# Final Project EDA

## Fall 2025

```
libraries <- read_csv("../data/libraries.csv")
```

```
## Rows: 6832 Columns: 28
## -- Column specification ---------------------------------------------------------
## Delimiter: ","
## chr  (2): state, Locale
## dbl (26): Service Area Population, Total Circulation, Percentage of Children...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
numeric_libaries <- libraries %>%
  select(where(is.numeric))
full_correlation_matrix <- cor(numeric_libaries)
```

Looking at the correlation plot, we see that there are many variables, with some of them having very little correlation to our variable of interest, Total Circulation. Let's drop some.

```
y_cor_matrix<- full_correlation_matrix[, "Total Circulation", drop = FALSE]
y_cor_matrix[order(abs(y_cor_matrix[, 1]), decreasing = TRUE), , drop = FALSE]
```

```
##                                   Total Circulation
## Total Circulation                        1.00000000
## Physical Visits                          0.89025206
## Local Revenue ($)                        0.85594139
## Hours/Year                               0.83766165
## Service Area Population                  0.81046353
## Registered Users                         0.78177823
## Branch Library                           0.77990997
## Internet Computers                       0.71332256
## Children's Program Attendance            0.70207707
## Computer Uses                            0.69813313
## Children's Programs                      0.66688666
## Reference Transactions                   0.62098179
## Young Adult Programs                     0.56534446
## Adult Program Attendance                 0.55176687
## Adult Programs                           0.51878890
## State Revenue ($)                        0.49280567
## Federal Revenue ($)                      0.36914870
## Bookmobiles                              0.36633538
## Other Revenue ($)                        0.32529456
## Central Library                         -0.30214050
## General Interest Programs                0.28204319
```

```
## Inter-library Loans from Other Library                        0.27570674
## General Interest Program Attendance                           0.25989250
## Wireless Sessions                                             0.24455640
## Inter-library Loans to Other Library                          0.24005347
## Percentage of Children's Material Circulation                 0.04351038
```

There are too many variables for a simple correlation matrix, and looking at this list, it seems nearly all of the numeric variables have some correlation to Total Circulation.

However, this does not show us potential collinearity, for that we can fit a simple model and check the VIF:

```r
model <- lm(`Total Circulation` ~ ., data = numeric_libaries)
vif(model)
```

```
##                       `Service Area Population`
##                                        6.594328
## `Percentage of Children's Material Circulation`
##                                        1.034565
##                              `Central Library`
##                                        1.183601
##                               `Branch Library`
##                                       21.020578
##                                     Bookmobiles
##                                        1.278856
##                             `Internet Computers`
##                                        4.810188
##                                  `Computer Uses`
##                                        2.862216
##                               `Wireless Sessions`
##                                        1.088883
##                             `Children's Programs`
##                                       11.928597
##                             `Young Adult Programs`
##                                        6.529996
##                                  `Adult Programs`
##                                       11.766650
##                   `Children's Program Attendance`
##                                        8.796485
##                        `Adult Program Attendance`
##                                        6.874101
##           `General Interest Program Attendance`
##                                        1.681073
##                   `General Interest Programs`
##                                        2.413284
##                              `Local Revenue ($)`
##                                        7.826613
##                              `State Revenue ($)`
##                                        1.711105
##                            `Federal Revenue ($)`
##                                        1.745681
##                              `Other Revenue ($)`
##                                        3.957116
##                                    `Hours/Year`
##                                       31.103517
```

```
##                             `Physical Visits`
##                                   15.167099
##                        `Reference Transactions`
##                                    4.545317
##                            `Registered Users`
##                                    5.614272
##           `Inter-library Loans to Other Library`
##                                    5.856800
##        `Inter-library Loans from Other Library`
##                                    6.043276
```
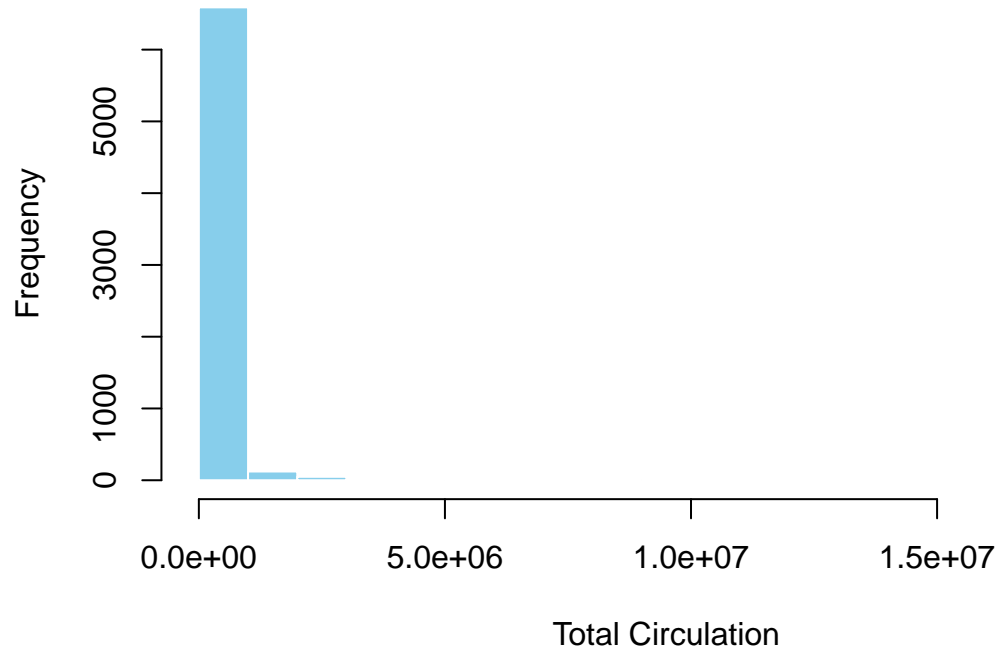
This shows that there is a generally high level of collinearity, and we should keep this in mind for the analysis.

```r
for (col in names(numeric_libaries)) {
  hist(
    numeric_libaries[[col]],
    main = paste("Histogram of", col),
    xlab = col,
    col = "skyblue",
    border = "white"
  )
}
```

