

# Diabetes - (Penalized) Logistic Regression Presentation

Group 6 - EBA1

2023-09-24

## Functions

```
# Count NA's per column
count_na_per_column <- function(df) {
  sapply(df, function(x) sum(is.na(x)))
}

# Create a function to calculate metrics
calc_metrics <- function(pred, true) {
  confusion <- table(pred, true)
  TP <- confusion[2, 2]
  FP <- confusion[2, 1]
  TN <- confusion[1, 1]
  FN <- confusion[1, 2]
  Sensitivity <- TP / (TP + FN)
  Specificity <- TN / (TN + FP)
  Accuracy <- (TP + TN) / (TP + FP + TN + FN)
  Precision <- TP / (TP + FP)
  F1 <- 2 * (Precision * Sensitivity) / (Precision + Sensitivity)
  pred_obj <- prediction(as.numeric(pred), as.numeric(true))
  perf <- performance(pred_obj, "auc")
  AUC <- as.numeric(perf@y.values)
  return(c(AUC, Sensitivity, Precision, F1, Accuracy))
}

### A string like "Logistic" has to be put in. <--- from the colnames of df_coef
importance_plot <- function(model_string) {

  # Create a data frame for plotting
  coef_df <- data.frame(
    Variable = rownames(df_coef)[-1], # Exclude the intercept
    Importance = abs(df_coef[-1, model_string]) # Exclude the intercept
  )

  # Order the variables by importance
  coef_df <- coef_df[order(coef_df$Importance), ]

  # Plot
  ggplot(coef_df, aes(x = reorder(Variable, Importance), y = Importance, fill = Importance)) +
    geom_bar(stat = "identity") +
    scale_fill_gradient(low = "lightblue2", high = "lightblue3") +
    theme_minimal() +
```

```

theme(axis.text.y = element_text(size = 12),
      axis.text.x = element_text(size = 10),
      title = element_text(size = 15),
      axis.title.y = element_blank(),
      axis.title.x = element_blank(),
      legend.position = "none") +
coord_flip() +
ggtitle("Variable Importance")
}

```

## 5. Data Processing

```

data("PimaIndiansDiabetes2", package = "mlbench")
diabetes <- PimaIndiansDiabetes2
diabetes$diabetes <- as.factor(ifelse(diabetes$diabetes == "pos", 1, 0))

```

### Removing NA's

```

diabetes <- na.omit(diabetes)

```

### Data for models

```

# Split
set.seed(222)
n <- nrow(diabetes)

# Training
training.samples <- sample(1:n, size = 0.75 * n)
train.data <- diabetes[training.samples, ]
scaled_train.data <- scale(train.data[, 1:8])
train.data[, 1:8] <- scaled_train.data

# Testing
test.data <- diabetes[-training.samples, ]
scaled_test.data <- scale(test.data[, 1:8])
test.data[, 1:8] <- scaled_test.data

# Train
X.train <- model.matrix(diabetes ~ ., data = train.data)[,-1]
y.train <- train.data$diabetes

# Test
X.test <- model.matrix(diabetes ~ ., data = test.data)[,-1]
y.test <- test.data$diabetes

```

### Logistic Regression

```

set.seed(123)
log.model <- glm(diabetes ~ ., data = train.data, family = "binomial")

# Make predictions

```

```

probabilities <- predict(log.model, newdata = test.data, type = "response")
predicted.classes.log <- as.factor(ifelse(probabilities > 0.5, 1, 0))

# Accuracy
# mean(predicted.classes == y.test)
log.conf <- confusionMatrix(predicted.classes.log, y.test, positive = "1")
log.conf

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0  1
##           0 61 10
##           1  6 21
##
##           Accuracy : 0.8367
##           95% CI : (0.7484, 0.9037)
##       No Information Rate : 0.6837
##       P-Value [Acc > NIR] : 0.0004538
##
##           Kappa : 0.609
##
##  McNemar's Test P-Value : 0.4532547
##
##           Sensitivity : 0.6774
##           Specificity : 0.9104
##       Pos Pred Value : 0.7778
##       Neg Pred Value : 0.8592
##           Prevalence : 0.3163
##       Detection Rate : 0.2143
##   Detection Prevalence : 0.2755
##       Balanced Accuracy : 0.7939
##
##       'Positive' Class : 1
##

```

## LASSO

```

set.seed(123)
cv.lasso.model <- cv.glmnet(X.train, y.train, alpha = 1, family = "binomial",
                           intercept = T)
#plot(cv.lasso.model)
cbind(coef(cv.lasso.model, s = cv.lasso.model$lambda.min), coef(cv.lasso.model, s = cv.lasso.model$lambda.1se))

## 9 x 2 sparse Matrix of class "dgCMatrix"
##           s1      s1
## (Intercept) -0.9087548 -0.79367862
## pregnant    0.2500270  0.04713706
## glucose     1.0571306  0.80577102
## pressure    .          .
## triceps     0.0626410  .
## insulin     .          .
## mass        0.4184043  0.19426613
## pedigree    0.3622318  0.10361089

```

```
## age          0.2235214  0.17759602
```

```
# Make predictions
```

```
probabilities <- predict(cv.lasso.model, newx = X.test, s = cv.lasso.model$lambda.min, type = "response")
predicted.classes.lasso.min <- as.factor(ifelse(probabilities > 0.5, 1, 0))
```

```
probabilities <- predict(cv.lasso.model, newx = X.test, s = cv.lasso.model$lambda.1se, type = "response")
predicted.classes.lasso.1se <- as.factor(ifelse(probabilities > 0.5, 1, 0))
```

```
# Accuracy Min
```

```
lasso.min.conf <- confusionMatrix(predicted.classes.lasso.min, y.test, positive = "1")
lasso.min.conf
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction  0  1
```

```
##           0 61 10
```

```
##           1  6 21
```

```
##
```

```
##           Accuracy : 0.8367
```

```
##           95% CI : (0.7484, 0.9037)
```

```
## No Information Rate : 0.6837
```

```
## P-Value [Acc > NIR] : 0.0004538
```

```
##
```

```
##           Kappa : 0.609
```

```
##
```

```
## Mcnemar's Test P-Value : 0.4532547
```

```
##
```

```
##           Sensitivity : 0.6774
```

```
##           Specificity : 0.9104
```

```
## Pos Pred Value : 0.7778
```

```
## Neg Pred Value : 0.8592
```

```
## Prevalence : 0.3163
```

```
## Detection Rate : 0.2143
```

```
## Detection Prevalence : 0.2755
```

```
## Balanced Accuracy : 0.7939
```

```
##
```

```
## 'Positive' Class : 1
```

```
##
```

```
# Accuracy 1se
```

```
lasso.1se.conf <- confusionMatrix(predicted.classes.lasso.1se, y.test, positive = "1")
```

```
lasso.1se.conf
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction  0  1
```

```
##           0 62 12
```

```
##           1  5 19
```

```
##
```

```
##           Accuracy : 0.8265
```

```
##           95% CI : (0.7369, 0.8956)
```

```
## No Information Rate : 0.6837
```

```
## P-Value [Acc > NIR] : 0.001064
```

```
##
##           Kappa : 0.573
##
## Mcnemar's Test P-Value : 0.145610
##
##           Sensitivity : 0.6129
##           Specificity : 0.9254
##           Pos Pred Value : 0.7917
##           Neg Pred Value : 0.8378
##           Prevalence : 0.3163
##           Detection Rate : 0.1939
##           Detection Prevalence : 0.2449
##           Balanced Accuracy : 0.7691
##
##           'Positive' Class : 1
##
```

## Ridge

```
set.seed(123)
cv.ridge.model <- cv.glmnet(X.train, y.train, alpha = 0, family = "binomial", intercept = T)

# Make predictions
probabilities <- predict(cv.ridge.model, newx = X.test, s = cv.ridge.model$lambda.min, type = "response")
predicted.classes.ridge <- as.factor(ifelse(probabilities > 0.5, 1, 0))

# Accuracy
ridge.conf <- confusionMatrix(as.factor(predicted.classes.ridge), y.test, positive = "1")
ridge.conf
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0  1
##           0 62 10
##           1  5 21
##
##           Accuracy : 0.8469
##           95% CI : (0.7601, 0.9117)
##           No Information Rate : 0.6837
##           P-Value [Acc > NIR] : 0.0001806
##
##           Kappa : 0.6301
##
## Mcnemar's Test P-Value : 0.3016996
##
##           Sensitivity : 0.6774
##           Specificity : 0.9254
##           Pos Pred Value : 0.8077
##           Neg Pred Value : 0.8611
##           Prevalence : 0.3163
##           Detection Rate : 0.2143
##           Detection Prevalence : 0.2653
##           Balanced Accuracy : 0.8014
```

```
##
##      'Positive' Class : 1
##
```

