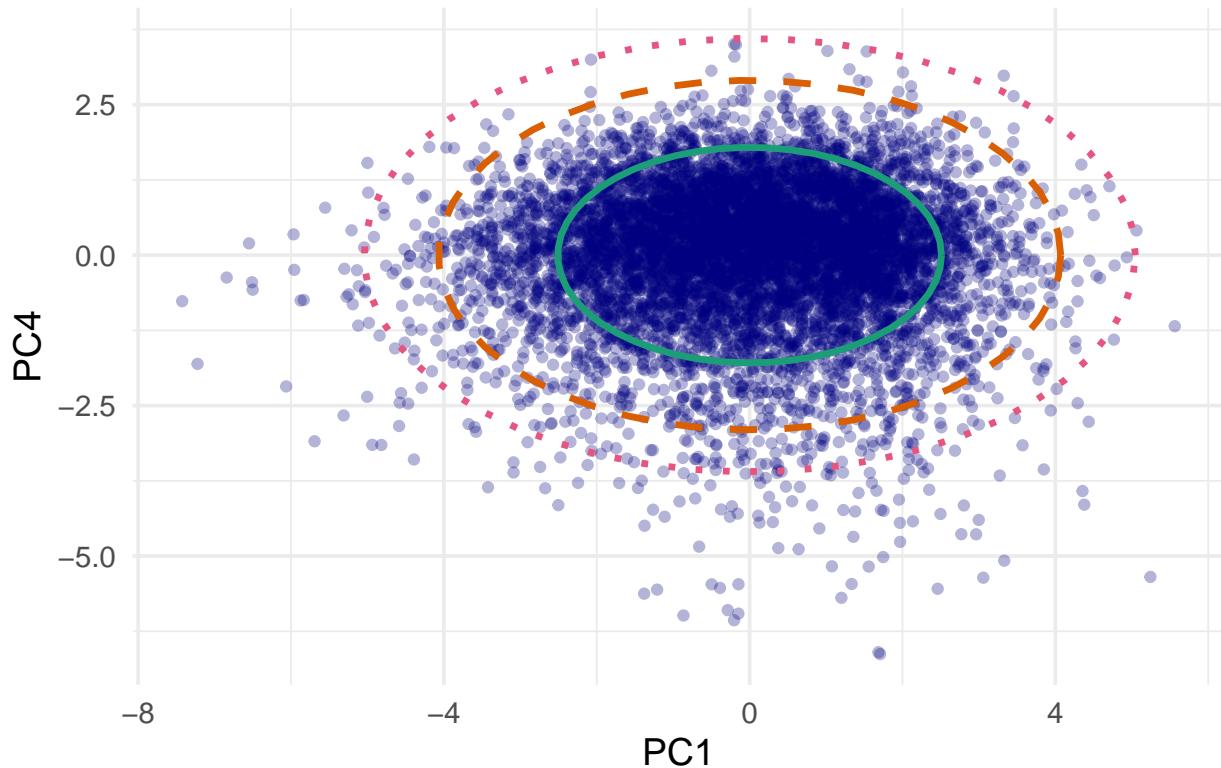
```
plot_pca_ellipse(pc_df, "PC1", "PC3")
```

PC1 vs PC3 with Confidence Ellipses



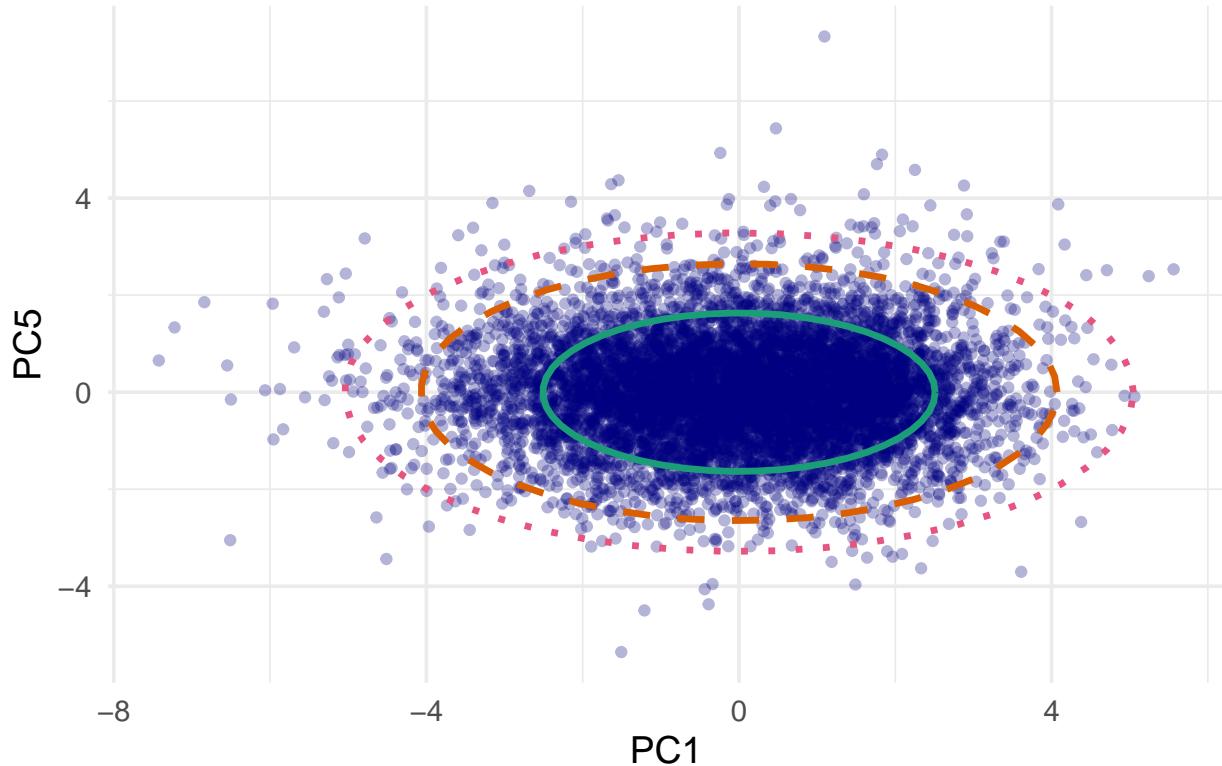
```
plot_pca_ellipse(pc_df, "PC1", "PC4")
```

