

D621 - Assignment 3

Marco Castro, Mubashira Qari, Puja Roy, Zach Rose & Erick Hadi

2025-03-07

DATA EXPLORATION

The training dataset contains 466 observations and 13 variables, including 12 predictor variables and one binary response variable (target). The target variable indicates whether a neighborhood's crime rate is above the median (1) or not (0). The dataset includes a mix of continuous and categorical features, such as housing characteristics, pollution levels, property taxes, and proximity to employment centers. Understanding the distributions, relationships, and correlations between these features is an essential first step in building an effective predictive model.

```
dim(df_training)
```

```
## [1] 466 13
```

```
str(df_training)
```

```
## spc_tbl_ [466 x 13] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ zn      : num [1:466] 0 0 0 30 0 0 0 0 0 80 ...
## $ indus   : num [1:466] 19.58 19.58 18.1 4.93 2.46 ...
## $ chas    : num [1:466] 0 1 0 0 0 0 0 0 0 0 ...
## $ nox     : num [1:466] 0.605 0.871 0.74 0.428 0.488 0.52 0.693 0.693 0.515 0.392 ...
## $ rm      : num [1:466] 7.93 5.4 6.49 6.39 7.16 ...
## $ age     : num [1:466] 96.2 100 100 7.8 92.2 71.3 100 100 38.1 19.1 ...
## $ dis     : num [1:466] 2.05 1.32 1.98 7.04 2.7 ...
## $ rad     : num [1:466] 5 5 24 6 3 5 24 24 5 1 ...
## $ tax     : num [1:466] 403 403 666 300 193 384 666 666 224 315 ...
## $ ptratio : num [1:466] 14.7 14.7 20.2 16.6 17.8 20.9 20.2 20.2 20.2 16.4 ...
## $ lstat   : num [1:466] 3.7 26.82 18.85 5.19 4.82 ...
## $ medv    : num [1:466] 50 13.4 15.4 23.7 37.9 26.5 5 7 22.2 20.9 ...
## $ target  : num [1:466] 1 1 1 0 0 0 1 1 0 0 ...
## - attr(*, "spec")=
## .. cols(
## ..   zn = col_double(),
## ..   indus = col_double(),
## ..   chas = col_double(),
## ..   nox = col_double(),
## ..   rm = col_double(),
## ..   age = col_double(),
## ..   dis = col_double(),
## ..   rad = col_double(),
## ..   tax = col_double(),
```

```
## .. ptratio = col_double(),
## .. lstat = col_double(),
## .. medv = col_double(),
## .. target = col_double()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
column_types <- supply(df_training, class)
print(column_types)
```

```
##      zn      indus      chas      nox      rm      age      dis      rad
## "numeric" "numeric" "numeric" "numeric" "numeric" "numeric" "numeric" "numeric"
##      tax      ptratio      lstat      medv      target
## "numeric" "numeric" "numeric" "numeric" "numeric"
```

The following three columns were imported as numerical but should be considered for converting to factors:
 - chas: binomial - rad: ordinal - target: binomial

```
# convert to factor
df_training <- df_training |>
  mutate(
    chas = as.factor(chas),
    rad = as.factor(rad),
    target = as.factor(target),
  )

numeric_cols <- c('zn', 'indus', 'nox', 'rm', 'age', 'dis', 'tax', 'ptratio', 'lstat', 'medv')

factor_cols <- c('chas', 'rad', 'target')

glimpse(df_training)
```

```
## Rows: 466
## Columns: 13
## $ zn      <dbl> 0, 0, 0, 30, 0, 0, 0, 0, 0, 80, 22, 0, 0, 22, 0, 0, 100, 20, 0~
## $ indus   <dbl> 19.58, 19.58, 18.10, 4.93, 2.46, 8.56, 18.10, 18.10, 5.19, 3.6~
## $ chas    <fct> 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,~
## $ nox     <dbl> 0.605, 0.871, 0.740, 0.428, 0.488, 0.520, 0.693, 0.693, 0.515,~
## $ rm      <dbl> 7.929, 5.403, 6.485, 6.393, 7.155, 6.781, 5.453, 4.519, 6.316,~
## $ age     <dbl> 96.2, 100.0, 100.0, 7.8, 92.2, 71.3, 100.0, 100.0, 38.1, 19.1,~
## $ dis     <dbl> 2.0459, 1.3216, 1.9784, 7.0355, 2.7006, 2.8561, 1.4896, 1.6582~
## $ rad     <fct> 5, 5, 24, 6, 3, 5, 24, 24, 5, 1, 7, 5, 24, 7, 3, 3, 5, 5, 24, ~
## $ tax     <dbl> 403, 403, 666, 300, 193, 384, 666, 666, 224, 315, 330, 398, 66~
## $ ptratio <dbl> 14.7, 14.7, 20.2, 16.6, 17.8, 20.9, 20.2, 20.2, 20.2, 16.4, 19~
## $ lstat   <dbl> 3.70, 26.82, 18.85, 5.19, 4.82, 7.67, 30.59, 36.98, 5.68, 9.25~
## $ medv    <dbl> 50.0, 13.4, 15.4, 23.7, 37.9, 26.5, 5.0, 7.0, 22.2, 20.9, 24.8~
## $ target  <fct> 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0,~
```

A closer examination of the rad data shows that that our observations have a rad index value of 1-8 or 24 in this column. Below are the counts:

```
##
##  1  2  3  4  5  6  7  8  24
## 17 20 36 103 109 25 15 20 121
```

Summary Statistics

Summary statistics reveal substantial variability in several features. For instance, the variable tax (full-value property tax rate per \$10,000) has a mean of 409.5 but a maximum value of 711, indicating a right-skewed distribution with significant outliers. Similar patterns are observed in zn (zoned residential land), age, and dis (distance to employment centers). These skewed distributions may require transformation to reduce leverage effects during modeling.

```
# only show summary stats for numeric values
for (param in numeric_cols) {
  cat("\nSummary for", param, ":\n")
  print(describe(df_training[[param]]))
}

##
## Summary for zn :
##   vars  n mean   sd median trimmed mad min max range skew kurtosis  se
## X1    1 466 11.58 23.36      0   5.35  0  0 100   100 2.18      3.81 1.08
##
## Summary for indus :
##   vars  n mean   sd median trimmed mad min  max range skew kurtosis  se
## X1    1 466 11.11 6.85   9.69   10.91 9.34 0.46 27.74 27.28 0.29   -1.24 0.32
##
## Summary for nox :
##   vars  n mean   sd median trimmed mad min  max range skew kurtosis  se
## X1    1 466 0.55 0.12   0.54   0.54 0.13 0.39 0.87  0.48 0.75   -0.04 0.01
##
## Summary for rm :
##   vars  n mean   sd median trimmed mad min  max range skew kurtosis  se
## X1    1 466 6.29 0.7   6.21   6.26 0.52 3.86 8.78  4.92 0.48    1.54 0.03
##
## Summary for age :
##   vars  n mean   sd median trimmed mad min max range skew kurtosis  se
## X1    1 466 68.37 28.32 77.15   70.96 30.02 2.9 100  97.1 -0.58   -1.01 1.31
##
## Summary for dis :
##   vars  n mean   sd median trimmed mad min  max range skew kurtosis  se
## X1    1 466  3.8 2.11   3.19   3.54 1.91 1.13 12.13   11  1    0.47 0.1
##
## Summary for tax :
##   vars  n mean   sd median trimmed mad min max range skew kurtosis  se
## X1    1 466 409.5 167.9 334.5  401.51 104.52 187 711   524 0.66   -1.15 7.78
##
## Summary for ptratio :
##   vars  n mean   sd median trimmed mad min max range skew kurtosis  se
## X1    1 466 18.4 2.2   18.9   18.6 1.93 12.6  22   9.4 -0.75   -0.4 0.1
##
## Summary for lstat :
##   vars  n mean   sd median trimmed mad min  max range skew kurtosis  se
## X1    1 466 12.63 7.1   11.35   11.88 7.07 1.73 37.97 36.24 0.91    0.5 0.33
##
## Summary for medv :
##   vars  n mean   sd median trimmed mad min max range skew kurtosis  se
## X1    1 466 22.59 9.24   21.2   21.63  6  5  50   45 1.08    1.37 0.43
```

Key Observations: Crime Rate Target (target)

The median (p50) is 0, indicating that more than half of the data points fall in the low-crime category (target = 0).

Median Home Value (medv) Mean = 22.59 (\$22,590 in \$1000s), Median = 21.2. The range (p0 = 5, p75 = 25) suggests that most homes are valued between \$5,000 and \$25,000 (in \$1000s). The standard deviation (9.23) indicates a relatively high spread in home values.

Lower Status Population (lstat) Mean = 12.63%, Median = 11.93%. A positively skewed distribution (p0 = 1.73, p75 = 16.93), meaning some areas have much higher lower-status populations than others.

Property Tax Rate (tax) High variance (Mean = 409.5, SD = 167.9). Large difference between the 25th percentile (281) and 75th percentile (666), suggesting significant variability in tax rates among neighborhoods.

Average Number of Rooms (rm) Mean = 6.29, Median = 6.21, with a relatively small spread (SD = 0.70). Indicates most homes have around 6 rooms.

Distance to Employment Centers (dis) Median = 3.19, but the 25th percentile is quite low (2.10), meaning some neighborhoods are much closer to employment centers than others. Higher standard deviation (2.1) suggests some neighborhoods are much more remote.

Industrial Land Proportion (indus) Mean = 11.10, Median = 9.69, and right-skewed distribution (p0 = 0.46, p75 = 18.1). Some areas have much higher proportions of industrial land, potentially influencing crime.

Highway Accessibility (rad) Highly right-skewed: The median is 5, but the 75th percentile is 24, meaning some neighborhoods have much greater access to highways than others. This might be an important predictor for crime.

Potential Data Transformations zn, indus, tax, rad, lstat, and medv show skewness, so applying a log transformation might improve normality. age can be categorized into bins (e.g., young, middle-aged, old) since it ranges from 2.9 to 94.1. dis has a wide range, so normalization might be needed.

Missing data

No missing values were found in the training dataset. Therefore, no imputation or flagging was necessary at this stage. Several continuous variables demonstrated right-skewed distributions and a high number of extreme values. To reduce the influence of outliers and better align the data with the assumptions of logistic regression, log-transformations were applied to tax, zn, dis, and lstat. This transformation helps normalize the data, reduce variance, and enhance model interpretability. A small constant was added to zn before the transformation to account for zero values.

```
#introduce(df_training, echo=FALSE)
missing_values_count <- sapply(data, function(x) sum(is.na(x)))
print(missing_values_count)
```

```
##      ...      list  package  lib.loc  verbose  enviro overwrite
##      0         0         0         0         0         0         0         0
```

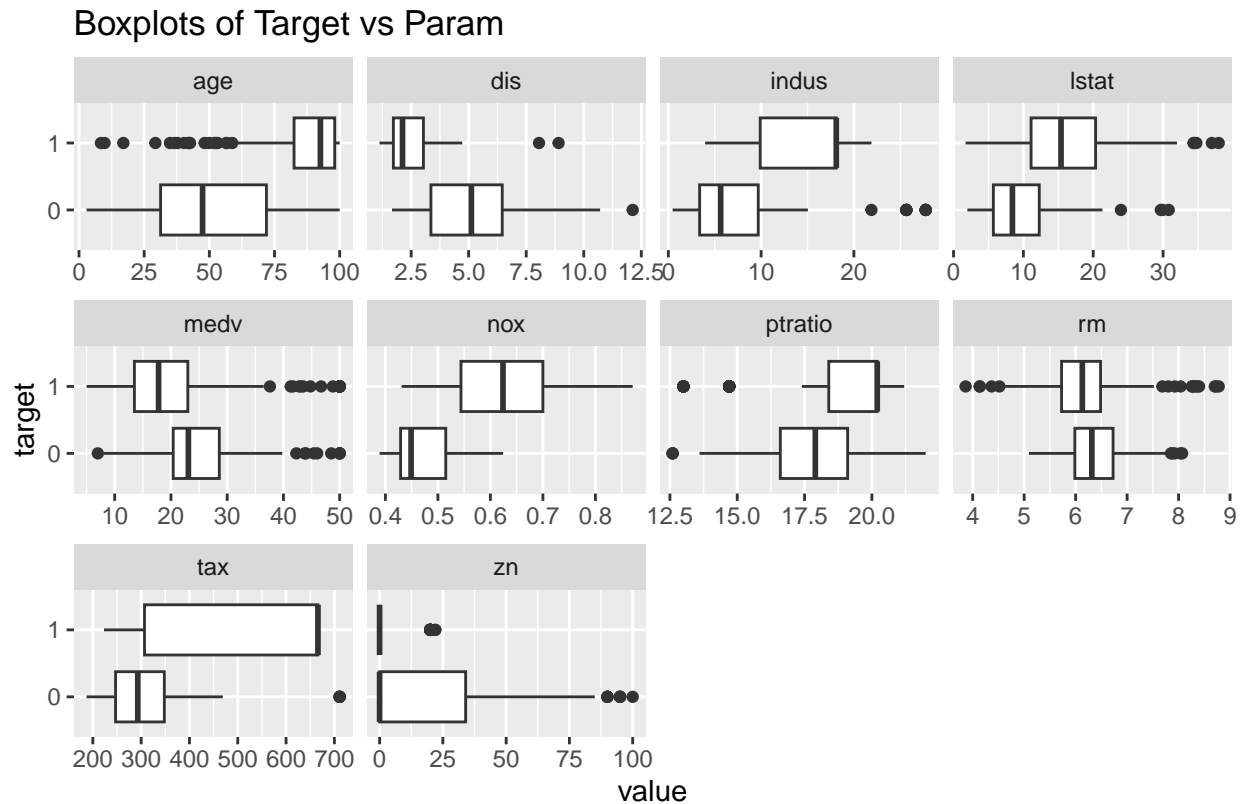
Plots of data

Boxplots A boxplot of the numeric features further highlights the presence of skewness and outliers in variables such as tax, zn, and age. Many variables, such as chas (a binary indicator for bordering the Charles River), show relatively limited spread, while others, like nox and indus, display more variability across neighborhoods. These differences in scale and distribution suggest that transformation or normalization may be beneficial in the modeling stage. A correlation heatmap was constructed to examine multicollinearity and the relationships between predictors. Strong positive correlations were observed between nox, age, tax,

indus, and rad, suggesting that these variables may capture related structural or geographic aspects of the neighborhoods. Several variables also show moderate to strong correlation with the target variable — in particular, nox, age, rad, and tax were positively correlated with higher crime risk, while dis and rm were negatively correlated.

Below is series of boxplots for all numeric parameters where target is our dependent variable.

```
plot_boxplot(df_training, by = "target", title="Boxplots of Target vs Param")
```



```
df_training_hi_crime <- df_training |>
  filter(target == 1) |>
  subset(select = -c(chas, rad, target))

df_training_lo_crime <- df_training |>
  filter(target == 0) |>
  subset(select = -c(chas, rad, target))

cat("\nIQR for High Crime Neighborhoods\n")
```

```
##
## IQR for High Crime Neighborhoods
```

```
summary(df_training_hi_crime)
```

```
##          zn          indus          nox          rm
```

```
## Min. : 0.000 Min. : 3.97 Min. :0.4310 Min. :3.863
## 1st Qu.: 0.000 1st Qu.: 9.90 1st Qu.:0.5440 1st Qu.:5.727
## Median : 0.000 Median :18.10 Median :0.6240 Median :6.130
## Mean : 1.328 Mean :15.31 Mean :0.6404 Mean :6.181
## 3rd Qu.: 0.000 3rd Qu.:18.10 3rd Qu.:0.7000 3rd Qu.:6.484
## Max. :22.000 Max. :21.89 Max. :0.8710 Max. :8.780
## age dis tax ptratio
## Min. : 8.4 Min. :1.130 Min. :223.0 Min. :13.00
## 1st Qu.: 82.5 1st Qu.:1.728 1st Qu.:307.0 1st Qu.:18.40
## Median : 92.6 Median :2.125 Median :666.0 Median :20.20
## Mean : 86.5 Mean :2.471 Mean :513.8 Mean :18.96
## 3rd Qu.: 98.1 3rd Qu.:3.033 3rd Qu.:666.0 3rd Qu.:20.20
## Max. :100.0 Max. :8.907 Max. :666.0 Max. :21.20
## lstat medv
## Min. : 1.73 Min. : 5.00
## 1st Qu.:11.10 1st Qu.:13.50
## Median :15.39 Median :17.80
## Mean :16.02 Mean :20.05
## 3rd Qu.:20.34 3rd Qu.:23.00
## Max. :37.97 Max. :50.00
```

```
cat("\nIQR for Low Crime Neighborhoods\n")
```

```
##
## IQR for Low Crime Neighborhoods
```

```
summary(df_training_lo_crime)
```

```
## zn indus nox rm
## Min. : 0.00 Min. : 0.460 Min. :0.3890 Min. :5.093
## 1st Qu.: 0.00 1st Qu.: 3.370 1st Qu.:0.4290 1st Qu.:5.985
## Median : 0.00 Median : 5.640 Median :0.4490 Median :6.315
## Mean : 21.48 Mean : 7.039 Mean :0.4711 Mean :6.396
## 3rd Qu.: 34.00 3rd Qu.: 9.690 3rd Qu.:0.5150 3rd Qu.:6.727
## Max. :100.00 Max. :27.740 Max. :0.6240 Max. :8.069
## age dis tax ptratio
## Min. : 2.90 Min. : 1.669 Min. :187.0 Min. :12.60
## 1st Qu.: 31.30 1st Qu.: 3.360 1st Qu.:247.0 1st Qu.:16.60
## Median : 47.40 Median : 5.118 Median :293.0 Median :17.90
## Mean : 50.84 Mean : 5.076 Mean :308.8 Mean :17.86
## 3rd Qu.: 71.90 3rd Qu.: 6.458 3rd Qu.:348.0 3rd Qu.:19.10
## Max. :100.00 Max. :12.127 Max. :711.0 Max. :22.00
## lstat medv
## Min. : 1.98 Min. : 7.00
## 1st Qu.: 5.70 1st Qu.:20.40
## Median : 8.43 Median :23.10
## Mean : 9.36 Mean :25.04
## 3rd Qu.:12.27 3rd Qu.:28.60
## Max. :30.81 Max. :50.00
```

The boxplots show the distribution numerical parameters grouped by the dependent variable **target**. The plots are useful for getting a sense as to which parameters may be good predictors based on how different

the parameter's IQRs are. Conversely, similar IQRs may provide insight into which may not add much information to our model. Based on these box plots, we see that the IQR for `rm` are very similar where `target` is 0 and 1 and should be flagged for potential removal of our plot. `ptratio` and `medv` have some overlap. All other variables appear somewhat

Further, we see that param `zn` has a median value around zero, suggesting that few neighborhoods have residential areas zoned for large plots as shown below. We should also consider omitting this variable from our model down the line

```
count_zeros <- sum(df_training_hi_crime$zn == 0)
cat("\nAbove Median Crime Rate Neighborhoods have ", count_zeros, " rows with a value of 0 for param zn\n",
(count_zeros / nrow(df_training_hi_crime)), "%\n")
```

```
##
## Above Median Crime Rate Neighborhoods have 214 rows with a value of 0 for param zn out of 229 obs
```

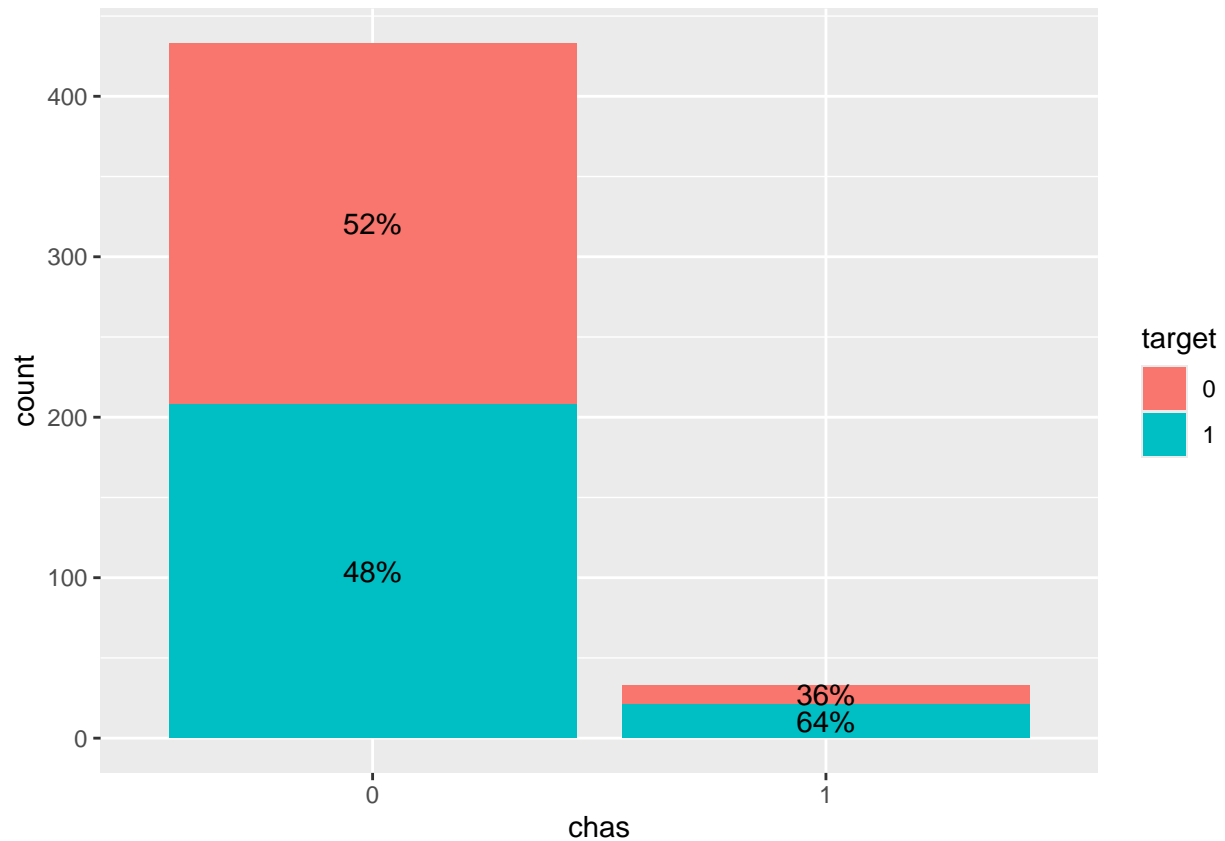
```
count_zeros <- sum(df_training_lo_crime$zn == 0)
cat("\nBelow Median Crime Rate Neighborhoods have ", count_zeros, " rows with a value of 0 for param zn\n",
(count_zeros / nrow(df_training_lo_crime)), "%\n")
```

```
##
## Below Median Crime Rate Neighborhoods have 125 rows with a value of 0 for param zn out of 237 obs
```

