

COMPAS Data Wrangling and Analysis

Dasha Asienaga

2024-02-27

Contents

Reading in the Data	1
The Data Set	2
Data Wrangling	3
Descriptive Statistics	3
Univariate Analysis	4
Bivariate Analysis	16
Multivariate Analysis	26
Demographic Group Analysis	31
Bivariate Analysis by Race	31
COMPAS Analysis	38
Logistic Regression	43
Check for Missing Data	43
Train and Test Split	43
Modeling	44
Seldonian Classification	44
Results	44

The thesis body will have more in-depth descriptions of the data analysis as well as select output and results from this file. This file is intended for general preliminary analysis of the COMPAS data set. Note that the results can only be generalized to Broward County, Florida, but there are important findings about the United States judicial system nevertheless.

Reading in the Data

```
#read in the data
compas_path <- "/home/dasienga24/Statistics-Senior-Honors-Thesis/Data Sets/COMPAS/compas_data.csv"
compasdata <- read.csv(compas_path)
```

The Data Set

The COMPAS data set has 12076 observations of defendants that were evaluated for the risk of recidivism by the COMPAS tool. There are 29 variables of interest as described below:

- **id**: unique person identifier.
- **compas_person_id**: unique COMPAS case identifier.
- **name**: full name.
- **first**: first name.
- **last**: last name.
- **sex**: sex categorized as male or female.
- **race**: race categorized as African-American, Asian, Caucasian, Hispanic, Native American, or Other.
- **age**: numeric age, ranging from 18 to 96.
- **age_cat**: age categorized as Less than 25, 25 - 45, or Greater than 45.
- **marital_status**: marital status categorized as Single, Significant Other, Married, Widowed, Separated, Divorced, or Unknown.
- **custody_status**: custody status categorized as Jail Inmate, Prison Inmate, Pretrial Defendant, Parole, Residential Program, or Probation.
- **juv_fel_count**: number of prior juvenile felonies, ranging from 0 to 20.
- **juv_misd_count**: number of prior juvenile misdemeanors, ranging from 0 to 13.
- **juv_other_count**: number of other prior juvenile offenses, ranging from 0 to 17.
- **priors_count**: number of non-juvenile prior offenses, ranging from 0 to 43.
- **days_b_screening_arrest**: number of days between COMPAS screening and arrest.
- **c_days_from_compas**: the number of days since COMPAS screening.
- **c_charge_degree**: the charge degree according to the appropriate laws.
- **c_charge_desc**: the charge description in words.
- **type_of_assessment**: the type of assessment, in this case, the assessment is 'Risk of Recidivism'.
- **raw_score**: COMPAS tool raw score on risk of recidivism.
- **decile_score**: decile rank on a scale of 1 - 10 based on the COMPAS raw score.
- **score_text**: COMPAS risk of recidivism based on the decile scores and categorized as High, Medium, or Low.
- **is_violent_recid**: categorical variable recording whether a defendant was accused of a violent crime within 2 years (0 = N, 1 = Y).
- **num_vr_cases**: number of times a defendant was accused of a violent crime within 2 years.
- **is_recid**: categorical variable recording whether a defendant was accused of a crime within 2 years (0 = N, 1 = Y).
- **num_r_cases**: number of times a defendant was accused of a crime within 2 years.
- **days_in_jail**: number of days spent in jail.
- **days_in_prison**: number of days spent in prison.

```
colnames(compasdata)
```

```
## [1] "id"                "compas_person_id"
## [3] "name"              "first"
## [5] "last"              "sex"
## [7] "race"              "age"
## [9] "age_cat"           "marital_status"
## [11] "custody_status"    "juv_fel_count"
## [13] "juv_misd_count"    "juv_other_count"
## [15] "priors_count"      "days_b_screening_arrest"
## [17] "c_days_from_compas" "c_charge_degree"
## [19] "c_charge_desc"     "type_of_assessment"
## [21] "raw_score"         "decile_score"
```

```
## [23] "score_text"           "is_violent_recid"
## [25] "num_vr_cases"        "is_recid"
## [27] "num_r_cases"         "days_in_jail"
## [29] "days_in_prison"
```

Data Wrangling

Before proceeding with the data analysis, we first need to handle some data anomalies. We'll also only consider COMPAS cases within 30 days of arrest to improve the data quality. This resulted in 9638 total observations.

```
compasdata <- compasdata %>%
  filter(decile_score > 0 & is_recid != -1 & days_b_screening_arrest >= -30 &
    days_b_screening_arrest <= 30) %>%
  mutate(days_b_screening_arrest = abs(days_b_screening_arrest))

count(compasdata)
```

```
##      n
## 1 9638
```

Next, let's also make sure that there are no duplicate defendants.

```
clean_compasdata <- compasdata[-which(duplicated(compasdata$id)), ]
```

We'll proceed with this data set and 9387 observations total.

Descriptive Statistics

Now that the data is clean, let's generate some descriptive statistics to understand the distribution of the variables in the data set and their relationships with each other.

First, below is a glimpse of the data as described above. Notice that there is a lot of missing data for `num_vr_cases` and `num_r_cases` because that information is only recorded for defendants that recommit a crime in the next 2 years.

```
glimpse(clean_compasdata)

## Rows: 9,387
## Columns: 29
## $ id                <int> 1, 3, 4, 7, 8, 9, 10, 11, 12, 14, 15, 16, 17, ~
## $ compas_person_id  <int> 56418, 51601, 38864, 59301, 61330, 56890, 6199~
## $ name              <chr> "miguel hernandez", "kevon dixon", "ed philo", ~
## $ first             <chr> "miguel", "kevon", "ed", "marsha", "edward", "~
## $ last              <chr> "hernandez", "dixon", "philo", "miles", "riddl~
## $ sex               <chr> "Male", "Male", "Male", "Male", "Male", "Male"~
## $ race              <chr> "Other", "African-American", "African-American~
## $ age               <int> 69, 34, 24, 44, 41, 43, 39, 20, 26, 27, 23, 37~
## $ age_cat           <chr> "Greater than 45", "25 - 45", "Less than 25", ~
## $ marital_status    <chr> "Single", "Single", "Single", "Separated", "Si~
## $ custody_status    <chr> "Jail Inmate", "Jail Inmate", "Jail Inmate", "~
```

```
## $ juv_fel_count      <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ juv_misd_count     <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ juv_other_count    <int> 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ priors_count       <int> 0, 0, 4, 0, 14, 3, 0, 0, 0, 0, 3, 0, 0, 0, 1, ~
## $ days_b_screening_arrest <int> 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 20, ~
## $ c_days_from_compas  <int> 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 490, ~
## $ c_charge_degree     <chr> "(F3)", "(F3)", "(F3)", "(M1)", "(F3)", "(F3)"~
## $ c_charge_desc       <chr> "Aggravated Assault w/Firearm", "Felony Batter~
## $ type_of_assessment  <chr> "Risk of Recidivism", "Risk of Recidivism", "R~
## $ raw_score           <dbl> -2.78, -0.76, -0.66, -1.93, -0.16, -0.72, -1.7~
## $ decile_score        <int> 1, 3, 4, 1, 6, 4, 1, 10, 5, 4, 6, 1, 3, 4, 1, ~
## $ score_text          <chr> "Low", "Low", "Low", "Low", "Medium", "Low", "~
## $ is_violent_recid    <int> 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ num_vr_cases        <lgl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA~
## $ is_recid            <int> 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1~
## $ num_r_cases         <int> NA, 3, 1, NA, 3, NA, NA, NA, NA, NA, 1, NA, NA~
## $ days_in_jail        <dbl> 8, 10, 139, 1, 48, 17, 3, 46, 87, 1, 4, 1, 0, ~
## $ days_in_prison      <dbl> 0, 53, 0, 0, 2130, 0, 0, 3948, 0, 0, 0, 0, 0, ~
```

Next, we will perform some univariate analysis for the variables in the data set before proceeding to conduct some bivariate and multivariate analysis.