PC1 vs PC4 with Confidence Ellipses



```
plot_pca_ellipse(pc_df, "PC1", "PC5")
```

PC1 vs PC5 with Confidence Ellipses



```

library(ggplot2)

# Combine PC scores with Mahalanobis distances
pc_df <- as.data.frame(pca_result$pc_scores)
pc_df$mahalanobis_dist <- dfICU_transformed$MahalanobisDist
threshold_98 <- qchisq(0.98, df = 13) # 98% threshold using 13 PCs

# Add a binary flag for outliers > 98%
pc_df$outlier_98 <- pc_df$mahalanobis_dist > threshold_98

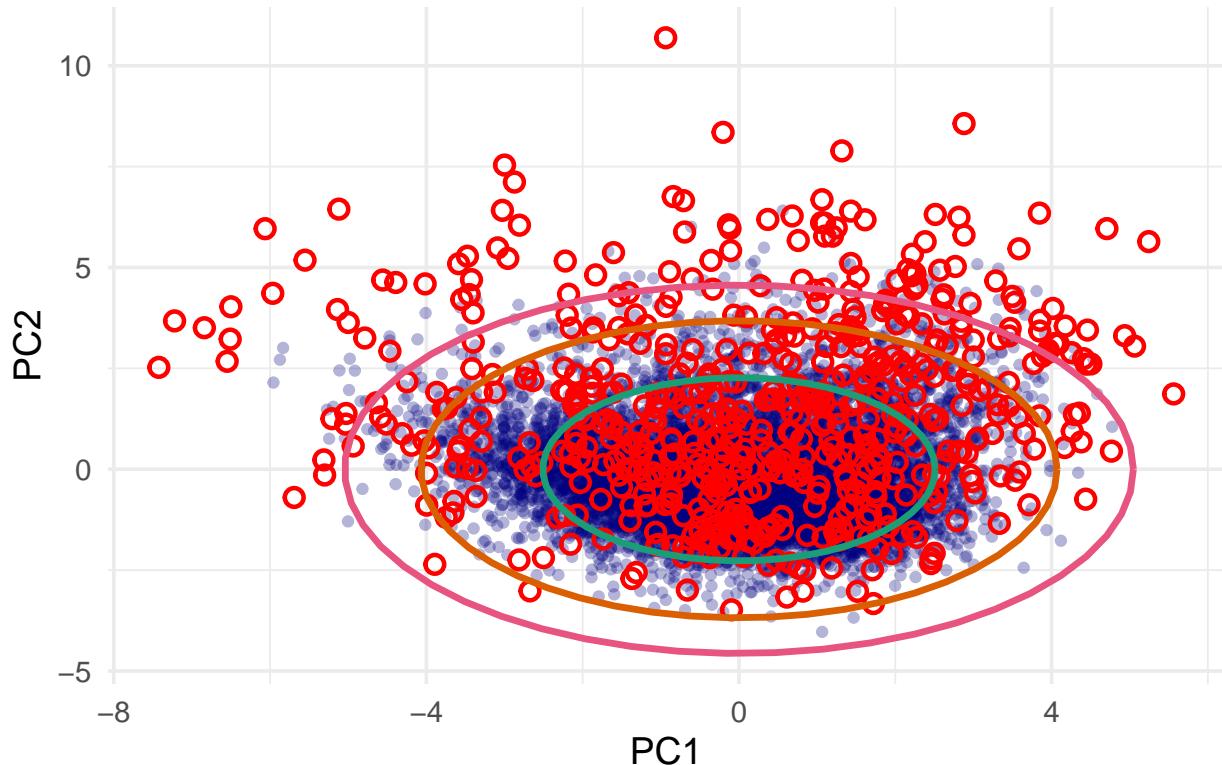
# Updated plotting function
plot_pca_ellipse <- function(pc_df, x = "PC1", y = "PC2") {
  ggplot(pc_df, aes_string(x = x, y = y)) +
    # Plot all non-outlier points first
    geom_point(data = subset(pc_df, !outlier_98), aes_string(x = x, y = y),
               alpha = 0.3, color = "navy") +
    # Overlay outliers > 98% in red
    geom_point(data = subset(pc_df, outlier_98), aes_string(x = x, y = y),
               color = "red", size = 2.5, shape = 21, stroke = 1.2) +
    # Confidence ellipses
    stat_ellipse(type = "norm", level = 0.68, color = "#1b9e77", size = 1.2) +      # Green
    stat_ellipse(type = "norm", level = 0.95, color = "#d95f02", size = 1.2) + # Orange
    stat_ellipse(type = "norm", level = 0.99, color = "#E75480", size = 1.2) + # Pink
    labs(title = paste(x, "vs", y, "with Confidence Ellipses and 98% Outliers"),
         x = x,
         y = y) +
}

