

## NYPD Project

Ryan Ruff

2024-11-13

```
## — Attaching go
```

```
## ✓ ggplot2 3.5.1 ✓ tibble 3.2.1
## ✓ lubridate 1.9.3 ✓ tidyr 1.3.1
## ✓ purrr 1.0.2
## --- Conflicts (use tidyverse_conflicts() --
## # dplyr::filter() masks stats::filter()
## # dplyr::lag() masks stats::lag()
## # I use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors

library(lubridate)
library(ggplot2)
library(caret)
```

```
## Loading require
##
## Attaching packa
##
## The following o
```

```
library(rsample)

# Import the dataset
nypd_data <- read_csv("NYPD_Shooting_Incident_Data__His

## Rows: 28562 Columns: 21
```

```
## Delimiter: ",
## chr (12): 0C
```

```
## time (1): OCCUR_TIME
##
```

```
# Inspect the data
glimpse(npyd_data)

## Rows: 28,562
## Columns: 21
## $ INCIDENT_KEY      <dbl> 244608249, 247542571, 84967535, 202853370, 270...
## $ OCCUR_DATE         <chr> "5/5/2022", "7/4/2022", "5/27/2012", "9/24/201...
## $ OCCUR_TIME         <time> 00:10:00, 22:20:00, 19:35:00, 21:00:00, 21:00:...
## $ BORO               <chr> "MANHATTAN", "BROOKX", "QUEENS", "BROOKX", "BROO...
## $ LOC OF OCCUR DESC  <chr> "INSIDE", "OUTSIDE", "NA_NA", "NA_NA", "NA_NA", ...
```

```
## $ LOC_CLASSFCTN
## $ LOCATION_DESC
## $ STATISTICAL_M
```

```
## $ PERP_RACE      <chr> "BLACK", "(null)", NA, "UNKNOWN", "BLACK", NA,...
## $ VIC_AGE_GROUP  <chr> "25-44", "18-24", "18-24", "25-44", "25-44", "...
## $ VIC_SEX        <chr> "M", "M", "M", "M", "M", "M", "M", "M", "...
## $ VIC_RACE        <chr> "BLACK", "BLACK", "BLACK", "BLACK", "BLACK", "...
## $ X_COORD_CD      <dbl> 986500, 1016802, 1048632, 1014493, 1009149, 99,...
## $ X_COORD_CD      <dbl> 214231.0, 250588.0, 190262.0, 242565.0, 190104...
## $ Latitude        <dbl> 40.75469, 40.85440, 40.71063, 40.83242, 40.688...
## $ Longitude       <dbl> -73.99350, -73.88233, -73.76777, -73.89071, -7...
## $ LonLat          <chr> "POINT (-73.9935 40.754692)", "POINT (-73.8823...
```

```
# Convert date and time columns to appropriate formats
npyd.data <- npyd.data %>%
  mutate(OCURR_DATE = mdy(OCURR_DATE),
         OCURR_TIME = hms(OCURR_TIME),
         YEAR = year(OCURR_DATE))
```