## Elastic Net

```
set.seed(123)

# CV and tuning grid
myFolds <- createFolds(y.train, k = 10, list = TRUE)
myControl <- trainControl(index = myFolds)
tuneGrid <- expand.grid(alpha = seq(0, 1, by = 0.05), lambda = 10^seq(1, -3, length=100))

# Train model
elasticnet.model <- train(X.train, y.train, method = "glmnet", trControl = myControl,
                          tuneGrid = tuneGrid, intercept = T, family = "binomial")

# Make predictions
probabilities <- predict(elasticnet.model, newdata = X.test, type = "prob")
predicted.classes.elasticnet <- ifelse(probabilities$`1` > 0.5, 1, 0)

# Accuracy
elasticnet.conf <- confusionMatrix(as.factor(predicted.classes.elasticnet), y.test, positive = "1")
elasticnet.conf

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0   1
##           0 63 11
##           1   4 20
##
##           Accuracy : 0.8469
##           95% CI : (0.7601, 0.9117)
##           No Information Rate : 0.6837
##           P-Value [Acc > NIR] : 0.0001806
##
##           Kappa : 0.6233
##
## Mcnemar's Test P-Value : 0.1213353
##
##           Sensitivity : 0.6452
##           Specificity : 0.9403
##           Pos Pred Value : 0.8333
##           Neg Pred Value : 0.8514
##           Prevalence : 0.3163
##           Detection Rate : 0.2041
##           Detection Prevalence : 0.2449
##           Balanced Accuracy : 0.7927
##
##           'Positive' Class : 1
##

# Coefficients
coef(elasticnet.model$finalModel, s = elasticnet.model$bestTune$lambda)
```

```
## 9 x 1 sparse Matrix of class "dgCMatrix"
##              s1
## (Intercept) -0.783644314
## pregnant    0.124519132
## glucose     0.616288277
## pressure    0.003625147
## triceps     0.073458775
## insulin     0.067361284
## mass        0.200698201
## pedigree    0.158043114
## age         0.188599082
```

---

```
# Starting here is all the code for the generation of Plots, Latex Tables etc.
```

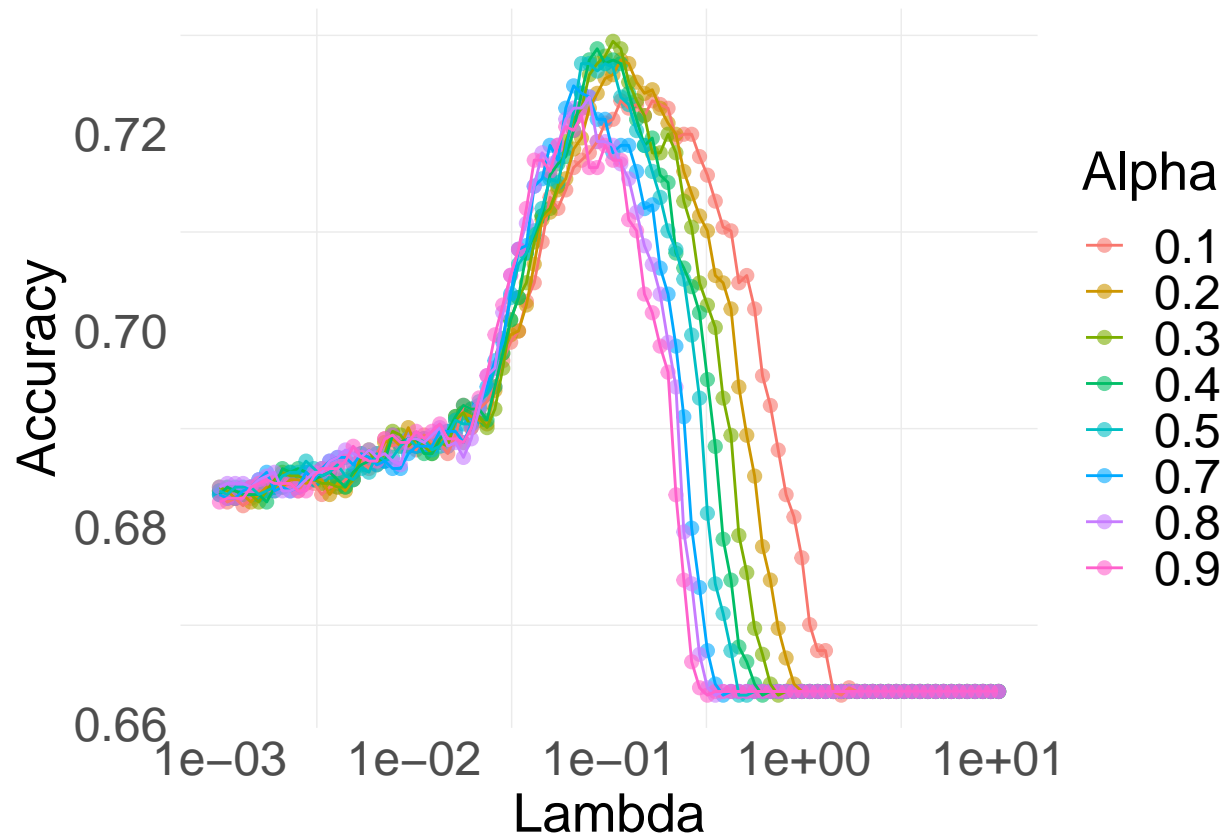
---

### Plotting the Parameter Grid for ElasticNet

```
# Only using some alphas
filtered_data <- subset(elasticnet.model$results, alpha %in% seq(0.1, 0.9, by=0.1))

# Create the plot
ggplot(filtered_data, aes(x=lambda, y=Accuracy, color=factor(alpha))) +
  geom_point(size = 2, alpha = 0.6) +
  geom_line(size = 0.5) +
  scale_x_log10() +
  scale_color_discrete(name = "Alpha") +
  labs(x = "Lambda", y = "Accuracy") +
  theme_minimal() +
  theme(
    axis.text = element_text(size = 18),
    axis.title = element_text(size = 20),
    legend.text = element_text(size = 18),
    legend.title = element_text(size = 20),
    panel.grid.major = element_blank()
  )
)
```

```
## 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.
```



### Correlation Matrix Plot

```
numeric_data <- subset(diabetes, select = -c(diabetes))
corr <- round(cor(numeric_data), 1)
ggcorrplot(corr,
  type = "lower",
  lab = TRUE,
  lab_size = 3,
  colors = c("red", "white", "cyan4"),
  title = "Correlogram of Diabetes Dataset",
  ggtheme = theme_bw()
)
```



Correlogram of Diabetes Dataset

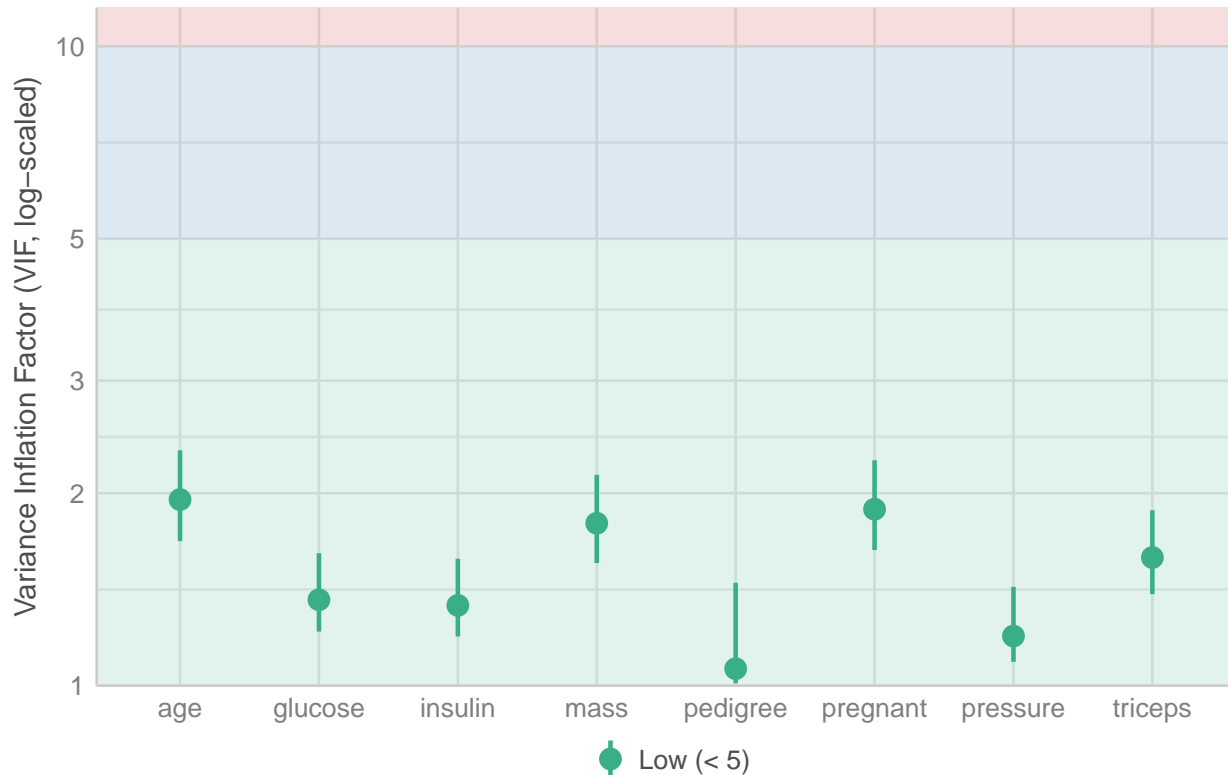


### Collinearity Plot

```
result <- check_collinearity(log.model)
plot(result) + ggtitle("")
```

## Variable 'Component' is not in your data frame :/

High collinearity (VIF) may inflate parameter uncertainty



### Calculate model metrics

*# Calculate metrics for each model*

```
metrics.log <- calc_metrics(predicted.classes.log, y.test)
metrics.lasso.min <- calc_metrics(predicted.classes.lasso.min, y.test)
metrics.lasso.1se <- calc_metrics(predicted.classes.lasso.1se, y.test)
metrics.ridge <- calc_metrics(predicted.classes.ridge, y.test)
metrics.elasticnet <- calc_metrics(predicted.classes.elasticnet, y.test)
```