```

```
    theme_minimal(base_size = 14)
}

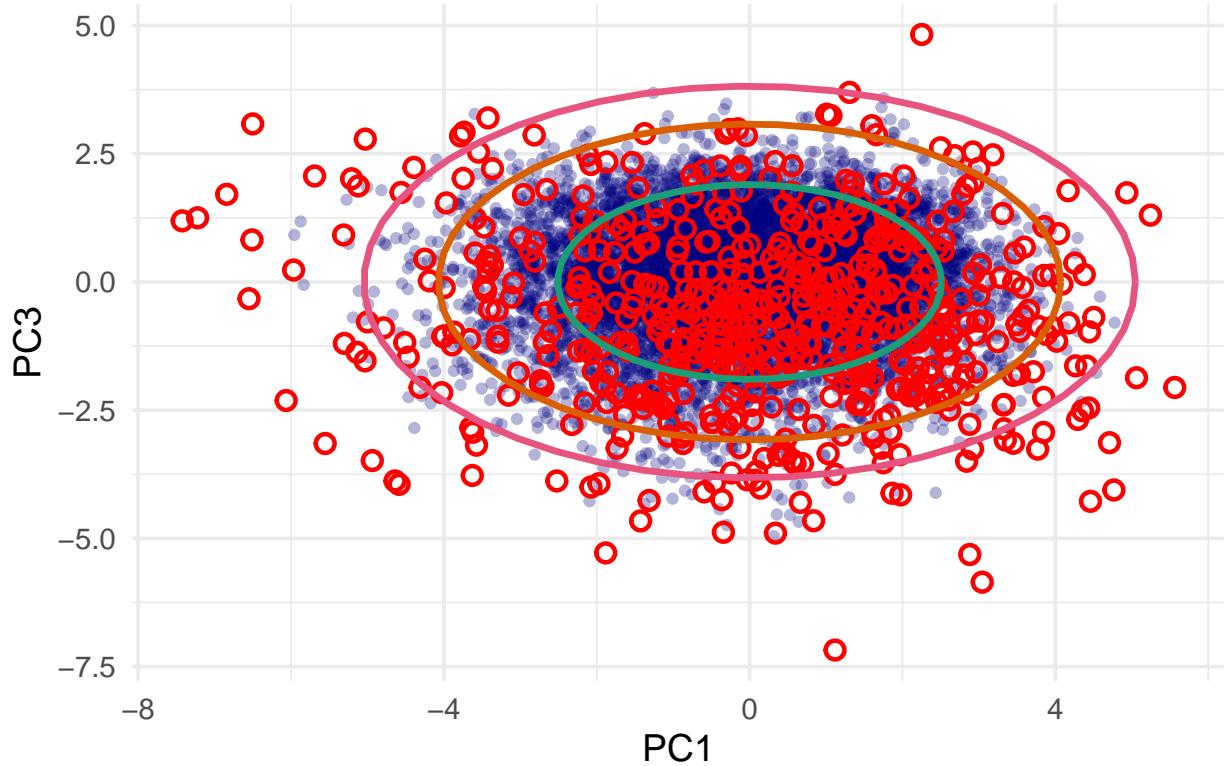
# Draw individual plots
plot_pca_ellipse(pc_df, "PC1", "PC2")
```

PC1 vs PC2 with Confidence Ellipses and 98% Outliers



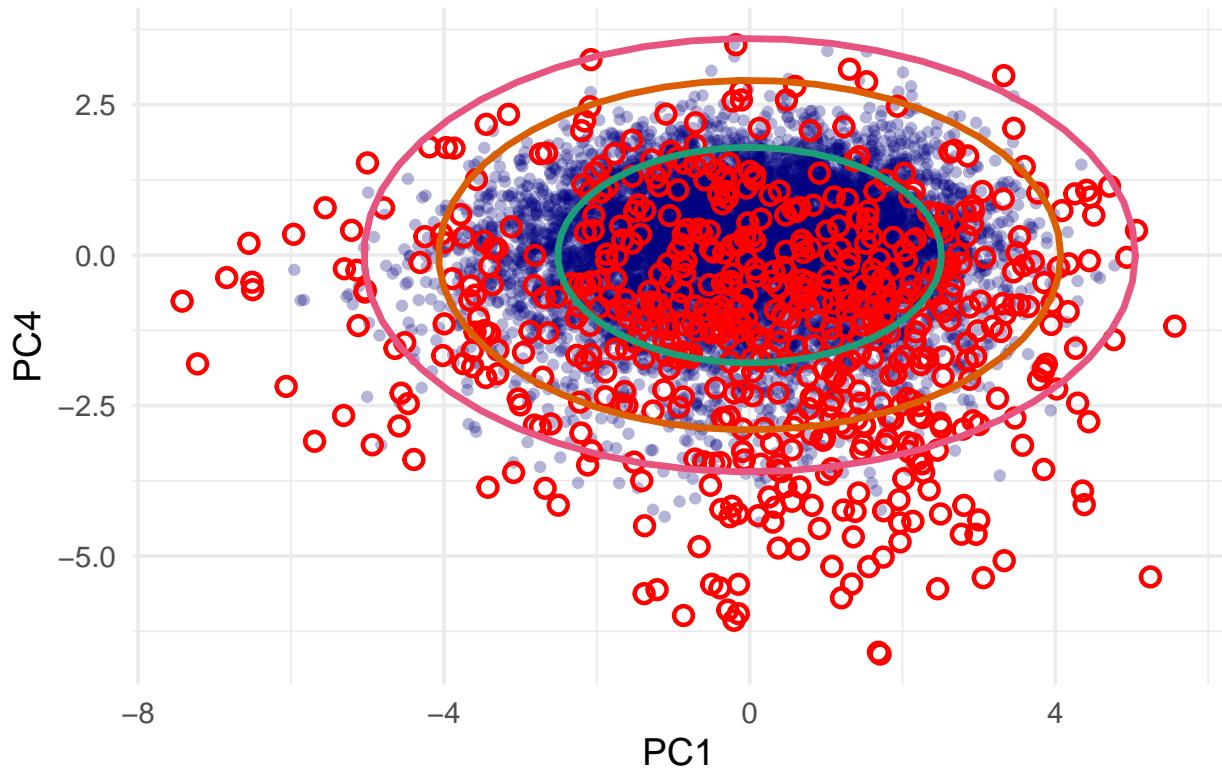
```
plot_pca_ellipse(pc_df, "PC1", "PC3")
```