Categorical Variables

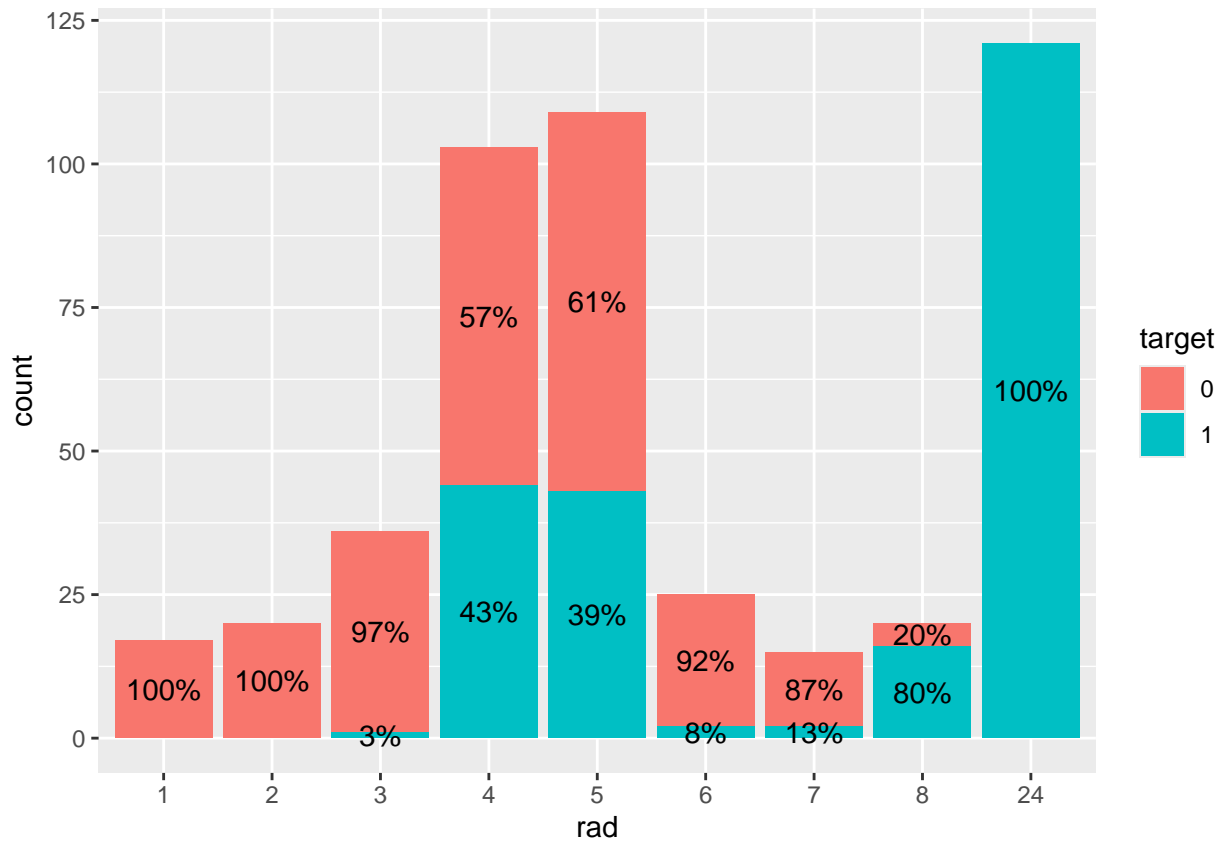
For our categorical variables, we can use barglaphs to get a sense of the parameter's impact on `target`.

```
df_training |>
  group_by(
    target, chas
  ) |>
  dplyr::summarise(
    count = n()
  ) |>
  ungroup() |>
  group_by(chas) |>
  mutate(
    percent = 100 * count / sum(count),
    label = paste0(round(percent), "%")
  ) |>
  ggplot() +
  aes(x = chas, y = count, label = label, fill=target) +
  geom_col() +
  geom_text(position = position_stack(0.5))
```



The bargraph for `chas` shows fairly equal values for 0 and 1 accross the `chas` values. This suggests that the variable will may have low impact on our model and we should consider removing it.

```
### rad
df_training |>
  group_by(
    target, rad
  ) |>
  dplyr::summarise(
    count = n()
  ) |>
  ungroup() |>
  group_by(rad) |>
  mutate(
    percent = 100 * count / sum(count),
    label = paste0(round(percent), "%")
  ) |>
  ggplot() +
  aes(x = rad, y = count, label = label, fill=target) +
  geom_col() +
  geom_text(position = position_stack(0.5))
```

The bargraphs for **rad** are somewhat more revealing. They suggest a strong relationship between low **rad** index values of 1-3 and below median crime rate, while an index value of 24 (the highest **rad** index) has a strong relationship with above median crime rate.

Pairs Using the pair function, we can print scatterplots comparing each of the variables to the others.

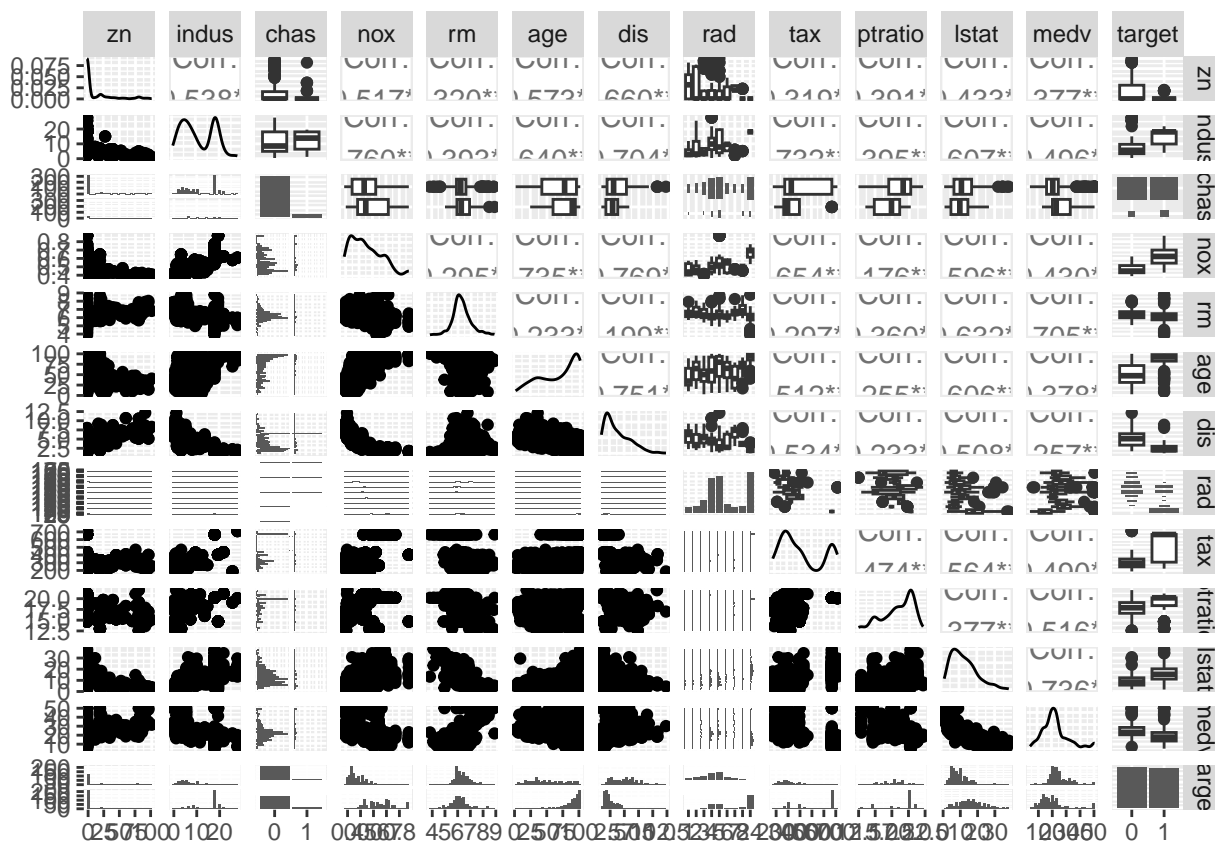
```
png("scatterplot_matrix.png", width = 800, height = 800)
```

```
pairs(df_training, main="")
dev.off()
```

```
## pdf
## 2
```

GGpairs plots take this a step further and show normal distribution and boxplots to get a fuller sense of how the data parameters relate to one another.

```
ggpairs(df_training)
```



Checking Binary Logistic Regression Assumptions

Before interpreting results from a binary logistic regression, we must verify three key assumptions:

- Independence of Observations
- Linearity of the Logit
- No Multicollinearity

Spearman's correlation measures monotonic relationships between variables, making it suitable when we have a mix of ordinal and continuous predictors. Unlike the Pearson test, it does not assume normality. Additionally, it is more robust against outliers than Pearson.

Checking Independence Assumption The independence assumption in binary logistic regression states that each observation (row) in the dataset should be independent of the others. This means:

No duplicated data points (e.g., same neighborhood appearing multiple times). No clustered observations (e.g., observations grouped by region, time, or other factors). No strong correlations between residuals of observations, meaning observations do not systematically affect each other.

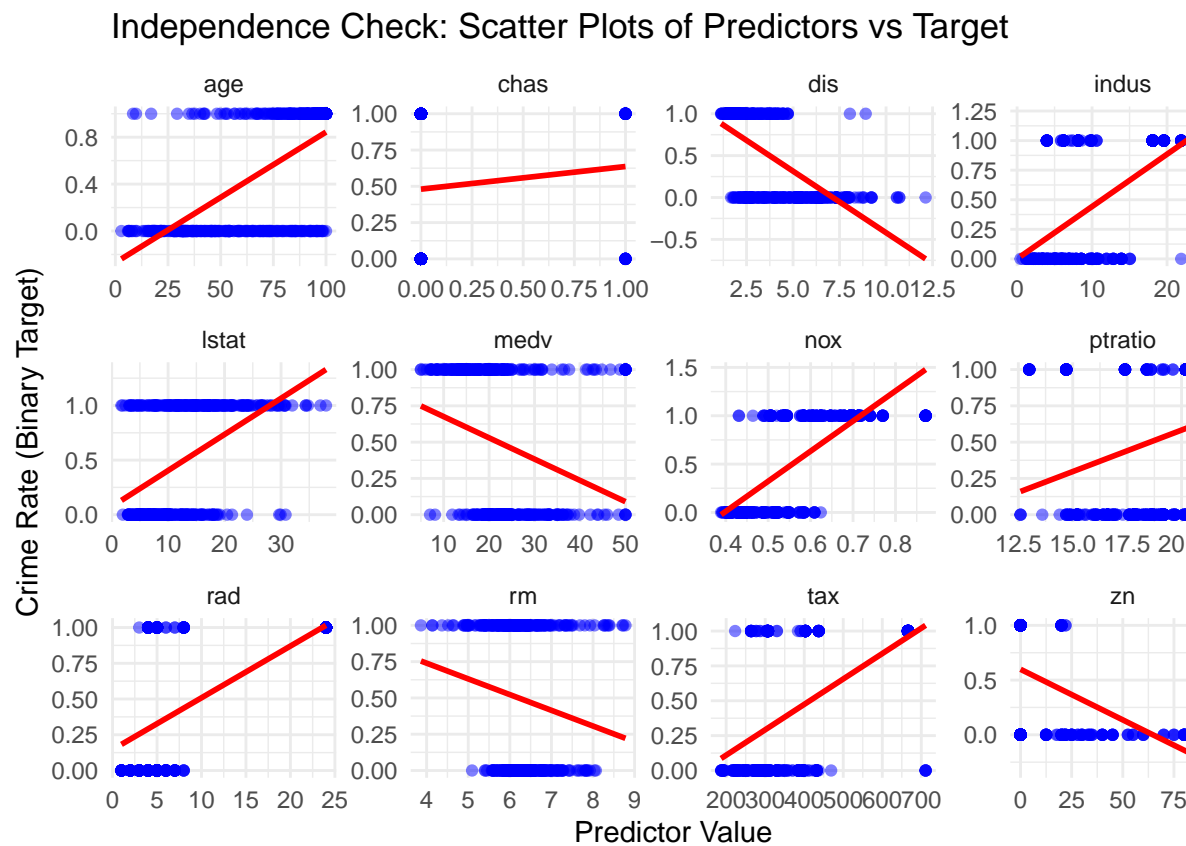
```
# Exclude categorical columns before pivoting
crime_long <- crime_training_df %>%
  dplyr::select(where(is.numeric)) %>% # Keep only numeric columns
```

```

pivot_longer(cols = -target, names_to = "Variable", values_to = "Value")

# Create scatter plots for each predictor vs target with a fitted line
ggplot(crime_long, aes(x = Value, y = target)) +
  geom_point(alpha = 0.5, color = "blue") + # Scatter plot
  geom_smooth(method = "lm", color = "red", se = FALSE) + # Linear fit
  facet_wrap(~Variable, scales = "free") + # Multiple plots for each variable
  theme_minimal() +
  labs(title = "Independence Check: Scatter Plots of Predictors vs Target",
       x = "Predictor Value",
       y = "Crime Rate (Binary Target)")

```



Visual Inspection

```

# Initialize an empty data frame to store results
independence_results <- data.frame(Variable = character(),
                                   Correlation = numeric(),
                                   P_Value = numeric(),
                                   stringsAsFactors = FALSE)

# Loop through each predictor variable (excluding the target)
for (var in colnames(crime_training_df)[colnames(crime_training_df) != "target"]) {

  # Ensure the variable is numeric before computing Spearman correlation
  if (is.numeric(crime_training_df[[var]])) {

    # Perform Spearman correlation test

```

```

test_result <- cor.test(crime_training_df[[var]], crime_training_df$target, method = "spearman")

# Store results in a data frame
independence_results <- rbind(independence_results,
                              data.frame(Variable = var,
                                          Correlation = test_result$estimate,
                                          P_Value = test_result$p.value))
}
}

# View results in tabular format
print(independence_results)

```

```

##      Variable Correlation      P_Value
## rho      zn -0.47299691 2.365331e-27
## rho1     indus 0.61915270 1.159150e-50
## rho2     chas 0.08004187 8.434811e-02
## rho3     nox 0.75471235 5.629874e-87
## rho4     rm -0.17719123 1.204401e-04
## rho5     age 0.64569520 2.544840e-56
## rho6     dis -0.65908410 2.147507e-59
## rho7     rad 0.57842556 5.796797e-43
## rho8     tax 0.59604354 3.693938e-46
## rho9     ptratio 0.35733899 1.756795e-15
## rho10    lstat 0.47946598 3.674288e-28
## rho11    medv -0.40154322 1.755130e-19

```

Checking for Independence Using Spearman’s Correlation In Spearman’s method, we check the correlation between:

Each independent variable and the dependent variable (target) If correlation is too low ($|p| < 0.1$) and p-value > 0.05 , the variable might not be useful in predicting the target.

Variables to Consider Removing:

chas ($p = 0.0800$, $p = 0.0843$) → No meaningful correlation with crime. Possibly rm ($p = -0.1772$, $p = 0.00012$) → Weak correlation but could check its importance in the model.

Variables to Keep (for now, but monitor multicollinearity):

nox ($p = 0.7547$) age ($p = 0.6457$) dis ($p = -0.6591$) indus ($p = 0.6192$) rad, tax, ptratio (moderate correlation)
Transform / Create Interaction Terms:

Log transform: dis (since it has a strong negative correlation) Categorize: age into “Young”, “Middle-aged”, “Old” Interaction term: tax * rad (both impact crime)

```

crime_training_df <- crime_training_df %>%
  mutate(across(where(is.character), as.numeric))

crime_evaluation_df <- crime_evaluation_df %>%
  mutate(across(where(is.character), as.numeric))

# Verify again
str(crime_training_df)

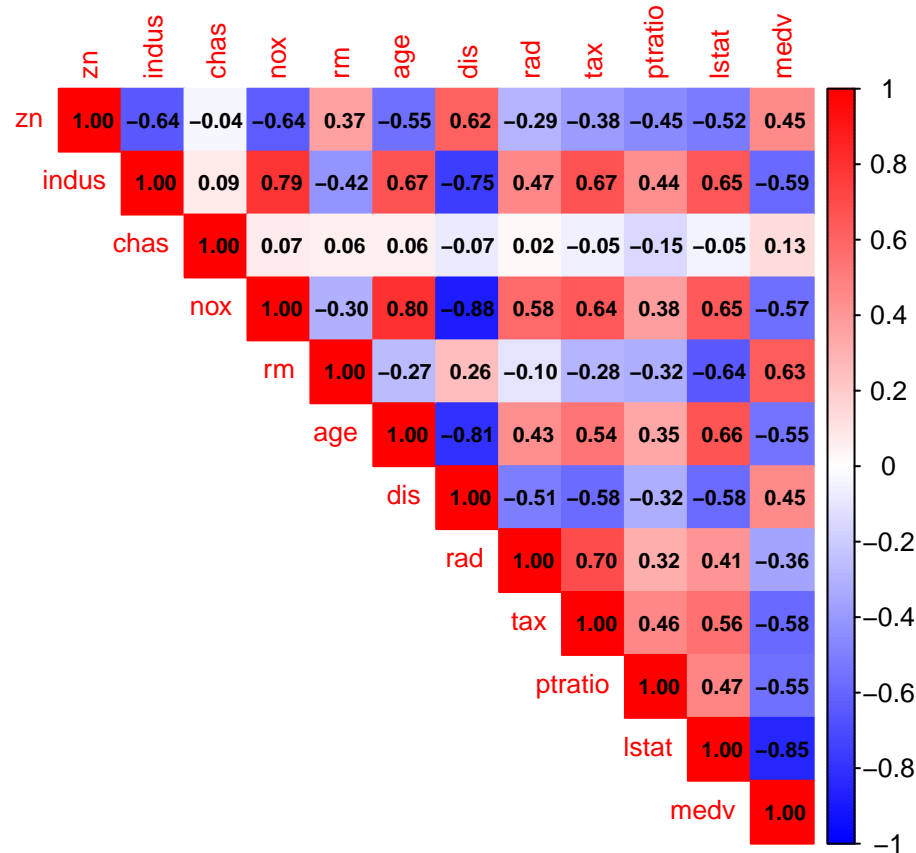
```

```
## tibble [466 x 13] (S3: tbl_df/tbl/data.frame)
## $ zn      : num [1:466] 0 0 0 30 0 0 0 0 0 80 ...
## $ indus   : num [1:466] 19.58 19.58 18.1 4.93 2.46 ...
## $ chas    : num [1:466] 0 1 0 0 0 0 0 0 0 0 ...
## $ nox     : num [1:466] 0.605 0.871 0.74 0.428 0.488 0.52 0.693 0.693 0.515 0.392 ...
## $ rm      : num [1:466] 7.93 5.4 6.49 6.39 7.16 ...
## $ age     : num [1:466] 96.2 100 100 7.8 92.2 71.3 100 100 38.1 19.1 ...
## $ dis     : num [1:466] 2.05 1.32 1.98 7.04 2.7 ...
## $ rad     : num [1:466] 5 5 24 6 3 5 24 24 5 1 ...
## $ tax     : num [1:466] 403 403 666 300 193 384 666 666 224 315 ...
## $ ptratio : num [1:466] 14.7 14.7 20.2 16.6 17.8 20.9 20.2 20.2 20.2 16.4 ...
## $ lstat   : num [1:466] 3.7 26.82 18.85 5.19 4.82 ...
## $ medv    : num [1:466] 50 13.4 15.4 23.7 37.9 26.5 5 7 22.2 20.9 ...
## $ target  : num [1:466] 1 1 1 0 0 0 1 1 0 0 ...
```

Visualize Correlation Matrix

```
# Compute Spearman correlation matrix
numeric_data <- crime_training_df %>% dplyr::select(where(is.numeric), -target)
spearman_cor <- cor(numeric_data, method = "spearman", use = "pairwise.complete.obs")

# Plot the Spearman correlation matrix with labels
corrplot::corrplot(spearman_cor,
  method = "color",
  type = "upper",
  tl.cex = 0.8,
  addCoef.col = "black",
  number.cex = 0.7,
  col = colorRampPalette(c("blue", "white", "red"))(200)
)
```



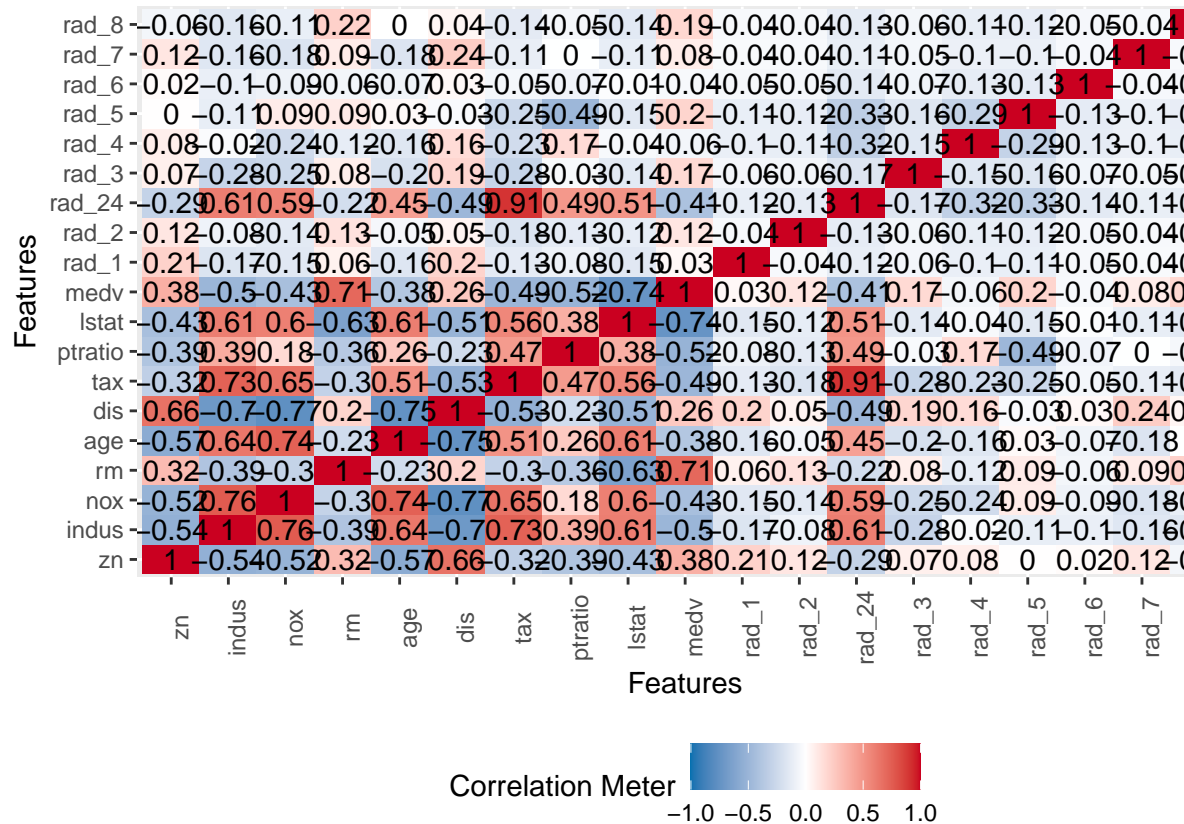
Checking for Multicollinearity (High Correlation Between Predictors) A general rule of thumb is that if $|p| > 0.7$, it indicates strong correlation between variables, which can lead to multicollinearity in the logistic regression model.