*# Store in a dataframe*

```
df_metrics <- as.data.frame(matrix(c(metrics.log, metrics.lasso.min, metrics.lasso.1se, metrics.ridge, metrics.elasticnet),
  nrow = 5, ncol = 8))
colnames(df_metrics) <- c("Logistic", "Lasso_min", "Lasso_1se", "Ridge", "ElasticNet")
rownames(df_metrics) <- c("AUC", "Recall", "Precision", "F1", "Accuracy")
df_metrics
```

	Logistic	Lasso_min	Lasso_1se	Ridge	ElasticNet
AUC	0.7939336	0.7939336	0.7691382	0.8013962	0.7927299
Recall	0.6774194	0.6774194	0.6129032	0.6774194	0.6451613
Precision	0.7777778	0.7777778	0.7916667	0.8076923	0.8333333
F1	0.7241379	0.7241379	0.6909091	0.7368421	0.7272727
Accuracy	0.8367347	0.8367347	0.8265306	0.8469388	0.8469388

Latex Table Metrics

```
metrics_latex <- df_metrics[-1, ] %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), digits = 3, escape = FALSE,
    col.names = c("Logistic", "Lasso  $\lambda_{\min}$ ",
```

```

      "Lasso  $\lambda_{1se}$ ", "Ridge", "ElasticNet")) %>%
kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
row_spec(0, bold = TRUE) %>%
column_spec(1, bold = TRUE) %>%
add_header_above(c(" " = 1, "Models" = 5), bold = TRUE) %>%
row_spec(nrow(df_metrics) - 2, extra_latex_after = "\\midrule[.08em]") %>%
row_spec(nrow(df_metrics) - 1, bold = TRUE)

print(metrics_latex)

```

	Models				
	Logistic	Lasso $\lambda_{\min}$	Lasso $\lambda_{1se}$	Ridge	ElasticNet
<b>Recall</b>	0.677	0.677	0.613	0.677	0.645
<b>Precision</b>	0.778	0.778	0.792	0.808	0.833
<b>F1</b>	0.724	0.724	0.691	0.737	0.727
<b>Accuracy</b>	<b>0.837</b>	<b>0.837</b>	<b>0.827</b>	<b>0.847</b>	<b>0.847</b>

#### Data Frame of all model coefficients

```

coef.log <- coef(log.model)
coef.lasso.min <- (coef(cv.lasso.model, s = cv.lasso.model$lambda.min))
coef.lasso.1se <- coef(cv.lasso.model, s = cv.lasso.model$lambda.1se)
coef.ridge <- coef(cv.ridge.model)
coef.elasticnet <- coef(elasticnet.model$finalModel, s = elasticnet.model$bestTune$lambda)

# Create a data frame to store coefficients
df_coef <- data.frame(
  Logistic = as.vector(coef.log),
  Lasso_min = as.vector(as.matrix(coef.lasso.min)),
  Lasso_1se = as.vector(as.matrix(coef.lasso.1se)),
  Ridge = as.vector(as.matrix(coef.ridge)),
  ElasticNet = as.vector(as.matrix(coef.elasticnet))
)
rownames(df_coef) <- names(coef.log)

options(scipen = 999)
df_coef_sparse <- round(df_coef, 3)
df_coef_sparse[] <- apply(df_coef_sparse, 2, function(x) ifelse(x == 0, '.', x))

df_coef_sparse

```

```

##           Logistic Lasso_min Lasso_1se Ridge ElasticNet
## (Intercept) -0.976   -0.909   -0.794 -0.777   -0.784
## pregnant    0.348     0.25    0.047  0.157     0.125
## glucose     1.197     1.057    0.806  0.437     0.616
## pressure    -0.008     .        .    0.084     0.004
## triceps      0.099     0.063    .    0.121     0.073
## insulin     -0.011     .        .    0.157     0.067

```

```
## mass      0.527    0.418    0.194  0.181    0.201
## pedigree  0.489    0.362    0.104  0.177    0.158
## age       0.24     0.224    0.178  0.186    0.189
```

## Latex Table Coef

```
coef_latex <- df_coef_sparse %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), escape = FALSE,
        col.names = c("Logistic", "Lasso  $\lambda_{\min}$ ",
                      "Lasso  $\lambda_{1se}$ ", "Ridge", "ElasticNet")) %>%
  kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
  row_spec(0, bold = TRUE) %>%
  column_spec(1, bold = TRUE) %>%
  add_header_above(c(" " = 1, "Models" = 5), bold = TRUE)

print(coef_latex)
```

	Models				
	Logistic	Lasso $\lambda_{\min}$	Lasso $\lambda_{1se}$	Ridge	ElasticNet
(Intercept)	-0.976	-0.909	-0.794	-0.777	-0.784
pregnant	0.348	0.25	0.047	0.157	0.125
glucose	1.197	1.057	0.806	0.437	0.616
pressure	-0.008	.	.	0.084	0.004
triceps	0.099	0.063	.	0.121	0.073
insulin	-0.011	.	.	0.157	0.067
mass	0.527	0.418	0.194	0.181	0.201
pedigree	0.489	0.362	0.104	0.177	0.158
age	0.24	0.224	0.178	0.186	0.189

## Confusion matrix latex - log.model

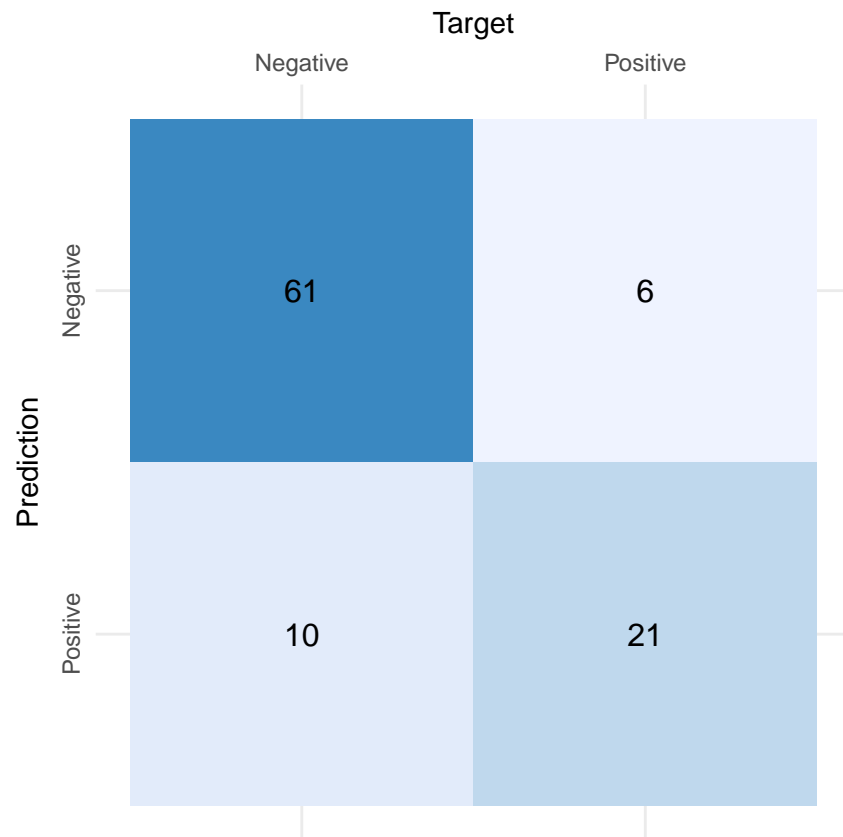
```
# Convert the table to a data frame
log.conf.df <- as.data.frame(log.conf$table)

# Replace the numeric labels with text labels
log.conf.df$Prediction <- factor(log.conf.df$Prediction, levels = c("0", "1"), labels = c("Negative", "Positive"))
log.conf.df$Reference <- factor(log.conf.df$Reference, levels = c("0", "1"), labels = c("Negative", "Positive"))

# Rename columns as required by cums
names(log.conf.df) <- c('Target', 'Prediction', 'N')
# Plot the confusion matrix using cums
plot_confusion_matrix(log.conf.df, add_normalized = FALSE,
  add_row_percentages = FALSE, add_col_percentages = FALSE, rotate_y_text = TRUE,
  place_x_axis_above = TRUE, class_order = c("Positive", "Negative"))

## Warning in plot_confusion_matrix(log.conf.df, add_normalized = FALSE,
## add_row_percentages = FALSE, : 'ggimage' is missing. Will not plot arrows and
## zero-shading.
```

```
## Warning in plot_confusion_matrix(log.conf.df, add_normalized = FALSE,
## add_row_percentages = FALSE, : 'rsvg' is missing. Will not plot arrows and
## zero-shading.
```