## Histogram of Service Area Population



3

# Histogram of Total Circulation

# Histogram of Percentage of Children's Material Circulation



Percentage of Children's Material Circulation

# Histogram of Central Library

# Histogram of Branch Library

**Histogram of Bookmobiles**

# Histogram of Internet Computers

# Histogram of Computer Uses

**Histogram of Wireless Sessions**

Frequency

Wireless Sessions

# Histogram of Children's Programs



Children's Programs

**Histogram of Young Adult Programs**

# Histogram of Adult Programs



Frequency

Adult Programs

# Histogram of Children's Program Attendance



Children's Program Attendance

# Histogram of Adult Program Attendance

# Histogram of General Interest Program Attendance



General Interest Program Attendance

# Histogram of General Interest Programs

# Histogram of Local Revenue ($)

# Histogram of State Revenue ($)

# Histogram of Federal Revenue ($)

# Histogram of Other Revenue ($)

# Histogram of Hours/Year

# Histogram of Physical Visits

# Histogram of Reference Transactions

# Histogram of Registered Users

# Histogram of Inter−library Loans to Other Library

## Histogram of Inter−library Loans from Other Library



The histograms of the numeric data show that many of these have extreme outliers, so let's remove them:

library(dplyr)

```r
# Function to remove outliers using the IQR rule
remove_outliers <- function(x) {
  if (is.numeric(x) && length(unique(na.omit(x))) > 2) {  # skip binary
    q1 <- quantile(x, 0.25, na.rm = TRUE)
    q3 <- quantile(x, 0.75, na.rm = TRUE)
    iqr <- q3 - q1
    lower <- q1 - 1.5 * iqr
    upper <- q3 + 1.5 * iqr
    x[x < lower | x > upper] <- NA  # mark outliers
  }
  return(x)
}

# Columns to ignore
ignore_cols <- c("Bookmobiles", "Branch Library")

# Apply only to numeric, non-binary, and not ignored columns
libraries_no_outliers <- libraries %>%
  mutate(across(
    .cols = where(is.numeric) & !all_of(ignore_cols),
    .fns = remove_outliers
  )) %>%
  drop_na()
```

```
libraries_no_outliers_numeric <- libraries_no_outliers %>% select(where(is.numeric))
for (col in names(libraries_no_outliers_numeric)) {
  hist(
    libraries_no_outliers_numeric[[col]],
    main = paste("Histogram of", col),
    xlab = col,
    col = "skyblue",
    border = "white"
  )
}
```

**Histogram of Service Area Population**

## Histogram of Total Circulation

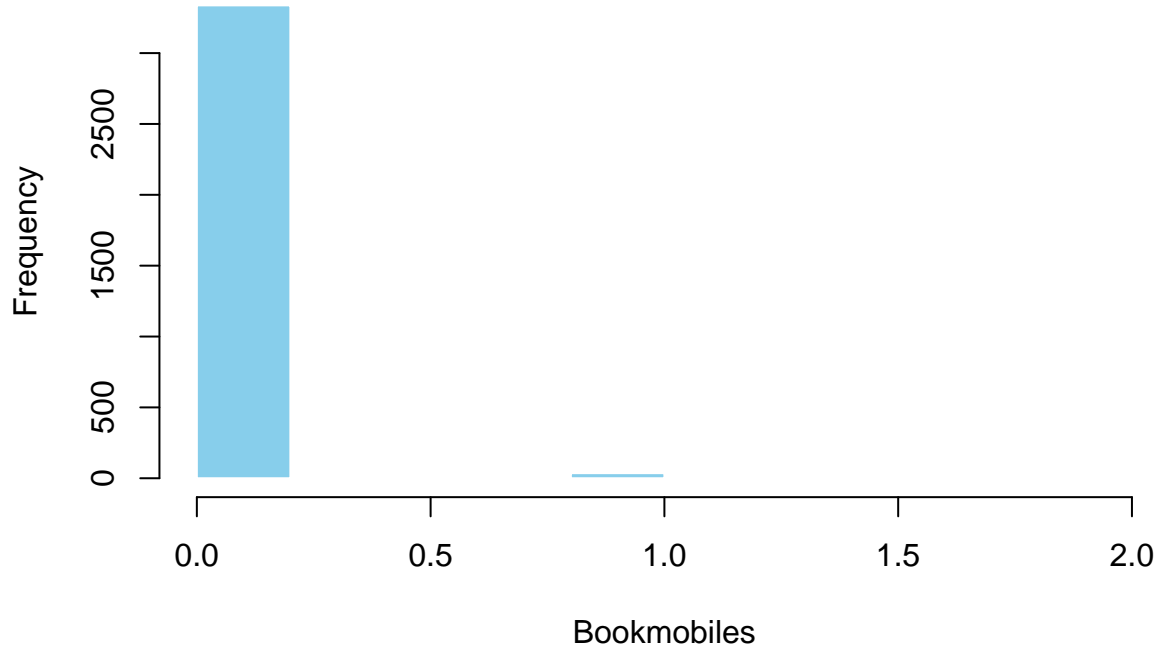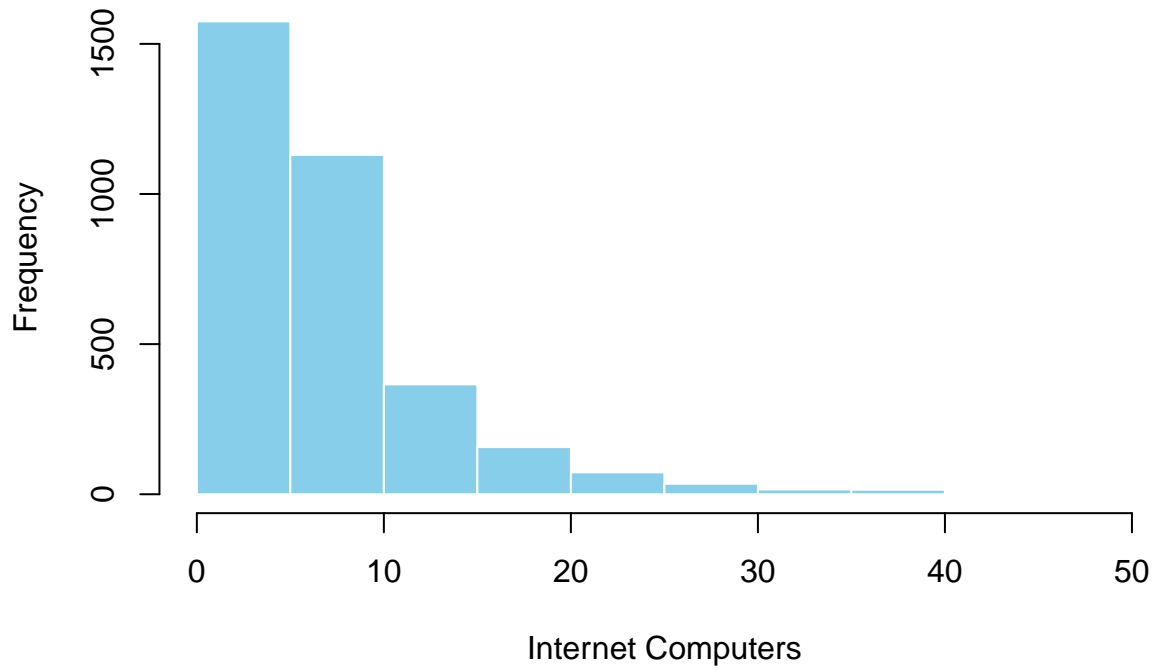# Histogram of Percentage of Children's Material Circulation