From the matrix: indus & nox ($p = 0.79$) → Strong positive correlation, meaning they provide redundant information. tax & rad ($p = 0.70$) → These variables are highly correlated, indicating one may be removed. lstat & medv ($p = -0.85$) → Very strong negative correlation; keeping both might be problematic. log_medv & medv ($p = 1.00$) → Perfect correlation (since log transformation was applied), meaning one must be removed to prevent redundancy.

Final Takeaways:

High correlation between predictors ($|p| > 0.7$) indicates potential multicollinearity. Consider removing indus, rad, or medv to improve model stability. No evidence of entire rows/columns being highly correlated, suggesting no major independence violations. Use VIF for confirmation and decide on variable selection accordingly.

```
df_training |>
  subset(select=-c(target, chas)) |>
  plot_correlation(type = "all")
```



Visual Inspection

Variance Inflation Factor (VIF) Variance Inflation Factor (VIF) helps quantify multicollinearity by measuring how much the variance of regression coefficients is inflated due to correlation among predictors. A VIF > 5 (or more conservatively, VIF > 10) suggests severe multicollinearity.

```
car::vif(glm(target ~ nox + age + rad + tax + dis + zn + medv, family = binomial, data = crime_training
```

```
##      nox      age      rad      tax      dis      zn      medv
## 2.859311 1.548306 1.421341 1.766596 3.576500 1.667181 1.784929
```

Conclusion There is NO severe multicollinearity (all VIF values are below 5). No immediate need to drop variables based on VIF. The variable dis (VIF = 3.58) shows moderate correlation with other predictors, but it's not problematic.

Keep all predictors in the model. If we suspect redundancy, check pairwise correlations again or test removing dis to see if model performance improves.

Checking the Assumption of Linearity of the Logit for Binary Logistic Regression In logistic regression, we assume that each continuous predictor has a linear relationship with the log-odds (logit) of the target variable. If this assumption is violated, the model may be misleading or inaccurate.

```
# Load necessary libraries

# Fit the logistic regression model
model <- glm(target ~ nox + age + rad + tax + dis + zn + medv,
```

```

family = binomial, data = crime_training_df)

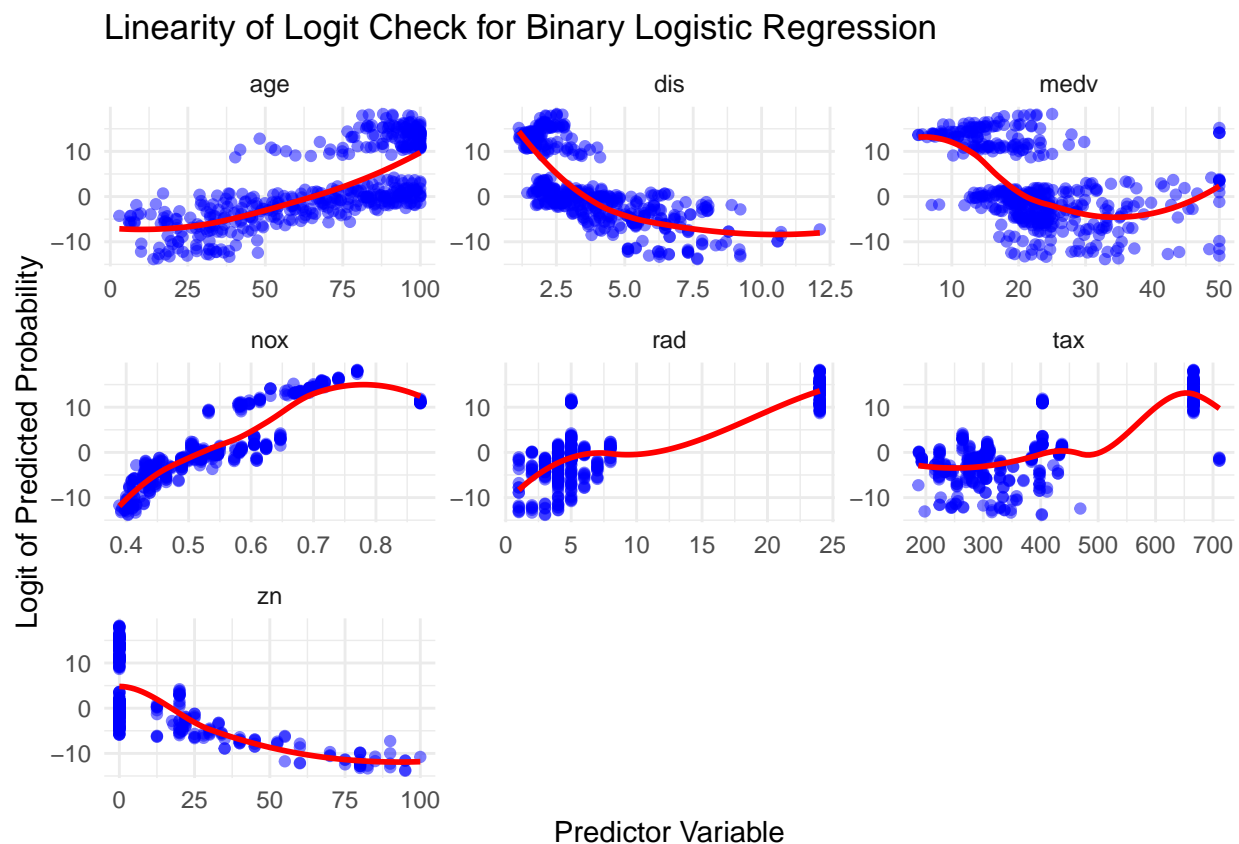
# Compute predicted probabilities
crime_training_df$predicted_prob <- predict(model, type = "response")

# Compute logit (log-odds) transformation
crime_training_df$logit <- log(crime_training_df$predicted_prob /
                              (1 - crime_training_df$predicted_prob))

# Reshape data for faceting
plot_data <- crime_training_df %>%
  dplyr::select(logit, nox, age, rad, tax, dis, zn, medv) %>%
  pivot_longer(cols = -logit, names_to = "Predictor", values_to = "Value")

# Create faceted scatter plot
ggplot(plot_data, aes(x = Value, y = logit)) +
  geom_point(alpha = 0.5, color = "blue") +
  geom_smooth(method = "loess", color = "red", se = FALSE) +
  facet_wrap(~ Predictor, scales = "free") +
  theme_minimal() +
  labs(title = "Linearity of Logit Check for Binary Logistic Regression",
       x = "Predictor Variable",
       y = "Logit of Predicted Probability")

```



Interpretation of the Linearity of Logit Check for Binary Logistic Regression This faceted scatter plot assesses the assumption of linearity of the logit for binary logistic regression. The blue dots represent the relationship between each predictor variable and the logit of the predicted probability, while the red LOESS (Locally Estimated Scatterplot Smoothing) curve helps visualize patterns.

Overall Conclusion Several predictors (e.g., dis, medv, tax, zn) violate the linearity of logit assumption.

DATA PREPARATION

This step involves cleaning and transforming data to improve model performance.

Fixing missing values

Luckily, there are no missing values in the training set.

Transforming data by bucketing and combining variables

The variable `rad` contains an ordinal factor that represents an index of accessibility to radial highways with values ranging from 1-24. A count of the `rad` values reveals that the `rad` column contains only values 1-8 and 24. This data set does not include any rows with a `rad` value of 9-23.

Since this column contains values for an index value where 1 is assigned to neighborhoods with the poorest accessibility to a highway and 24 is assigned to neighborhoods with the most accessibility, we can simplify our variables by binning our `rad` values. Here we are using quantiles to bin the values into three buckets of nearly equal sizes for low, moderate and high accessibility. This method ensures a more balanced distribution of rows across the bins over using equal sized bins (1-8, 9-16, 17-24). This is especially useful when the data is not uniformly distributed across the range such as in our case where we do not have any `rad` values of 1-23.

```
rad_counts
```

```
##
##  1  2  3  4  5  6  7  8 24
## 17 20 36 103 109 25 15 20 121
```

```
quantile_breaks <- quantile(as.numeric(df_training$rad), probs = c(0, 1/3, 2/3, 1))

df_training$radq <- cut(as.numeric(df_training$rad),
  breaks = quantile_breaks,
  labels = c('_low', '_mid', '_hi'),
  include.lowest = TRUE,
  right = TRUE)

table(df_training$radq)
```

```
##
## _low _mid _hi
## 176 149 141
```

While the `glm` function should automatically perform one-hot encoding to factors, we should consider one-hot encoding on the `rad_quantile` parameter to perform other operations, such as calculating correlation using the spearman test.

We will drop one of the one-hot encoded params as the presence of this additional param will result in correlation issues down the line. `radq_mid` was selected, as it seemed to have the most mixed results in our plots above.

```
# one-hot encode rad values
rad_one_hot <- model.matrix(~ radq - 1, data = df_training)

# combine new columns
df_training_one_hot <- cbind(df_training[, !names(df_training) %in% "rad"], rad_one_hot) |>
  subset(select=-c(radq, radq_mid))

glimpse(df_training_one_hot)
```

```
## Rows: 466
## Columns: 14
## $ zn      <dbl> 0, 0, 0, 30, 0, 0, 0, 0, 0, 80, 22, 0, 0, 22, 0, 0, 100, 20, ~
## $ indus   <dbl> 19.58, 19.58, 18.10, 4.93, 2.46, 8.56, 18.10, 18.10, 5.19, 3.~
## $ chas    <fct> 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ nox     <dbl> 0.605, 0.871, 0.740, 0.428, 0.488, 0.520, 0.693, 0.693, 0.515~
## $ rm      <dbl> 7.929, 5.403, 6.485, 6.393, 7.155, 6.781, 5.453, 4.519, 6.316~
## $ age     <dbl> 96.2, 100.0, 100.0, 7.8, 92.2, 71.3, 100.0, 100.0, 38.1, 19.1~
## $ dis     <dbl> 2.0459, 1.3216, 1.9784, 7.0355, 2.7006, 2.8561, 1.4896, 1.658~
## $ tax     <dbl> 403, 403, 666, 300, 193, 384, 666, 666, 224, 315, 330, 398, 6~
## $ ptratio <dbl> 14.7, 14.7, 20.2, 16.6, 17.8, 20.9, 20.2, 20.2, 20.2, 16.4, 1~
## $ lstat   <dbl> 3.70, 26.82, 18.85, 5.19, 4.82, 7.67, 30.59, 36.98, 5.68, 9.2~
## $ medv    <dbl> 50.0, 13.4, 15.4, 23.7, 37.9, 26.5, 5.0, 7.0, 22.2, 20.9, 24.~
## $ target  <fct> 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0~
## $ radq_low <dbl> 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1~
## $ radq_hi  <dbl> 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0~
```

We will use the one-hot encoded dataframe to diagnose a preliminary model with all of the predictors.

```
model_full <- glm(target ~., binomial(link = "logit"), data=df_training_one_hot)
summary(model_full)
```

```
##
## Call:
## glm(formula = target ~ ., family = binomial(link = "logit"),
##      data = df_training_one_hot)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -3.736e+01  6.617e+00  -5.646 1.65e-08 ***
## zn          -8.026e-02  4.105e-02  -1.955 0.050551 .
## indus       -1.724e-01  4.986e-02  -3.457 0.000547 ***
## chas1        1.179e+00  8.097e-01   1.456 0.145311
## nox          5.946e+01  9.038e+00   6.579 4.74e-11 ***
## rm          -9.564e-01  6.992e-01  -1.368 0.171345
## age          2.030e-02  1.322e-02   1.536 0.124585
```

```
## dis      8.131e-01  2.456e-01  3.310 0.000931 ***
## tax      9.619e-04  2.291e-03  0.420 0.674583
## ptratio  1.190e-01  1.333e-01  0.893 0.372038
## lstat    5.447e-02  5.231e-02  1.041 0.297705
## medv     1.760e-01  5.907e-02  2.980 0.002887 **
## radq_low 1.563e+00  5.254e-01  2.975 0.002931 **
## radq_hi  5.118e+00  9.416e-01  5.436 5.46e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 180.97  on 452  degrees of freedom
## AIC: 208.97
##
## Number of Fisher Scoring iterations: 8
```

Reviewing the summary statistics for full model indicates that the variable `indus`, `nox`, `dis`, and `radq_hi` has very strong statistically signification. Two additional variables, `medv`, `dis` and `radq_mid`, have high statistical significance while `zn` has weak statistical significance. `chas1`, `rm`, `age`, `tax`, `ptratio` and `lstat` have weak statistical significance values.

Note: had we one-hot encoded all of the values for `rad` instead of binning them first, all `rad` params would have very weak statistical significance, as their p-values are nearly 1.0.

Multicollinearity

To test if correlation exists between the dependent and independent variables, we used a Pearson's Correlation test. The function below loops through each of our columns and prints out the correlation of the dependent variable `target` with each of the predictors. For predictors where Pearson's Correlation coefficient is close to zero, we can determine that collinearity does not exist.

```
# is above .7 would be too highly correlated
cor_results <- data.frame(name = character(0), value = numeric(0))
for (param in colnames(df_training_one_hot)) {
  cat("\nPearson Test score for", param, ":\n")
  x <- as.numeric(df_training_one_hot$target)
  y <- as.numeric(df_training_one_hot[[param]])
  pears <- cor.test(x, y, method = "pearson")
  print(pears)
  # calc pearson cor value only
  cor_object <- data.frame(name = param, value = cor(x, y))
  assign("cor_results", rbind(cor_results, cor_object), envir = .GlobalEnv)
}
```

```
##
## Pearson Test score for zn :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = -10.309, df = 464, p-value < 2.2e-16
```

```

## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.5028019 -0.3547564
## sample estimates:
##      cor
## -0.4316818
##
##
## Pearson Test score for indus :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 16.361, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.5438976 0.6594549
## sample estimates:
##      cor
## 0.6048507
##
##
## Pearson Test score for chas :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 1.7297, df = 464, p-value = 0.08435
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.01087336 0.16964461
## sample estimates:
##      cor
## 0.08004187
##
##
## Pearson Test score for nox :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 22.748, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.6801291 0.7663936
## sample estimates:
##      cor
## 0.7261062
##
##
## Pearson Test score for rm :
##
## Pearson's product-moment correlation
##

```

```

## data:  x and y
## t = -3.325, df = 464, p-value = 0.0009542
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.24006288 -0.06258443
## sample estimates:
##      cor
## -0.1525533
##
##
## Pearson Test score for age :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 17.479, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.5720099 0.6819122
## sample estimates:
##      cor
## 0.6301062
##
##
## Pearson Test score for dis :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = -16.963, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.6717579 -0.5592666
## sample estimates:
##      cor
## -0.6186731
##
##
## Pearson Test score for tax :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 16.631, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.5508558 0.6650327
## sample estimates:
##      cor
## 0.6111133
##
##
## Pearson Test score for ptratio :
##

```

```

## Pearson's product-moment correlation
##
## data:  x and y
## t = 5.5819, df = 464, p-value = 4.053e-08
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.1637438 0.3340729
## sample estimates:
##      cor
## 0.2508489
##
##
## Pearson Test score for lstat :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 11.443, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.3951287 0.5370764
## sample estimates:
##      cor
## 0.469127
##
##
## Pearson Test score for medv :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = -6.0536, df = 464, p-value = 2.925e-09
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.3527185 -0.1842424
## sample estimates:
##      cor
## -0.2705507
##
##
## Pearson Test score for target :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = Inf, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  1 1
## sample estimates:
## cor
##  1
##
##

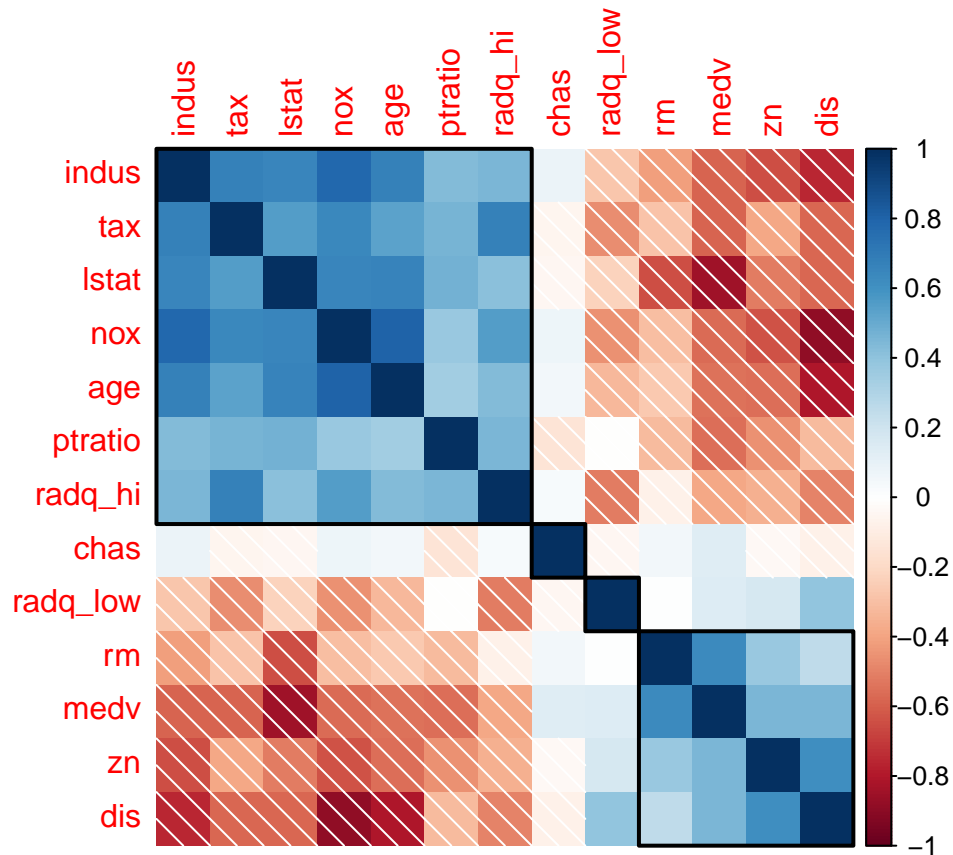
```

```
## Pearson Test score for radq_low :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = -8.5077, df = 464, p-value = 2.466e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.4433864 -0.2860542
## sample estimates:
##      cor
## -0.3673453
##
##
## Pearson Test score for radq_hi :
##
## Pearson's product-moment correlation
##
## data:  x and y
## t = 17.599, df = 464, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.5749042 0.6842126
## sample estimates:
##      cor
## 0.6326995
```

```
print(cor_results)
```

```
##      name      value
## 1      zn -0.43168176
## 2    indus  0.60485074
## 3     chas  0.08004187
## 4      nox  0.72610622
## 5       rm -0.15255334
## 6      age  0.63010625
## 7      dis -0.61867312
## 8      tax  0.61111331
## 9  ptratio  0.25084892
## 10    lstat  0.46912702
## 11    medv -0.27055071
## 12   target  1.00000000
## 13 radq_low -0.36734528
## 14 radq_hi  0.63269952
```

Correlation Clusters Next we can visualize the correlations in clusters.



Using “hclust”, our corplot shows four distinct groups, each with strong correlation between the parameters within each group. This suggests that we may want to select specific parameters from within these groups or conduct principal component analysis on each of these groups.

```
car::vif(model_full) |> sort()
```

Variance Inflation Factor

```
##      chas radq_low      zn radq_hi      age      tax ptratio      lstat
## 1.208878 1.939100 1.979366 2.126242 2.146346 2.205786 2.314164 2.561351
##      indus      dis      rm      nox      medv
## 3.388295 4.138425 4.732821 5.329898 5.983892
```

A VIF test suggests that we should remove nox and medv.

Presence of Outliers

We can examine our diagnostic plots to find potential outliers and leverage. First we will examine the Cook’s Distance and Cook’s Distance vs Leverage plots. Cook’s Distance measures the influence of an observation on the fitted values of the model.