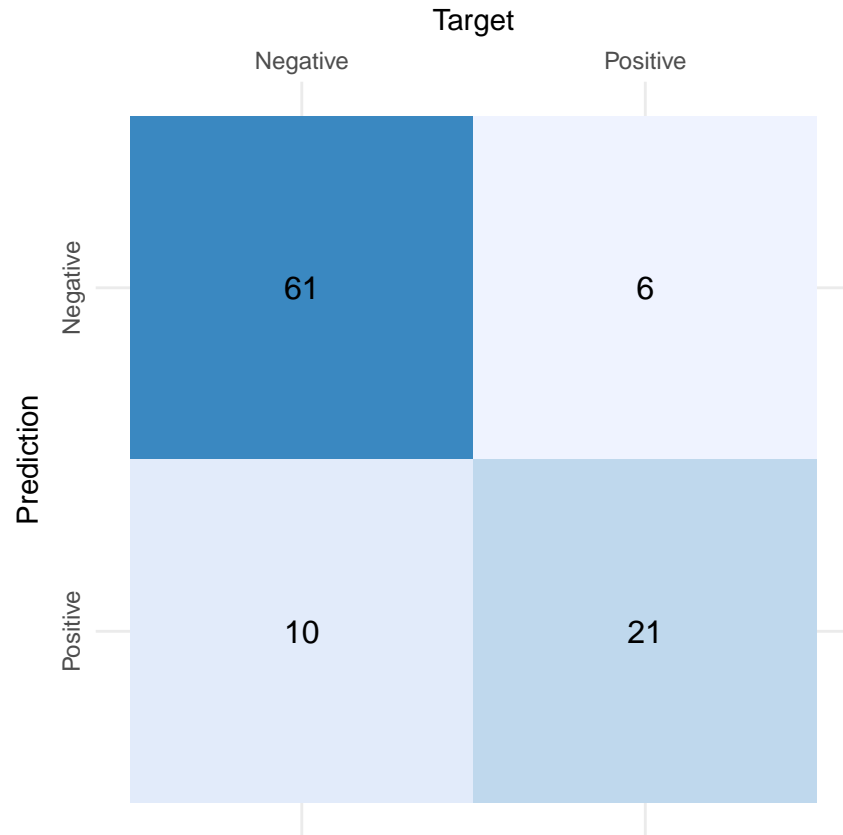


```
### Confusion matrix latex - lasso.min.model
```

```
# For lasso.min.conf
lasso.min.conf.df <- as.data.frame(lasso.min.conf$table)
lasso.min.conf.df$Prediction <- factor(lasso.min.conf.df$Prediction, levels = c("0", "1"), labels = c("0", "1"))
lasso.min.conf.df$Reference <- factor(lasso.min.conf.df$Reference, levels = c("0", "1"), labels = c("0", "1"))
names(lasso.min.conf.df) <- c('Target', 'Prediction', 'N')
plot_confusion_matrix(lasso.min.conf.df, add_normalized = FALSE, add_row_percentages = FALSE, add_col_p
```

```
## Warning in plot_confusion_matrix(lasso.min.conf.df, add_normalized = FALSE, :
## 'ggimage' is missing. Will not plot arrows and zero-shading.

## Warning in plot_confusion_matrix(lasso.min.conf.df, add_normalized = FALSE, :
## 'rsvg' is missing. Will not plot arrows and zero-shading.
```

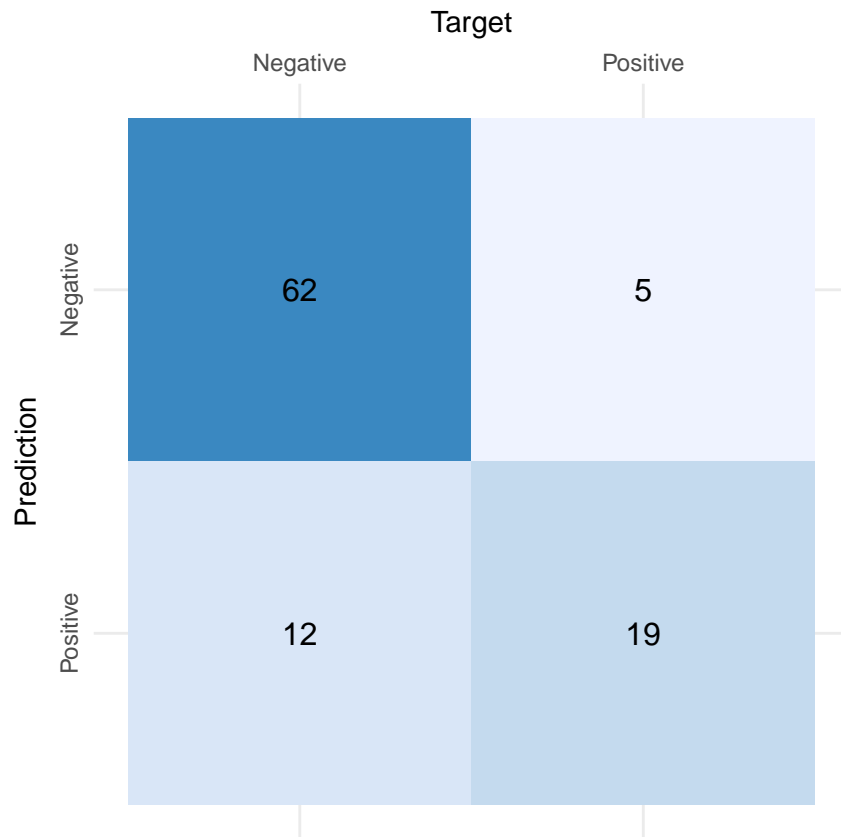


Confusion matrix latex - lasso.1se.model

```
lasso.1se.conf.df <- as.data.frame(lasso.1se.conf$table)
lasso.1se.conf.df$Prediction <- factor(lasso.1se.conf.df$Prediction, levels = c("0", "1"), labels = c("Neg", "Pos"))
lasso.1se.conf.df$Reference <- factor(lasso.1se.conf.df$Reference, levels = c("0", "1"), labels = c("Neg", "Pos"))
names(lasso.1se.conf.df) <- c('Target', 'Prediction', 'N')
plot_confusion_matrix(lasso.1se.conf.df, add_normalized = FALSE, add_row_percentages = FALSE, add_col_percentages = FALSE)

## Warning in plot_confusion_matrix(lasso.1se.conf.df, add_normalized = FALSE, :
## 'ggimage' is missing. Will not plot arrows and zero-shading.

## Warning in plot_confusion_matrix(lasso.1se.conf.df, add_normalized = FALSE, :
## 'rsvg' is missing. Will not plot arrows and zero-shading.
```