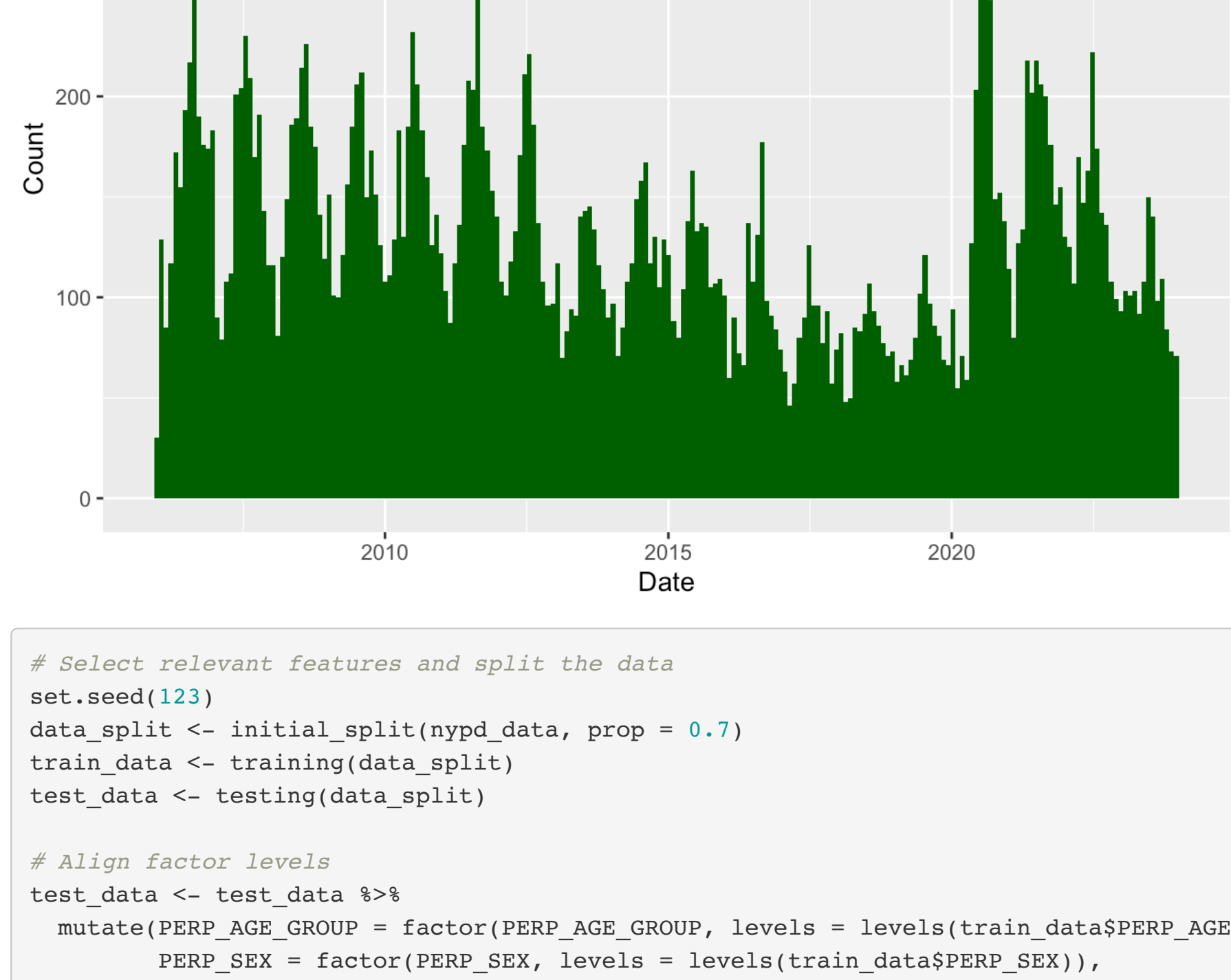
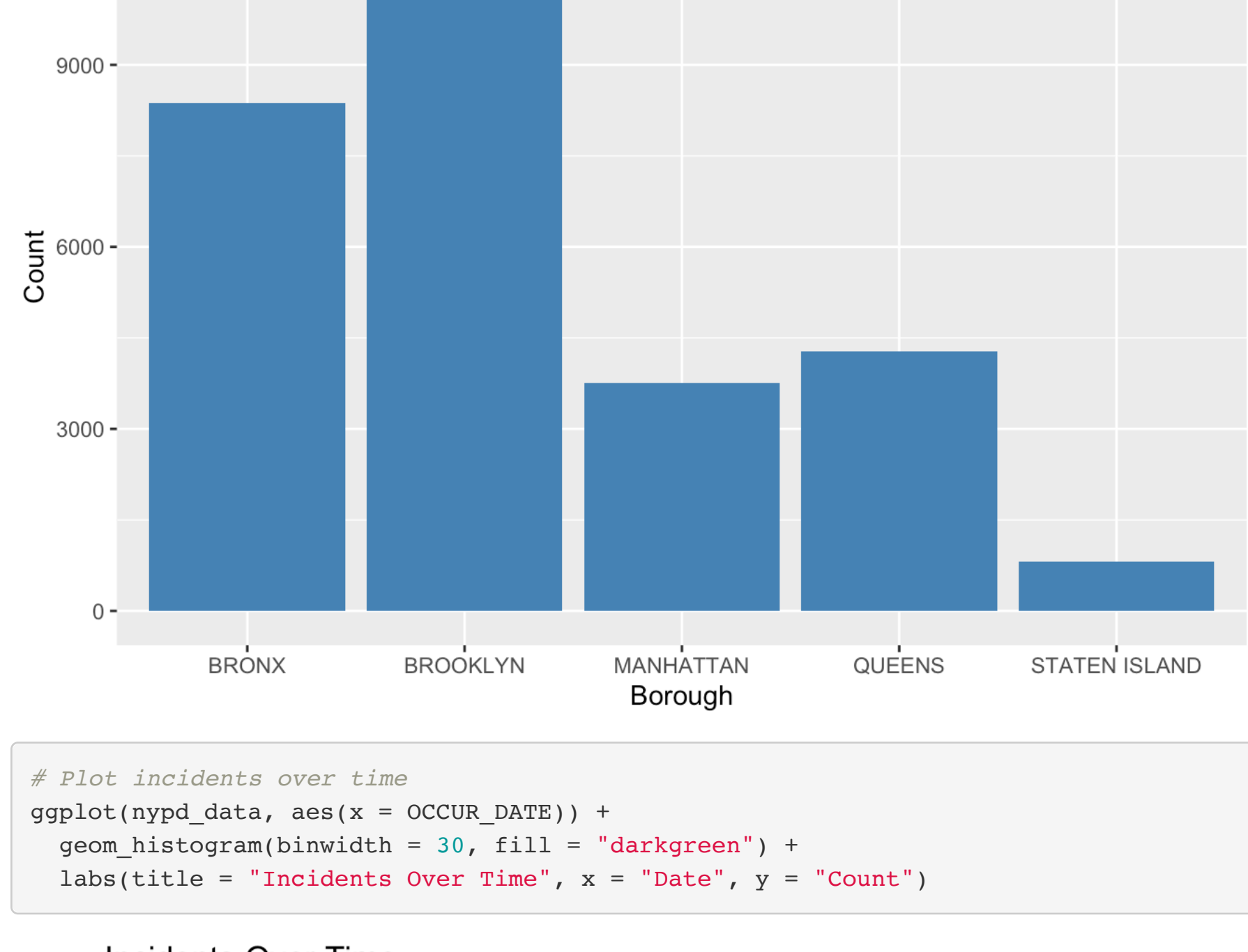
```
nypd_data <- nypd_data %>%
  filter(PERP_AGE_GROUP, PERP_SEX, PERP_RACE, .direction = "downup") %>%
  filter(VIC_AGE_GROUP, VIC_SEX, VIC_RACE, .direction = "downup")

# Ensure the relevant features are factors with the same levels
nypd_data <- nypd_data %>%
  mutate(PERP_AGE_GROUP = factor(PERP_AGE_GROUP),
```

```
VIC_SEX = factor(VIC_SEX),
FATAL = ifelse(STATISTICAL_MURDER_FLAG == "TRUE", 1, 0))

# Plot the number of incidents by borough
ggplot(nypd_data, aes(x = BORO)) +
  geom_bar(fill = "#e69d00") +
```