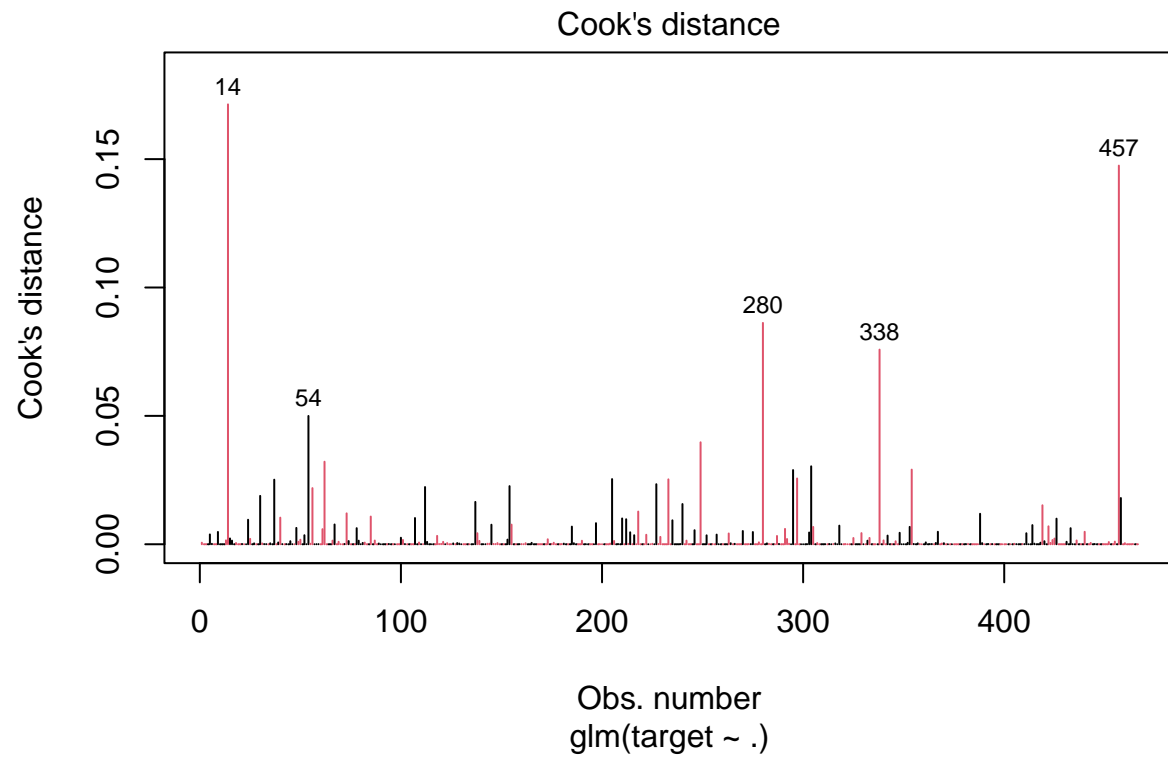
```
## Rows: 5
```

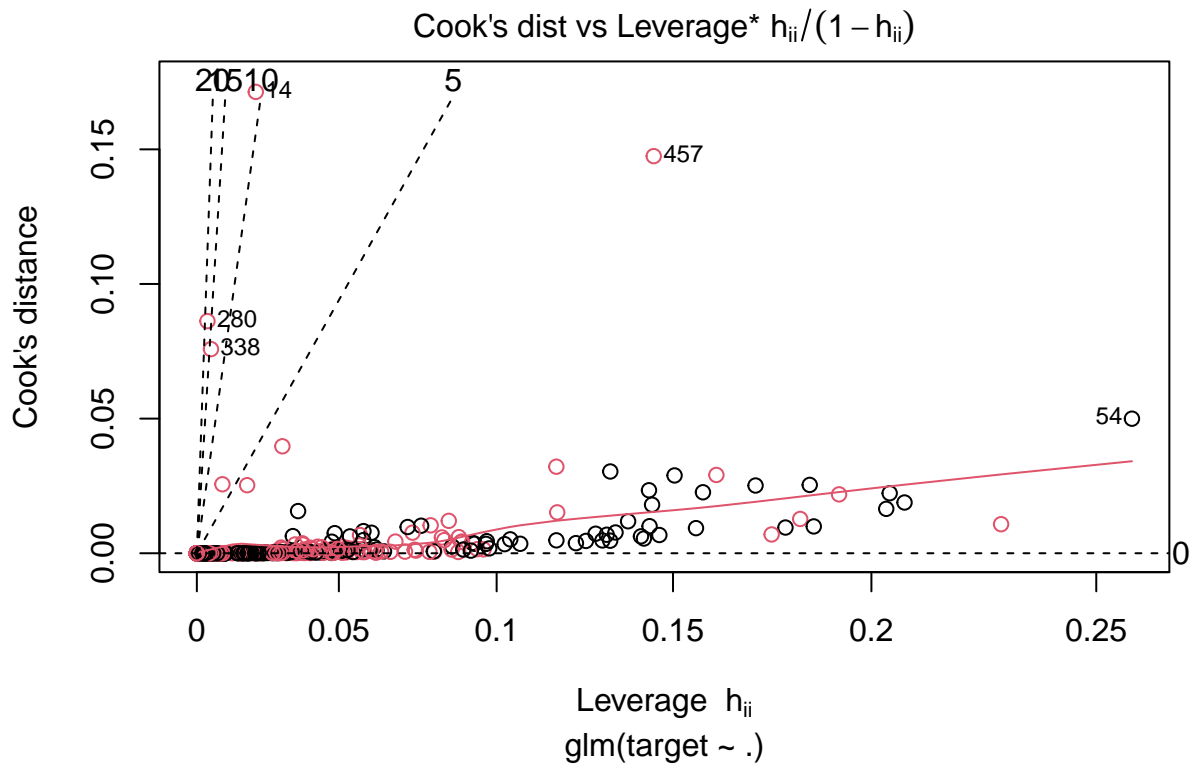


```
## Columns: 21
## $ index      <int> 14, 457, 280, 338, 54
## $ target     <chr> "Hi", "Hi", "Hi", "Hi", "Lo"
## $ zn         <dbl> 22, 0, 22, 20, 0
## $ indus      <dbl> 5.86, 10.59, 5.86, 6.96, 1.89
## $ chas       <fct> 0, 0, 0, 0, 0
## $ nox        <dbl> 0.431, 0.489, 0.431, 0.464, 0.518
## $ rm         <dbl> 8.259, 5.412, 6.108, 5.856, 6.540
## $ age        <dbl> 8.4, 9.8, 34.9, 42.1, 59.7
## $ dis        <dbl> 8.9067, 3.5875, 8.0555, 4.4290, 6.2669
## $ tax        <dbl> 330, 277, 330, 223, 422
## $ ptratio    <dbl> 19.1, 18.6, 19.1, 18.6, 15.9
## $ lstat      <dbl> 3.54, 29.55, 9.16, 13.00, 8.65
## $ medv       <dbl> 42.8, 23.7, 24.3, 21.1, 16.5
## $ radq_low   <dbl> 0, 1, 0, 1, 1
## $ radq_hi    <dbl> 0, 0, 0, 0, 0
## $ .fitted    <dbl> -4.6773398, -2.3444828, -5.7239448, -5.3060581, 0.4048362
## $ .resid     <dbl> 3.061568, 2.207286, 3.384437, 3.259143, -1.353450
## $ .hat       <dbl> 0.021373959, 0.144810967, 0.003911828, 0.005210891, 0.25736~
## $ .sigma     <dbl> 0.6164625, 0.6234026, 0.6129971, 0.6144815, 0.6291214
## $ .cooks     <dbl> 0.17134297, 0.14748407, 0.08620525, 0.07580775, 0.04996881
## $ .std.resid <dbl> 3.094821, 2.386863, 3.391076, 3.267668, -1.570564
```

The calculation above shows that points 14, 457, 280, 338, and 54 have the highest Cook's distance values (ordered from highest to lowest) and should be investigated as potential outliers.

```
par(mar = c(5, 4, 4, 2) + 0.1)
plot(model_full, which = c(4, 6), col=df_training_one_hot$target, id.n = 5)
```



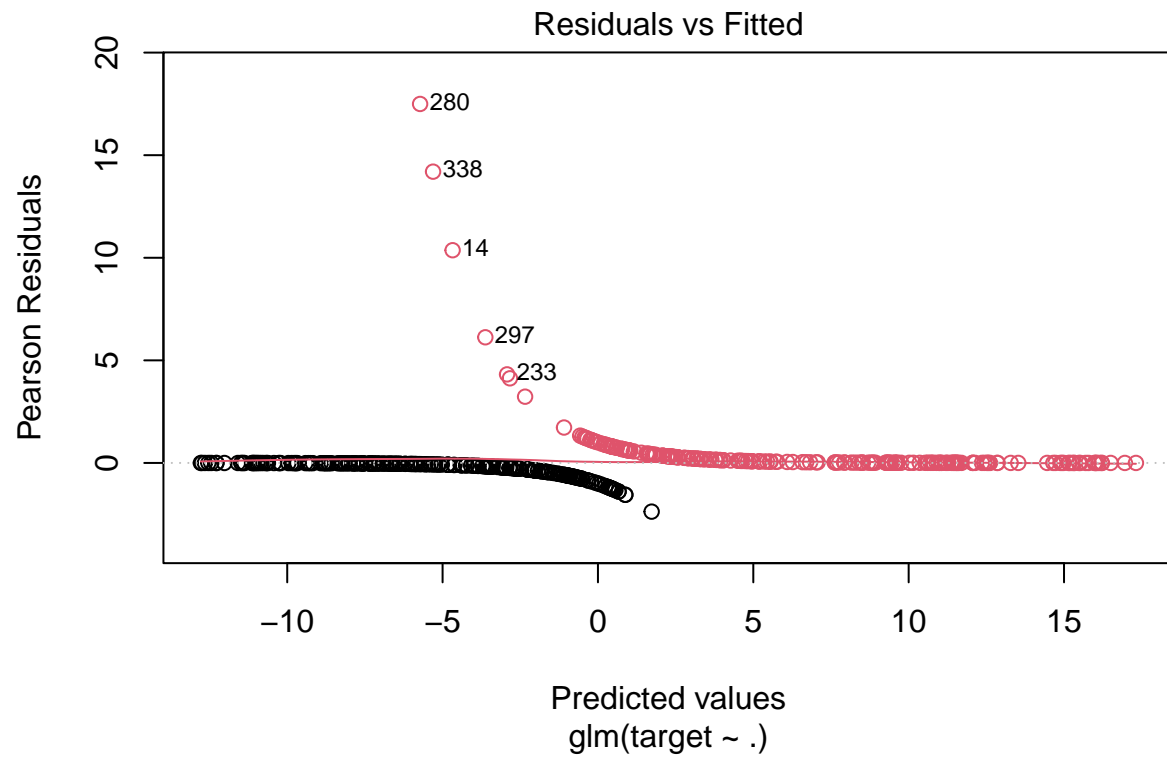


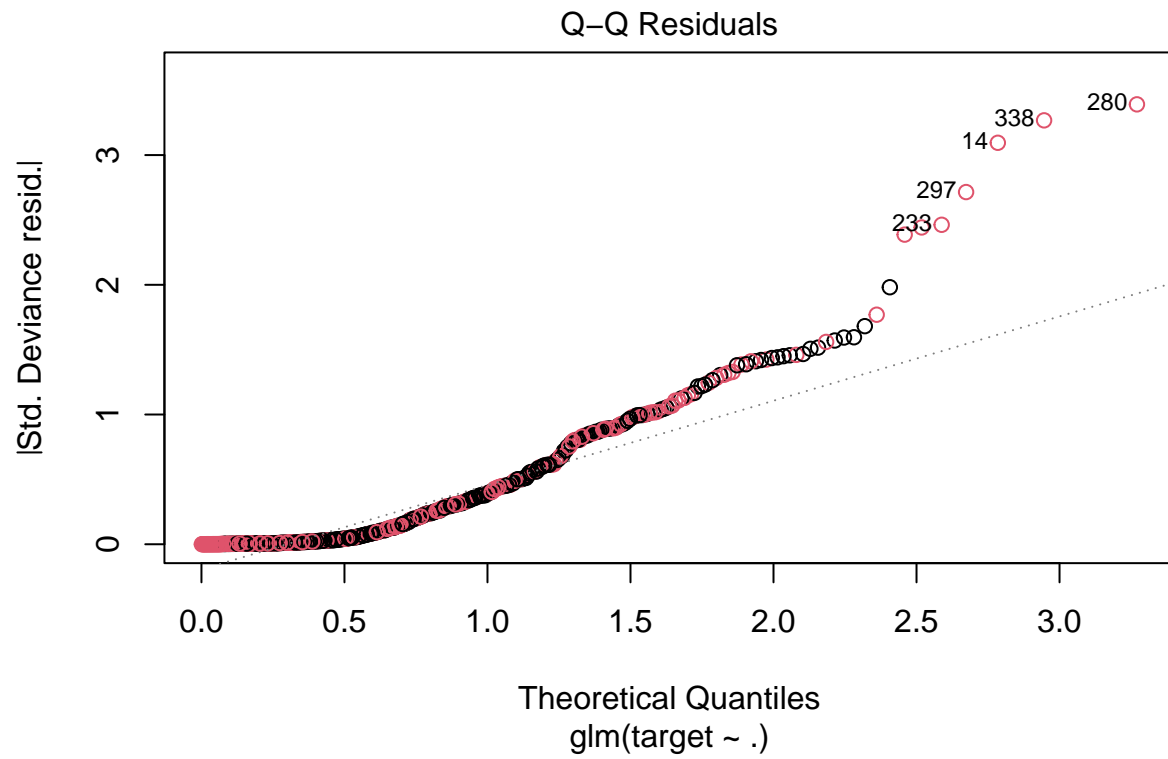
We see on the Cook's dist vs Leverage plot that points 280 and 338 may have very high leverage on our model, followed by point 14. Point 457 also stand out and should be investigated but appears to have less leverage.

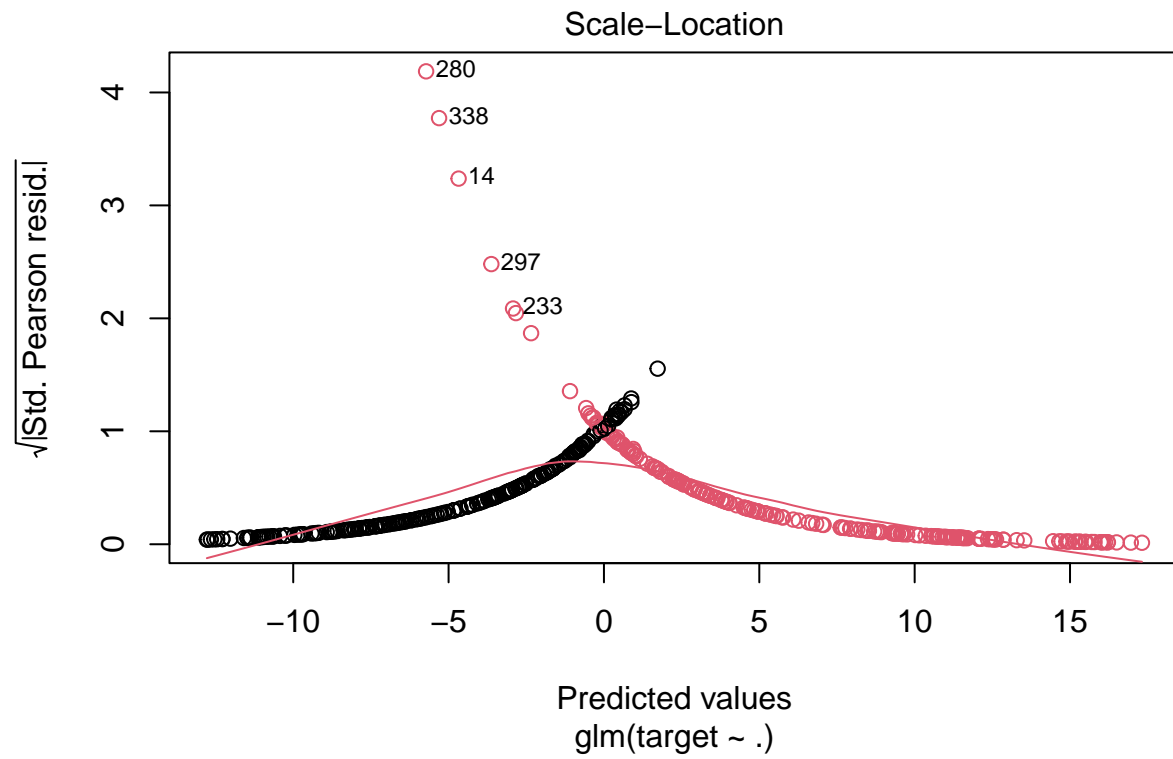
```
# print influential points using cooks-distance
cooksdist <- cooks.distance(model_full)
influential <- which(cooksdist > (4 / length(cooksdist)))
print(influential)
```

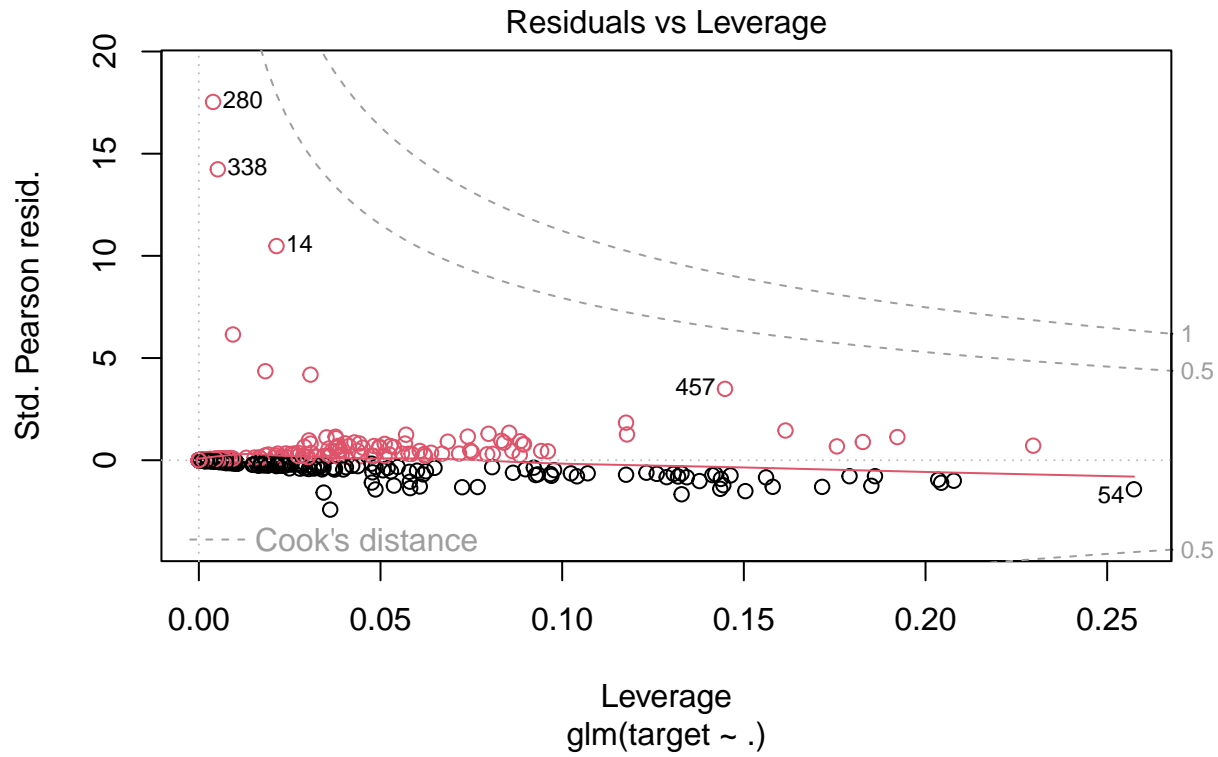
```
## 14 24 30 37 40 54 56 62 73 85 107 112 137 154 205 210 212 218 227 233
## 14 24 30 37 40 54 56 62 73 85 107 112 137 154 205 210 212 218 227 233
## 235 240 249 280 295 297 304 338 354 388 419 426 457 458
## 235 240 249 280 295 297 304 338 354 388 419 426 457 458
```

The formula above is used to identify influential points defined as points Cook's Distance value is greater than $4 / \text{length of cooksdist}$. This contains all three points (280, 338, and 14) as being influential.









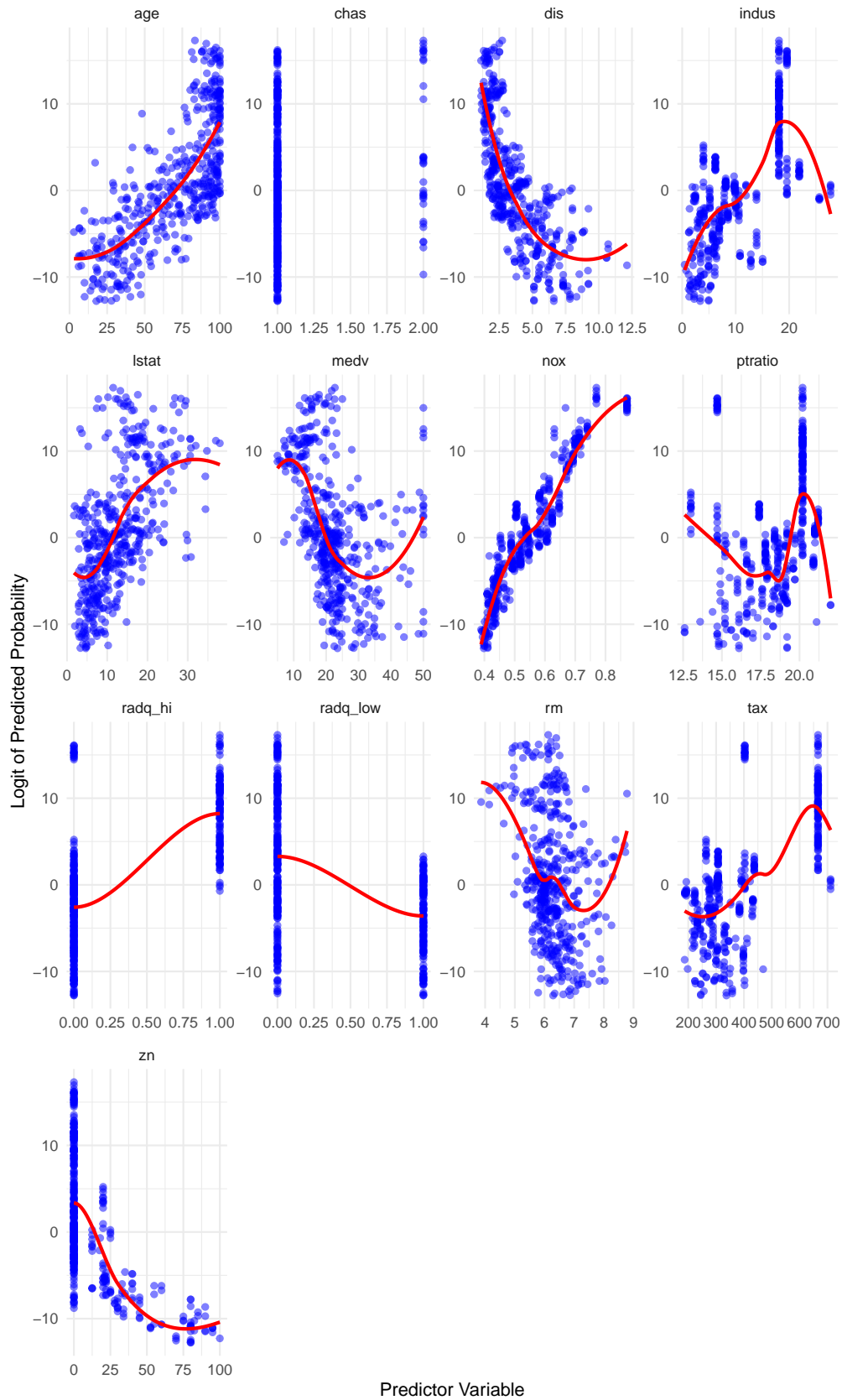
Our residual vs fitted, QQ, Scale-Location and Residual vs Leverage plots all confirm that points 280, 338, and 14 should be investigated and could be outliers with high influence. Points 457 appear to have less leverage and does not stand out in these plots.

Below is the output for the three points identified as potential outliers in our diagnostic plots. A quick review of the data doesn't reveal anything that stands out as being out of the ordinary.

Linearity

To check this condition, I created a scatterplot with a loess line to check that there is a linear relationship between the logit of the dependent variable and the independent variables.

