

# Chicago Crime

Xin Guan, Vera Hudak, Yuqi Zhang

2023-11-24

```
# Packages required
```

```
library(lubridate)
```

```
##
```

```
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##     date, intersect, setdiff, union
```

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##     filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##     intersect, setdiff, setequal, union
```

```
library(magrittr)
```

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
```

```
## v forcats 1.0.0      v stringr 1.5.0
```

```
## v ggplot2 3.4.3      v tibble  3.2.1
```

```
## v purrr   1.0.2      v tidyr   1.3.0
```

```
## v readr   2.1.4
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x tidyr::extract()   masks magrittr::extract()
```

```
## x dplyr::filter()    masks stats::filter()
```

```
## x dplyr::lag()       masks stats::lag()
```

```
## x purrr::set_names() masks magrittr::set_names()
```

```
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

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

```
## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##     lift
```

```
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
##
## Attaching package: 'pROC'
##
## The following objects are masked from 'package:stats':
##
##     cov, smooth, var
```

```
library(MLmetrics)
```

```
##
## Attaching package: 'MLmetrics'
##
## The following objects are masked from 'package:caret':
##
##     MAE, RMSE
##
## The following object is masked from 'package:base':
##
##     Recall
```

```
library(ROSE)
```

```
## Loaded ROSE 0.0-4
```

## Data Processing

```
# Read data
Crime_data <- read.csv("OriginalData.csv")
```

Processing variable 'Date'

```

# We want to firstly convert the format of date for the `lubridate` package, then extract information f
Crime_data$newDate <- mdy_hms(Crime_data$Date)
Crime_data$Hour <- hour(Crime_data$newDate)
Crime_data$WeekDay <- weekdays(Crime_data$newDate)
Crime_data$DayOfMonth <- day(Crime_data$newDate)
Crime_data$DayOfYear <- yday(Crime_data$newDate)
Crime_data$Month <- month(Crime_data$newDate, label = TRUE, abbr = FALSE)
Crime_data$Time <- hour(Crime_data$newDate)*100 + minute(Crime_data$newDate) #This format of `Time` var
Crime_data$TimeOfDay <- cut(
  hour(Crime_data$newDate),
  breaks= c(-Inf, 5, 12, 17, 20, Inf),
  labels = c("Night", "Early Morning", "Morning", "Afternoon", "Evening"),
  include.lowest = TRUE
)

```

## Processing missing data

```
colSums(is.na(Crime_data))
```

```

##          ID          Case.Number          Date
##          0              0              0
##          Block          IUCR          Primary.Type
##          0              0              0
##          Description Location.Description          Arrest
##          0              0              0
##          Domestic          Beat          District
##          0              0              0
##          Ward          Community.Area          FBI.Code
##          15              0              0
##          X.Coordinate          Y.Coordinate          Year
##          2205              2205              0
##          Updated.On          Latitude          Longitude
##          0              2205              2205
##          Location          newDate          Hour
##          0              0              0
##          WeekDay          DayOfMonth          DayOfYear
##          0              0              0
##          Month          Time          TimeOfDay
##          0              0              0

```

We can see that the number of missing values for the variables `X.Coordinate`, `Y.Coordinate`, `Latitude`, and `Longitude` are exactly the same, we can deduce that the coordinates are calculated from the latitude and longitude, so we don't need to include both pairs of location information. Also since the number of missing value is small compare to the number of data size, we will just eliminate the rows with missing values.

```
Crime_data %<>% na.omit()
```

## Expalnatory Data Analysis

### Feature Selection

Our first task is to predict whether arrest or not given time, location, and the crime type.