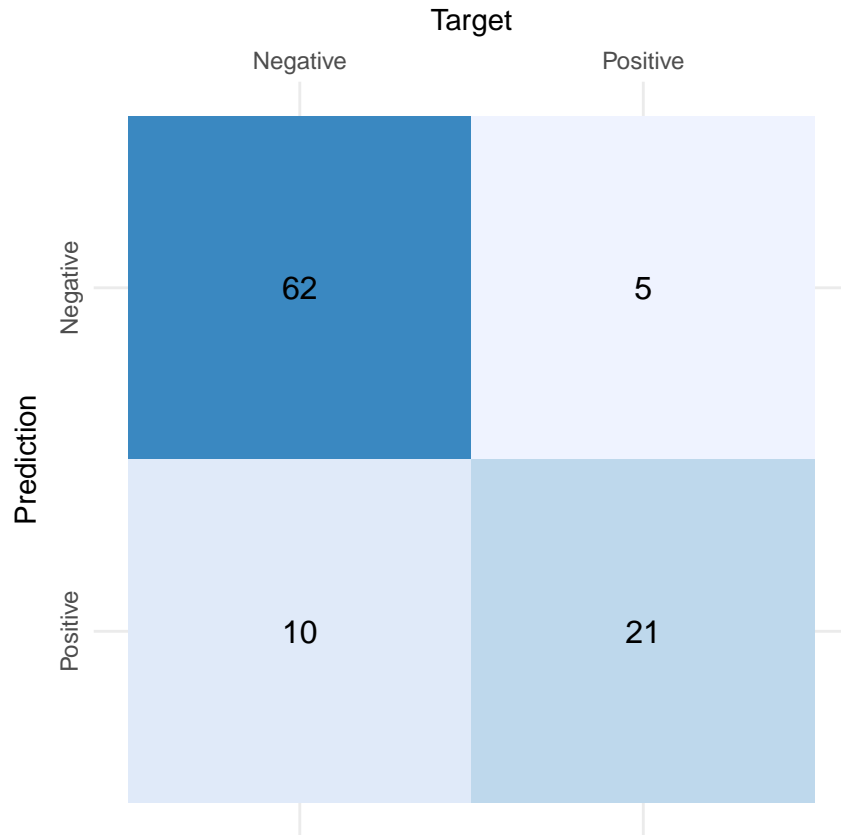


Confusion matrix latex - ridge.model

```
# For ridge.conf
ridge.conf.df <- as.data.frame(ridge.conf$table)
ridge.conf.df$Prediction <- factor(ridge.conf.df$Prediction, levels = c("0", "1"), labels = c("Negative", "Positive"))
ridge.conf.df$Reference <- factor(ridge.conf.df$Reference, levels = c("0", "1"), labels = c("Negative", "Positive"))
names(ridge.conf.df) <- c('Target', 'Prediction', 'N')
plot_confusion_matrix(ridge.conf.df, add_normalized = FALSE, add_row_percentages = FALSE, add_col_percentages = FALSE)

## Warning in plot_confusion_matrix(ridge.conf.df, add_normalized = FALSE, :
## 'ggimage' is missing. Will not plot arrows and zero-shading.

## Warning in plot_confusion_matrix(ridge.conf.df, add_normalized = FALSE, :
## 'rsvg' is missing. Will not plot arrows and zero-shading.
```



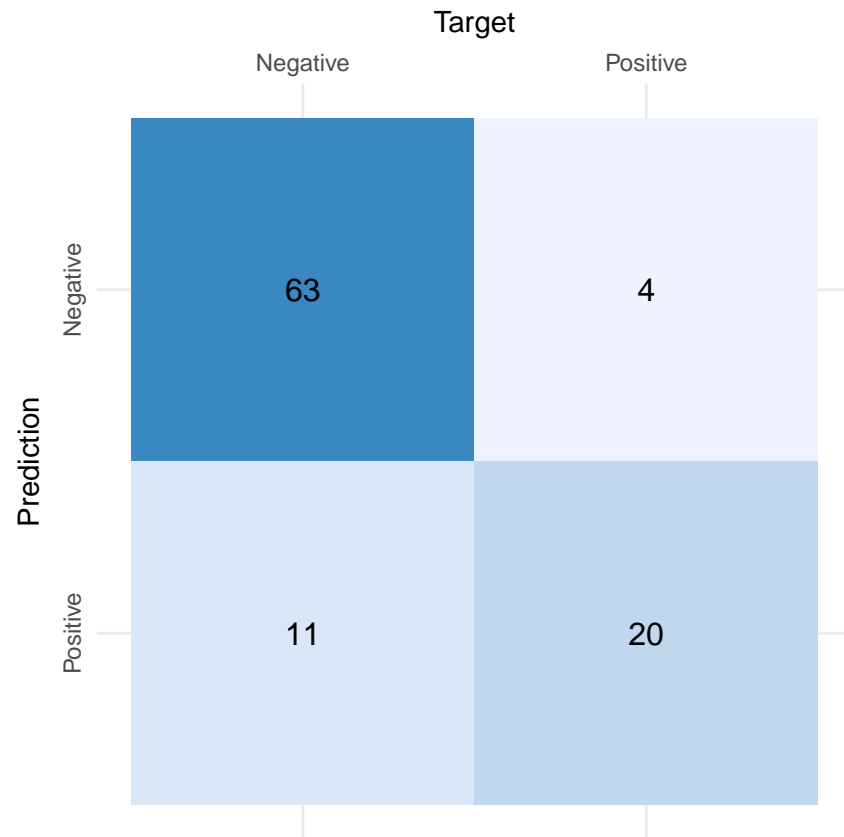
#### Confusion matrix latex - elasticnet.model

```
# For elasticnet.conf
elasticnet.conf.df <- as.data.frame(elasticnet.conf$table)
elasticnet.conf.df$Prediction <- factor(elasticnet.conf.df$Prediction, levels = c("0", "1"), labels = c("Negative", "Positive"))
elasticnet.conf.df$Reference <- factor(elasticnet.conf.df$Reference, levels = c("0", "1"), labels = c("Negative", "Positive"))
names(elasticnet.conf.df) <- c('Target', 'Prediction', 'N')
plot_confusion_matrix(elasticnet.conf.df, add_normalized = FALSE, add_row_percentages = FALSE, add_col_percentages = FALSE)
```

```
## Warning in plot_confusion_matrix(elasticnet.conf.df, add_normalized = FALSE, :
## 'ggimage' is missing. Will not plot arrows and zero-shading.

## Warning in plot_confusion_matrix(elasticnet.conf.df, add_normalized = FALSE, :
## 'rsvg' is missing. Will not plot arrows and zero-shading.
```