Number of Incidents by Borough



```
FATAL = factor(FATAL, levels = c(0, 1))

# Train a logistic regression model
model <- train(FATAL ~ BORO + PERP_AGE_GROUP + PERP_AGE_GROUP^2,
               data = train_data, method = "glm",
```

```
## regression and your outcome only has two possible values Are you trying to do
## classification? If so, use a 2 level factor as your outcome column.

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type = if (type == :
## prediction from rank-deficient fit; attr(,"non-estim") has doubtful cases
## Warning in predict.lm(object, newdata, se.fit, scale = 1, type = if (type == :
```

[illegible]

~~##~~ Call:

```
## Coefficients: (2 not due to overall scale)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -2.114e+00  1.234e-01 -17.125 < 2e-16
## BOROBOROOKALYN   3.167e-03  4.418e-02  0.072  0.942848
## BOROMANNHAYLN    -1.320e-01  6.131e-02 -2.154  0.031276
## BOROQUEENS        3.957e-03  5.734e-02  0.069  0.944979
## "BOROSTATENS ISLAND"
##    1.189e-01  1.078e-01  1.104  0.269768
##    "PERP_AGE_GROUP18"
##    6.927e-01  1.583e-01  4.375  0.000151
## PERP_AGE_GROUP1020
##    -1.044e+01  2.292e+02 -0.046  0.963668
## PERP_AGE_GROUP1028
##    NA NA NA NA
## PERP_AGE_GROUP18+24"
##    7.748e-01  1.481e-01  5.230  0.69e-07
## PERP_AGE_GROUP224
##    -1.030e+01  3.247e+02 -0.032  0.974669
## PERP_AGE_GROUP25+44"
##    9.294e-01  1.679e-01  6.285  0.27e-10
## PERP_AGE_GROUP25+64"
##    1.194e+01  1.495e-01  7.140  9.36e-13
## PERP_AGE_GROUP65+"
##    1.265e+00  2.872e-01  4.076  0.05e-05
## PERP_AGE_GROUP940
##    1.061e+01  1.451e+02 -0.403  0.941698
## PERP_AGE_GROUPUNKNOWN
##    -4.256e-01  1.205e-01 -3.790  0.000151
## PERP_SEXF
##    -2.620e-01  1.582e-01 -1.403  0.160506
## PERP_SEXM
##    -3.776e-01  1.161e-01 -3.251  0.001151
## PERP_SEXU
##    NA NA NA NA
## VIC_AGE_GROUP1022
##    -1.072e+01  3.247e+02 -0.633  0.973669
## VIC_AGE_GROUP18+24"
##    2.234e-01  7.231e-02  3.048  0.002303
## VIC_AGE_GROUP25+44"
##    4.855e-01  7.160e-02  6.781  1.19e-11
## VIC_AGE_GROUP25+64"
##    7.088e-01  9.191e-02  7.774  6.73e-15
## VIC_AGE_GROUP65+"
##    1.220e+00  1.951e-01  6.254  3.99e-10
## VIC_AGE_GROUPUNKNOWN
##    8.975e-01  3.510e-01  2.466  0.013654
## VIC_SEXM
##    9.674e-04  6.158e-02  0.016  0.987078
## VIC_SEXU
##    -1.274e+00  1.083e+00 -1.176  0.239751
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 19679 on 1992 degrees of freedom
## Residual deviance: 19053 on 1969 degrees of freedom
## AIC: 19101
##
## Number of Fisher Scoring iterations: 11
##
## Ensure there are no NA values in the test data
test_na( test_data, %>%
drop_na(PERP_AGE_GROUP, PERP_SEXF, VIC_AGE_GROUP, VIC_SEXF)
)
##
## Make predictions
predictions <- predict(model, test_data)
##
## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## prediction from rank-deficient fit: attr(,"non-estim")=