```
Binary_pred_df <- Crime_data %>% select(c("ID", "X.Coordinate", "Y.Coordinate", "Hour", "Time", "WeekDay", "Primary.Type"))
#location variables
```

### Imbalance Data Experiments

#### Baseline Model

```
#make this example reproducible
set.seed(1)

#use 80% of dataset as training set and 30% as test set
train <- Binary_pred_df %>% dplyr::sample_frac(0.80)
train$Arrest <- as.factor(train$Arrest)
test <- dplyr::anti_join(Binary_pred_df, train, by = 'ID')
test$Arrest <- as.factor(test$Arrest)

# Define a 5-fold cross validation
ctrl <- trainControl(method = "cv", number = 5, summaryFunction = twoClassSummary, classProbs = TRUE)

base_m <- train(Arrest ~ X.Coordinate + Y.Coordinate + DayOfYear + Hour +
  Primary.Type, data = train, method = "glm", family = "binomial", trControl = ctrl, metric = "ROC")

summary(base_m)
```

```
##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value
## (Intercept) -3.422e+00  9.502e-01 -3.601
## X.Coordinate  3.328e-06  4.963e-07  6.706
## Y.Coordinate -1.313e-06  2.568e-07 -5.114
## DayOfYear    -2.520e-04  6.583e-05 -3.828
## Hour         9.616e-03  1.016e-03  9.463
## Primary.TypeASSAULT  4.199e-01  1.793e-01  2.342
## Primary.TypeBATTERY  6.061e-01  1.786e-01  3.395
## Primary.TypeBURGLARY -8.733e-01  1.848e-01 -4.725
## 'Primary.TypeCONCEALED CARRY LICENSE VIOLATION'  6.068e+00  6.090e-01  9.963
## 'Primary.TypeCRIM SEXUAL ASSAULT' -5.924e-01  2.291e-01 -2.586
## 'Primary.TypeCRIMINAL DAMAGE' -7.697e-01  1.804e-01 -4.267
## 'Primary.TypeCRIMINAL SEXUAL ASSAULT' -3.619e-01  2.337e-01 -1.549
## 'Primary.TypeCRIMINAL TRESPASS'  2.154e+00  1.802e-01 11.951
## 'Primary.TypeDECEPTIVE PRACTICE' -1.015e+00  1.823e-01 -5.568
```

## Primary.TypeGAMBLING	1.650e+01	8.373e+01	0.197
## Primary.TypeHOMICIDE	1.429e+00	2.055e-01	6.954
## 'Primary.TypeHUMAN TRAFFICKING'	-1.260e+01	2.788e+02	-0.045
## 'Primary.TypeINTERFERENCE WITH PUBLIC OFFICER'	5.147e+00	2.331e-01	22.080
## Primary.TypeINTIMIDATION	-1.305e+00	4.895e-01	-2.665
## Primary.TypeKIDNAPPING	-7.778e-01	3.878e-01	-2.006
## 'Primary.TypeLIQUOR LAW VIOLATION'	1.651e+01	6.706e+01	0.246
## 'Primary.TypeMOTOR VEHICLE THEFT'	-8.693e-01	1.853e-01	-4.691
## Primary.TypeNARCOTICS	9.728e+00	4.815e-01	20.205
## 'Primary.TypeNON-CRIMINAL'	1.212e+00	1.238e+00	0.979
## Primary.TypeOBSCENITY	3.024e+00	3.939e-01	7.678
## 'Primary.TypeOFFENSE INVOLVING CHILDREN'	1.006e-01	1.903e-01	0.528
## 'Primary.TypeOTHER NARCOTIC VIOLATION'	2.636e+00	8.850e-01	2.979
## 'Primary.TypeOTHER OFFENSE'	6.890e-01	1.793e-01	3.842
## Primary.TypePROSTITUTION	1.651e+01	3.741e+01	0.441
## 'Primary.TypePUBLIC INDECENCY'	1.654e+01	3.120e+02	0.053
## 'Primary.TypePUBLIC PEACE VIOLATION'	2.822e+00	1.889e-01	14.938
## Primary.TypeROBBERY	-3.863e-01	1.835e-01	-2.105
## 'Primary.TypeSEX OFFENSE'	7.572e-02	1.998e-01	0.379
## Primary.TypeSTALKING	3.947e-01	2.663e-01	1.482
## Primary.TypeTHEFT	-2.664e-01	1.788e-01	-1.490
## 'Primary.TypeWEAPONS VIOLATION'	2.627e+00	1.806e-01	14.544
##	Pr(> z )		
## (Intercept)	0.000317	***	
## X.Coordinate	2.01e-11	***	
## Y.Coordinate	3.16e-07	***	
## DayOfYear	0.000129	***	
## Hour	< 2e-16	***	
## Primary.TypeASSAULT	0.019162	*	
## Primary.TypeBATTERY	0.000687	***	
## Primary.TypeBURGLARY	2.30e-06	***	
## 'Primary.TypeCONCEALED CARRY LICENSE VIOLATION'	< 2e-16	***	
## 'Primary.TypeCRIM SEXUAL ASSAULT'	0.009720	**	
## 'Primary.TypeCRIMINAL DAMAGE'	1.98e-05	***	
## 'Primary.TypeCRIMINAL SEXUAL ASSAULT'	0.121424		
## 'Primary.TypeCRIMINAL TRESPASS'	< 2e-16	***	
## 'Primary.TypeDECEPTIVE PRACTICE'	2.57e-08	***	
## Primary.TypeGAMBLING	0.843792		
## Primary.TypeHOMICIDE	3.56e-12	***	
## 'Primary.TypeHUMAN TRAFFICKING'	0.963939		
## 'Primary.TypeINTERFERENCE WITH PUBLIC OFFICER'	< 2e-16	***	
## Primary.TypeINTIMIDATION	0.007693	**	
## Primary.TypeKIDNAPPING	0.044899	*	
## 'Primary.TypeLIQUOR LAW VIOLATION'	0.805587		
## 'Primary.TypeMOTOR VEHICLE THEFT'	2.72e-06	***	
## Primary.TypeNARCOTICS	< 2e-16	***	
## 'Primary.TypeNON-CRIMINAL'	0.327808		
## Primary.TypeOBSCENITY	1.61e-14	***	
## 'Primary.TypeOFFENSE INVOLVING CHILDREN'	0.597272		
## 'Primary.TypeOTHER NARCOTIC VIOLATION'	0.002896	**	
## 'Primary.TypeOTHER OFFENSE'	0.000122	***	
## Primary.TypePROSTITUTION	0.658993		
## 'Primary.TypePUBLIC INDECENCY'	0.957728		
## 'Primary.TypePUBLIC PEACE VIOLATION'	< 2e-16	***	