Percentage of Children's Material Circulation

# Histogram of Central Library

# Histogram of Branch Library

# Histogram of Bookmobiles

# Histogram of Internet Computers

# Histogram of Computer Uses

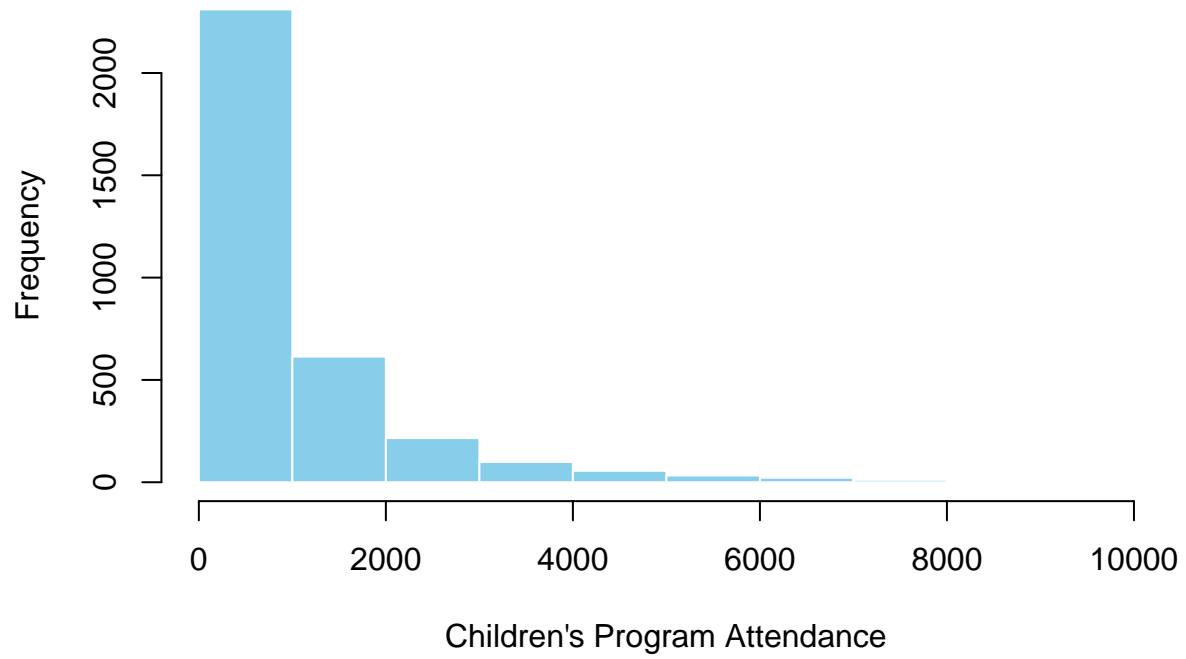# Histogram of Wireless Sessions

# Histogram of Children's Programs

**Histogram of Young Adult Programs**