PC1 vs PC3 with Confidence Ellipses and 98% Outliers



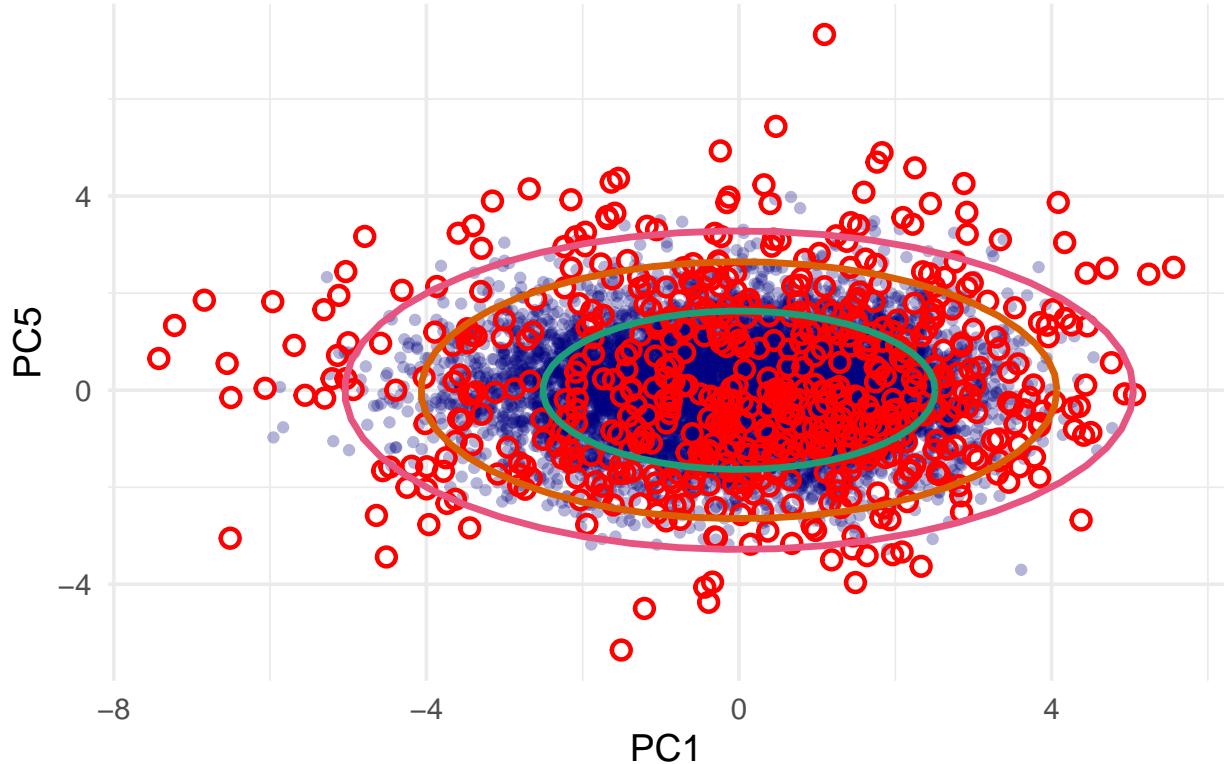
```
plot_pca_ellipse(pc_df, "PC1", "PC4")
```

PC1 vs PC4 with Confidence Ellipses and 98% Outliers



```
plot_pca_ellipse(pc_df, "PC1", "PC5")
```