Univariate Analysis

Univariate analysis will involve looking at some summary statistics and visualizations of the different variables in the data set.

Categorical Variables

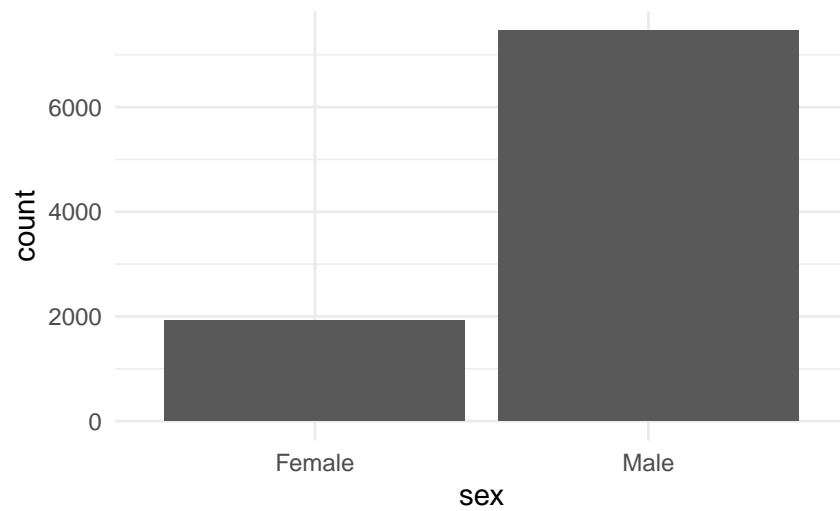
There 7457 males and 1930 females in the data set.

```
tally(clean_compasdata$sex)
```

```
## X
## Female   Male
##   1930    7457
```

```
ggplot(data = clean_compasdata, mapping = aes(x = sex)) +
  geom_bar() +
  theme_minimal() +
  labs(title = "Sex in the COMPAS Data Set")
```

Sex in the COMPAS Data Set



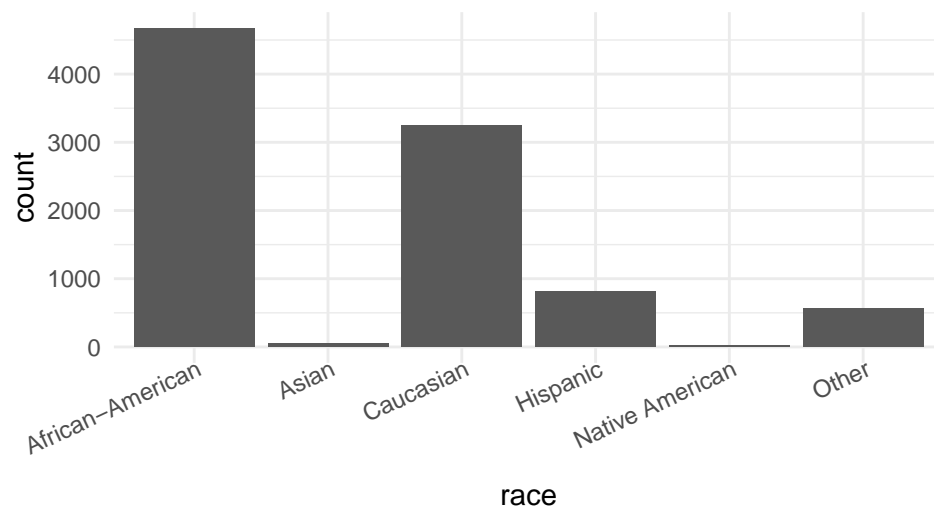
Most of the defendants are African-American and Caucasian, with only 27 Native Americans and 48 Asians.

```
tally(clean_compasdata$race)
```

```
## X
## African-American      Asian      Caucasian      Hispanic
##           4674           48           3250           818
## Native American      Other
##           27           570
```

```
ggplot(data = clean_compasdata, mapping = aes(x = race)) +
  geom_bar() +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  labs(title = "Race in the COMPAS Data Set")
```

Race in the COMPAS Data Set



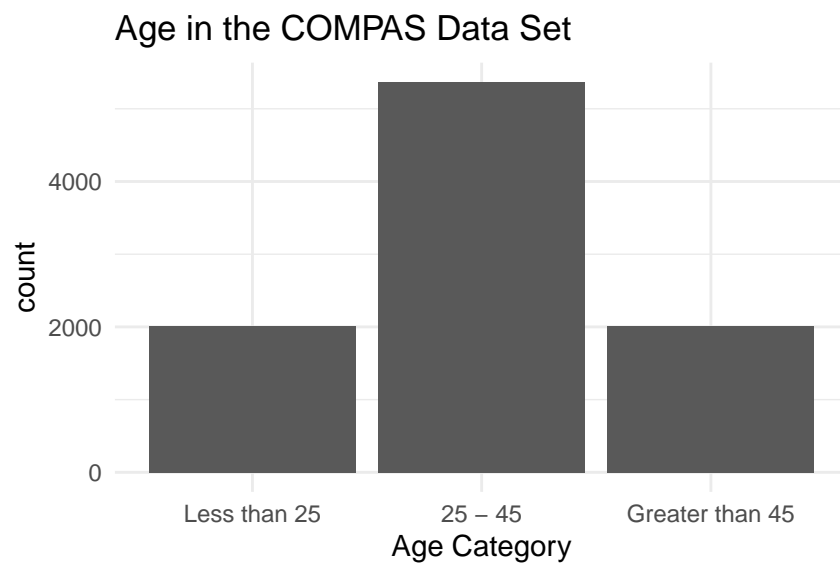
Majority of the defendants are between the age of 25 and 45, with about the same number of defendants less than 25 and greater than 25.

```
tally(clean_compasdata$age_cat)
```

```
## X
##      25 - 45 Greater than 45    Less than 25
##      5366          2012          2009

order <- c("Less than 25", "25 - 45", "Greater than 45")
```

```
ggplot(data = clean_compasdata, mapping = aes(x = age_cat)) +
  geom_bar() +
  theme_minimal() +
  scale_x_discrete(limits = order) +
  labs(x = "Age Category",
       title = "Age in the COMPAS Data Set")
```