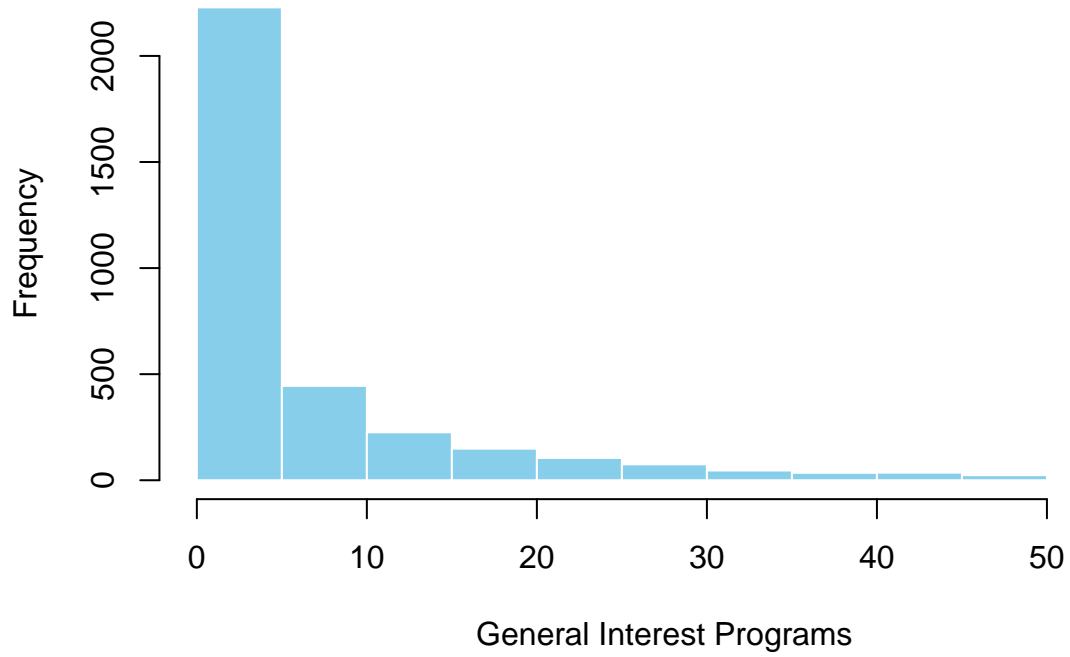Histogram of Adult Programs

# Histogram of Children's Program Attendance

# Histogram of Adult Program Attendance

# Histogram of General Interest Program Attendance



General Interest Program Attendance

# Histogram of General Interest Programs
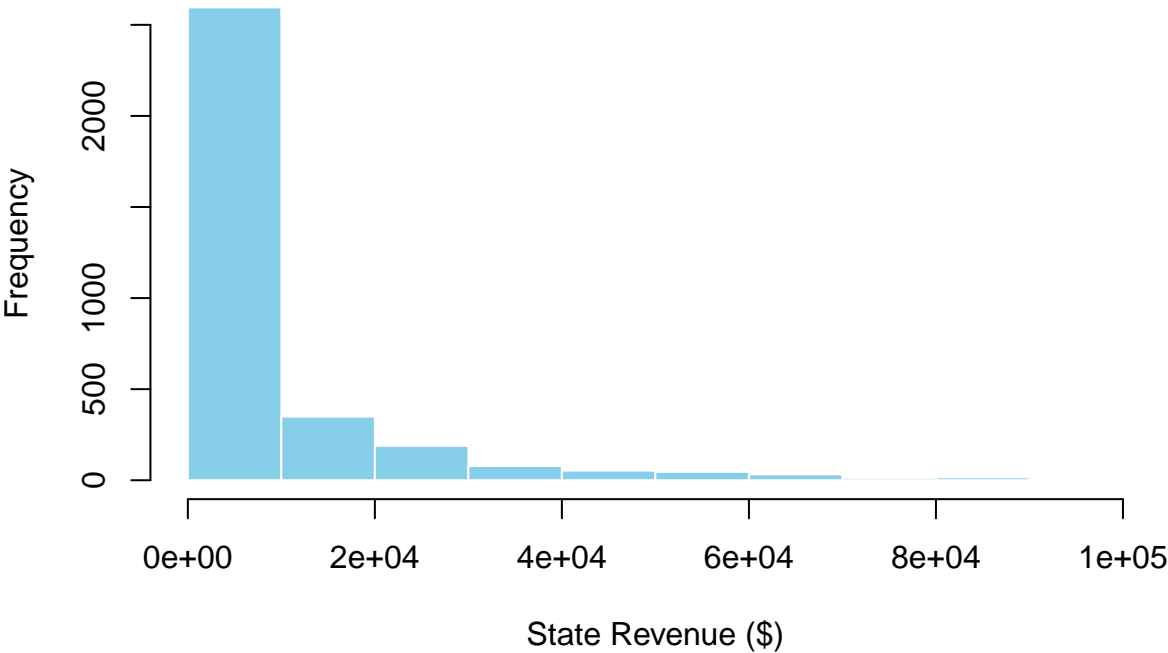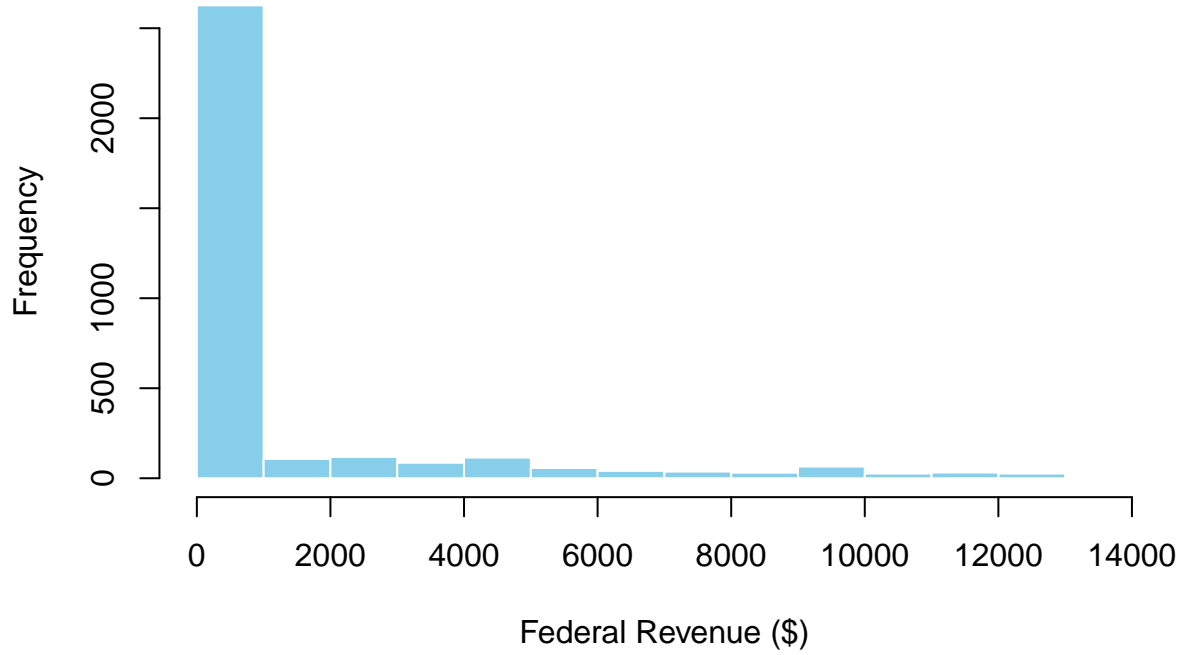