PC1 vs PC5 with Confidence Ellipses and 98% Outliers



```

library(ggplot2)

# Combine PC scores with Mahalanobis distances
pc_df <- as.data.frame(pca_result$pc_scores)
pc_df$mahalanobis_dist <- dfICU_transformed$MahalanobisDist
threshold_99 <- qchisq(0.99, df = 13) # 99% threshold using 13 PCs

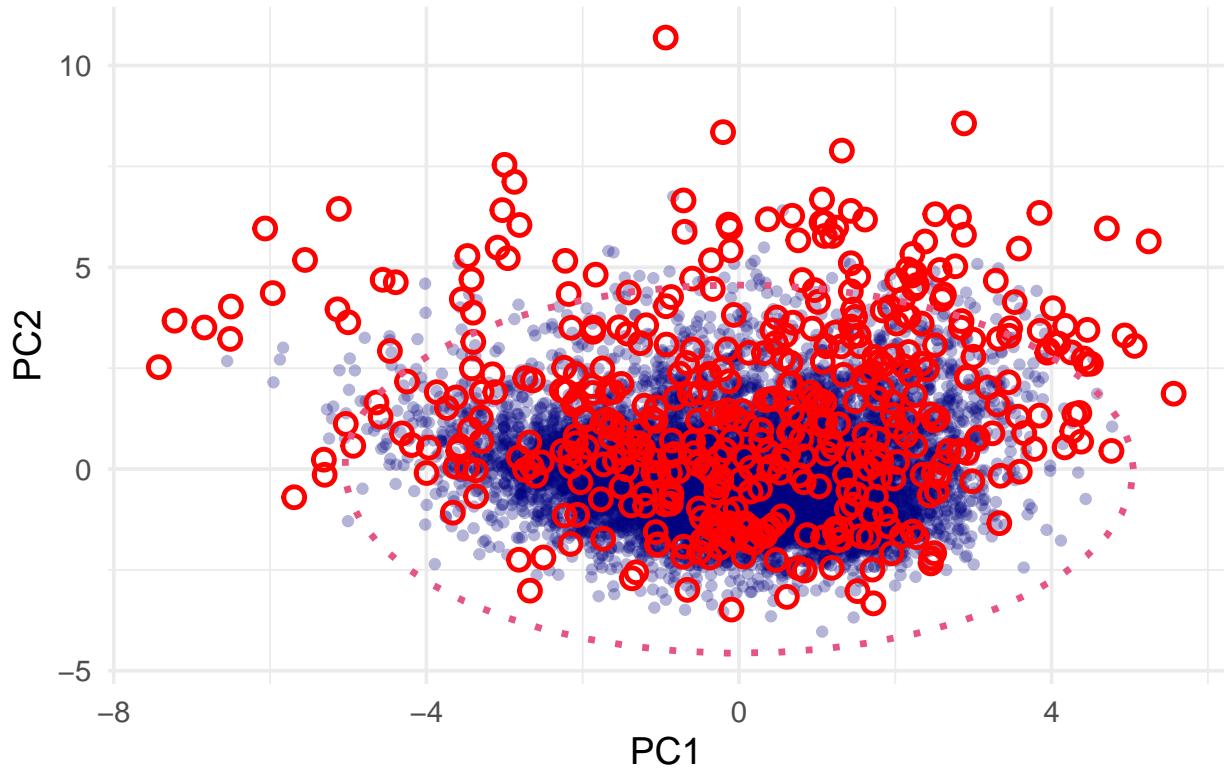
# Add a binary flag for outliers > 99%
pc_df$outlier_99 <- pc_df$mahalanobis_dist > threshold_99

# Updated plotting function for 99% outliers
plot_pca_ellipse <- function(pc_df, x = "PC1", y = "PC2") {
  ggplot(pc_df, aes_string(x = x, y = y)) +
    # Plot normal points
    geom_point(data = subset(pc_df, !outlier_99), aes_string(x = x, y = y),
               alpha = 0.3, color = "navy") +
    # Highlight 99% outliers
    geom_point(data = subset(pc_df, outlier_99), aes_string(x = x, y = y),
               color = "red", size = 2.8, shape = 21, stroke = 1.3) +
    # Confidence ellipse for 99%
    stat_ellipse(type = "norm", level = 0.99, color = "#E75480", size = 1.2, linetype = "dotted") + # P
    labs(title = paste(x, "vs", y, "with 99% Outliers"),
         x = x,
         y = y) +
    theme_minimal(base_size = 14)
}