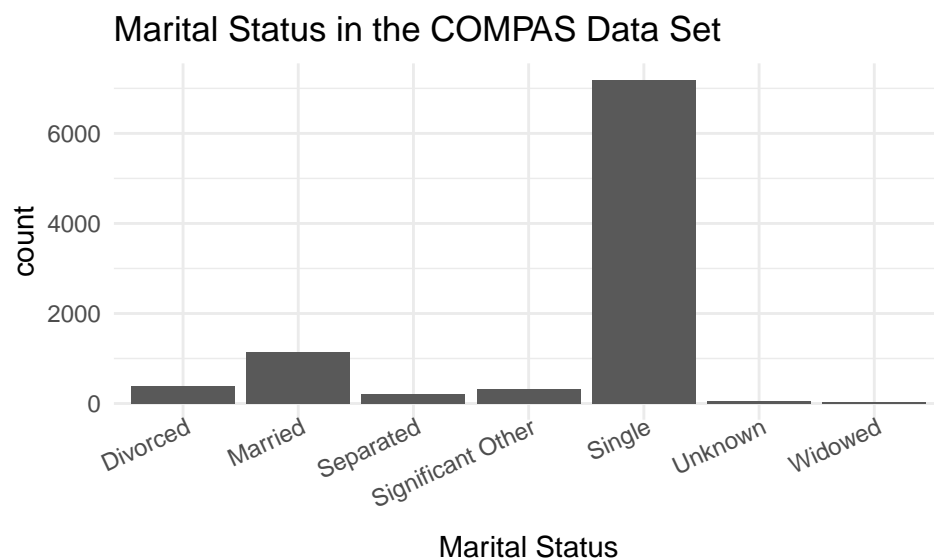


Most of the defendants are single, followed by married.

```
tally(clean_compasdata$marital_status)
```

```
## X
##      Divorced      Married      Separated Significant Other
##          398          1145          219             333
##      Single      Unknown      Widowed
##      7195          57          40
```

```
ggplot(data = clean_compasdata, mapping = aes(x = marital_status)) +
  geom_bar() +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  labs(x = "Marital Status",
       title = "Marital Status in the COMPAS Data Set")
```

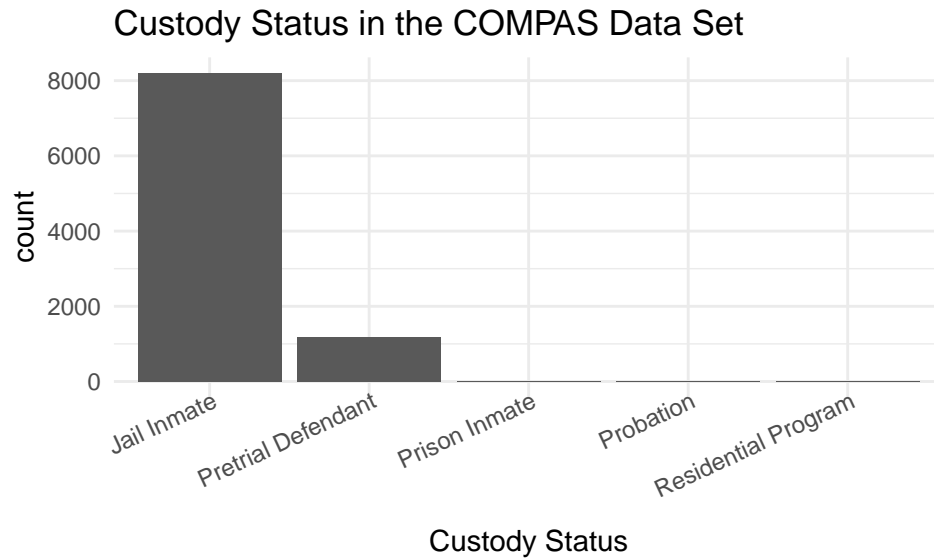


Most of the defendants are jail inmates, with only a handful of prison inmates, probationers, and defendants of the residential program.

```
tally(clean_compasdata$custody_status)
```

```
## X
##      Jail Inmate  Pretrial Defendant      Prison Inmate      Probation
##           8208             1170              4              3
## Residential Program
##              2
```

```
ggplot(data = clean_compasdata, mapping = aes(x = custody_status)) +
  geom_bar() +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  labs(x = "Custody Status",
       title = "Custody Status in the COMPAS Data Set")
```



As a data check, all the assessments are for risk of recidivism.

```
tally(clean_compasdata$type_of_assessment)
```

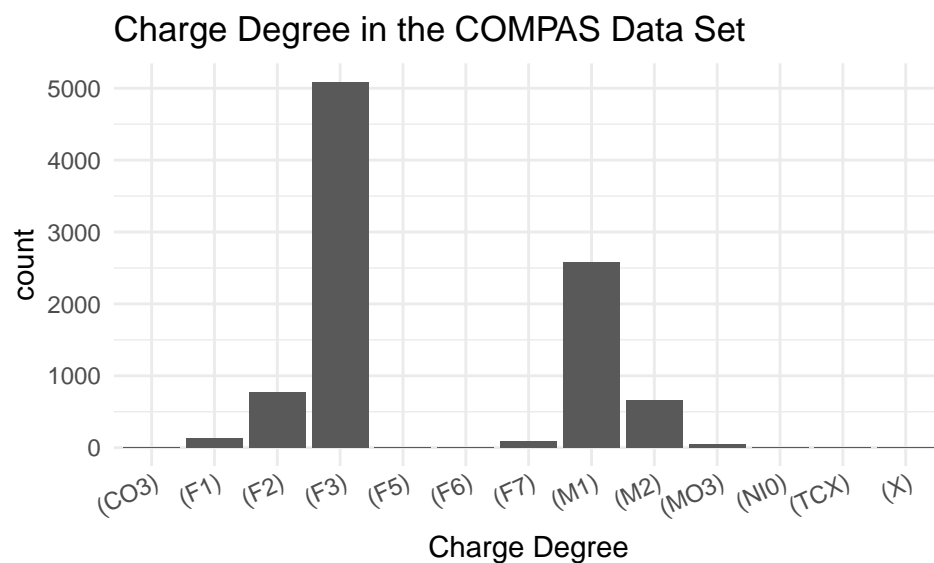
```
## X
## Risk of Recidivism
##          9387
```

There are 13 different charge degrees present in the data set. Most defendants were charged with (F3), which are felonies of the third degree. These are the least serious felonies in Florida and typically include crimes like breaking and entering, collecting and keeping stolen property, fraud, and petty theft. Many other defendants were also charged with (M1), which are a first-degree misdemeanors and can be punished by up to one year in jail. These include simple battery, disorderly conduct, DUI, indecent exposure, marijuana possession, shoplifting, prostitution, and vandalism, among others.

```
tally(clean_compasdata$c_charge_degree)
```

```
## X
## (C03) (F1) (F2) (F3) (F5) (F6) (F7) (M1) (M2) (M03) (NI0) (TCX) (X)
##      1  129  774 5091      5    3   85 2584  658   51    4    1    1
```

```
ggplot(data = clean_compasdata, mapping = aes(x = c_charge_degree)) +
  geom_bar() +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  labs(x = "Charge Degree",
       title = "Charge Degree in the COMPAS Data Set")
```

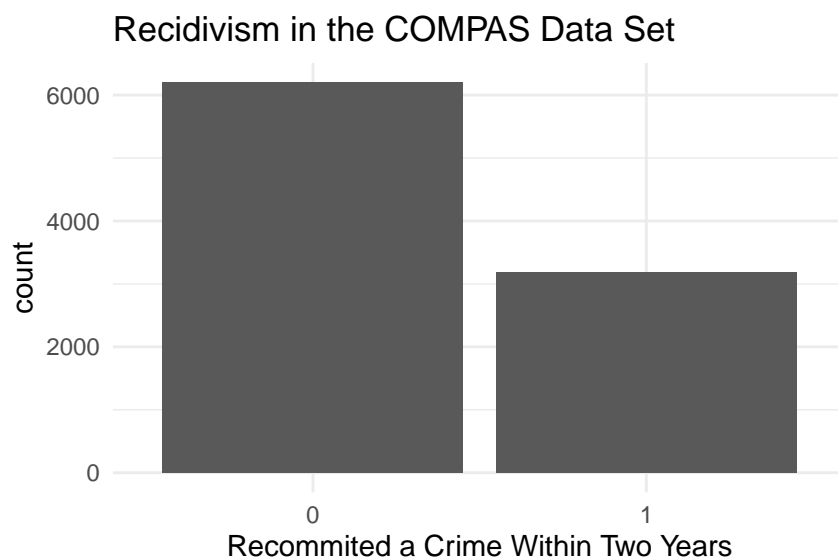



About two-thirds of the defendants did not recommit a crime within two years, while one-thirds did. This is our response variable and is indicative of class imbalance, which can affect the performance of machine learning classification algorithms. This is important to keep in mind when assessing model performance later on.

```
tally(clean_compasdata$is_recid)
```

```
## X
##   0   1
## 6199 3188
```

```
ggplot(data = clean_compasdata, mapping = aes(x = as.factor(is_recid))) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Recommitted a Crime Within Two Years",
       title = "Recidivism in the COMPAS Data Set")
```