General Interest Programs

# Histogram of Local Revenue ($)

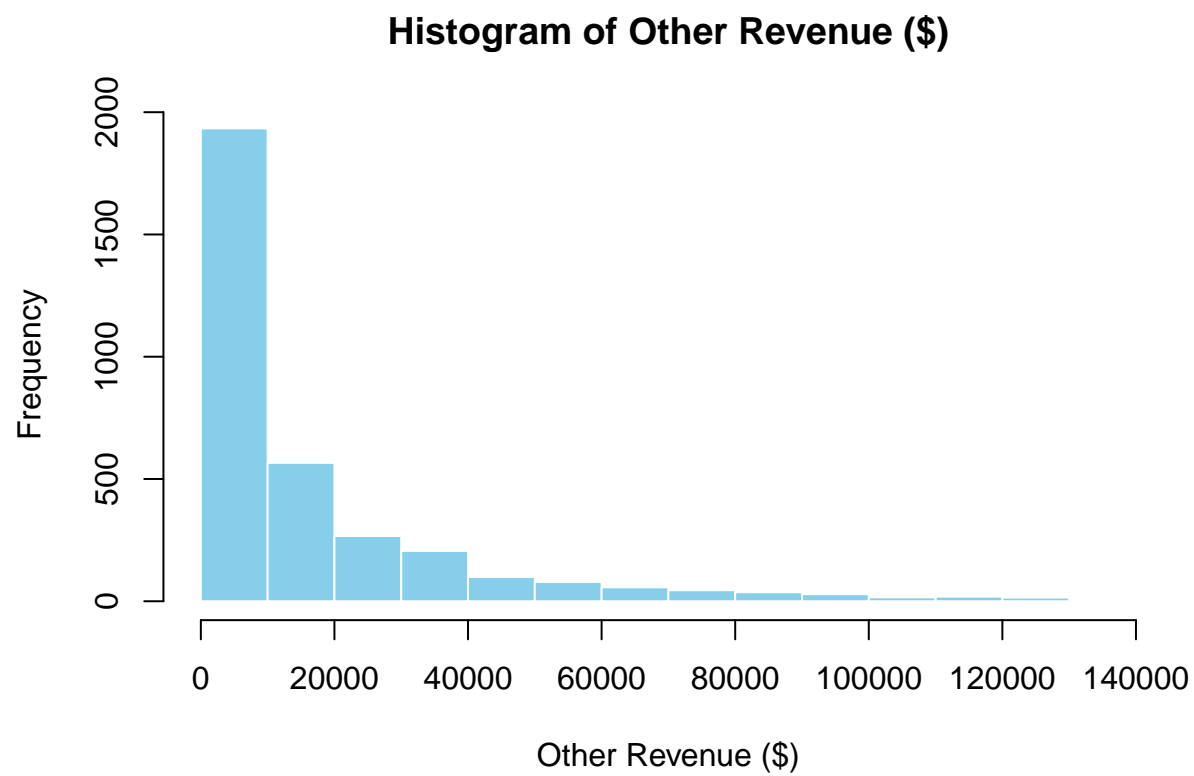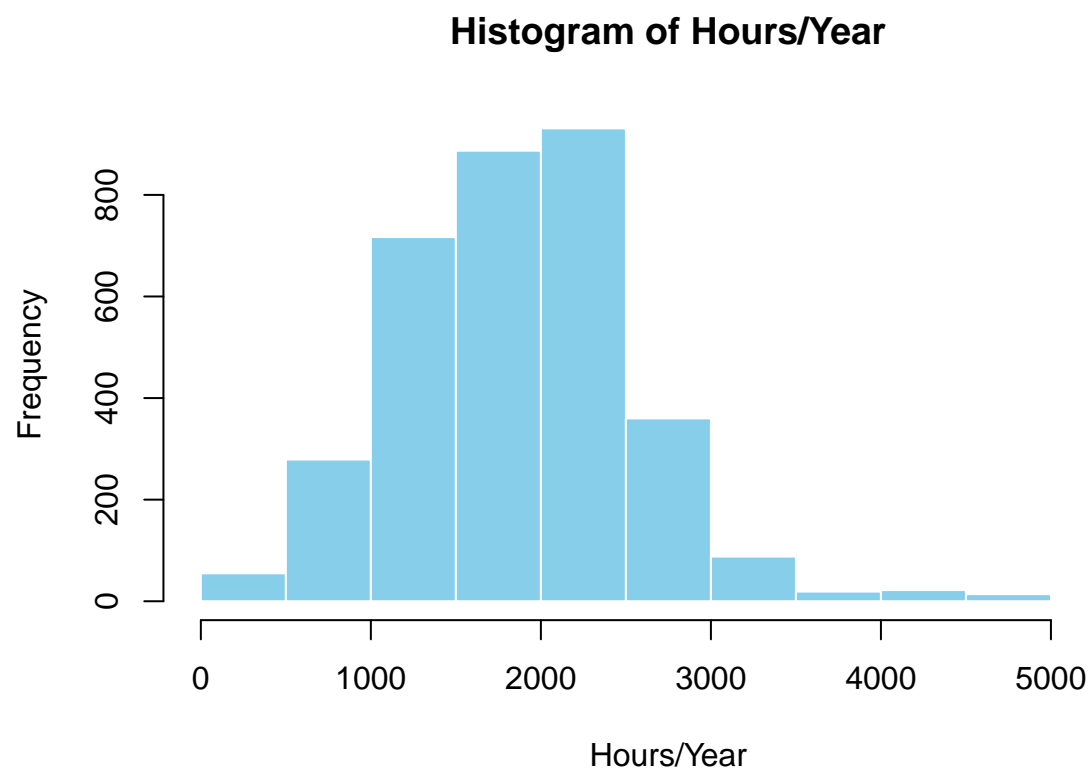