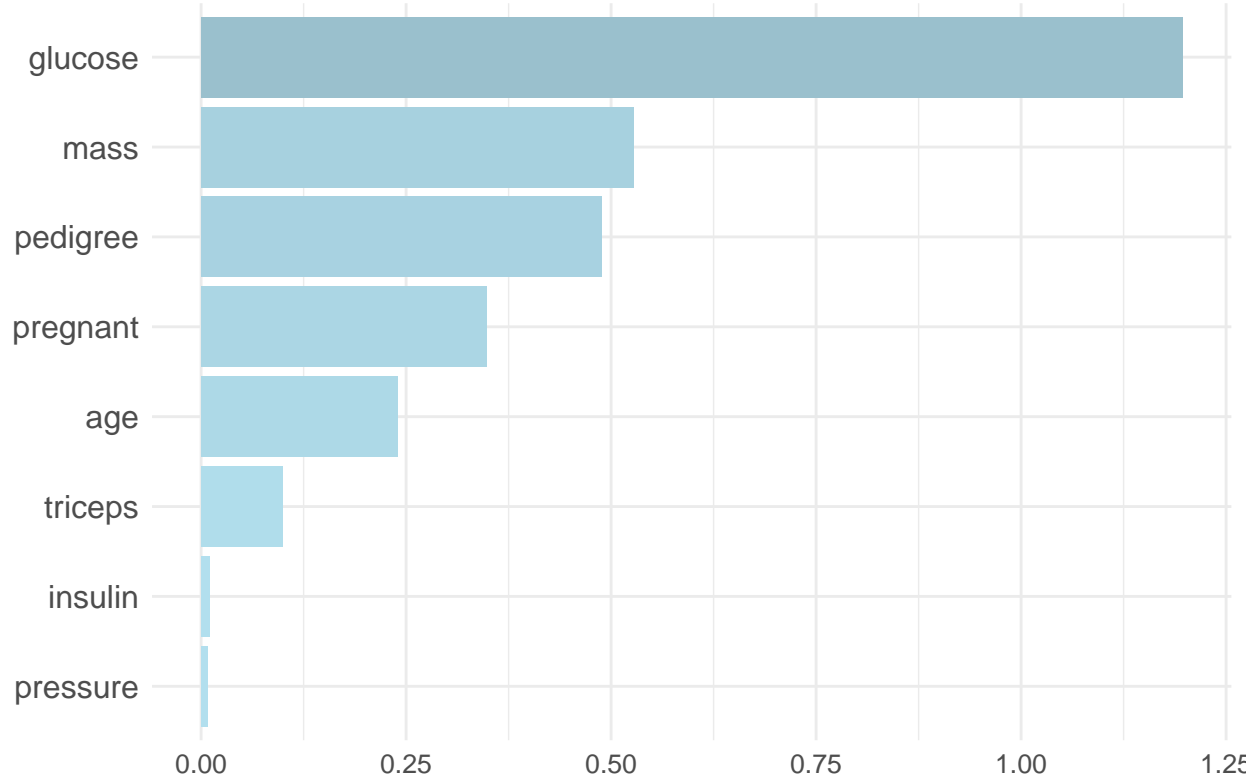




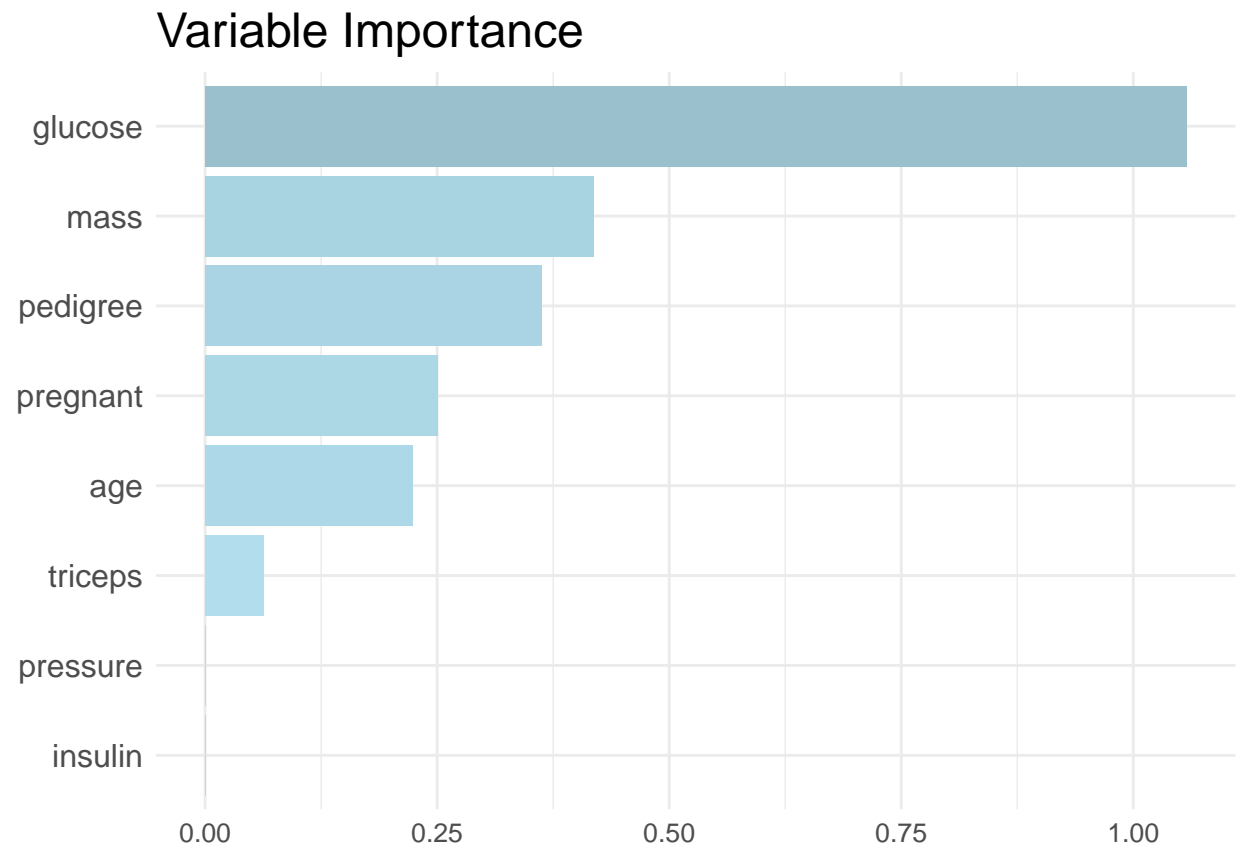
### Variable Importance Plots

```
importance_plot("Logistic")
```