Only 745 defendants recommitted a violent crime.

```
tally(clean_compasdata$is_violent_recid)
```

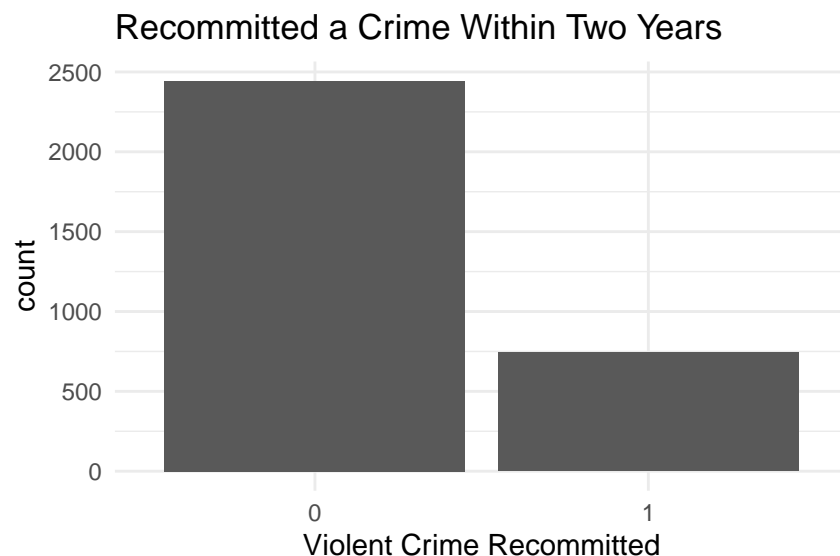
```
## X
##   0    1
## 8642  745
```

Out of the 3188 who recommitted a crime, 2443 re-committed a non-violent crime,

```
tally(clean_compasdata[clean_compasdata$is_recid == 1, ]$is_violent_recid,
      margins = TRUE)
```

```
## X
##   0    1 Total
## 2443  745 3188
```

```
ggplot(data = clean_compasdata[clean_compasdata$is_recid == 1, ],
       mapping = aes(x = as.factor(is_violent_recid))) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Violent Crime Recommited",
       title = "Recommitted a Crime Within Two Years")
```



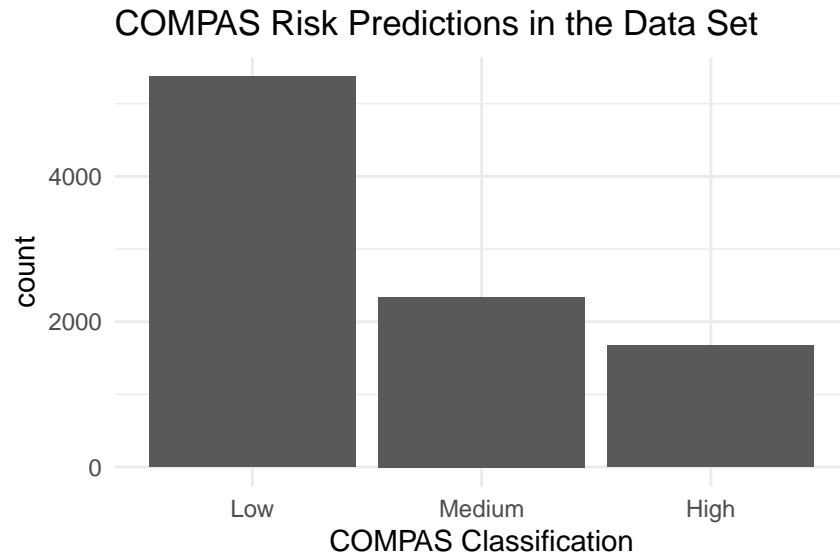
Finally, the COMPAS tool classified more than half of the defendants as low risk. In particular, 5370 were classified as low risk and 1677 as high risk, with the remaining 2340 as medium risk. This is expected since most of the defendants did not recommit a crime within the two year time window.

```
tally(clean_compasdata$score_text)
```

```
## X
##   High   Low Medium
##   1677  5370  2340
```

```
order <- c("Low", "Medium", "High")

ggplot(data = clean_compasdata, mapping = aes(x = score_text)) +
  geom_bar() +
  theme_minimal() +
  scale_x_discrete(limits = order) +
  labs(x = "COMPAS Classification",
       title = "COMPAS Risk Predictions in the Data Set")
```



This wraps up our univariate analysis of the categorical variables. Next, let's examine the univariate distribution of the continuous variables.

Continuous Variables

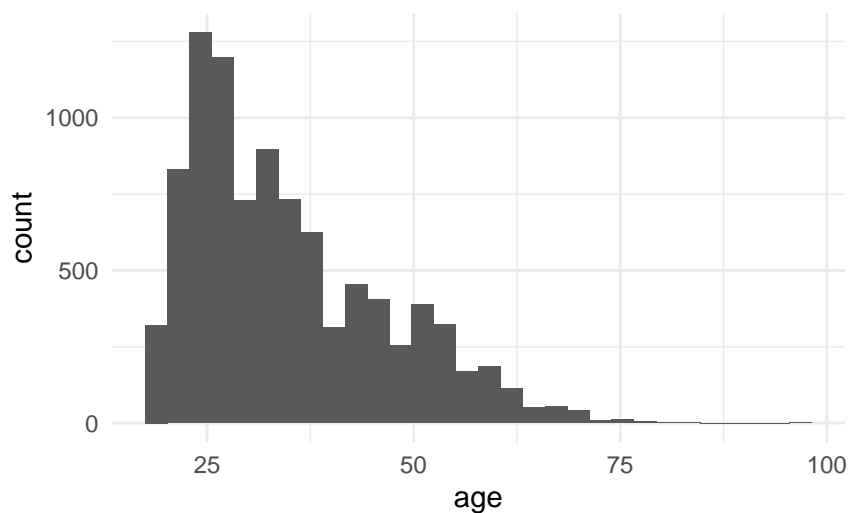
The age of the defendants ranges from 18 to 96 with a mean of 34 and a median of 32. There is no missing data. There's a right-skew in the distribution because of the few really old defendants.

```
favstats(clean_compasdata$age)
```

```
##  min Q1 median Q3 max    mean    sd    n missing
##   18 25     32 42  96 34.75413 11.80854 9387      0
```

```
ggplot(data = clean_compasdata, mapping = aes(x = age)) +
  geom_histogram() +
  theme_minimal() +
  labs(title = "Age in the COMPAS Data Set")
```

Age in the COMPAS Data Set



Most of the defendants had no juvenile felony accounts. The maximum juvenile felony count is 20. There is not enough variation in this variable.

```
favstats(clean_compasdata$juv_fel_count)
```

```
##  min Q1 median Q3 max      mean      sd    n missing
##    0  0      0  0  20 0.05837861 0.4518127 9387      0
```