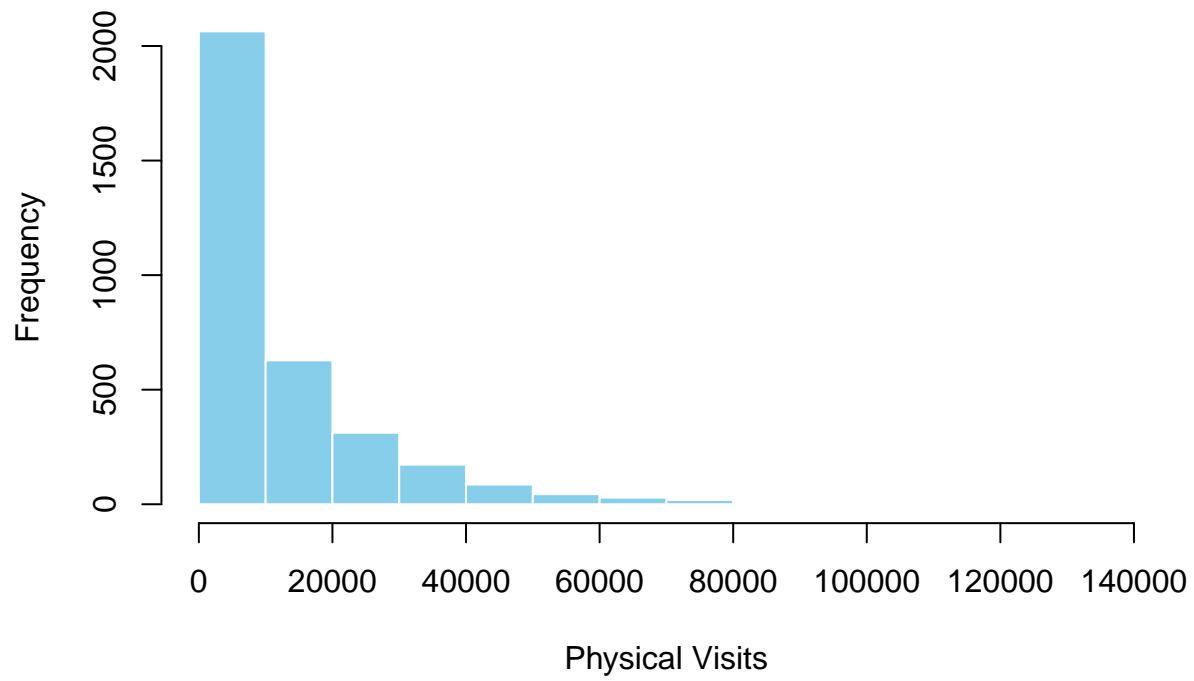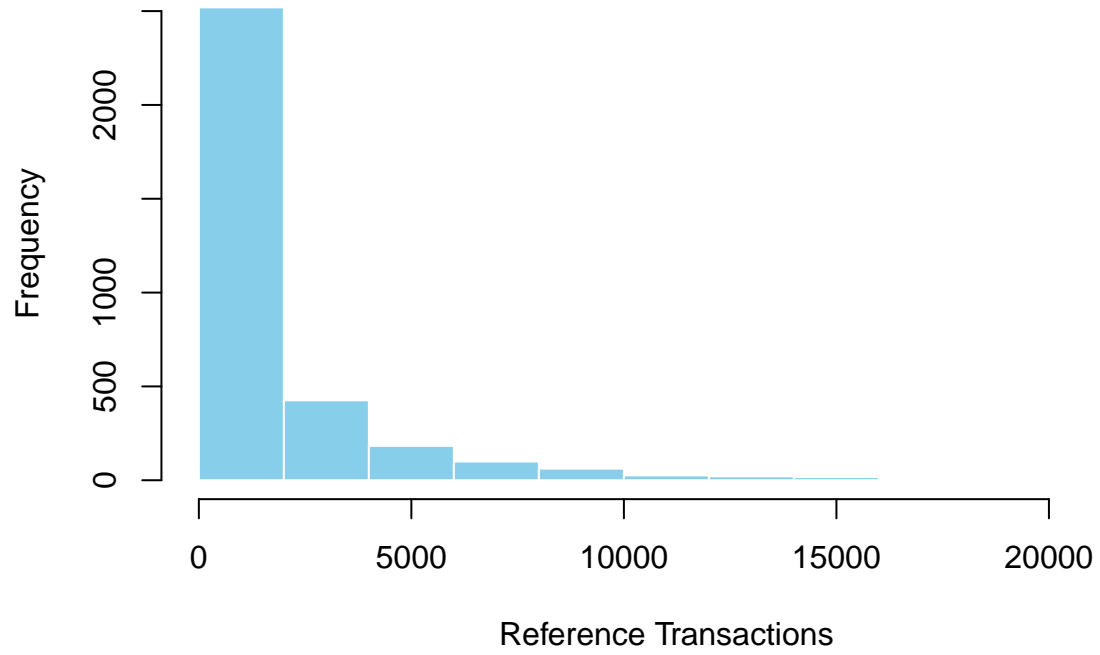# Histogram of State Revenue ($)