## Variable Importance

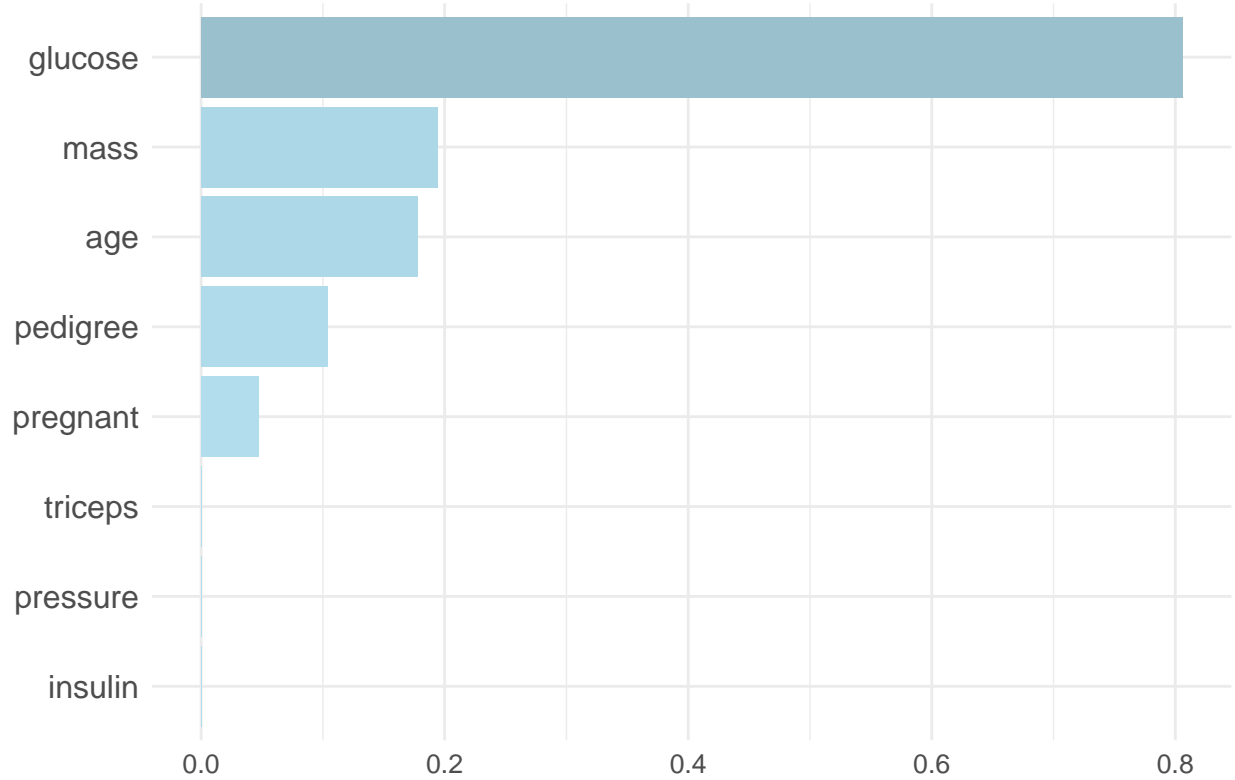


```
importance_plot("Lasso_min")
```

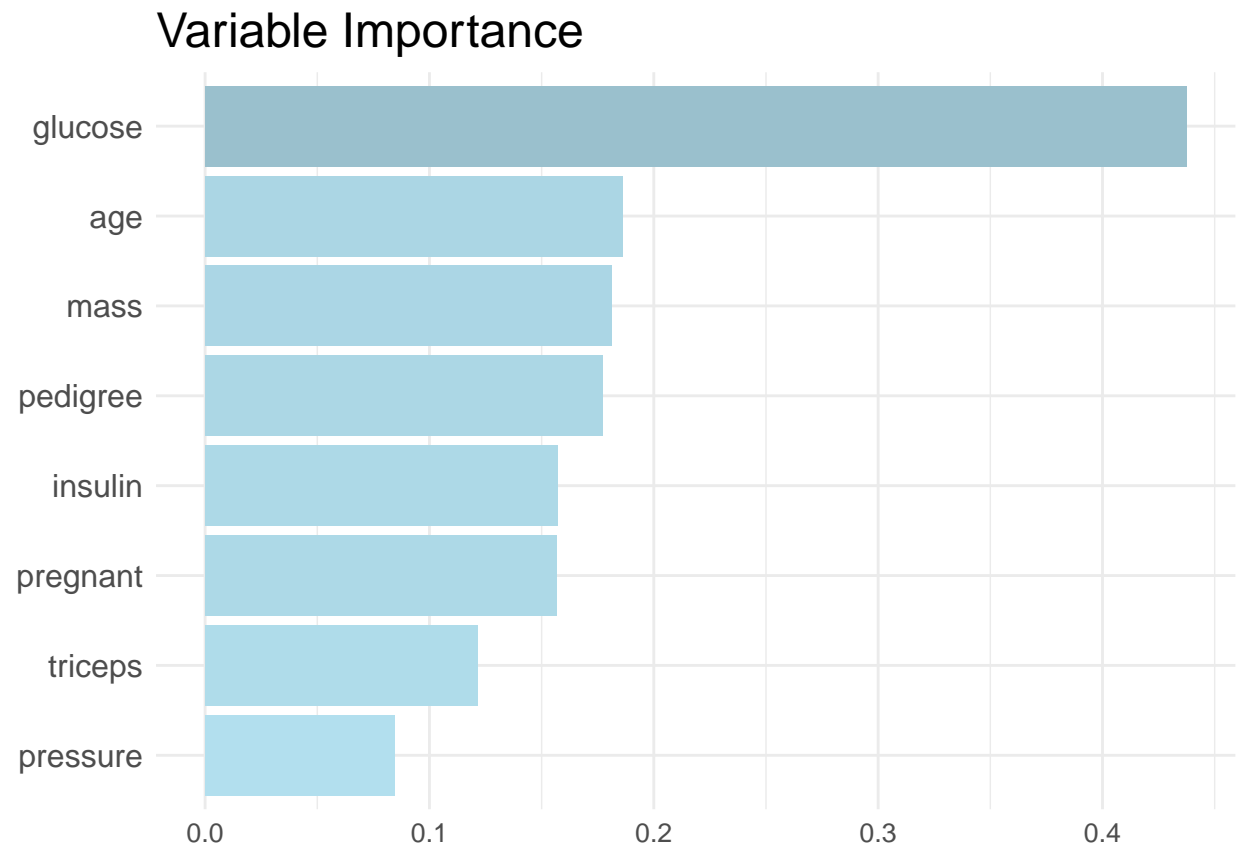


```
importance_plot("Lasso_1se")
```

## Variable Importance

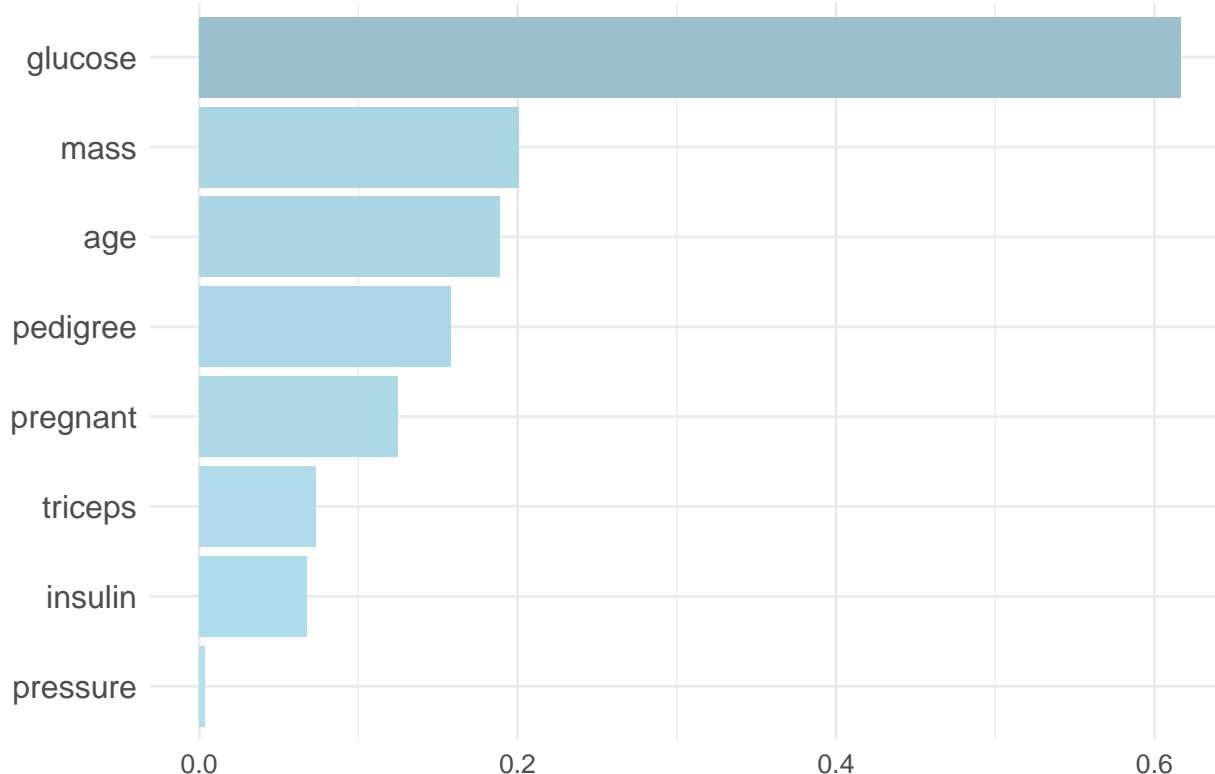


```
importance_plot("Ridge")
```



```
importance_plot("ElasticNet")
```

## Variable Importance



### logistic model summary table latex

```
# Extract coefficients from the model summary
coef_summary <- summary(log.model)$coefficients

# Convert the matrix to a data frame for kable()
coef_summary <- round(as.data.frame(coef_summary)[, c(1, 4)], 3)
colnames(coef_summary) <- c("Estimate", "p-value")

# coef_summary$p-value` <- cell_spec(coef_summary$p-value`, "latex",
#                                     bold = ifelse(coef_summary$p-value` < 0.05, TRUE, FALSE))

# Create the LaTeX table
# kable(coef_summary, "latex", booktabs = TRUE, align = "c", escape = F) %>%
#   kable_styling(latex_options = c("striped", "scale_down", "hold_position")) %>%
#   row_spec(0, bold = TRUE)

# Add stars based on significance level and make significant p-values bold
coef_summary$p-value` <- sapply(coef_summary$p-value`, function(p) {
  stars <- ifelse(p < 0.001, "***", ifelse(p < 0.01, "**", ifelse(p < 0.05, "*", "")))
  value_with_stars <- paste0(round(p, 3), stars)
  cell_spec(value_with_stars, "latex", bold = p < 0.05)
})

# Create the LaTeX table
kable(coef_summary, "latex", booktabs = TRUE, align = "c", escape = F) %>%
  kable_styling(latex_options = c("striped", "scale_down", "hold_position")) %>%
  row_spec(0, bold = TRUE)
```

	Estimate	p-value
(Intercept)	-0.976	<b>0***</b>
pregnant	0.348	0.087
glucose	1.197	<b>0***</b>
pressure	-0.008	0.963
triceps	0.099	0.622
insulin	-0.011	0.95
mass	0.527	<b>0.016*</b>
pedigree	0.489	<b>0.004**</b>
age	0.240	0.251

continued log model latex table

```
coef_latex <- coef_summary %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), escape = FALSE,
        col.names = c("Estimate", "p value")) %>%
  kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
  row_spec(0, bold = TRUE) %>%
  column_spec(1, bold = TRUE) %>%
  #add_header_above(c(" " = 1, "Logistic Regression" = 2), bold = TRUE)

print(coef_latex)
```

coef latex per model - lasso both

```
coef_latex <- df_coef_sparse[, 2:3] %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), escape = FALSE,
        col.names = c("Lasso  $\\lambda_{\\text{min}}$ ", "Lasso  $\\lambda_{\\text{1se}}$ ")) %>%
  kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
  row_spec(0, bold = TRUE) %>%
  column_spec(1, bold = TRUE) %>%
```

	Estimate	p value
(Intercept)	-0.976	<b>0***</b>
pregnant	0.348	0.087
glucose	1.197	<b>0***</b>
pressure	-0.008	0.963
triceps	0.099	0.622
insulin	-0.011	0.95
mass	0.527	<b>0.016*</b>
pedigree	0.489	<b>0.004**</b>
age	0.240	0.251

```
add_header_above(c(" " = 1, "Estimates" = 2), bold = TRUE)

print(coef_latex)
```

coef latex per model - log, ridge

```
coef_latex <- df_coef_sparse[, c(1,4)] %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), escape = FALSE,
        col.names = c("Logistic", "Ridge")) %>%
  kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
  row_spec(0, bold = TRUE) %>%
  column_spec(1, bold = TRUE) %>%
  add_header_above(c(" " = 1, "Estimates" = 2), bold = TRUE)

print(coef_latex)
```

coef latex per model - log, elastic net

```
coef_latex <- df_coef_sparse[, c(1,5)] %>%
  kable(format = "latex", booktabs = TRUE, align = c('c'), escape = FALSE,
        col.names = c("Logistic", "ElasticNet")) %>%
```