```

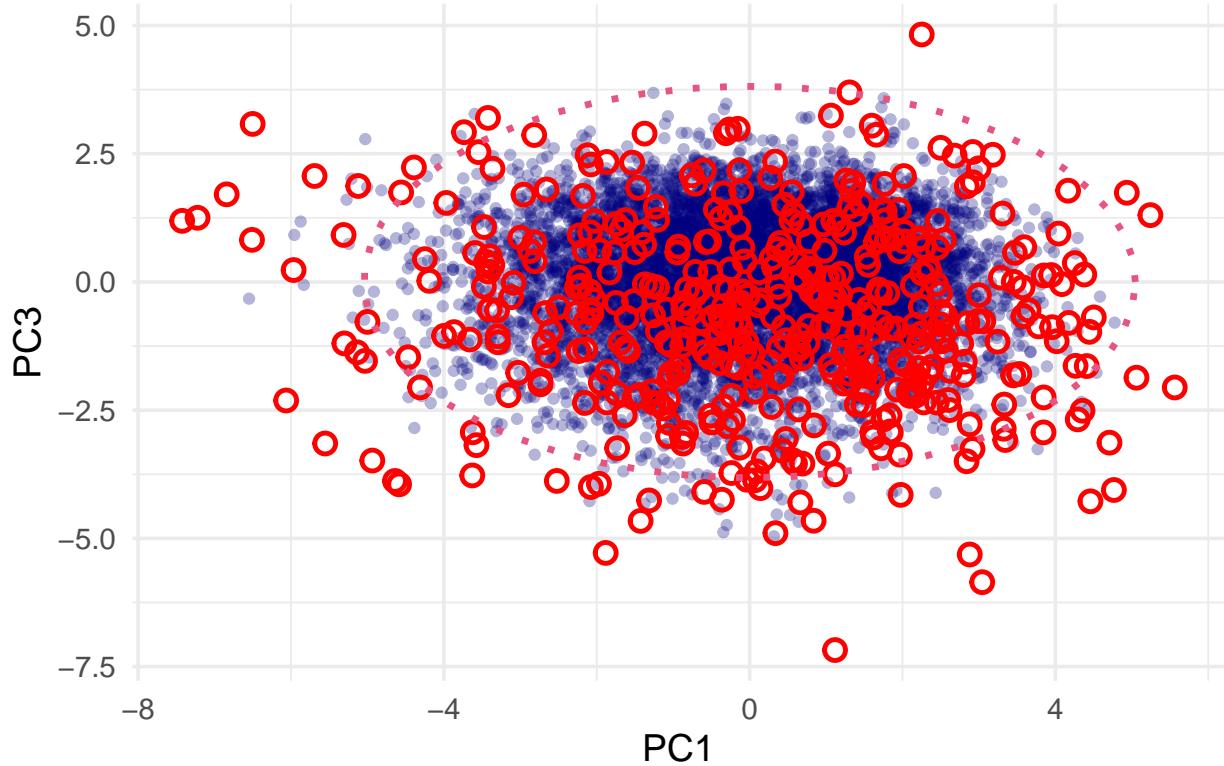
```
# Draw each plot individually for PC1 vs other PCs  
plot_pca_ellipse(pc_df, "PC1", "PC2")
```

PC1 vs PC2 with 99% Outliers



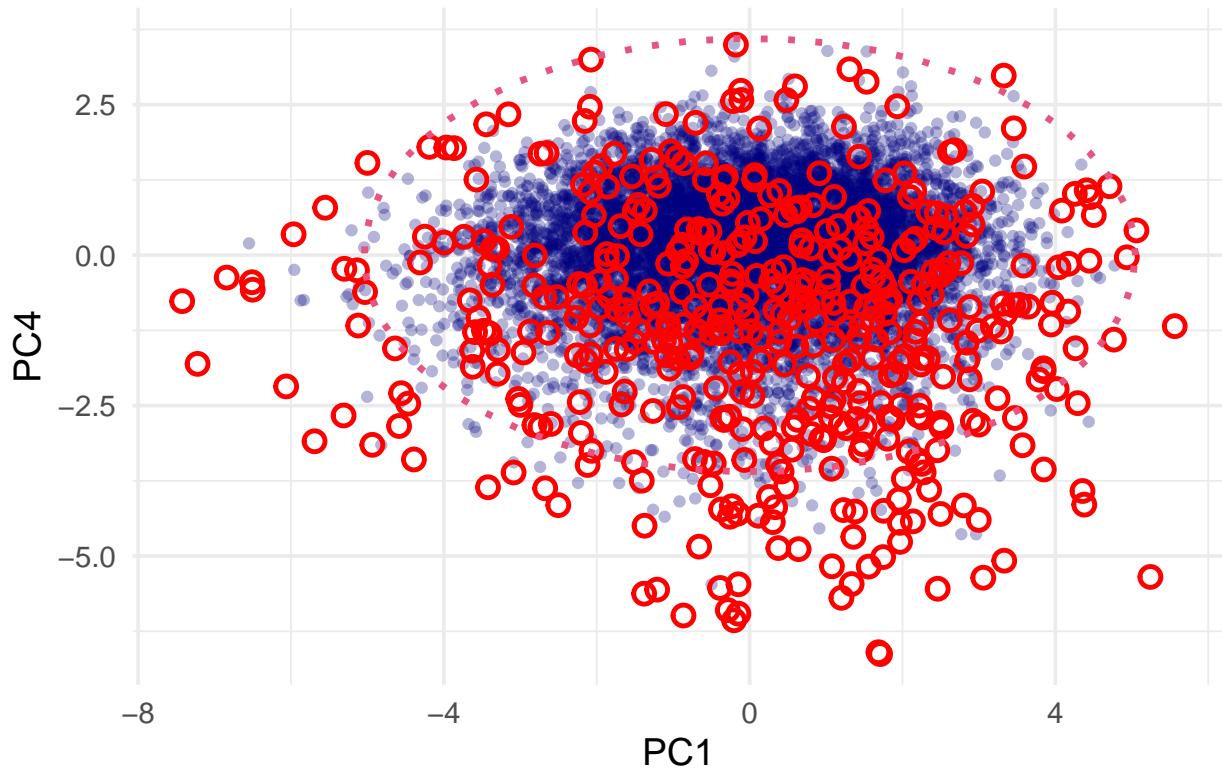
```
plot_pca_ellipse(pc_df, "PC1", "PC3")
```