Similarly, most defendants had no juvenile misdemeanor counts, which are less serious crimes than felonies. The maximum was 13, but there is not enough variation in this variable.

```
favstats(clean_compasdata$juv_misd_count)
```

```
##  min Q1 median Q3 max      mean      sd    n missing
##    0  0      0  0  13 0.0787259 0.4640061 9387      0
```

Similarly, most defendants had no other juvenile counts, excluding misdemeanors and felonies. The maximum was 11, but there is not enough variation in this variable.

```
favstats(clean_compasdata$juv_other_count)
```

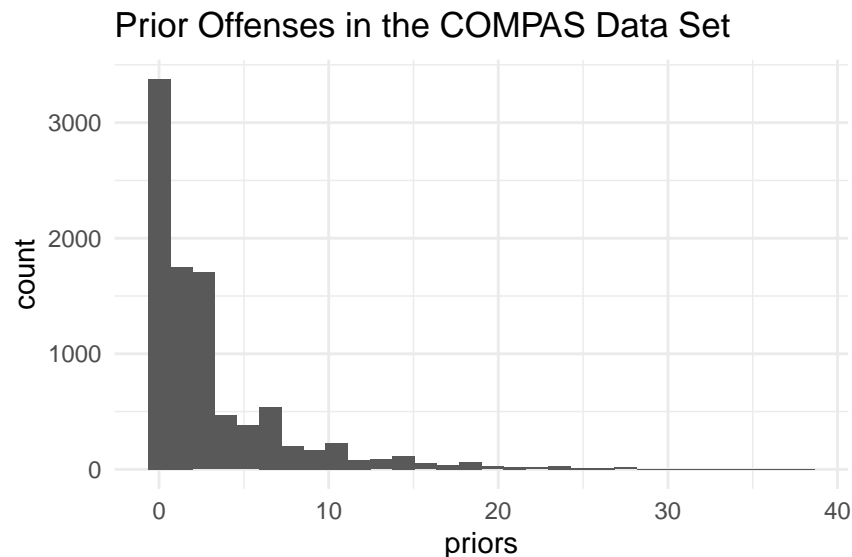
```
##  min Q1 median Q3 max      mean      sd    n missing
##    0  0      0  0  11 0.09917972 0.4683305 9387      0
```

There is slightly more variation in the `priors_count` variable which records the number of non-juvenile prior offenses for each defendant. It ranges from 0 to 38, with a median of 1 and a mean of 3.02, indicating a right skew as visualized in the histogram below. There is no missing data and the standard deviation is 4.586, suggesting that this may be a more informative variable when modeling.

```
favstats(clean_compasdata$priors_count)
```

```
## min Q1 median Q3 max      mean      sd    n missing
##   0  0       1  4  38 3.023863 4.586441 9387      0
```

```
ggplot(data = clean_compasdata, mapping = aes(x = priors_count)) +
  geom_histogram() +
  theme_minimal() +
  labs(x = "priors",
       title = "Prior Offenses in the COMPAS Data Set")
```



The `days_b_screening_arrest` variable indicates how many days passed between arrest and COMPAS screening. It may not be indicative of recidivism, however. We will evaluate this when performing bivariate analysis.

```
favstats(clean_compasdata$days_b_screening_arrest)
```

```
## min Q1 median Q3 max      mean      sd    n missing
##   0  1       1  1  30 2.140194 4.89312 9387      0
```

The interpretation of this variable is not clear – it seems to indicate the number of days since COMPAS screening to date. We will not include this in the analysis.

```
favstats(clean_compasdata$c_days_from_compas)
```

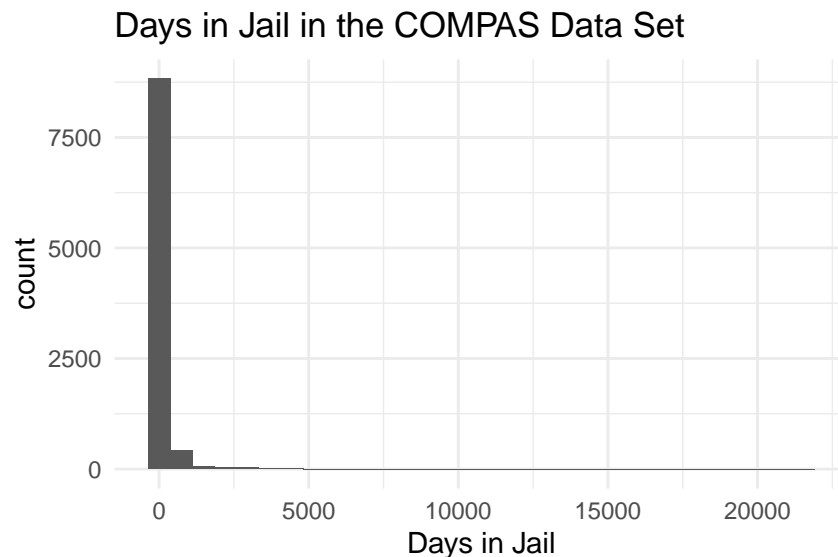
```
## min Q1 median Q3 max      mean      sd    n missing
##   0  1       1  1 9485 24.92436 263.4065 9387      0
```

The number of days spent in jail ranges from 0 to 21540, with a median of 4 days and a mean of 100 days. This variable is extremely right skewed, as visualized in the histogram. The standard deviation is also 393, indicating a lot of variation that may potentially be useful for predicting the risk of recidivism.

```
favstats(clean_compasdata$days_in_jail)
```

```
## min Q1 median Q3 max mean sd n missing
## 0 1 4 60 21540 100.1712 393.2173 9387 0
```

```
ggplot(data = clean_compasdata, mapping = aes(x = days_in_jail)) +
  geom_histogram() +
  theme_minimal() +
  labs(x = "Days in Jail",
       title = "Days in Jail in the COMPAS Data Set")
```



The days spent in prison is not as variable as the days spent in jail. The minimum 0 and the maximum is 190739. This skews the mean to 784.7951, but the median is 0. The distinction between jail and prison is still unclear.

```
favstats(clean_compasdata$days_in_prison)
```

```
## min Q1 median Q3 max mean sd n missing
## 0 0 0 0 190739 784.7951 3473.352 9387 0
```

The number of crimes recommitted by the defendants who re-committed a crime within two years ranges from 1 to 55, with a median of 1 and a mean of 1.73.

```
favstats(clean_compasdata$num_r_cases)
```

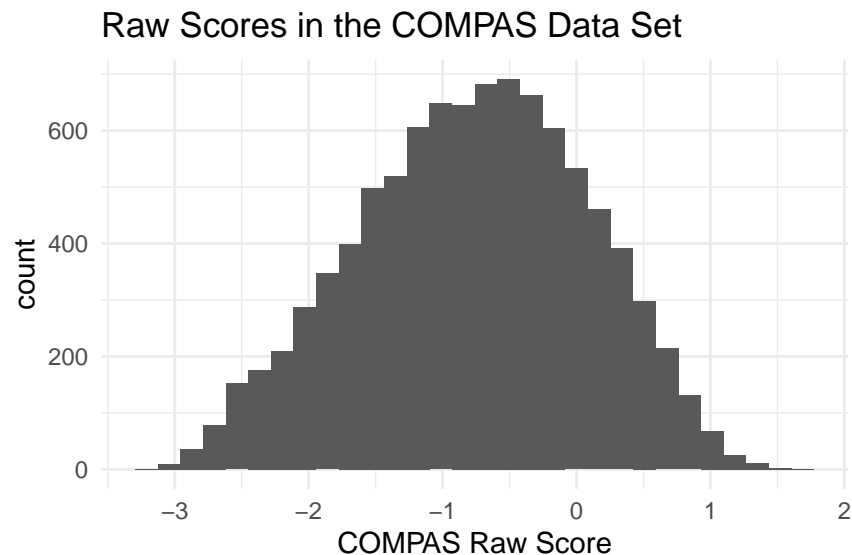
```
## min Q1 median Q3 max mean sd n missing
## 1 1 1 2 55 1.736512 1.629916 3188 6199
```