	Estimates	
	Lasso $\lambda_{\min}$	Lasso $\lambda_{1se}$
(Intercept)	-0.909	-0.794
pregnant	0.25	0.047
glucose	1.057	0.806
pressure	.	.
triceps	0.063	.
insulin	.	.
mass	0.418	0.194
pedigree	0.362	0.104
age	0.224	0.178

```
kable_styling(latex_options = c("striped", "scale_down", "hold_position"), position = "center") %>%
row_spec(0, bold = TRUE) %>%
column_spec(1, bold = TRUE) %>%
add_header_above(c(" " = 1, "Estimates" = 2), bold = TRUE)

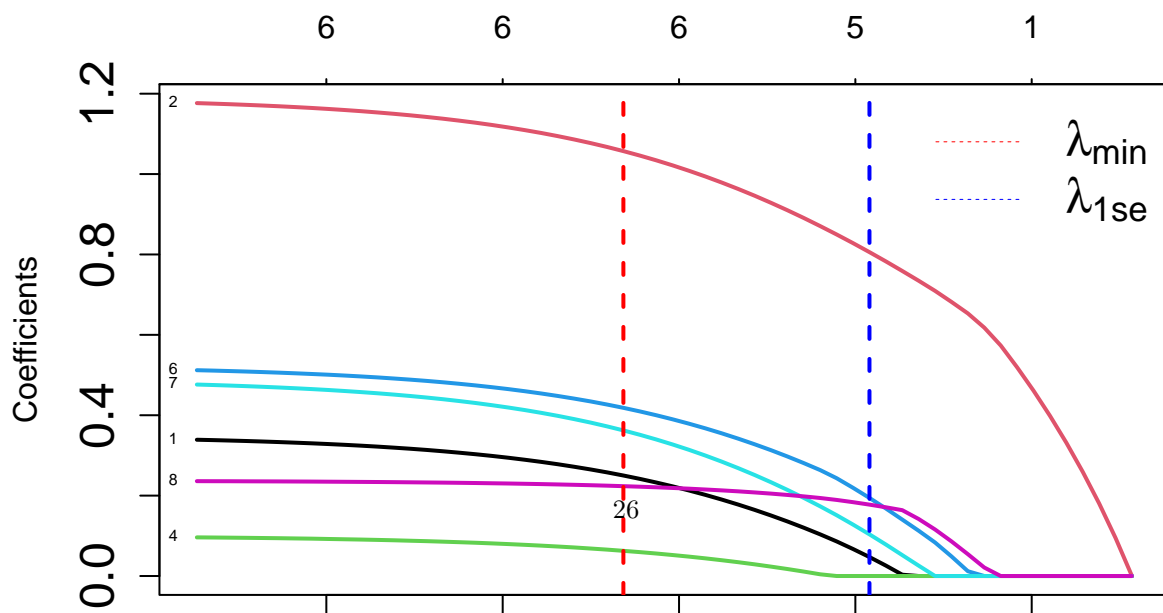
print(coef_latex)
```

Coefficient path plot - not included in slides

```
plot(cv.lasso.model$glmnet.fit, xvar = "lambda", label = TRUE, cex.axis = 1.5, lwd = 2)
abline(v = log(cv.lasso.model$lambda.min), col = "red", lty = 2, lwd = 2)
abline(v = log(cv.lasso.model$lambda.1se), col = "blue", lty = 2, lwd = 2)

# Add a legend with LaTeX-style text and increased font size, and thicker lines
legend("topright",
      legend = expression(lambda["min"], lambda["1se"]),
      col = c("red", "blue"),
      lty = 2,
      cex = 1.5,
      lwd = 0.5, bty = "n")
```

	Estimates	
	Logistic	Ridge
(Intercept)	-0.976	-0.777
pregnant	0.348	0.157
glucose	1.197	0.437
pressure	-0.008	0.084
triceps	0.099	0.121
insulin	-0.011	0.157
mass	0.527	0.181
pedigree	0.489	0.177
age	0.24	0.186



	Estimates	
	Logistic	ElasticNet
(Intercept)	-0.976	-0.784
pregnant	0.348	0.125
glucose	1.197	0.616
pressure	-0.008	0.004
triceps	0.099	0.073
insulin	-0.011	0.067
mass	0.527	0.201
pedigree	0.489	0.158
age	0.24	0.189

### Manual calculation of p-value - Acc > NIR – Logistic Model Example -> identical to result from confusionMatrix()

```
# Number of trials in test data
n <- 98

# Number of successes (correctly classified instances)
# Accuracy is 0.8367, so the number of successes would be 0.8367 * 392
n_success <- round(0.8367 * n)

# No Information Rate
nir <- 0.6837

# Perform binomial test
binom_test_result <- binom.test(n_success, n, p = nir, alternative = "greater")

# Print the result
print(binom_test_result)
```

```
##
## Exact binomial test
```

```
##
## data:  n_success and n
## number of successes = 82, number of trials = 98, p-value = 0.0004546
## alternative hypothesis: true probability of success is greater than 0.6837
## 95 percent confidence interval:
##  0.7626479 1.0000000
## sample estimates:
## probability of success
##          0.8367347
```

Accuracy and p value table

```
log.conf$overall["Accuracy"]
```

Accuracy 0.8367347

```
log.conf$overall["AccuracyPValue"]
```

AccuracyPValue 0.0004537699

```
tibble(Accuracy = log.conf$overall["Accuracy"],
       `p-value [Acc > NIR]` = log.conf$overall["AccuracyPValue"])
```

A tibble: 1 x 2

Accuracy p-value [Acc > NIR] 1 0.837 0.000454

```
result_table <- tibble(Accuracy = log.conf$overall["Accuracy"],
                      `p-value [Acc > NIR]` = log.conf$overall["AccuracyPValue"])
```

*# Round values for better formatting, if needed*

```
result_table$Accuracy <- round(result_table$Accuracy, 4)
```

```
result_table$`p-value [Acc > NIR]` <- round(result_table$`p-value [Acc > NIR]`, 4)
```

*# Generate LaTeX table*

```
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

---

Accuracy	p-value [Acc > NIR]
0.8367	0.0005

---

Accuracy and p value table

```
result_table <- tibble(Accuracy = log.conf$overall["Accuracy"],
                      `p-value [Acc > NIR]` = log.conf$overall["AccuracyPValue"])
```

*# Round values for better formatting, if needed*

```
result_table$Accuracy <- round(result_table$Accuracy, 4)
result_table$p-value [Acc > NIR]` <- round(result_table$p-value [Acc > NIR]`, 4)

# Generate LaTeX table
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

---

Accuracy	p-value [Acc > NIR]
0.8367	0.0005

---

Accuracy and p value table

```
result_table <- tibble(Accuracy = ridge.conf$overall["Accuracy"],
  p-value [Acc > NIR]` = ridge.conf$overall["AccuracyPValue"])

# Round values for better formatting, if needed
result_table$Accuracy <- round(result_table$Accuracy, 4)
result_table$p-value [Acc > NIR]` <- round(result_table$p-value [Acc > NIR]`, 4)

# Generate LaTeX table
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

---

Accuracy	p-value [Acc > NIR]
0.8469	0.0002

---

Accuracy and p value table

```
#c("Lasso $\lambda_{\text{min}}$", "Lasso $\lambda_{\text{1se}}$")
result_table <- tibble(Accuracy = lasso.min.conf$overall["Accuracy"],
  p-value [Acc > NIR]` = lasso.min.conf$overall["AccuracyPValue"])

# Round values for better formatting, if needed
result_table$Accuracy <- round(result_table$Accuracy, 4)
result_table$p-value [Acc > NIR]` <- round(result_table$p-value [Acc > NIR]`, 4)

# Generate LaTeX table
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
```

```
kable_styling(latex_options = c("scale_down", "hold_position"))
```

Accuracy	p-value [Acc > NIR]
0.8367	0.0005

Accuracy and p value table

```
result_table <- tibble(Accuracy = lasso.1se.conf$overall["Accuracy"],
  `p-value [Acc > NIR]` = lasso.1se.conf$overall["AccuracyPValue"])

# Round values for better formatting, if needed
result_table$Accuracy <- round(result_table$Accuracy, 4)
result_table$`p-value [Acc > NIR]` <- round(result_table$`p-value [Acc > NIR]`, 4)

# Generate LaTeX table
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

Accuracy	p-value [Acc > NIR]
0.8265	0.0011

Accuracy and p value table

```
result_table <- tibble(Accuracy = elasticnet.conf$overall["Accuracy"],
  `p-value [Acc > NIR]` = elasticnet.conf$overall["AccuracyPValue"])

# Round values for better formatting, if needed
result_table$Accuracy <- round(result_table$Accuracy, 4)
result_table$`p-value [Acc > NIR]` <- round(result_table$`p-value [Acc > NIR]`, 4)

# Generate LaTeX table
kable(result_table, "latex", booktabs = TRUE, align = "c") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

Plotting the distributions of the 8 predictor variables

```
par(mfrow=c(2, 4))
hist(diabetes$glucose, breaks = 5, col = "sky blue", main = "Glucose", xlab = NULL)
```

---

Accuracy	p-value [Acc > NIR]
----------	---------------------

---

0.8469

0.0002

```
hist(diabetes$pressure, breaks = 5, col = "sky blue", main = "Blood Pressure", xlab = NULL)
hist(diabetes$triceps, breaks = 10, col = "sky blue", main = "Skin Thickness", xlab = NULL)
hist(diabetes$mass, breaks = 10, col = "sky blue", main = "BMI", xlab = NULL)
hist(diabetes$pregnant, breaks = 10, col = "sky blue", main = "No. of Pregnancies", xlab = NULL)
hist(diabetes$insulin, breaks = 10, col = "sky blue", main = "Insulin", xlab = NULL)
hist(diabetes$pedigree, breaks = 10, col = "sky blue", main = "Pedigree", xlab = NULL)
hist(diabetes$age, breaks = 10, col = "sky blue", main = "Age", xlab = NULL)
```