PC1 vs PC3 with 99% Outliers



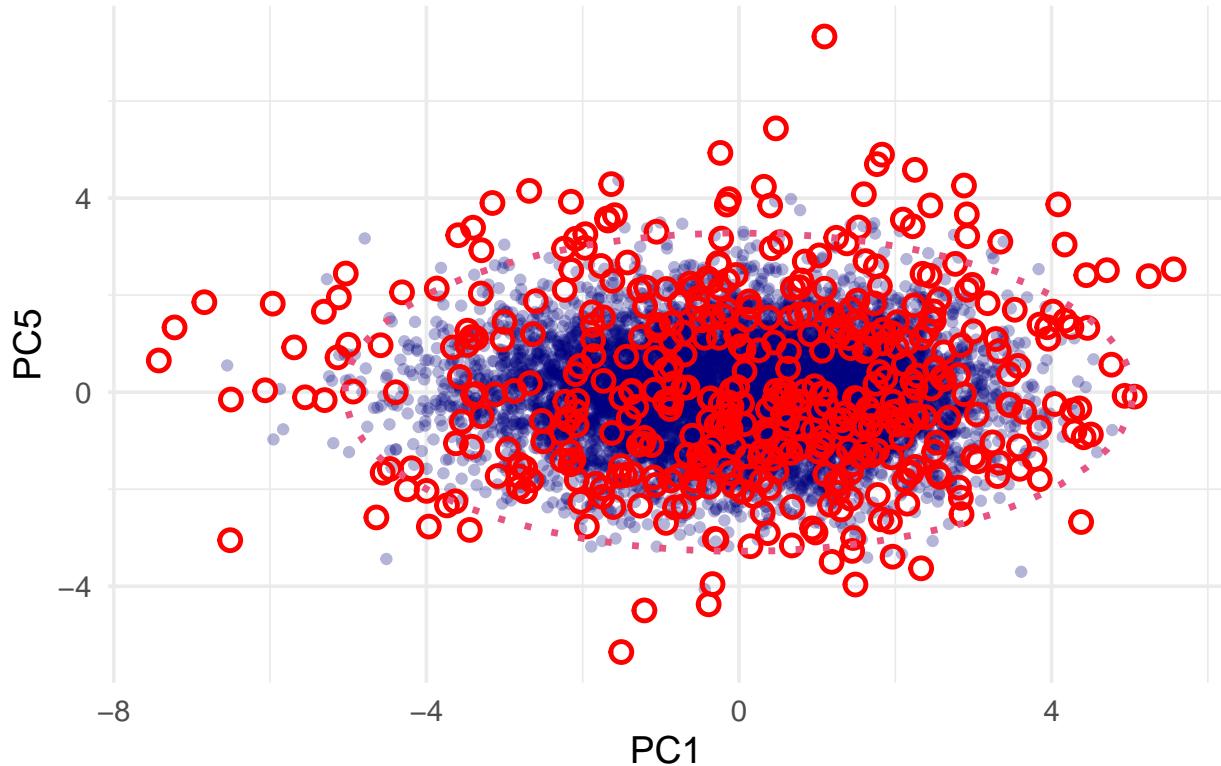
```
plot_pca_ellipse(pc_df, "PC1", "PC4")
```