Linearity of Logit Check for Binary Logistic Regression



Using mathematical transformations

To reduce the influence of outliers and better align the data with the assumptions of logistic regression, log-transformations were applied to tax, zn, dis, and lstat. This transformation helps normalize the data, reduce variance, and enhance model interpretability. A small constant was added to zn before the transformation to account for zero values.

```
df_training_1h_log <- df_training_one_hot |>
  mutate(
    log_tax = log(tax),
    log_dis = log(dis),
    log_zn = log(zn + 1),
    log_lstat = log(lstat),
  ) |>
  subset(select = -c(tax, dis, zn, lstat))

model_full_log <- glm(target ~ ., binomial(link = "logit"), data=df_training_1h_log)
summary(model_full_log)
```

```
##
## Call:
## glm(formula = target ~ ., family = binomial(link = "logit"),
##      data = df_training_1h_log)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -51.47583    9.47210  -5.434 5.50e-08 ***
## indus       -0.15299    0.05221  -2.930 0.003387 **
## chas1        1.09580    0.80693   1.358 0.174472
## nox         62.12852    8.98196   6.917 4.61e-12 ***
## rm          -1.28605    0.71865  -1.790 0.073529 .
## age          0.02785    0.01352   2.060 0.039431 *
## ptratio      0.14295    0.13969   1.023 0.306125
## medv         0.22933    0.06474   3.542 0.000397 ***
## radq_low     1.65334    0.53895   3.068 0.002157 **
## radq_hi      4.85868    0.90038   5.396 6.80e-08 ***
## log_tax      1.73235    1.00185   1.729 0.083782 .
## log_dis      4.25555    1.04849   4.059 4.93e-05 ***
## log_zn      -0.49631    0.25373  -1.956 0.050460 .
## log_lstat    0.46418    0.67418   0.689 0.491133
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 172.97  on 452  degrees of freedom
## AIC: 200.97
##
## Number of Fisher Scoring iterations: 8
```

Alternate Transformation

Log transformations: medv, lstat, and dis might benefit from log transformations to reduce skewness. Binning: Convert age into categorical buckets (e.g., young, middle-aged, old). Interactions: Create new features (e.g., lstat*medv to capture relationships between income and home values). Standardization: Normalize numerical variables to bring them to the same scale.

```
# Log transformation
crime_training_df <- crime_training_df %>%
  mutate(
    log_medv = log(medv + 1), # Avoid log(0)
    log_lstat = log(lstat + 1),
    log_dis = log(dis + 1)
  )

#Standardization (Normalize Continuous Variables)
#Standardization ensures that variables are on the same scale, improving model stability.
crime_training_df <- crime_training_df %>%
  mutate(
    zn_scaled = as.numeric(scale(zn)),
    indus_scaled = as.numeric(scale(indus)),
    nox_scaled = as.numeric(scale(nox)),
    rm_scaled = as.numeric(scale(rm)),
    age_scaled = as.numeric(scale(age)),
    dis_scaled = as.numeric(scale(dis)),
    rad_scaled = as.numeric(scale(rad)),
    tax_scaled = as.numeric(scale(tax)),
    ptratio_scaled = as.numeric(scale(ptratio))
  )

# Create categorical age groups bins
crime_training_df$age_group <- cut(crime_training_df$age, breaks=c(0, 30, 60, 100), labels=c("Young", "Middle-aged", "Old"))

# Create Interaction Terms
crime_training_df <- crime_training_df %>%
  mutate(
    lstat_medv_interact = log_lstat * log_medv, # Income & housing price interaction
    tax_rad_interact = tax * rad # Tax burden & highway accessibility
  )

colnames(crime_training_df)
```

## [1]	"zn"	"indus"	"chas"
## [4]	"nox"	"rm"	"age"
## [7]	"dis"	"rad"	"tax"
## [10]	"ptratio"	"lstat"	"medv"
## [13]	"target"	"predicted_prob"	"logit"
## [16]	"log_medv"	"log_lstat"	"log_dis"
## [19]	"zn_scaled"	"indus_scaled"	"nox_scaled"
## [22]	"rm_scaled"	"age_scaled"	"dis_scaled"
## [25]	"rad_scaled"	"tax_scaled"	"ptratio_scaled"
## [28]	"age_group"	"lstat_medv_interact"	"tax_rad_interact"

Alternate Transformation

```
### Puja Modified below:
pj_crime_training_df <- read_csv("https://raw.githubusercontent.com/uzmabb182/Data_621/refs/heads/main/

### Data Preparation
# Log Transformations
pj_crime_training_df <- pj_crime_training_df %>%
  mutate(
    log_zn = log1p(zn),
    log_indus = log1p(indus),
    log_tax = log1p(tax),
    log_rad = log1p(rad),
    log_lstat = log1p(lstat),
    log_medv = log1p(medv)
  )

# Binning 'age' into Categories
pj_crime_training_df <- pj_crime_training_df %>%
  mutate(age_group = cut(age, breaks = c(0, 40, 70, 100), labels = c("Young", "Middle-aged", "Old")))

# Normalize 'dis'
pj_crime_training_df <- pj_crime_training_df %>%
  mutate(dis_scaled = scale(dis))
```

Log Transformation

```
crime_training_zr <- crime_training_df
```

To address skewness and improve model interpretability, the following features were log transformed: tax, dis, zn, and lstat. These variables exhibited non normal distributions with long tails and outliers. Applying a log transformation reduces variance, brings extreme values closer to the center, and better satisfies the assumptions of logistic regression.

```
crime_training_zr <- crime_training_zr %>%
  mutate(
    zr_log_tax = log(tax),
    zr_log_dis = log(dis),
    zr_log_zn = log(zn + 1),
    zr_log_lstat = log(lstat)
  )
```

Binned Transformation

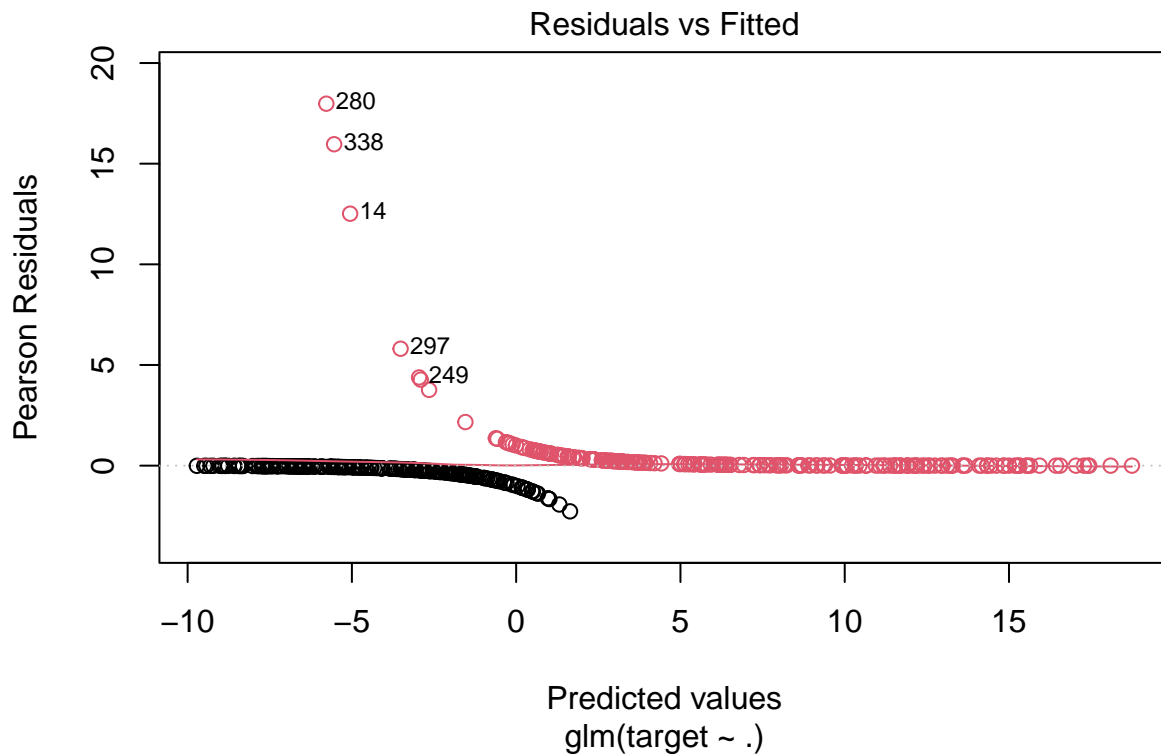
To simplify the effects of age and distance while preserving interpretability, these two variables were transformed into categorical bins using quantiles. This can help reduce the influence of extreme values and better capture non-linear effects in logistic regression.

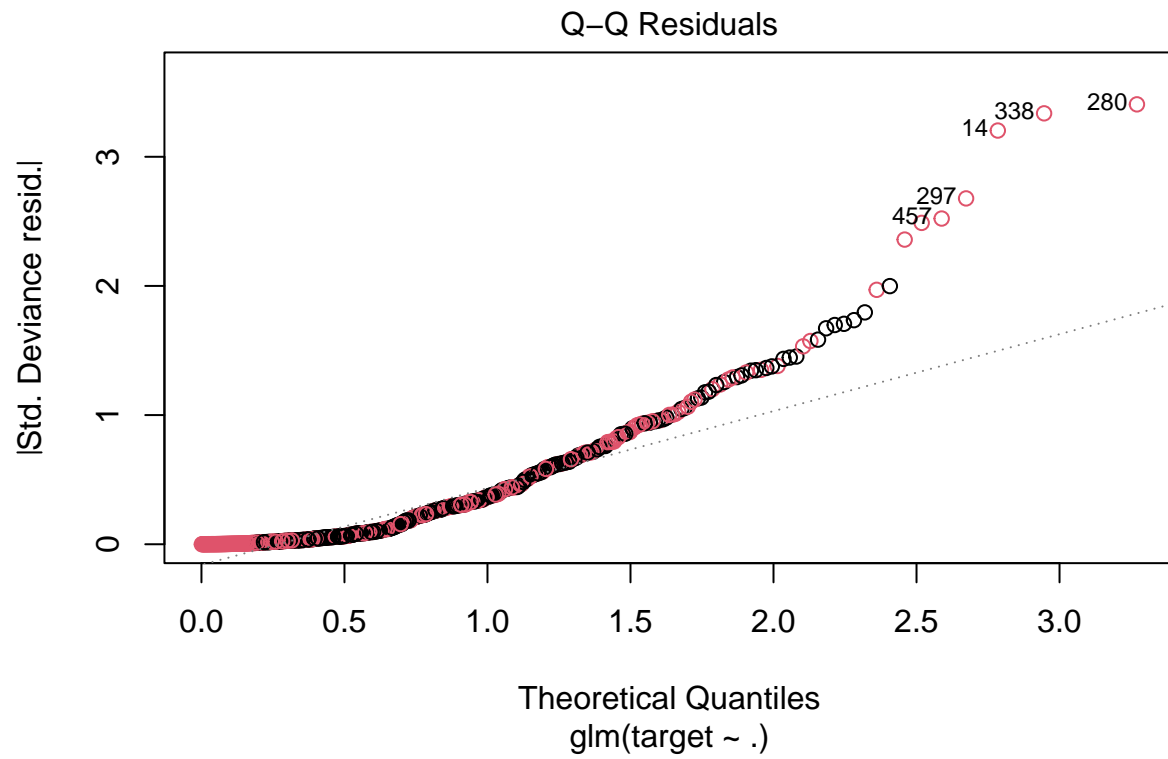
```
# Create age and dis bins
crime_training_zr$zr_age_bin <- cut(crime_training_zr$age,
                                   breaks = quantile(crime_training_zr$age, probs = c(0, 0.33, 0.66, 1),
                                   labels = c("Young", "Middle-aged", "Old"),
                                   include.lowest = TRUE)

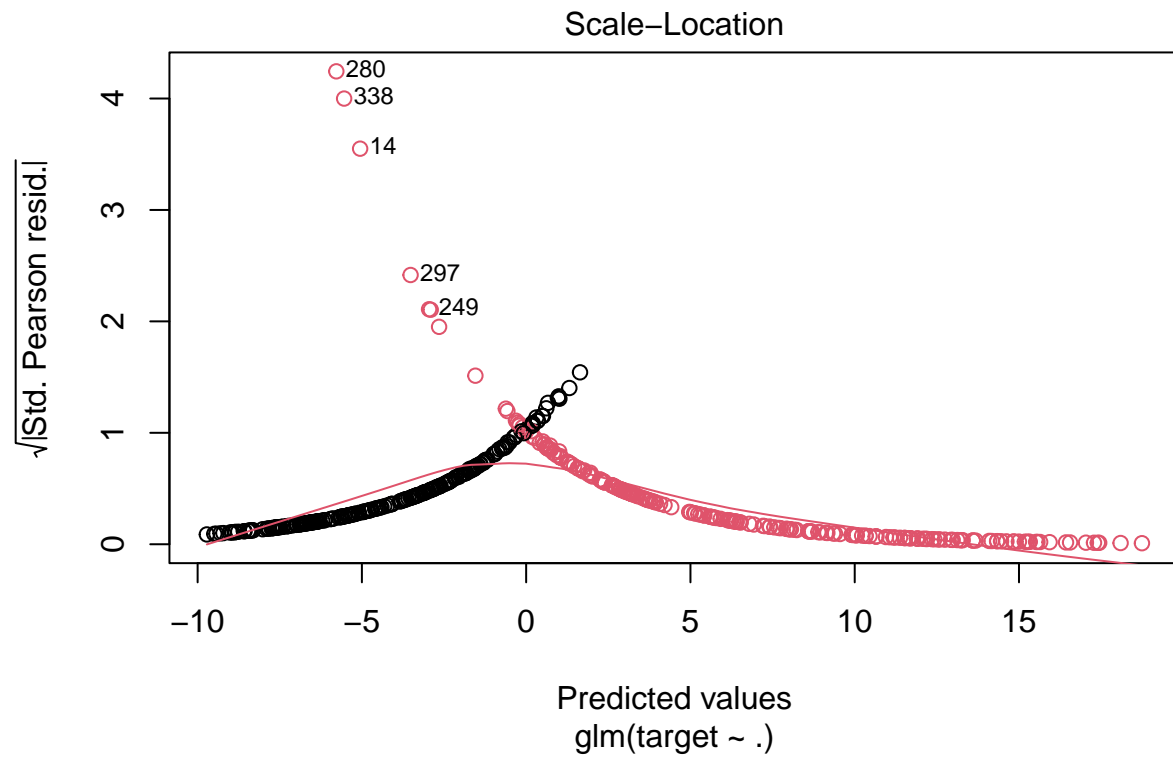
crime_training_zr$zr_dis_bin <- cut(crime_training_zr$dis,
                                   breaks = quantile(crime_training_zr$dis, probs = c(0, 0.33, 0.66, 1),
                                   labels = c("Near", "Mid-range", "Far"),
                                   include.lowest = TRUE)
```

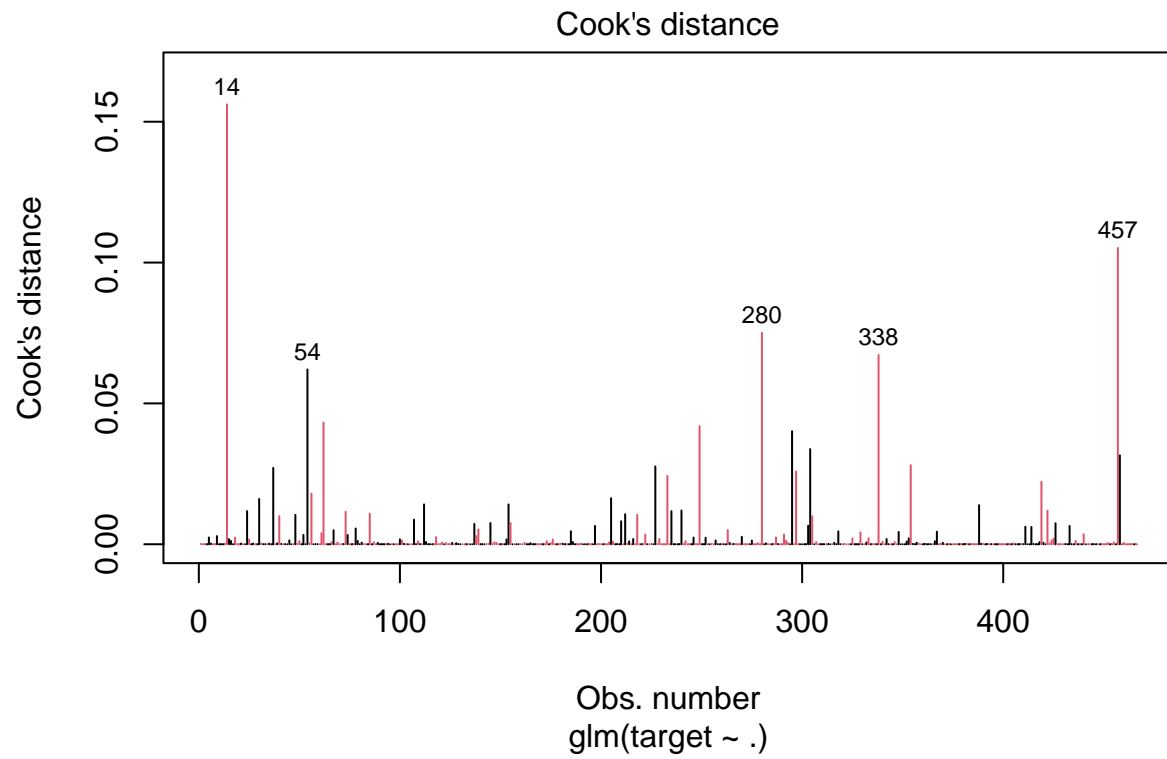
Outliers Applying the log transformation didn't make too much of a difference with our questionable points (280, 338, and 14)

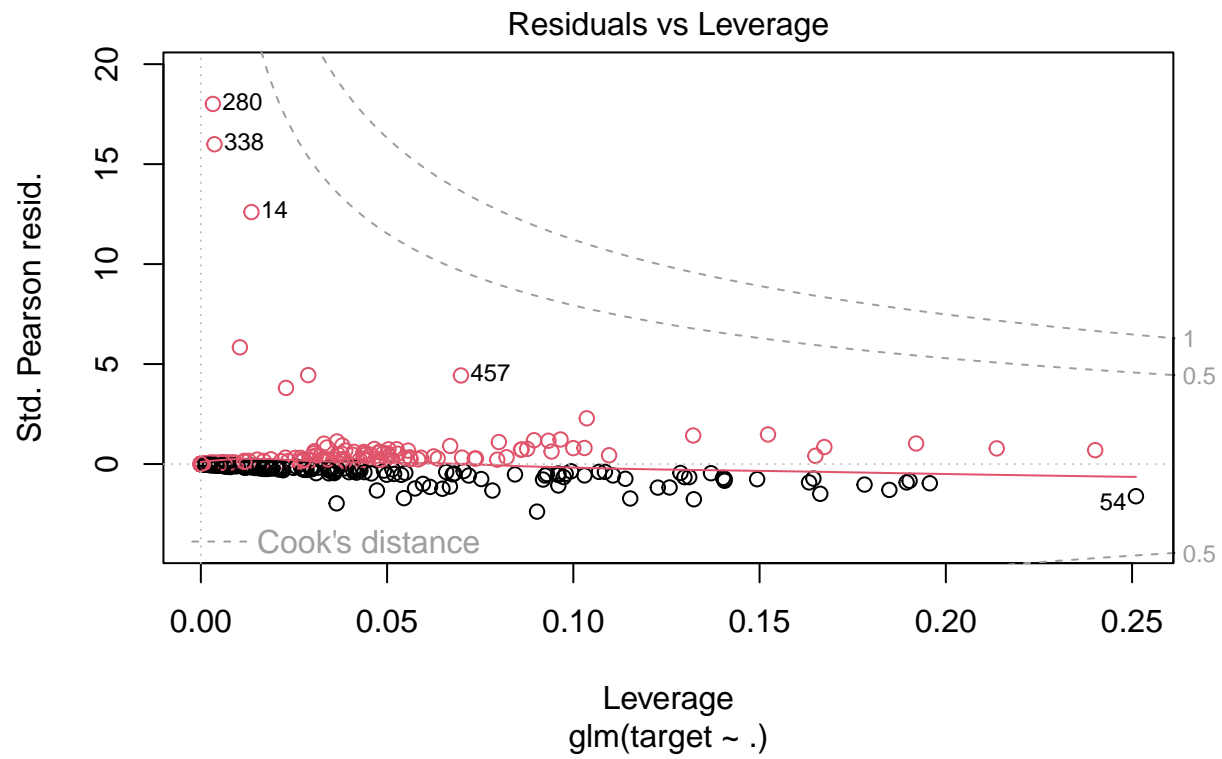
```
par(mar = c(5, 4, 4, 2) + 0.1)
plot(model_full_log, which = c(4, 6, 1, 2, 3, 5), col=df_training_1h_log$target, id.n = 5)
```