```

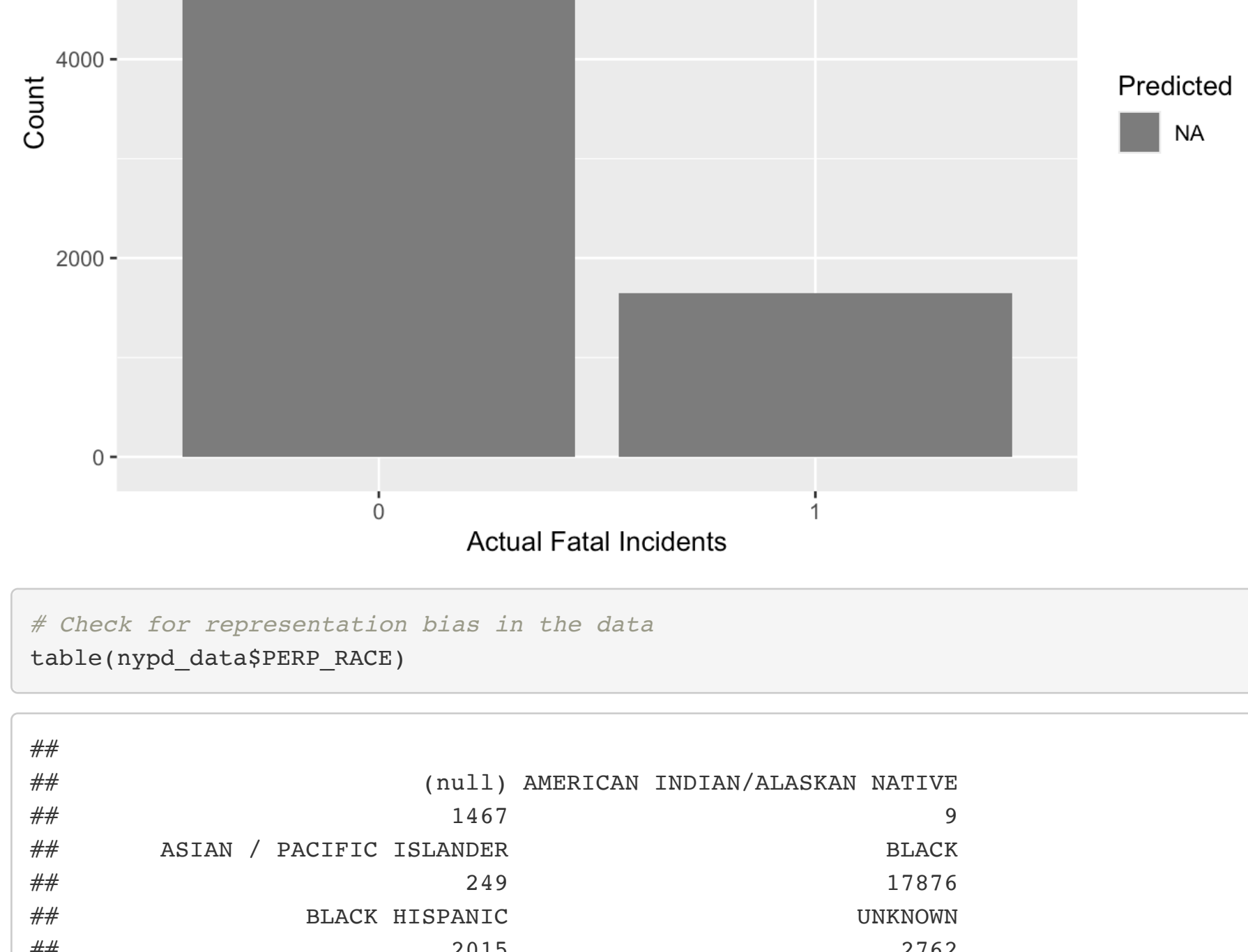
```
# Ensure predictions and reference are factors with the same levels
predictions <- factor(predictions, levels = levels(test_data$PRIORITY))
test_data$PRIORITY = factor(test_data$PRIORITY, levels = levels(predictions))
```

```
# Evaluate the model
confusionMatrix(predictions, test_data$FATAL)
```

```
## CONCLUSION Matrix and Statistics
##
## Reference
```

```
##      1 0 0
##
##      Accuracy : NaN
##              95 % CI : NA (NA)
##      No Information Rate : NA
##      P-Value [Acc > NIR] : NA
##
##              Kappa : NaN
##
## Mcnemar's Test P-Value : NA
##
##      Sensitivity : NA
##      Specificity : NA
##      Pos Pred Value : NA
##      Neg Pred Value : NA
##      Prevalence : NaN
##      Detection Rate : NaN
##      Detection Prevalence : NaN
##      Balanced Accuracy : NA
##
##      'Positive' Class : 0
##
##
##
## # Create a visual model comparison
results <- data.frame(
  Actual = test_data$FATAL,
  Predicted = predictions
),
```

```
geom_bar(position = "dodge") +
  labs(title = "Predicted vs Actual Fatal Incidents", x = "Actual Fatal Incidents", y = "Count", fill = "Predicted")
```



```
## AMERICAN INDIAN/ALASKAN NATIVE      ASIAN / PACIFIC ISLANDER
##              11                      440
##              BLACK                    BLACK HISPANIC
##              20235                    2796
##              UNKNOWN                  WHITE
##              70                      728
##              WHITE HISPANIC
##              4283
```

---

*# Discuss potential biases and their impact on the analysis*

```
#as not all shooting incidents may be reported or recorded accurately,
#potentially leading to underreporting or selective reporting. Additionally,
#data collection bias might arise from the subjective nature of how data is
#categorized, including racial classifications and incident descriptions.
#Survivor bias is also a factor, as the dataset only includes recorded
#incidents, leaving out unreported cases where victims or witnesses did not
#inform authorities. Geographical bias may occur if certain boroughs or
#neighborhoods are overrepresented or underrepresented due to variations
```

#reflecting changes in law enforcement practices, policies, and socio-  
#factors that affect the number and type of reported incidents.  
#Recognizing these biases is essential for accurately interpreting the