```
## Primary.TypeROBBERY                0.035254 *
## 'Primary.TypeSEX OFFENSE'           0.704668
## Primary.TypeSTALKING                 0.138236
## Primary.TypeTHEFT                   0.136209
## 'Primary.TypeWEAPONS VIOLATION'     < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 216422  on 207323  degrees of freedom
## Residual deviance: 148083  on 207288  degrees of freedom
## AIC: 148155
##
## Number of Fisher Scoring iterations: 13
```

```
print(base_m)
```

```
## Generalized Linear Model
##
## 207324 samples
##      5 predictor
##      2 classes: 'false', 'true'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 165859, 165859, 165860, 165859, 165859
## Resampling results:
##
##      ROC          Sens          Spec
## 0.8131141 0.9728524 0.4757505
```

```
pred <- predict(base_m, test[,c("X.Coordinate", "Y.Coordinate", "DayOfYear", "Hour", "Primary.Type")],
```

```
# Create a ROC curve object
roc_curve <- roc(test$Arrest, pred[,2])
```

```
## Setting levels: control = false, case = true
```

```
## Setting direction: controls < cases
```

```
# Plot the ROC curve
plot(roc_curve, main = "ROC Curve", col = "blue", lwd = 2)

# Add AUC to the plot
auc_value <- auc(roc_curve)
text(0.8, 0.2, paste("AUC =", round(auc_value, 3)), col = "blue", cex = 1.2)
```