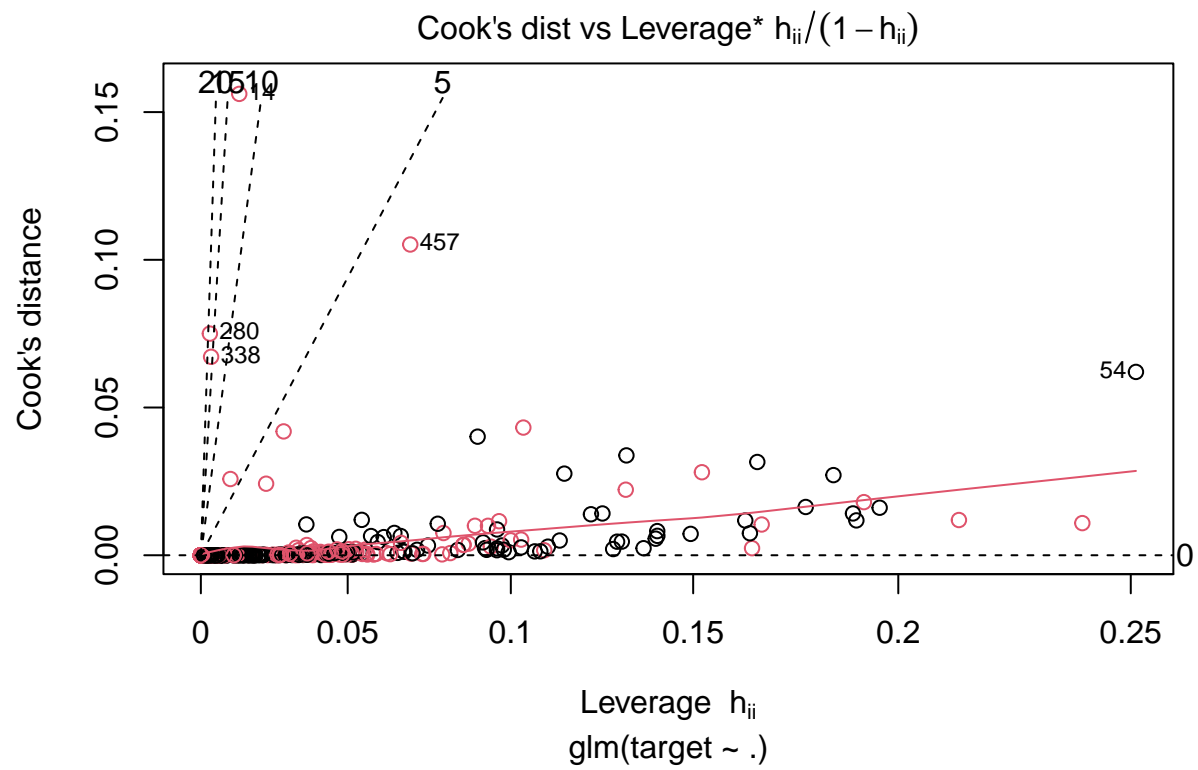






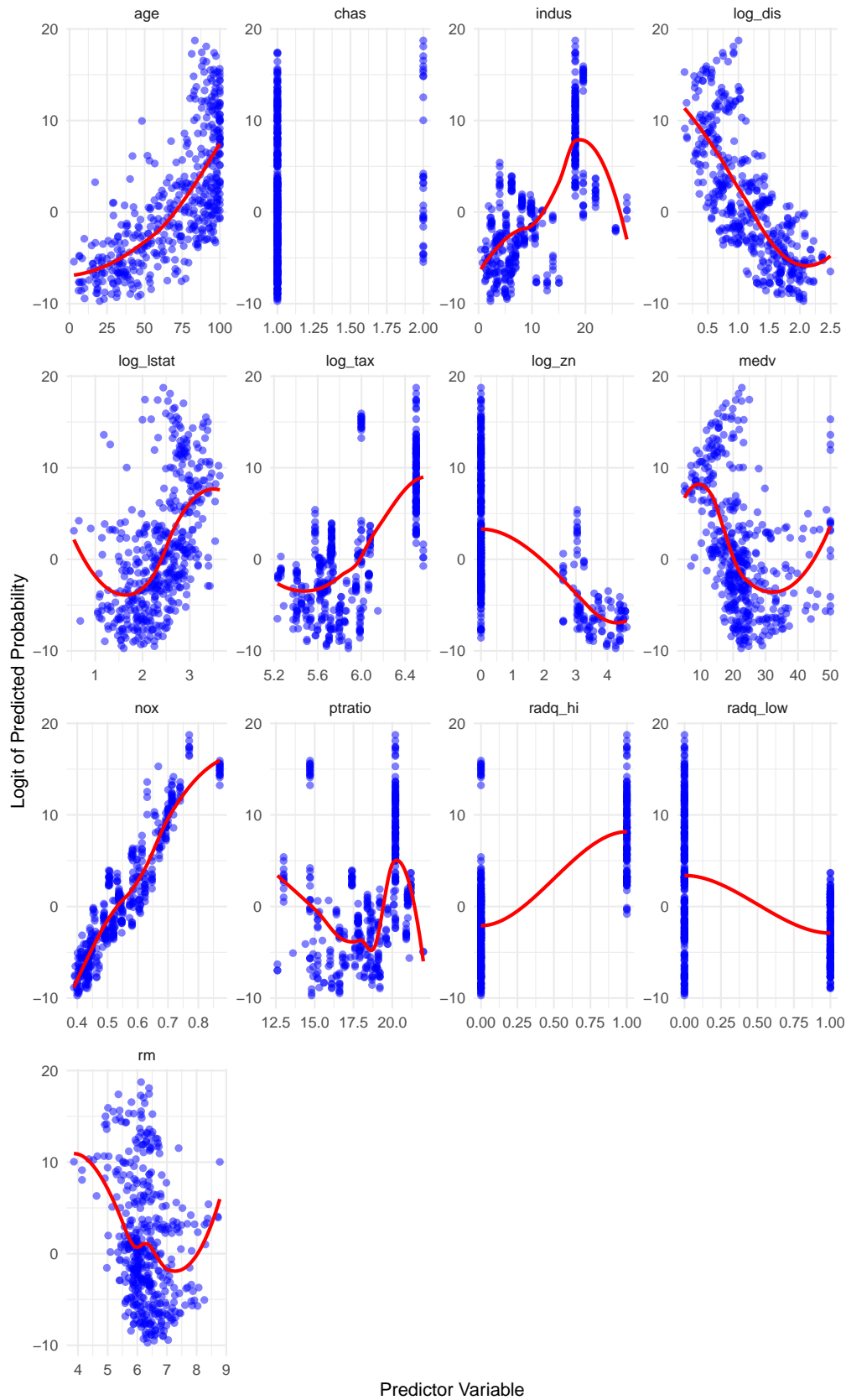






Linearity Applying the log transformation helped the linearity for some of the variables. It had less of an effect on medv, indus, and ptratio

Linearity of Logit Check for Binary Logistic Regression



Colinearity A Variance Inflation Factor test on our model with logged predictors shows that `nox` and `medv` should be considered for removal. `rm` and `log_dis` may also need to be considered.

```
car::vif(model_full_log) |> sort()
```

```
##      chas  radq_hi  log_zn radq_low  age  log_tax  ptratio log_lstat
## 1.227214 1.803965 1.834951 1.918943 2.121983 2.407600 2.516205 3.047489
##      indus      rm  log_dis      nox      medv
## 3.274047 4.731058 4.930674 5.318899 6.713174
```

MODEL BUILDING

Using a binomial (`target`) for our dependent variable would violate the common assumptions for linear regression. Specifically:

- the observations will not be normally distributed as they are binary
- the variance of error may be heteroskedastic instead of homoskedastic
- R-squared may not be a good fit

To account for these violations, we will use a Generalized Linear Model (GLM) to conduct logistic regression.

Approach 1:

Model 1: Baseline (All variables)

```
# Load library

# Logistic Regression Model 1 (Baseline)
modell1 <- glm(target ~ chas + lstat + medv + log_medv + log_lstat + log_dis +
              zn_scaled + indus_scaled + nox_scaled + rm_scaled + age_scaled + dis_scaled +
              rad_scaled + tax_scaled + ptratio_scaled +
              lstat_medv_interact + tax_rad_interact + age_group,
              data = crime_training_df, family = binomial)

# Summary of models
summary(modell1)

##
## Call:
## glm(formula = target ~ chas + lstat + medv + log_medv + log_lstat +
##      log_dis + zn_scaled + indus_scaled + nox_scaled + rm_scaled +
##      age_scaled + dis_scaled + rad_scaled + tax_scaled + ptratio_scaled +
##      lstat_medv_interact + tax_rad_interact + age_group, family = binomial,
##      data = crime_training_df)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -5.0602687 44.5969910  -0.113 0.909661
## chas             0.4993133  0.8073679   0.618 0.536282
## lstat           0.2652043  0.2277588   1.164 0.244258
```

```
## medv          0.2762597  0.2678575   1.031 0.302368
## log_medv      -3.8939721 14.6172280  -0.266 0.789934
## log_lstat     -5.6376391 13.5788530  -0.415 0.678012
## log_dis       17.0550235  5.6959404   2.994 0.002751 **
## zn_scaled     -0.6049684  0.7369038  -0.821 0.411669
## indus_scaled   0.1018175  0.3889005   0.262 0.793469
## nox_scaled     6.3323869  0.9839426   6.436 1.23e-10 ***
## rm_scaled     -0.7695076  0.5830524  -1.320 0.186905
## age_scaled     1.1937087  0.6715130   1.778 0.075463 .
## dis_scaled    -5.4654610  2.3511673  -2.325 0.020095 *
## rad_scaled    10.3465284  3.1747948   3.259 0.001118 **
## tax_scaled    -0.3596513  0.6158054  -0.584 0.559197
## ptratio_scaled 1.0016288  0.3022041   3.314 0.000918 ***
## lstat_medv_interact 0.6077222  3.5889380   0.169 0.865535
## tax_rad_interact -0.0011455  0.0006667  -1.718 0.085768 .
## age_groupMiddle-aged -1.8678051  1.2977251  -1.439 0.150068
## age_groupOld   -1.4112569  1.7240299  -0.819 0.413026
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 645.88 on 465 degrees of freedom
## Residual deviance: 174.22 on 446 degrees of freedom
## AIC: 214.22
##
## Number of Fisher Scoring iterations: 10
```

```
AIC(model1) # Compare AIC
```

```
## [1] 214.2244
```

Model 2: Stepwise Selection

```
# Logistic Regression Model 2 (Stepwise Selection)
model2 <- stepAIC(glm(target ~ ., data=crime_training_df, family=binomial), direction="both")

## Start: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
## ptratio + lstat + medv + predicted_prob + logit + log_medv +
## log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
## rm_scaled + age_scaled + dis_scaled + rad_scaled + tax_scaled +
## ptratio_scaled + age_group + lstat_medv_interact + tax_rad_interact
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
## ptratio + lstat + medv + predicted_prob + logit + log_medv +
## log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
## rm_scaled + age_scaled + dis_scaled + rad_scaled + tax_scaled +
## age_group + lstat_medv_interact + tax_rad_interact
```

```

##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
##   rm_scaled + age_scaled + dis_scaled + rad_scaled + age_group +
##   lstat_medv_interact + tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
##   rm_scaled + age_scaled + dis_scaled + age_group + lstat_medv_interact +
##   tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
##   rm_scaled + age_scaled + age_group + lstat_medv_interact +
##   tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
##   rm_scaled + age_group + lstat_medv_interact + tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + nox_scaled +
##   age_group + lstat_medv_interact + tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + indus_scaled + age_group +
##   lstat_medv_interact + tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##   ptratio + lstat + medv + predicted_prob + logit + log_medv +
##   log_lstat + log_dis + zn_scaled + age_group + lstat_medv_interact +
##   tax_rad_interact
##
##

```

```

## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##      ptratio + lstat + medv + predicted_prob + logit + log_medv +
##      log_lstat + log_dis + age_group + lstat_medv_interact + tax_rad_interact
##
##
## Step: AIC=215.97
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
##      ptratio + lstat + medv + predicted_prob + log_medv + log_lstat +
##      log_dis + age_group + lstat_medv_interact + tax_rad_interact
##
##
##      Df Deviance    AIC
## - indus      1    173.99 213.99
## - lstat_medv_interact 1    174.01 214.01
## - log_medv      1    174.07 214.07
## - tax          1    174.13 214.13
## - log_lstat     1    174.16 214.16
## - zn           1    174.19 214.19
## - predicted_prob 1    174.22 214.22
## - chas         1    174.31 214.31
## - age_group     2    176.44 214.44
## - tax_rad_interact 1    175.07 215.07
## - medv          1    175.09 215.09
## - lstat         1    175.27 215.27
## - rm           1    175.86 215.86
## <none>          173.97 215.97
## - age          1    176.11 216.11
## - rad          1    177.98 217.98
## - dis          1    180.14 220.14
## - log_dis      1    183.16 223.16
## - ptratio      1    186.62 226.62
## - nox          1    189.02 229.02
##
##
## Step: AIC=213.99
## target ~ zn + chas + nox + rm + age + dis + rad + tax + ptratio +
##      lstat + medv + predicted_prob + log_medv + log_lstat + log_dis +
##      age_group + lstat_medv_interact + tax_rad_interact
##
##
##      Df Deviance    AIC
## - lstat_medv_interact 1    174.03 212.03
## - log_medv          1    174.08 212.08
## - tax              1    174.14 212.14
## - log_lstat        1    174.17 212.17
## - zn              1    174.22 212.22
## - predicted_prob    1    174.29 212.29
## - chas            1    174.40 212.40
## - age_group        2    176.49 212.49
## - tax_rad_interact  1    175.07 213.07
## - medv            1    175.09 213.09
## - lstat           1    175.27 213.27
## - rm             1    175.90 213.90
## <none>            173.99 213.99
## - age            1    176.11 214.11
## + indus          1    173.97 215.97

```

```

## + indus_scaled      1    173.97 215.97
## - rad               1    178.36 216.36
## - dis              1    180.95 218.95
## - log_dis          1    185.17 223.17
## - ptratio          1    186.71 224.71
## - nox              1    189.36 227.36
##
## Step: AIC=212.03
## target ~ zn + chas + nox + rm + age + dis + rad + tax + ptratio +
##         lstat + medv + predicted_prob + log_medv + log_lstat + log_dis +
##         age_group + tax_rad_interact
##
##           Df Deviance    AIC
## - tax      1    174.22 210.22
## - log_medv  1    174.24 210.24
## - zn       1    174.29 210.29
## - predicted_prob 1    174.31 210.31
## - chas     1    174.44 210.44
## - age_group 2    176.65 210.65
## - tax_rad_interact 1    175.11 211.11
## - rm       1    175.91 211.91
## - lstat    1    175.97 211.97
## <none>      174.03 212.03
## - log_lstat 1    176.10 212.10
## - age      1    176.11 212.11
## - medv     1    176.32 212.32
## + lstat_medv_interact 1    173.99 213.99
## + indus_scaled 1    174.01 214.01
## + indus     1    174.01 214.01
## - rad      1    178.43 214.43
## - dis      1    180.97 216.97
## - log_dis  1    185.24 221.24
## - ptratio  1    186.71 222.71
## - nox      1    189.74 225.74
##
## Step: AIC=210.22
## target ~ zn + chas + nox + rm + age + dis + rad + ptratio + lstat +
##         medv + predicted_prob + log_medv + log_lstat + log_dis +
##         age_group + tax_rad_interact
##
##           Df Deviance    AIC
## - log_medv  1    174.40 208.40
## - zn       1    174.44 208.44
## - chas     1    174.59 208.59
## - age_group 2    176.70 208.70
## - predicted_prob 1    174.70 208.70
## - lstat    1    176.13 210.13
## - rm       1    176.19 210.19
## <none>      174.22 210.22
## - age      1    176.24 210.24
## - log_lstat 1    176.25 210.25
## - medv     1    176.48 210.48
## - tax_rad_interact 1    176.71 210.71
## + tax      1    174.03 212.03

```

```

## + logit          1    174.03 212.03
## + tax_scaled     1    174.03 212.03
## + lstat_medv_interact 1    174.14 212.14
## + indus          1    174.22 212.22
## + indus_scaled   1    174.22 212.22
## - rad            1    179.80 213.80
## - dis            1    182.17 216.17
## - log_dis        1    186.56 220.56
## - ptratio        1    187.37 221.37
## - nox            1    189.82 223.82
##
## Step:  AIC=208.4
## target ~ zn + chas + nox + rm + age + dis + rad + ptratio + lstat +
##         medv + predicted_prob + log_lstat + log_dis + age_group +
##         tax_rad_interact
##
##           Df Deviance    AIC
## - zn          1    174.62 206.62
## - age_group    2    176.72 206.72
## - chas         1    174.72 206.72
## - predicted_prob 1    174.78 206.78
## - rm           1    176.20 208.20
## <none>         174.40 208.40
## - age          1    176.55 208.55
## - tax_rad_interact 1    176.79 208.79
## - log_lstat    1    177.39 209.39
## - lstat        1    177.78 209.78
## + log_medv     1    174.22 210.22
## + tax           1    174.24 210.24
## + tax_scaled   1    174.24 210.24
## + logit        1    174.24 210.24
## + lstat_medv_interact 1    174.29 210.29
## + indus_scaled 1    174.40 210.40
## + indus        1    174.40 210.40
## - rad          1    179.93 211.93
## - medv         1    180.19 212.19
## - dis          1    182.56 214.56
## - log_dis      1    187.24 219.24
## - ptratio      1    187.90 219.90
## - nox          1    192.05 224.05
##
## Step:  AIC=206.62
## target ~ chas + nox + rm + age + dis + rad + ptratio + lstat +
##         medv + predicted_prob + log_lstat + log_dis + age_group +
##         tax_rad_interact
##
##           Df Deviance    AIC
## - age_group    2    176.75 204.75
## - chas         1    174.97 204.97
## - predicted_prob 1    175.65 205.65
## <none>         174.62 206.62
## - age          1    176.66 206.66
## - rm           1    176.72 206.72
## - tax_rad_interact 1    176.81 206.81

```



```

## - log_lstat      1    177.53 207.53
## - lstat          1    177.88 207.88
## + logit          1    174.25 208.25
## + zn             1    174.40 208.40
## + zn_scaled      1    174.40 208.40
## + log_medv       1    174.44 208.44
## + tax_scaled     1    174.51 208.51
## + tax            1    174.51 208.51
## + lstat_medv_interact 1    174.52 208.52
## + indus          1    174.62 208.62
## + indus_scaled   1    174.62 208.62
## - rad            1    179.93 209.93
## - medv           1    180.34 210.34
## - dis            1    185.08 215.08
## - log_dis        1    188.78 218.78
## - ptratio        1    190.24 220.24
## - nox            1    192.08 222.08
##
## Step:  AIC=204.75
## target ~ chas + nox + rm + age + dis + rad + ptratio + lstat +
##         medv + predicted_prob + log_lstat + log_dis + tax_rad_interact
##
##           Df Deviance    AIC
## - chas      1    177.01 203.01
## - predicted_prob 1    177.90 203.90
## - rm         1    178.59 204.59
## - tax_rad_interact 1    178.62 204.62
## <none>              176.75 204.75
## - log_lstat      1    180.45 206.45
## + age_group       2    174.62 206.62
## + logit           1    176.68 206.68
## + tax             1    176.72 206.72
## + tax_scaled      1    176.72 206.72
## + zn              1    176.72 206.72
## + zn_scaled       1    176.72 206.72
## + log_medv        1    176.73 206.73
## + indus_scaled    1    176.74 206.74
## + indus           1    176.74 206.74
## + lstat_medv_interact 1    176.75 206.75
## - lstat           1    181.07 207.07
## - age             1    181.25 207.25
## - rad             1    181.66 207.66
## - medv            1    181.92 207.92
## - dis             1    186.37 212.37
## - log_dis         1    189.72 215.72
## - ptratio         1    191.70 217.70
## - nox             1    193.03 219.03
##
## Step:  AIC=203.01
## target ~ nox + rm + age + dis + rad + ptratio + lstat + medv +
##         predicted_prob + log_lstat + log_dis + tax_rad_interact
##
##           Df Deviance    AIC
## - predicted_prob  1    178.45 202.45

```

```

## - tax_rad_interact      1   178.96 202.96
## - rm                    1   178.99 202.99
## <none>                  177.01 203.01
## + chas                  1   176.75 204.75
## - log_lstat             1   180.86 204.86
## + logit                 1   176.94 204.94
## + indus                 1   176.96 204.96
## + indus_scaled          1   176.96 204.96
## + age_group             2   174.97 204.97
## + zn_scaled             1   176.97 204.97
## + zn                    1   176.97 204.97
## + tax_scaled            1   176.99 204.99
## + tax                   1   176.99 204.99
## + log_medv              1   177.00 205.00
## + lstat_medv_interact   1   177.01 205.01
## - age                   1   181.59 205.59
## - lstat                 1   181.62 205.62
## - rad                   1   182.09 206.09
## - medv                  1   182.23 206.23
## - dis                   1   187.05 211.05
## - log_dis               1   190.08 214.08
## - ptratio               1   191.70 215.70
## - nox                   1   193.04 217.04
##
## Step:  AIC=202.45
## target ~ nox + rm + age + dis + rad + ptratio + lstat + medv +
##          log_lstat + log_dis + tax_rad_interact
##
##              Df Deviance    AIC
## <none>                178.45 202.45
## - rm                  1   180.84 202.84
## + predicted_prob      1   177.01 203.01
## + logit               1   177.56 203.56
## + zn                  1   177.73 203.73
## + zn_scaled           1   177.73 203.73
## + chas                1   177.90 203.90
## + indus               1   178.15 204.15
## + indus_scaled        1   178.15 204.15
## + tax_scaled          1   178.24 204.24
## + tax                 1   178.24 204.24
## + age_group           2   176.26 204.26
## + lstat_medv_interact 1   178.36 204.36
## + log_medv            1   178.41 204.41
## - log_lstat           1   182.72 204.72
## - lstat               1   183.31 205.31
## - tax_rad_interact    1   183.89 205.89
## - medv                1   185.16 207.16
## - age                 1   188.21 210.21
## - dis                 1   191.26 213.26
## - ptratio             1   193.88 215.88
## - log_dis             1   198.19 220.19
## - rad                 1   199.02 221.02
## - nox                 1   272.24 294.24

```

```
summary(model2)
```

```
##
## Call:
## glm(formula = target ~ nox + rm + age + dis + rad + ptratio +
##      lstat + medv + log_lstat + log_dis + tax_rad_interact, family = binomial,
##      data = crime_training_df)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -52.249633    9.248748  -5.649 1.61e-08 ***
## nox             54.605667    7.883747   6.926 4.32e-12 ***
## rm             -1.208141    0.790079  -1.529 0.126230
## age              0.042762    0.014541   2.941 0.003274 **
## dis            -2.916585    0.966434  -3.018 0.002545 **
## rad              1.207518    0.282327   4.277 1.89e-05 ***
## ptratio         0.468597    0.128001   3.661 0.000251 ***
## lstat           0.306041    0.144472   2.118 0.034146 *
## medv            0.187502    0.075197   2.493 0.012649 *
## log_lstat      -4.160910    2.078049  -2.002 0.045251 *
## log_dis        18.413523    4.919939   3.743 0.000182 ***
## tax_rad_interact -0.001236    0.000458  -2.698 0.006966 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 178.45  on 454  degrees of freedom
## AIC: 202.45
##
## Number of Fisher Scoring iterations: 9
```

```
AIC(model2) # Compare AIC
```

```
## [1] 202.4455
```

Model 3: Transformations & Interactions

```
# Logistic Regression Model 3 (With Transformations & Interactions)
model3 <- glm(target ~ log_medv + lstat + nox + ptratio + age_group, data=crime_training_df, family=binomial)

summary(model3)
```

```
##
## Call:
## glm(formula = target ~ log_medv + lstat + nox + ptratio + age_group,
##      family = binomial, data = crime_training_df)
##
## Coefficients:
```

```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -31.74902    4.89194  -6.490 8.58e-11 ***
## log_medv       2.82091    0.80743   3.494 0.000476 ***
## lstat          0.05754    0.03958   1.454 0.145995
## nox            30.75449    3.52525   8.724 < 2e-16 ***
## ptratio        0.28767    0.09449   3.044 0.002332 **
## age_groupMiddle-aged -0.17887    0.72673  -0.246 0.805579
## age_groupOld     0.38894    0.69726   0.558 0.576980
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 274.17  on 459  degrees of freedom
## AIC: 288.17
##
## Number of Fisher Scoring iterations: 6
```

Approach 2:

Model 1: Baseline Logistic Regression

```
### Model Building
# Model 1: Baseline Logistic Regression
pj_model1 <- glm(target ~ zn + indus + chas + nox + rm + age + dis + rad + tax + ptratio + lstat + medv,
                 data = pj_crime_training_df, family = binomial)
summary(pj_model1)
```

```
##
## Call:
## glm(formula = target ~ zn + indus + chas + nox + rm + age + dis +
##      rad + tax + ptratio + lstat + medv, family = binomial, data = pj_crime_training_df)
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -40.822934    6.632913  -6.155 7.53e-10 ***
## zn           -0.065946    0.034656  -1.903 0.05706 .
## indus        -0.064614    0.047622  -1.357 0.17485
## chas          0.910765    0.755546   1.205 0.22803
## nox           49.122297    7.931706   6.193 5.90e-10 ***
## rm           -0.587488    0.722847  -0.813 0.41637
## age           0.034189    0.013814   2.475 0.01333 *
## dis           0.738660    0.230275   3.208 0.00134 **
## rad           0.666366    0.163152   4.084 4.42e-05 ***
## tax          -0.006171    0.002955  -2.089 0.03674 *
## ptratio       0.402566    0.126627   3.179 0.00148 **
## lstat         0.045869    0.054049   0.849 0.39608
## medv          0.180824    0.068294   2.648 0.00810 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 645.88 on 465 degrees of freedom
## Residual deviance: 192.05 on 453 degrees of freedom
## AIC: 218.05
##
## Number of Fisher Scoring iterations: 9
```