Finally, the COMPAS tool outputs a raw score for each defendant. The raw score ranges from -3.21 to 1.69 with a median of -0.74 and a mean of -0.78. The distribution of the raw scores is visualized on the histogram below. The distribution is unimodal and symmetric with a slight left skew.

```
favstats(clean_compasdata$raw_score)
```

```
##   min   Q1 median   Q3  max      mean      sd    n missing  
## -3.21 -1.38 -0.74 -0.15 1.69 -0.7763417 0.856942 9387      0
```

```
ggplot(data = clean_compasdata, mapping = aes(x = raw_score)) +  
  geom_histogram() +  
  theme_minimal() +  
  labs(x = "COMPAS Raw Score",  
       title = "Raw Scores in the COMPAS Data Set")
```

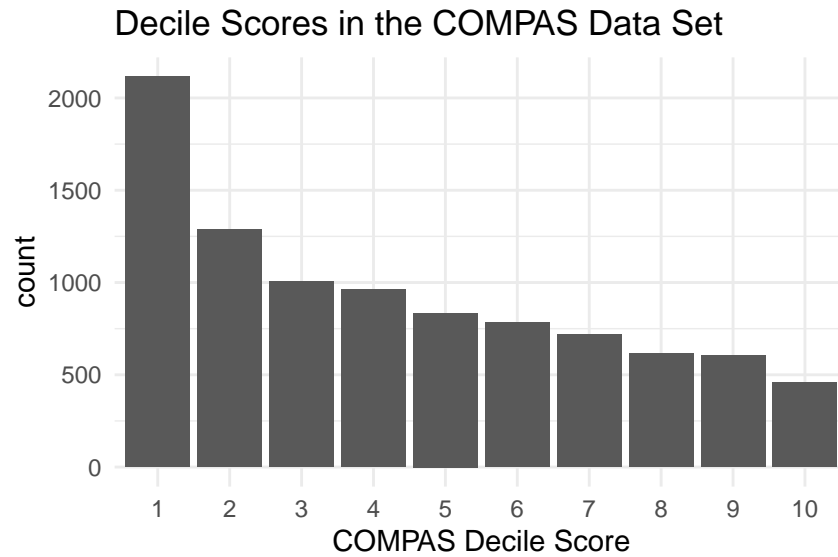


The raw scores are then converted into decile scores that determine the predicted risk of recidivism. The decile scores range from 1 to 10 with a median of 4 and a mean of 4.3. The histogram displays the distribution of the decile scores – it makes me wonder how, or whether, the decile scores are computed from the raw scores.

```
favstats(clean_compasdata$decile_score)
```

```
##   min Q1 median Q3 max      mean      sd    n missing  
##    1  2     4  7  10 4.305849 2.849011 9387      0
```

```
ggplot(data = clean_compasdata, mapping = aes(x = as.factor(decile_score))) +  
  geom_bar() +  
  theme_minimal() +  
  labs(x = "COMPAS Decile Score",  
       title = "Decile Scores in the COMPAS Data Set")
```



Note that the decile scores are mapped to ‘low’, ‘medium’, and ‘high’ risk as detailed in the table below.

```
clean_compasdata %>%
  dplyr::select(decile_score, score_text) %>%
  filter(score_text != 'N/A') %>%
  rename("Risk" = score_text) %>%
  group_by(Risk) %>%
  summarise("Min" = min(decile_score),
            "Max" = max(decile_score)) %>%
  arrange(Min) %>%
  kable(booktabs = TRUE)
```

Risk	Min	Max
Low	1	4
Medium	5	7
High	8	10

This concludes our univariate analysis of the variables in the COMPAS data set. Next, we will look at some of the bivariate relationships.

Bivariate Analysis

In this section, we will explore the relationships between our variables and the response variable, `is_recid`, which records whether or not a defendant recommitted a crime within 2 years.

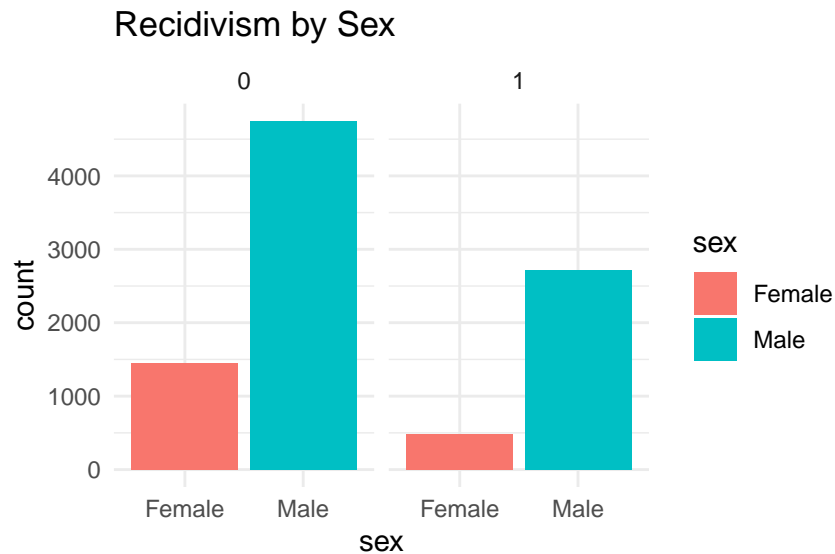
Categorical Variables

It doesn't appear as though there is much evident relationship between sex and recidivism.

```
ggplot(data = clean_compasdata, mapping = aes(x = sex, fill = sex)) +
  geom_bar() +
  theme_minimal() +
```



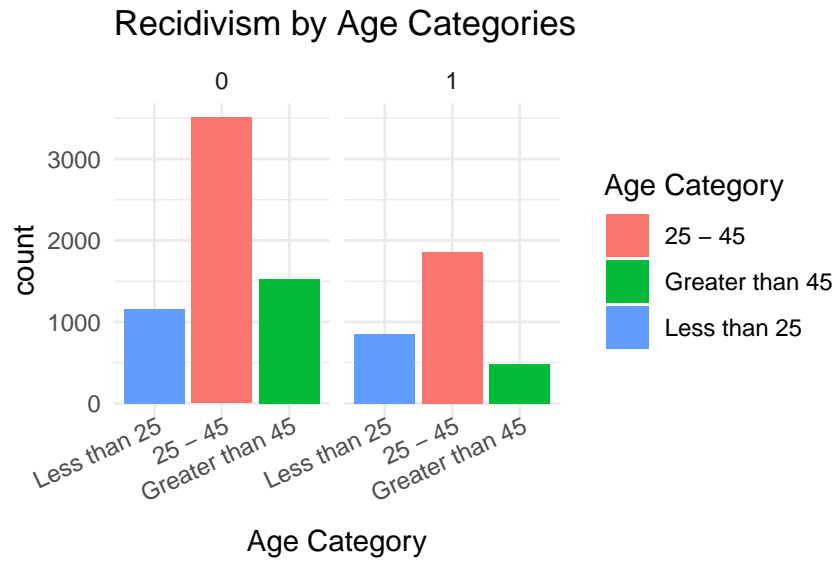
```
facet_wrap(~is_recid) +
labs(title = "Recidivism by Sex")
```



Among defendants who do not recidivate, there are more defendants that are aged 45 in comparison to those less than 25. However, among those that recidivated, there are more defendants that are less than 25 in comparison to those that are greater than 45. This indicates that age may hold some valuable information regarding a defendant's likelihood of recidivism.