## Histogram of Federal Revenue ($)

# Histogram of Other Revenue ($)
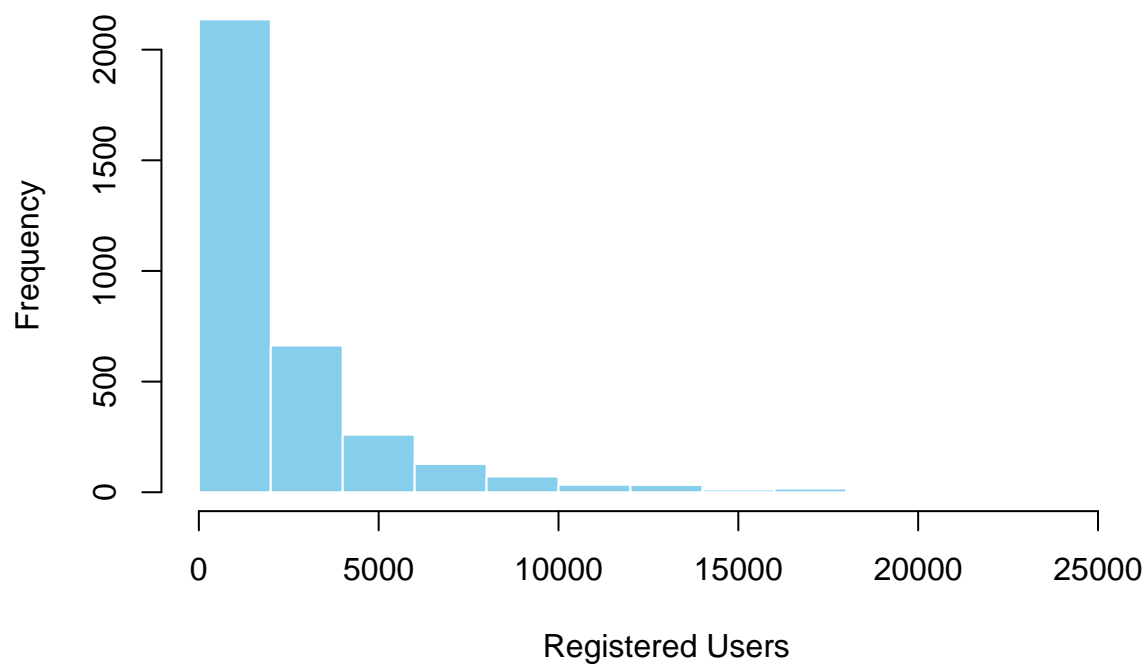
# Histogram of Hours/Year

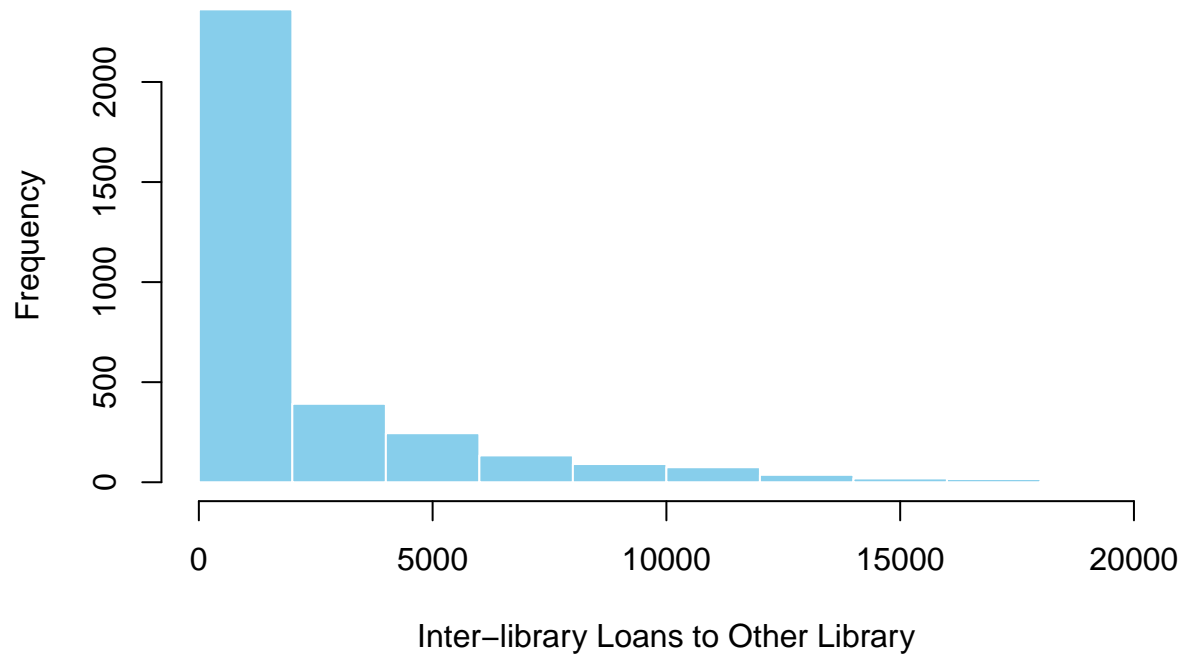**Histogram of Physical Visits**
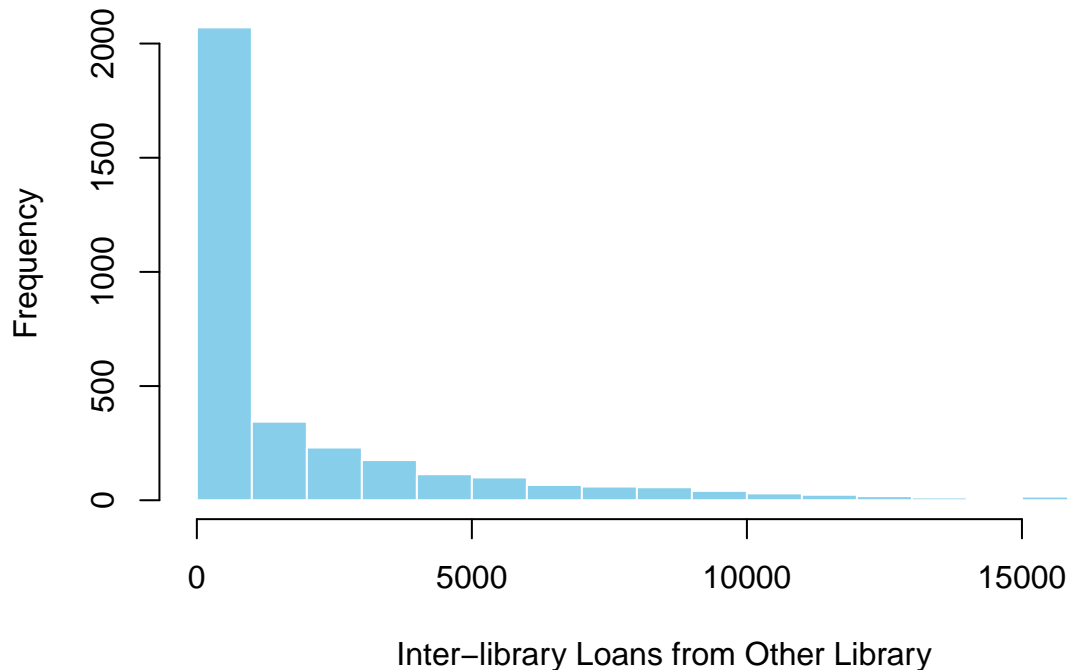
**Histogram of Reference Transactions**