Model 2: Stepwise Logistic Regression

```
#Model 2: Stepwise Logistic Regression
pj_model2 <- step(glm(target ~ ., data = pj_crime_training_df, family = binomial), direction = "both")
```

```
## Start: AIC=186.78
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
## ptratio + lstat + medv + log_zn + log_indus + log_tax + log_rad +
## log_lstat + log_medv + age_group + dis_scaled
##
##
## Step: AIC=186.78
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
## ptratio + lstat + medv + log_zn + log_indus + log_tax + log_rad +
## log_lstat + log_medv + age_group
##
##      Df Deviance    AIC
## - age_group  2   146.01 184.01
## - zn         1   144.78 184.78
## - log_indus  1   144.78 184.78
## - log_zn     1   145.21 185.21
## - chas       1   145.33 185.33
## - indus      1   145.43 185.43
## - log_rad    1   146.19 186.19
## - dis        1   146.43 186.43
## - log_lstat  1   146.45 186.45
## - lstat      1   146.76 186.76
## <none>       144.78 186.78
## - rm         1   148.14 188.14
## - log_medv   1   148.74 188.74
## - age        1   148.80 188.80
## - medv       1   153.46 193.46
## - ptratio    1   155.46 195.46
## - rad        1   166.66 206.66
## - log_tax    1   170.97 210.97
## - tax        1   175.73 215.73
## - nox        1   184.31 224.31
##
## Step: AIC=184.01
## target ~ zn + indus + chas + nox + rm + age + dis + rad + tax +
## ptratio + lstat + medv + log_zn + log_indus + log_tax + log_rad +
## log_lstat + log_medv
##
##      Df Deviance    AIC
```

```

## - zn      1    146.01 182.01
## - log_indus 1    146.05 182.05
## - log_zn   1    146.35 182.35
## - chas     1    146.43 182.43
## - indus    1    146.86 182.86
## - log_rad  1    147.20 183.20
## - dis      1    147.76 183.76
## <none>      146.01 184.01
## - log_lstat 1    148.30 184.30
## - lstat    1    148.74 184.74
## - rm       1    149.64 185.64
## - log_medv 1    150.12 186.12
## + age_group 2    144.78 186.78
## - age      1    150.98 186.98
## - medv     1    154.86 190.86
## - ptratio  1    157.24 193.24
## - rad      1    166.96 202.96
## - log_tax  1    171.68 207.68
## - tax      1    176.48 212.48
## - nox      1    186.30 222.30
##
## Step: AIC=182.01
## target ~ indus + chas + nox + rm + age + dis + rad + tax + ptratio +
##          lstat + medv + log_zn + log_indus + log_tax + log_rad + log_lstat +
##          log_medv
##
##           Df Deviance    AIC
## - log_indus 1    146.06 180.06
## - chas      1    146.44 180.44
## - indus     1    146.92 180.92
## - log_rad   1    147.21 181.21
## - dis       1    147.88 181.88
## <none>      146.01 182.01
## - log_lstat 1    148.30 182.30
## - log_zn    1    148.73 182.73
## - lstat     1    148.75 182.75
## - rm        1    149.65 183.65
## + zn        1    146.01 184.01
## - log_medv  1    150.16 184.16
## + age_group 2    144.78 184.78
## - age       1    151.04 185.04
## - medv      1    154.88 188.88
## - ptratio   1    157.51 191.51
## - rad       1    166.96 200.96
## - log_tax   1    171.79 205.79
## - tax       1    176.62 210.62
## - nox       1    192.85 226.85
##
## Step: AIC=180.06
## target ~ indus + chas + nox + rm + age + dis + rad + tax + ptratio +
##          lstat + medv + log_zn + log_tax + log_rad + log_lstat + log_medv
##
##           Df Deviance    AIC
## - chas      1    146.53 178.53

```

```

## - log_rad      1    147.35 179.35
## - dis          1    147.90 179.90
## <none>         146.06 180.06
## - log_lstat    1    148.39 180.39
## - lstat        1    148.78 180.78
## - log_zn       1    149.08 181.08
## - rm           1    149.72 181.72
## + log_indus    1    146.01 182.01
## + zn           1    146.05 182.05
## - log_medv     1    150.35 182.35
## + age_group    2    144.78 182.78
## - indus        1    151.19 183.19
## - age          1    151.24 183.24
## - medv         1    155.10 187.10
## - ptratio      1    157.72 189.72
## - rad          1    170.12 202.12
## - log_tax      1    187.07 219.07
## - tax          1    192.62 224.62
## - nox          1    192.90 224.90
##
## Step: AIC=178.53
## target ~ indus + nox + rm + age + dis + rad + tax + ptratio +
##          lstat + medv + log_zn + log_tax + log_rad + log_lstat + log_medv
##
##           Df Deviance    AIC
## - log_rad      1    147.77 177.77
## <none>         146.53 178.53
## - log_lstat    1    148.73 178.73
## - dis          1    148.79 178.79
## - lstat        1    149.01 179.01
## - log_zn       1    149.24 179.24
## - rm           1    150.03 180.03
## + chas         1    146.06 180.06
## + log_indus    1    146.44 180.44
## + zn           1    146.53 180.53
## - log_medv     1    150.98 180.98
## - indus        1    151.25 181.25
## - age          1    151.31 181.31
## + age_group    2    145.37 181.37
## - medv         1    155.57 185.57
## - ptratio      1    158.20 188.20
## - rad          1    170.28 200.28
## - log_tax      1    188.21 218.21
## - tax          1    194.54 224.54
## - nox          1    195.82 225.82
##
## Step: AIC=177.77
## target ~ indus + nox + rm + age + dis + rad + tax + ptratio +
##          lstat + medv + log_zn + log_tax + log_lstat + log_medv
##
##           Df Deviance    AIC
## - dis          1    149.74 177.74
## <none>         147.77 177.77
## - log_zn       1    150.43 178.43

```

```

## - log_lstat 1 150.50 178.50
## + log_rad 1 146.53 178.53
## - lstat 1 150.55 178.55
## - rm 1 150.96 178.96
## + chas 1 147.35 179.35
## + log_indus 1 147.69 179.69
## + zn 1 147.76 179.76
## - log_medv 1 152.48 180.48
## - age 1 152.56 180.56
## - indus 1 152.84 180.84
## + age_group 2 146.97 180.97
## - medv 1 156.74 184.74
## - ptratio 1 160.63 188.63
## - log_tax 1 189.82 217.82
## - tax 1 195.18 223.18
## - nox 1 199.38 227.38
## - rad 1 233.98 261.98
##
## Step: AIC=177.74
## target ~ indus + nox + rm + age + rad + tax + ptratio + lstat +
## medv + log_zn + log_tax + log_lstat + log_medv
##
## Df Deviance AIC
## - log_zn 1 150.81 176.81
## <none> 149.74 177.74
## + dis_scaled 1 147.77 177.77
## + dis 1 147.77 177.77
## - rm 1 152.11 178.11
## + log_rad 1 148.79 178.79
## + chas 1 148.98 178.98
## - log_lstat 1 153.13 179.13
## - age 1 153.15 179.15
## - lstat 1 153.26 179.26
## + zn 1 149.41 179.41
## + log_indus 1 149.51 179.51
## + age_group 2 148.79 180.79
## - log_medv 1 155.01 181.01
## - indus 1 157.37 183.37
## - medv 1 157.98 183.98
## - ptratio 1 164.13 190.13
## - log_tax 1 199.05 225.05
## - tax 1 207.11 233.11
## - nox 1 211.76 237.76
## - rad 1 249.42 275.42
##
## Step: AIC=176.81
## target ~ indus + nox + rm + age + rad + tax + ptratio + lstat +
## medv + log_tax + log_lstat + log_medv
##
## Df Deviance AIC
## <none> 150.81 176.81
## + zn 1 149.48 177.48
## + log_rad 1 149.69 177.69
## + log_zn 1 149.74 177.74

```



```
## + log_indus      1    150.04 178.04
## - lstat          1    154.10 178.10
## - log_lstat      1    154.15 178.15
## - rm             1    154.28 178.28
## + dis            1    150.43 178.43
## + dis_scaled     1    150.43 178.43
## + chas           1    150.57 178.57
## - age            1    154.99 178.99
## - log_medv       1    155.59 179.59
## + age_group      2    149.97 179.97
## - indus          1    158.63 182.63
## - medv           1    158.69 182.69
## - ptratio        1    171.16 195.16
## - log_tax        1    200.49 224.49
## - tax            1    208.54 232.54
## - nox            1    218.11 242.11
## - rad            1    251.50 275.50
```

```
summary(pj_model2)
```

```
##
## Call:
## glm(formula = target ~ indus + nox + rm + age + rad + tax + ptratio +
##       lstat + medv + log_tax + log_lstat + log_medv, family = binomial,
##       data = pj_crime_training_df)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -259.23712   49.39152  -5.249 1.53e-07 ***
## indus         0.19523    0.07065   2.763 0.00573 **
## nox          43.61151    7.29024   5.982 2.20e-09 ***
## rm          -1.61615    0.87896  -1.839 0.06596 .
## age           0.02858    0.01432   1.995 0.04600 *
## rad           1.26111    0.20382   6.187 6.12e-10 ***
## tax          -0.18471    0.03353  -5.509 3.62e-08 ***
## ptratio       0.61262    0.15168   4.039 5.37e-05 ***
## lstat         0.34249    0.18703   1.831 0.06707 .
## medv          0.62332    0.23615   2.640 0.00830 **
## log_tax       54.25065    9.94821   5.453 4.94e-08 ***
## log_lstat    -4.87315    2.67324  -1.823 0.06831 .
## log_medv     -11.13849    5.29557  -2.103 0.03543 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 150.81  on 453  degrees of freedom
## AIC: 176.81
##
## Number of Fisher Scoring iterations: 8
```

```
# Model 3: Logistic Regression with Transformed Variables
pj_model3 <- glm(target ~ log_zn + log_indus + chas + nox + rm + age_group + dis_scaled + log_rad + log
                data = pj_crime_training_df, family = binomial)
summary(pj_model3)
```

Model 3: Logistic Regression with Transformed Variables

```
##
## Call:
## glm(formula = target ~ log_zn + log_indus + chas + nox + rm +
##      age_group + dis_scaled + log_rad + log_tax + ptratio + log_lstat +
##      log_medv, family = binomial, data = pj_crime_training_df)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -34.7064     10.0151  -3.465 0.000529 ***
## log_zn          -0.4754      0.2365  -2.010 0.044456 *
## log_indus       -0.2371      0.5819  -0.408 0.683618
## chas             1.0171      0.7743   1.314 0.188965
## nox             44.4244      7.3505   6.044 1.51e-09 ***
## rm              0.3355      0.6008   0.559 0.576484
## age_groupMiddle-aged 0.2113      0.7028   0.301 0.763631
## age_groupOld       1.1467      0.7654   1.498 0.134096
## dis_scaled       1.3220      0.4593   2.878 0.004000 **
## log_rad          3.7247      0.8263   4.508 6.55e-06 ***
## log_tax         -2.0090      1.0548  -1.905 0.056829 .
## ptratio          0.2923      0.1235   2.366 0.017964 *
## log_lstat        0.4373      0.7558   0.579 0.562824
## log_medv         2.2550      1.5833   1.424 0.154370
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 205.35  on 452  degrees of freedom
## AIC: 233.35
##
## Number of Fisher Scoring iterations: 8
```

Approach 3:

Model 1: Baseline Predictions This model includes six variables identified through exploratory data analysis and correlation inspection: nox, dis, tax, rad, ptratio, and lstat. These were selected based on their relatively strong correlation with the binary crime outcome and theoretical reasoning. For example, higher pollution (nox) and tax rates (tax) may be indicative of urban density, which could correlate with higher crime, while greater distance to employment centers (dis) might have a protective effect. This model provides a straightforward process using untransformed, raw features.

```
zr_model_a <- glm(target ~ nox + dis + tax + rad + ptratio + lstat,
                  data = crime_training_zr, family = binomial)
```

Model 2: Log-Transformed Predictors This model includes a broader set of variables, with several of them log-transformed: log_tax, log_dis, log_zn, log_lstat, as well as nox, rm, ptratio, rad, chas, and age. This method retains the core features from Model 1 but adjusts for non-normality and skewness observed in variables like tax, zn, dis, and lstat. These features exhibited right-skewed distributions and extreme values, which could affect model stability. Applying log transformation helps to reduce the impact of outliers, normalize distributions, and potentially improve model performance.

```
zr_model_b <- glm(target ~ zr_log_tax + zr_log_dis + zr_log_zn + zr_log_lstat + nox + rm + ptratio + rad + age,
                  data = crime_training_zr, family = binomial)
```

Zach Model 3: Binned Variables Logistic Regression This model includes age and distance as categorical bins based on quantiles, plus tax, rad, ptratio, nox, and rm. These features were selected based on exploratory analysis and their correlation with the target. Binning helps reduce skewness and simplify interpretation.

```
zr_model_c <- glm(target ~ zr_age_bin + zr_dis_bin + tax + rad + ptratio + nox + rm,
                  data = crime_training_zr, family = binomial)
```

Model using Principal Components

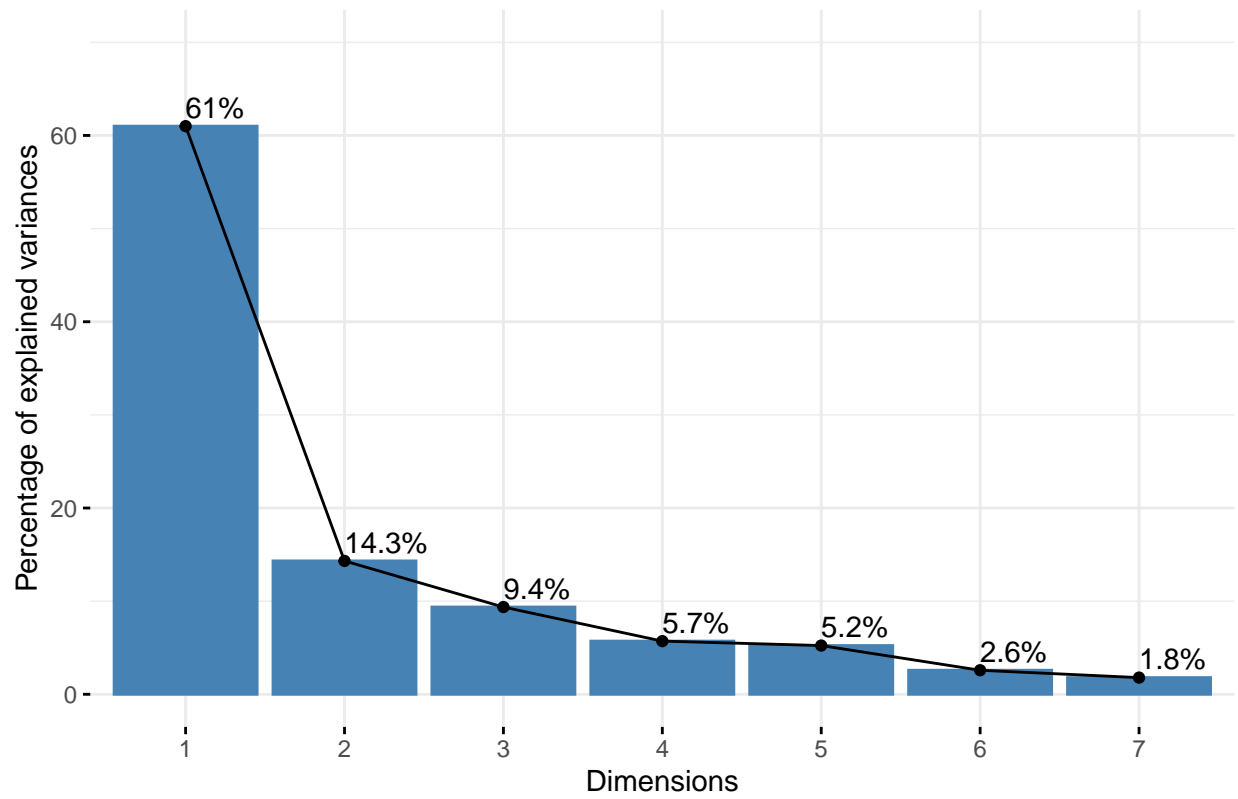
This section uses the correlation plot to perform Principal Component Analysis on the two large variable clusters shown in the plot. We will then substitute the variables in each of the two clusters with their respective PC scores in our model.

```
# Create PCA from first cluster in our correlation plot
df_pca_subset1 <- df_training_one_hot |>
  subset(select = c(indus, tax, lstat, nox, age, ptratio, radq_hi))

# calculate PCA
df_training_pca1 <- prcomp(df_pca_subset1, scale=TRUE)

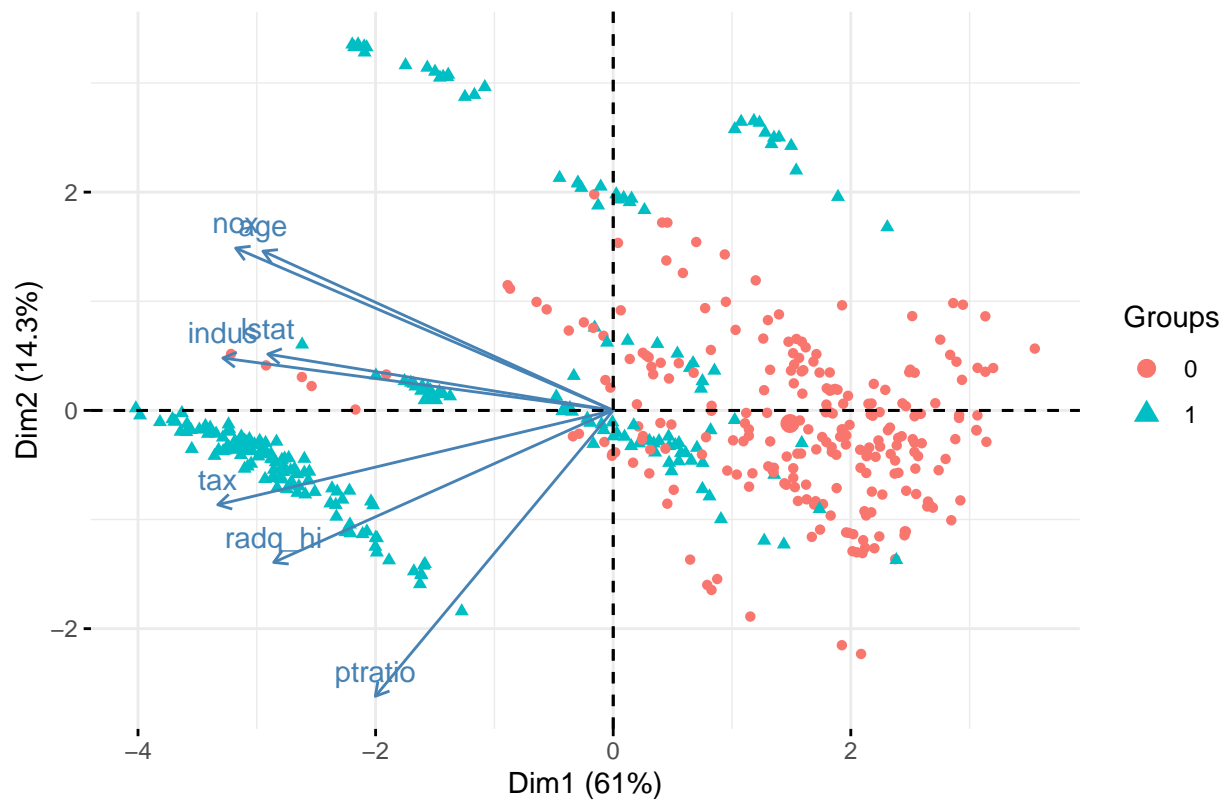
# use eigen vectors to plot % of data explained by PCA1
fviz_eig(df_training_pca1, addlabels=TRUE, ylim=c(0, 70))
```

Scree plot



```
# plot PCA biplot  
fviz_pca_biplot(df_training_pca1, label="var", habillage = df_training_one_hot$target)
```