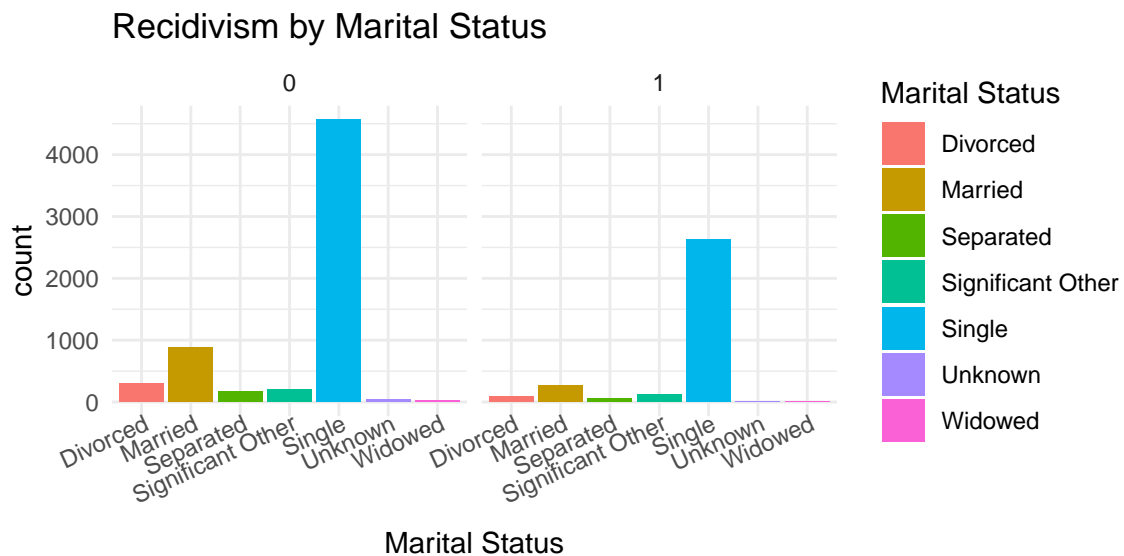
```
order <- c("Less than 25", "25 - 45", "Greater than 45")
```

```
ggplot(data = clean_compasdata,
       mapping = aes(x = age_cat, fill = age_cat)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~is_recid) +
  labs(title = "Recidivism by Age Categories",
       x = "Age Category",
       fill = "Age Category") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  scale_x_discrete(limits = order)
```



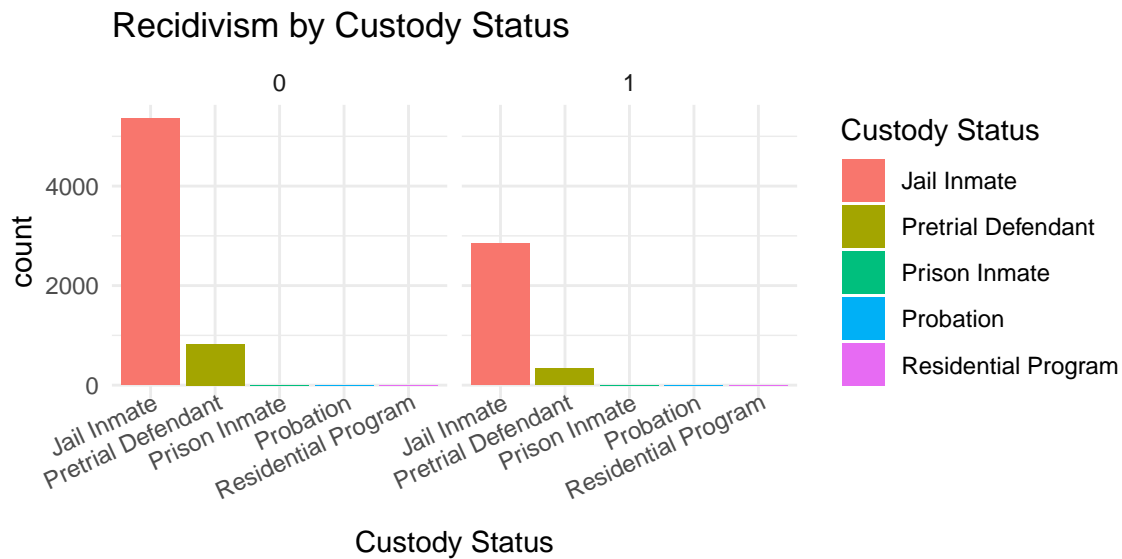
It doesn't appear as though there is much relationship between recidivism and marital status.

```
ggplot(data = clean_compasdata,
       mapping = aes(x = marital_status, fill = marital_status)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~is_recid) +
  labs(title = "Recidivism by Marital Status",
       x = "Marital Status",
       fill = "Marital Status") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



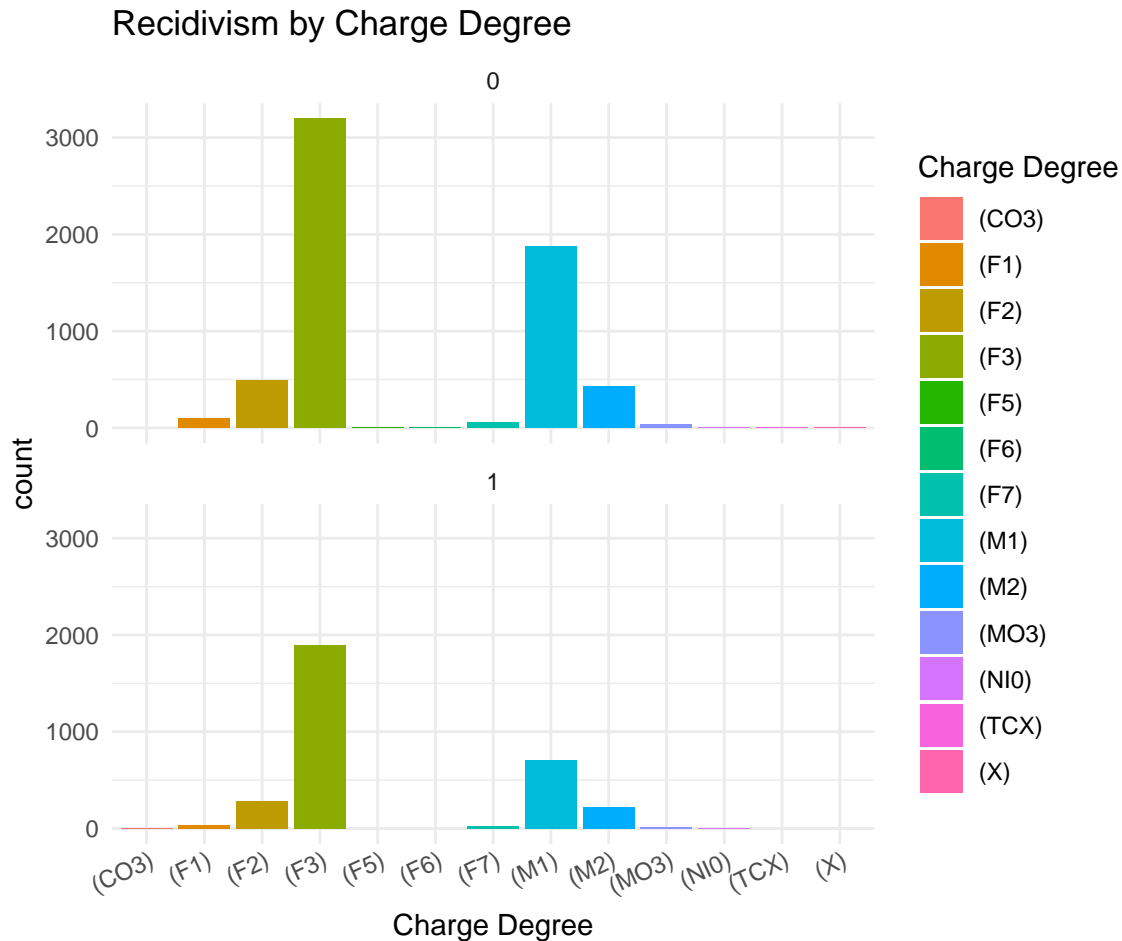
It doesn't appear as though there is much relationship between recidivism and custody status.