# Histogram of Registered Users

# Histogram of Inter−library Loans to Other Library

**Histogram of Inter−library Loans from Other Library**



With outliers removed, the data is more reasonable, but skewed to the right. With this being said, we should consider Box-Cox transformations when making the models.
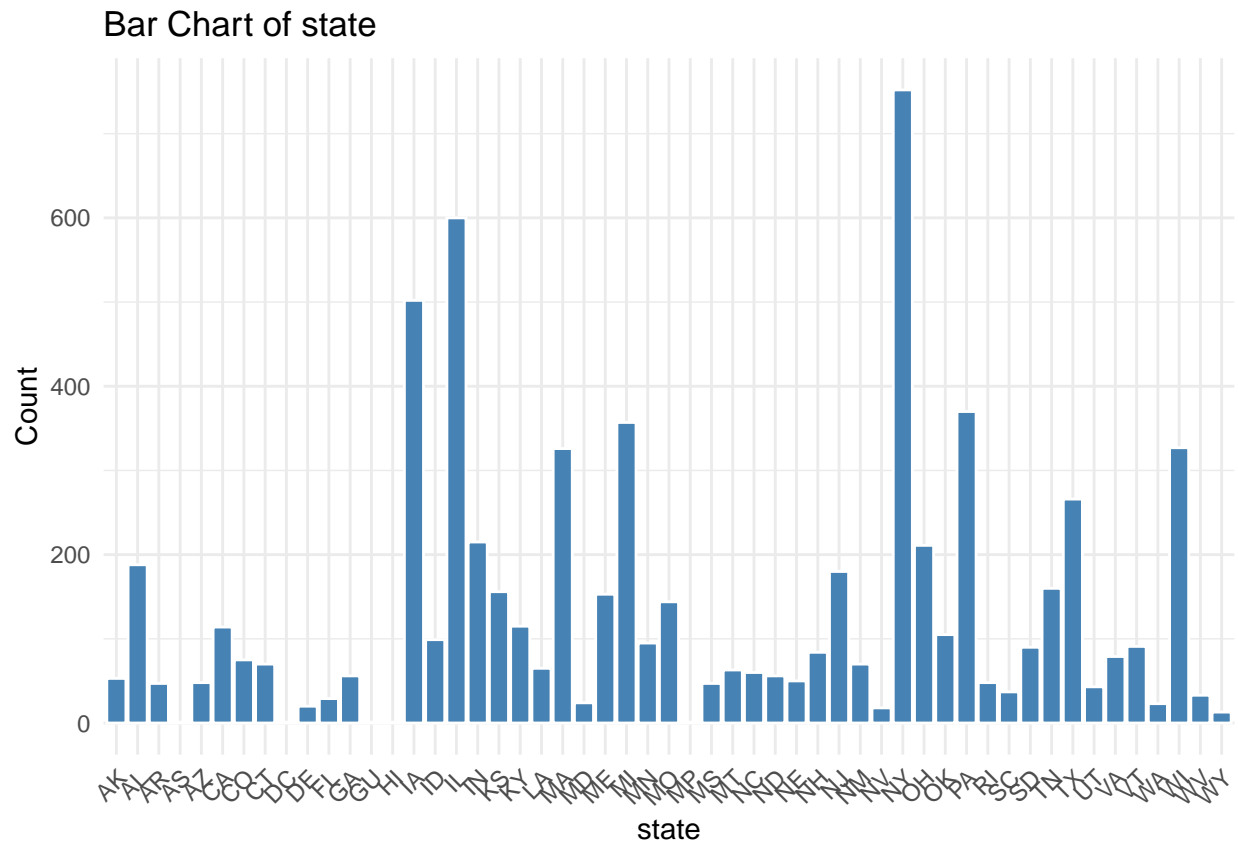
```r
categorical_data <- libraries %>% select(!where(is.numeric))

# Loop through each categorical column and plot
for (col in names(categorical_data)) {
  ggplot(categorical_data, aes_string(x = col)) +
    geom_bar(fill = "steelblue", color = "white") +
    theme_minimal() +
    labs(
      title = paste("Bar Chart of", col),
      x = col,
      y = "Count"
    ) +
    theme(
      axis.text.x = element_text(angle = 45, hjust = 1)
    ) -> p

  print(p)
}
```
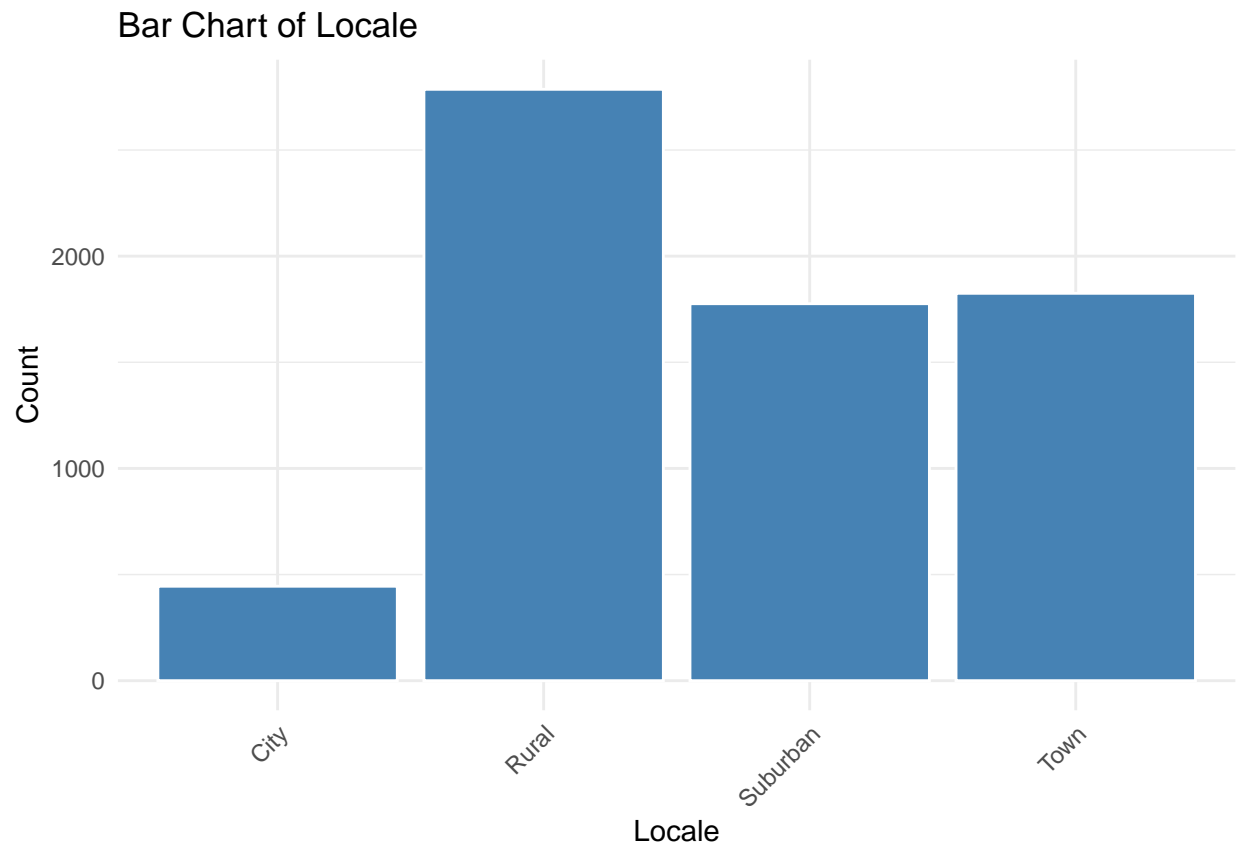
```
## Warning: `aes_string()` was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with `aes()`.
## i See also `vignette("ggplot2-in-packages")` for more information.
## This warning is displayed once every 8 hours.
```

```
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

Bar Chart of state

## Bar Chart of Locale



With outliers removed, we see an overrepresentation of NY, IL, and HI, which means we should be careful making generalizations about the data to the entire nation. The locale is not as surprising, but we should be aware of the fact that the city is less represented than other types of locale.

Because this is a regression task, we should be considering Ridge and Lasso multi variable regression models to minimize the effects of the high collinearity.