PC1 vs PC4 with 99% Outliers



```
plot_pca_ellipse(pc_df, "PC1", "PC5")
```

PC1 vs PC5 with 99% Outliers



```
library(Rtsne)
library(ggplot2)
library(dplyr)

# Remove NAs and scale data
df_clean <- na.omit(dfICU_transformed)
df_scaled <- scale(df_clean)

# Run t-SNE (2D)
set.seed(42)
tsne_result <- Rtsne(df_scaled, dims = 2, perplexity = 30, verbose = TRUE, max_iter = 1000)

## Performing PCA
## Read the 7886 x 18 data matrix successfully!
## OpenMP is working. 1 threads.
## Using no_dims = 2, perplexity = 30.000000, and theta = 0.500000
## Computing input similarities...
## Building tree...
## Done in 1.48 seconds (sparsity = 0.018015)!
## Learning embedding...
## Iteration 50: error is 93.005107 (50 iterations in 0.57 seconds)
## Iteration 100: error is 93.005107 (50 iterations in 0.71 seconds)
## Iteration 150: error is 93.005106 (50 iterations in 0.73 seconds)
## Iteration 200: error is 93.005074 (50 iterations in 0.67 seconds)
## Iteration 250: error is 93.003510 (50 iterations in 0.69 seconds)
```

```

## Iteration 300: error is 3.949140 (50 iterations in 0.79 seconds)
## Iteration 350: error is 3.643539 (50 iterations in 0.50 seconds)
## Iteration 400: error is 3.490864 (50 iterations in 0.48 seconds)
## Iteration 450: error is 3.396178 (50 iterations in 0.48 seconds)
## Iteration 500: error is 3.330883 (50 iterations in 0.49 seconds)
## Iteration 550: error is 3.280938 (50 iterations in 0.50 seconds)
## Iteration 600: error is 3.242335 (50 iterations in 0.49 seconds)
## Iteration 650: error is 3.212130 (50 iterations in 0.49 seconds)
## Iteration 700: error is 3.187815 (50 iterations in 0.49 seconds)
## Iteration 750: error is 3.168659 (50 iterations in 0.50 seconds)
## Iteration 800: error is 3.152891 (50 iterations in 0.50 seconds)
## Iteration 850: error is 3.140513 (50 iterations in 0.49 seconds)
## Iteration 900: error is 3.130909 (50 iterations in 0.50 seconds)
## Iteration 950: error is 3.122188 (50 iterations in 0.50 seconds)
## Iteration 1000: error is 3.115261 (50 iterations in 0.50 seconds)
## Fitting performed in 11.08 seconds.

# Create a dataframe of t-SNE results
tsne_df <- as.data.frame(tsne_result$Y)
colnames(tsne_df) <- c("TSNE1", "TSNE2")

# Mahalanobis on 2D t-SNE space
mean_vec <- colMeans(tsne_df)
cov_mat <- cov(tsne_df)
tsne_df$MahalanobisDist <- mahalanobis(tsne_df, center = mean_vec, cov = cov_mat)

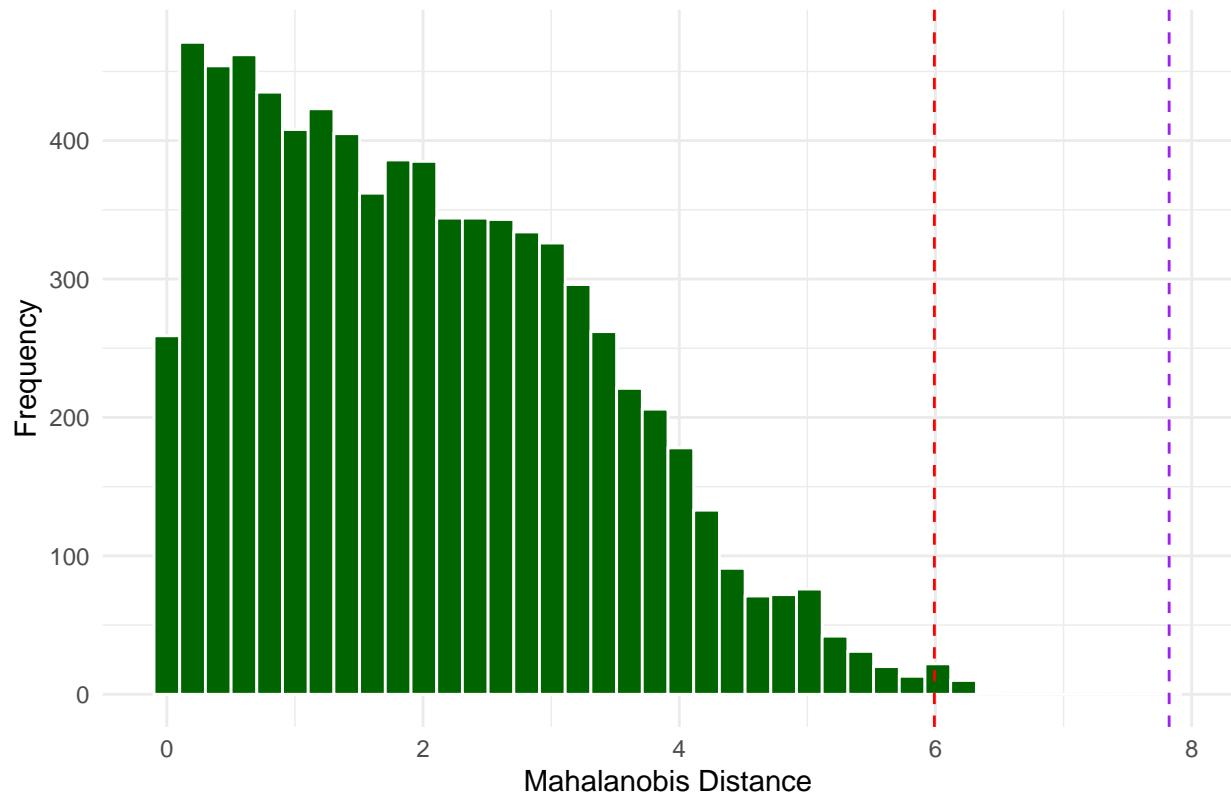
# Thresholds
threshold_95 <- qchisq(0.95, df = 2)
threshold_98 <- qchisq(0.98, df = 2)

# Flagging outliers
tsne_df$outlier_95 <- tsne_df$MahalanobisDist > threshold_95
tsne_df$outlier_98 <- tsne_df$MahalanobisDist > threshold_98

# Plot Mahalanobis Distance Histogram
ggplot(tsne_df, aes(x = MahalanobisDist)) +
  geom_histogram(bins = 40, fill = "darkgreen", color = "white") +
  geom_vline(xintercept = threshold_95, color = "red", linetype = "dashed") +
  geom_vline(xintercept = threshold_98, color = "purple", linetype = "dashed") +
  labs(title = "Mahalanobis Distance on t-SNE (2D)",
       x = "Mahalanobis Distance", y = "Frequency") +
  theme_minimal()

```

Mahalanobis Distance on t-SNE (2D)



```
# Visualize t-SNE points with outliers in red
ggplot(tsne_df, aes(x = TSNE1, y = TSNE2)) +
  geom_point(aes(color = outlier_98), alpha = 0.6, size = 2) +
  scale_color_manual(values = c("FALSE" = "blue", "TRUE" = "red")) +
  labs(title = "t-SNE 2D Plot with 98% Mahalanobis Outliers",
       color = "Outlier > 98%") +
  theme_minimal()
```

t-SNE 2D Plot with 98% Mahalanobis Outliers