PCA – Biplot

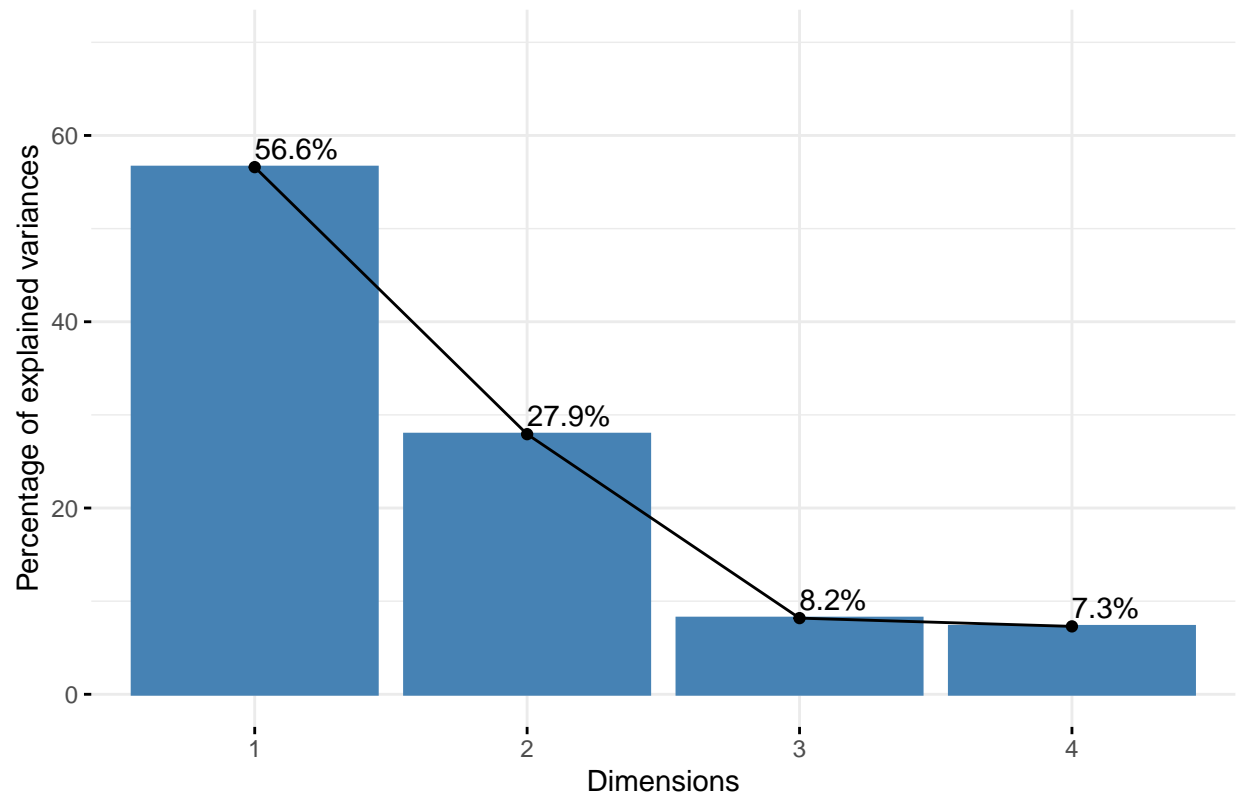


```
# Create PCA from second cluster in our correlation plot
df_pca_subset2 <- df_training_one_hot |>
  subset(select = c(rm, medv, zn, dis))

# calculate PCA
df_training_pca2 <- prcomp(df_pca_subset2, scale=TRUE)

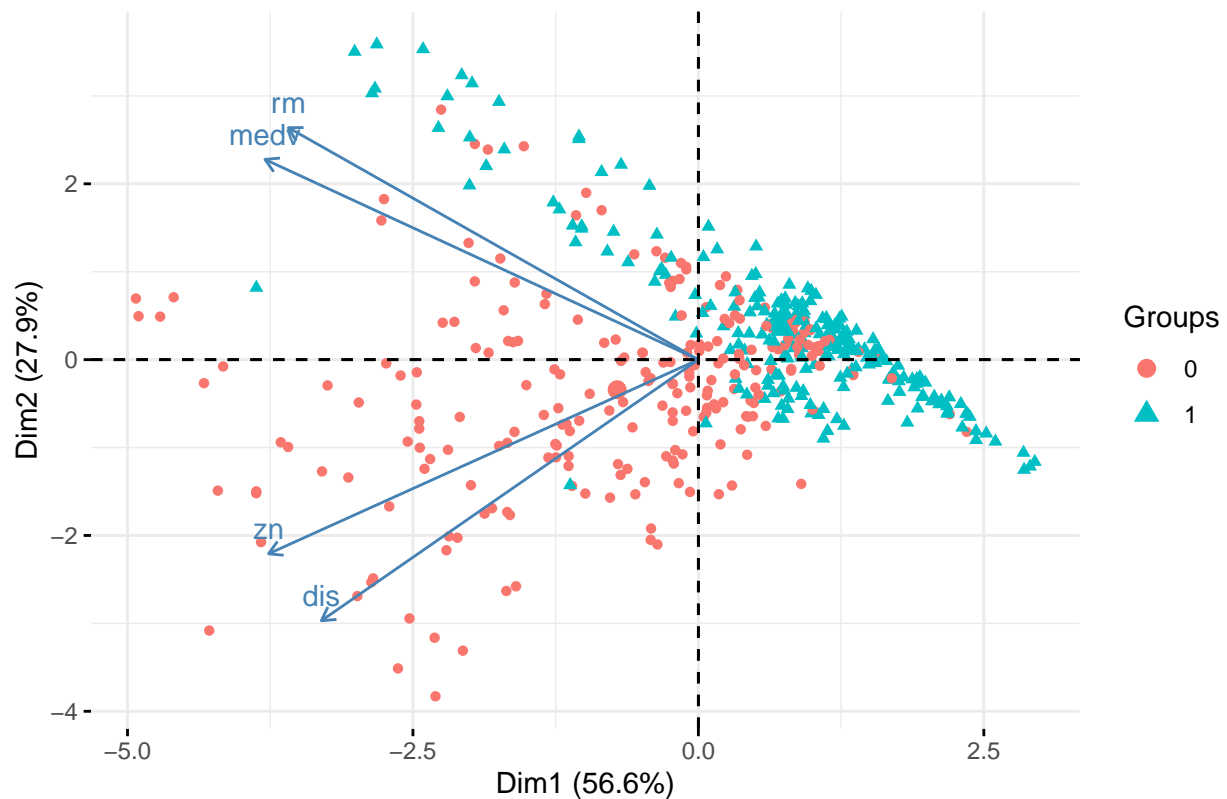
# use eigen vectors to plot % of data explained by PCA1
fviz_eig(df_training_pca2, addlabels=TRUE, ylim=c(0, 70))
```

Scree plot



```
# plot PCA biplot  
fviz_pca_biplot(df_training_pca2, label="var", habillage = df_training_one_hot$target)
```

PCA – Biplot



```
# add pca's to our dataset
df_training_one_hot_pca <- df_training_one_hot |>
  subset(select = c(target, chas, radq_low)) |>
  mutate(
    group1_pc1 = df_training_pca1$x[, "PC1"],
    group1_pc2 = df_training_pca1$x[, "PC2"],
    group2_pc1 = df_training_pca2$x[, "PC1"],
    group2_pc2 = df_training_pca2$x[, "PC2"],
  )

#ggpairs(df_training_one_hot_pca |> subset(select = -c(target)))

model_pca <- glm(target ~ ., binomial(link = "logit"), data=df_training_one_hot_pca)
summary(model_pca)
```

```
##
## Call:
## glm(formula = target ~ ., family = binomial(link = "logit"),
##      data = df_training_one_hot_pca)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.14990    0.22856   0.656 0.511908
## chas1        0.35486    0.52302   0.678 0.497468
## radq_low     -0.04425    0.32215  -0.137 0.890753
```

```
## group1_pc1 -1.45138 0.18342 -7.913 2.52e-15 ***
## group1_pc2 0.17208 0.15573 1.105 0.269166
## group2_pc1 -0.30550 0.20533 -1.488 0.136790
## group2_pc2 0.67822 0.18654 3.636 0.000277 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 645.88 on 465 degrees of freedom
## Residual deviance: 289.25 on 459 degrees of freedom
## AIC: 303.25
##
## Number of Fisher Scoring iterations: 6
```

Interestingly, only the primary principal component from group1 and the secondary principal component from group two have strong statistical significance. `radq_low` has a particularly high p-value and should be considered for removal.

```
model_pca2 <- update(model_pca, . ~ . - radq_low)
summary(model_pca2)
```

```
##
## Call:
## glm(formula = target ~ chas + group1_pc1 + group1_pc2 + group2_pc1 +
##      group2_pc2, family = binomial(link = "logit"), data = df_training_one_hot_pca)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.1331     0.1933   0.689 0.491070
## chas1        0.3545     0.5219   0.679 0.497018
## group1_pc1   -1.4576     0.1779  -8.194 2.53e-16 ***
## group1_pc2    0.1751     0.1541   1.136 0.256022
## group2_pc1   -0.3094     0.2032  -1.523 0.127809
## group2_pc2    0.6808     0.1856   3.668 0.000244 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 645.88 on 465 degrees of freedom
## Residual deviance: 289.27 on 460 degrees of freedom
## AIC: 301.27
##
## Number of Fisher Scoring iterations: 6
```

```
model_pca2 <- update(model_pca2, . ~ . - chas)
summary(model_pca2)
```

```
##
## Call:
## glm(formula = target ~ group1_pc1 + group1_pc2 + group2_pc1 +
```



```
##      group2_pc2, family = binomial(link = "logit"), data = df_training_one_hot_pca)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.1640     0.1881   0.871 0.383506
## group1_pc1   -1.4598     0.1780  -8.201 2.38e-16 ***
## group1_pc2    0.1848     0.1532   1.206 0.227758
## group2_pc1   -0.3124     0.2027  -1.542 0.123190
## group2_pc2    0.6823     0.1845   3.697 0.000218 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 289.73  on 461  degrees of freedom
## AIC: 299.73
##
## Number of Fisher Scoring iterations: 6
```

```
model_pca2 <- update(model_pca2, . ~ . - group1_pc2)
summary(model_pca2)
```

```
##
## Call:
## glm(formula = target ~ group1_pc1 + group2_pc1 + group2_pc2,
##      family = binomial(link = "logit"), data = df_training_one_hot_pca)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.1950     0.1905   1.024   0.306
## group1_pc1   -1.4594     0.1818  -8.027 9.99e-16 ***
## group2_pc1   -0.2843     0.2034  -1.398   0.162
## group2_pc2    0.7457     0.1783   4.181 2.90e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 291.22  on 462  degrees of freedom
## AIC: 299.22
##
## Number of Fisher Scoring iterations: 6
```

```
model_pca2 <- update(model_pca2, . ~ . - group2_pc1)
summary(model_pca2)
```

```
##
## Call:
## glm(formula = target ~ group1_pc1 + group2_pc2, family = binomial(link = "logit"),
##      data = df_training_one_hot_pca)
##
```

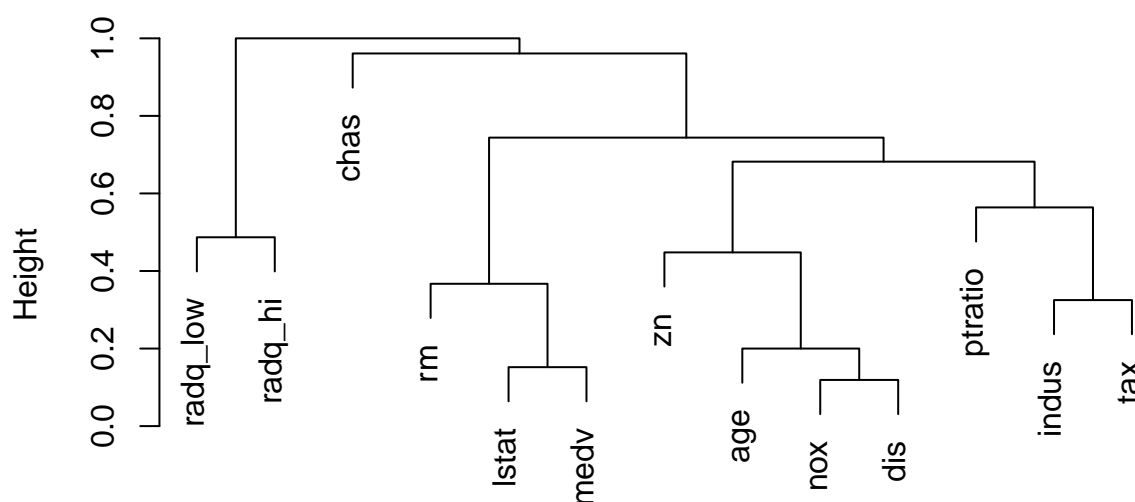
```
## Coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.05301    0.16131   0.329   0.742
## group1_pc1  -1.28804    0.12063 -10.678 < 2e-16 ***
## group2_pc2   0.87594    0.16084   5.446 5.15e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 293.06  on 463  degrees of freedom
## AIC: 299.06
##
## Number of Fisher Scoring iterations: 6
```

Model based on Variable Clustering

The dendrogram is a variable clustering technique that shows how the parameters progressively come together at different levels of similarity. It offers another way to visualize correlations between our parameters. In this model, we will use the dendrogram to prune parameters that are similar from the lower branches. In this model, we used the results from a T and Wilcoxon pairwise test to assist with the parameter selection.

```
dist_one_hot = as.dist(m = 1 - abs(df_training_cor))
par(mar = c(5, 4, 4, 2) + 0.1)
plot(hclust(dist_one_hot))
```

Cluster Dendrogram



dist_one_hot
hclust (*, "complete")

```
sapply(numeric_cols, function(param) {
  pairwise.t.test(
    x = df_training_one_hot[, param],
    g = df_training_one_hot$target,
    pool.sd = FALSE,
    paired = FALSE,
    alternative = "two.sided"
  )$p.value
}) |> sort()
```

```
##          nox          age          dis          indus          tax          lstat
## 1.486824e-70 3.953661e-52 1.762618e-48 7.522700e-48 2.028465e-45 4.663092e-26
##          zn          medv          ptratio          rm
## 1.545946e-21 3.868621e-09 4.851822e-08 1.036364e-03
```

```
sapply(numeric_cols, function(param) {
  pairwise.wilcox.test(
    x = df_training_one_hot[, param],
    g = df_training_one_hot$target,
    pool.sd = FALSE,
    paired = FALSE,
    alternative = "two.sided"
  )$p.value
}) |> sort()
```

```
##          nox          dis          age          indus          tax          lstat
```

```
## 1.505559e-59 7.713151e-46 4.570642e-44 1.169101e-40 8.311193e-38 4.704275e-25
##          zn          medv          ptratio          rm
## 1.999127e-24 4.781087e-18 1.305775e-14 1.331368e-04
```

```
model_dendo <- glm(target ~ radq_hi + chas + lstat + indus + age, binomial(link = "logit"), data=df_train)
summary(model_dendo)
```

```
##
## Call:
## glm(formula = target ~ radq_hi + chas + lstat + indus + age,
##      family = binomial(link = "logit"), data = df_training_one_hot)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.090764   0.560770  -9.078  < 2e-16 ***
## radq_hi      3.827787   0.576717   6.637 3.20e-11 ***
## chas1        0.266750   0.554497   0.481  0.6305
## lstat        0.003477   0.028225   0.123  0.9020
## indus        0.061046   0.025708   2.375  0.0176 *
## age          0.050076   0.008265   6.059 1.37e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 645.88  on 465  degrees of freedom
## Residual deviance: 297.52  on 460  degrees of freedom
## AIC: 309.52
##
## Number of Fisher Scoring iterations: 6
```

Model Using Quasi-Logit

Model comparison

```
library(vcdExtra)
library(pscl)
#models <- list(model_full, model_full_log, backward_model, backward_log_model, model_corr, model_pca, model_dendo)

stats <- LRstats(model_full, model_pca, model_dendo, model1, model2, model3, pj_model1, pj_model2, pj_model3)

stats$McFaddenR2 <- NA
stats$Accuracy <- NA
stats$Precision <- NA
#stats$Recall <- NA
stats$Sensitivity <- NA
stats$Specificity <- NA
stats$F1_score <- NA
stats$AUC <- NA
stats$CV_est_predict_err <- NA
stats$CV_adj_est <- NA
```

```

enhanceEvaluationMetrics <- function(df, model_name) {
  model <- get(model_name)

  if (model_name == "model_full_log" | model_name == "backward_log_model") {
    model_data <- df_training_1h_log
  } else if (model_name == "model_pca") {
    model_data <- df_training_one_hot_pca
  } else if (model_name == "model_full" | model_name == "backward_model" | model_name == "model_corr" |
  } else if (model_name == "pj_model1" | model_name == "pj_model2" | model_name == "pj_model3") {
    model_data <- pj_crime_training_df
  } else if (model_name == "zr_model_a" | model_name == "zr_model_b" | model_name == "zr_model_c") {
    model_data <- crime_training_zr
  } else {
    model_data <- crime_training_df
  }

  df[model_name, "McFaddenR2"] <- pR2(model)["McFadden"]

  pred_probs <- predict(model, type = "response")

  pred_probs_factor <- as.factor(ifelse(pred_probs > 0.5, 1, 0))
  conf_matrix <- confusionMatrix(pred_probs_factor, as.factor(model_data$target))
  df[model_name, "Accuracy"] <- conf_matrix$overall['Accuracy']
  df[model_name, "Precision"] <- conf_matrix$byClass['Precision']
  #df[model_name, "Recall"] <- conf_matrix$byClass['Recall']
  df[model_name, "F1_score"] <- conf_matrix$byClass['F1']
  df[model_name, "Sensitivity"] <- conf_matrix$byClass["Sensitivity"]
  df[model_name, "Specificity"] <- conf_matrix$byClass["Specificity"]

  #roc_model <- roc(as.factor(model_data$target), pred_probs)
  #plot(roc_model, main = "ROC Curve using pROC", col = "red", lwd = 2)
  # roc_auc not working, so use MLmetrics
  df[model_name, "AUC"] <- MLmetrics::AUC(y_true = model_data$target, y_pred = pred_probs)

  # Cross-Validation using 10 folds
  cv_result <- boot::cv.glm(model_data, model, K= 10)
  df[model_name, "CV_est_predict_err"] <- cv_result$delta[1]
  df[model_name, "CV_adj_est"] <- cv_result$delta[2]

  return(df)
}

# Loop through the list of models and update the dataframe for each
for (model_name in rownames(stats)) {

  stats <- enhanceEvaluationMetrics(stats, model_name)
}

## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2

```

```
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
```

```
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
## fitting null model for pseudo-r2
```

```
stats
```

```
## Likelihood summary table:
```

##	AIC	BIC	LR	Chisq	Df	Pr(>Chisq)	McFaddenR2	Accuracy	Precision						
## model_full	208.97	266.99	180.97	452		1	0.71981	0.92275	0.92405						
## model_pca	303.25	332.26	289.25	459		1	0.55216	0.85193	0.83871						
## model_dendo	309.52	334.39	297.52	460		1	0.53935	0.85837	0.84615						
## model1	214.22	297.11	174.22	446		1	0.73025	0.93348	0.92917						
## model2	202.45	252.18	178.45	454		1	0.72372	0.92704	0.92116						
## model3	288.17	317.17	274.17	459		1	0.57551	0.85193	0.84146						
## pj_model1	218.05	271.92	192.05	453		1	0.70266	0.91631	0.90909						
## pj_model2	176.81	230.68	150.81	453		1	0.76651	0.93562	0.92946						
## pj_model3	233.35	291.37	205.35	452		1	0.68206	0.89270	0.88163						
## zr_model_a	231.73	260.73	217.73	459		1	0.66290	0.87339	0.85317						
## zr_model_b	224.25	269.84	202.25	455		1	0.68686	0.91631	0.91250						
## zr_model_c	223.64	265.08	203.64	456		1	0.68470	0.91202	0.89837						
##	Sensitivity		Specificity		F1_score		AUC	CV_est_predict_err							
## model_full	0.92405		0.92140		0.92405		0.97771	0.066949							
## model_pca	0.87764		0.82533		0.85773		0.94054	0.099136							
## model_dendo	0.88186		0.83406		0.86364		0.93288	0.103471							
## model1	0.94093		0.92576		0.93501		0.97872	0.065748							
## model2	0.93671		0.91703		0.92887		0.97756	0.065880							
## model3	0.87342		0.82969		0.85714		0.94723	0.102254							
## pj_model1	0.92827		0.90393		0.91858		0.97376	0.071437							
## pj_model2	0.94515		0.92576		0.93724		0.98434	0.054832							
## pj_model3	0.91139		0.87336		0.89627		0.96901	0.080245							
## zr_model_a	0.90717		0.83843		0.87935		0.96814	0.081149							
## zr_model_b	0.92405		0.90830		0.91824		0.97059	0.072095							
## zr_model_c	0.93249		0.89083		0.91511		0.96984	0.075726							
##	CV_adj_est														
## model_full	0.066409														
## model_pca	0.098969														
## model_dendo	0.103286														
## model1	0.065230														
## model2	0.065130														
## model3	0.101942														
## pj_model1	0.070969														
## pj_model2	0.054372														
## pj_model3	0.079707														
## zr_model_a	0.080866														
## zr_model_b	0.071756														
## zr_model_c	0.075303														

Model Evaluation:

Evaluate model performance and select the best model based on multiple criteria.

Evaluation Metrics: Accuracy: $(TP + TN) / (TP + TN + FP + FN)$ Precision: $TP / (TP + FP)$ Recall (Sensitivity): $TP / (TP + FN)$ Specificity: $TN / (TN + FP)$ F1 Score: $2 * (Precision * Recall) / (Precision + Recall)$ AUC-ROC Curve: Evaluate model discrimination.

For logistic regression, the “prediction error” is the mean squared error (difference between the predicted probabilities and the actual outcomes).

MODEL SELECTION

Checking the Model’s Conditions

We will examine the following key conditions for fitting a logistic model:

1. dependent variable is binary
2. large enough sample
3. observations are independent, not matched
4. independent (predictor) variables do not correlate too strongly with each other
5. linearity of independent variables and log odds
6. no outliers in data

As a result, Model 2: Stepwise Logistic Regression was selected as the best binary logistic regression model due to achieving a trade-off between model simplicity and performance. The optimal model not only has to perform excellently in prediction but also be interpretable and extendable to new data. Although a complicated model will provide marginally better performance, it will overfit if too many extraneous parameters are introduced. Therefore, we selected Model 2 since it achieves a balance between parsimony and performance as it retains the strongest predictors only. Stepwise logistic regression (direction = “both”) reduced the model by selecting the optimal subset of features and hence giving a more efficient and stable model.

To compare Model 2, we employed a range of statistical measures that assess different aspects of performance. Akaike Information Criterion (AIC), on which the model fit will be judged, was 176.81, reflecting high performance compared to other models. Area Under the Curve (AUC) of 0.9843 reflects that the model was working very well to distinguish between the two classes. Accuracy (0.9372) also reflects that the model is correctly classifying the majority of the cases and classification error rate (0.0644) is low, reflecting high reliability.

Accuracy, recall or sensitivity, and specificity also act to define Model 2’s performance. The accuracy of the model at 0.9451 means that whenever it is positive, it is correct 94.51% of the time. The sensitivity at 0.9356 means that it identifies 93.56% of real positive cases correctly, and the specificity at 0.9295 means 92.95% of non-positive cases are identified correctly. The 0.9372 F1 score as a compromise between the recall and precision measures how good the model is at predicting things correctly. The confusion matrix also confirms that the false positives and false negatives are zero, yet again proving correct.

Lastly, Model 2 was chosen since it is the optimal compromise among predictiveness, interpretability, and model fit. It has very high AUC value, good specificity and sensitivity, and very low classification error, and hence very dependable in binary classification. Stepwise selection of removing extraneous predictors avoids overfitting but not super-predictiveness. Due to its low AIC, good performance on a range of measures of evaluation, and relatively well-scaled set of predictors, Model 2 is optimal to be utilized in this analysis.

Apply the Best Model to Evaluation Data

Once the best model is selected, we use it for prediction on `crime_evaluation_df`.

```

# Apply same transformations as before
# Apply the same transformations to the evaluation dataset as used in training
crime_evaluation_df$log_tax <- log(crime_evaluation_df$tax)
crime_evaluation_df$log_lstat <- log(crime_evaluation_df$lstat)
crime_evaluation_df$log_medv <- log(crime_evaluation_df$medv)

# If scaling or categorical transformations were applied, replicate them here
crime_evaluation_df$age_group <- cut(crime_evaluation_df$age, breaks = c(0, 35, 70, 100), labels = c("y
crime_evaluation_df$dis_scaled <- scale(crime_evaluation_df$dis)

# Predict probabilities using model1
eval_pred_prob <- predict(pj_model2, newdata = crime_evaluation_df, type = "response")

# Convert probabilities to binary class (0 or 1) using a threshold of 0.5
eval_pred_class <- ifelse(eval_pred_prob > 0.5, 1, 0)

# Add predictions to the evaluation data
crime_evaluation_df$predicted_prob <- eval_pred_prob
crime_evaluation_df$predicted_class <- eval_pred_class

# use the predicted_class and predicted_prob values as your final output

head(crime_evaluation_df[, c("predicted_prob", "predicted_class")])

```

```

## # A tibble: 6 x 2
##   predicted_prob predicted_class
##           <dbl>           <dbl>
## 1         0.0171             0
## 2         0.812             1
## 3         0.710             1
## 4         0.962             1
## 5         0.279             0
## 6         0.783             1

```

This ensures that Model 2: Stepwise Logistic Regression is applied consistently and provides accurate predictions on the evaluation dataset.