```
ggplot(data = clean_compasdata,
       mapping = aes(x = custody_status, fill = custody_status)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~is_recid) +
  labs(title = "Recidivism by Custody Status",
       x = "Custody Status",
       fill = "Custody Status") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



It doesn't appear as though there is much relationship between recidivism and charge degree.

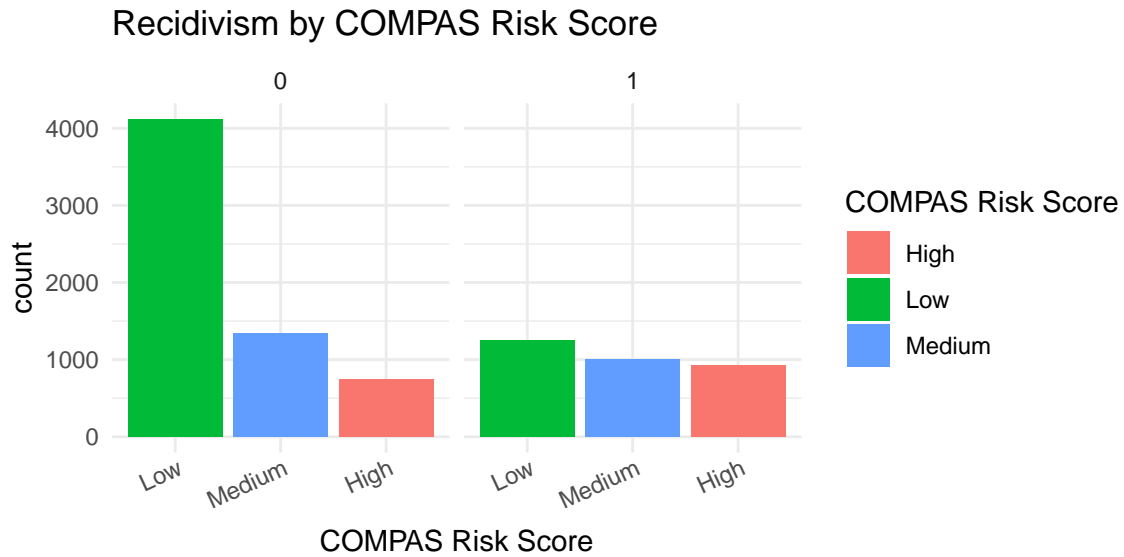
```
ggplot(data = clean_compasdata,
       mapping = aes(x = c_charge_degree, fill = c_charge_degree)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~is_recid, ncol = 1) +
  labs(title = "Recidivism by Charge Degree",
       x = "Charge Degree",
       fill = "Charge Degree") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



However, it appears as though the COMPAS tool classifies defendants who recommit a crime as almost as equally risky of recidivism – there is no significant distinction between ‘low’, ‘medium’, and ‘high’ risk for these defendants. For the defendants that don’t recommit a crime, most are predicted as ‘low’ risk, followed by ‘medium’, and then ‘high’ risk. Note, however, that this variable will not be included as a predictor in the model as the purpose of this analysis is to assess COMPAS performance, or more generally, standard ML approaches, in comparison to the Seldonian framework.

```
order <- c("Low", "Medium", "High")

ggplot(data = clean_compasdata,
       mapping = aes(x = score_text, fill = score_text)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~is_recid) +
  labs(title = "Recidivism by COMPAS Risk Score",
       x = "COMPAS Risk Score",
       fill = "COMPAS Risk Score") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  scale_x_discrete(limits = order)
```



Next, let's perform a similar analysis for the continuous variables.

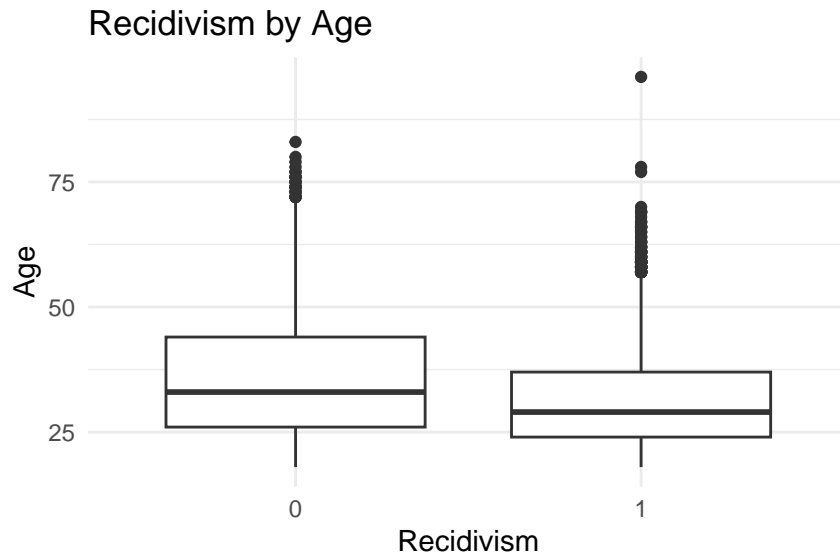
Continuous Variables

There is a difference in the mean and median ages for defendants who recommit a crime within two years versus those who don't. Those who recidivate tend to be younger than those who don't, indicating that this will be a useful variable in the model. This is in line with intuition from society.

```
favstats(data = clean_compasdata, age ~ is_recid)
```

	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	18	26	33	44	83	36.04646	12.17525	6199	0
## 2	1	18	24	29	37	96	32.24122	10.62147	3188	0

```
ggplot(data = clean_compasdata,
       mapping = aes(x = as.factor(is_recid), y = age)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Recidivism by Age",
       x = "Recidivism",
       y = "Age")
```



There is not much distributional difference in juvenile felony counts for defendants who recidivate versus those who don't.

```
favstats(data = clean_compasdata, juv_fel_count ~ is_recid)
```

##	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	0	0	0	0	13	0.03645749	0.3297061	6199	0
## 2	1	0	0	0	0	20	0.10100376	0.6221204	3188	0

There is not much distributional difference in juvenile misdemeanor counts for defendants who recidivate versus those who don't.

```
favstats(data = clean_compasdata, juv_misd_count ~ is_recid)
```

##	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	0	0	0	0	12	0.04387804	0.3356002	6199	0
## 2	1	0	0	0	0	13	0.14648683	0.6388211	3188	0

There is not much distributional difference in juvenile offenses for defendants who recidivate versus those who don't.

```
favstats(data = clean_compasdata, juv_other_count ~ is_recid)
```

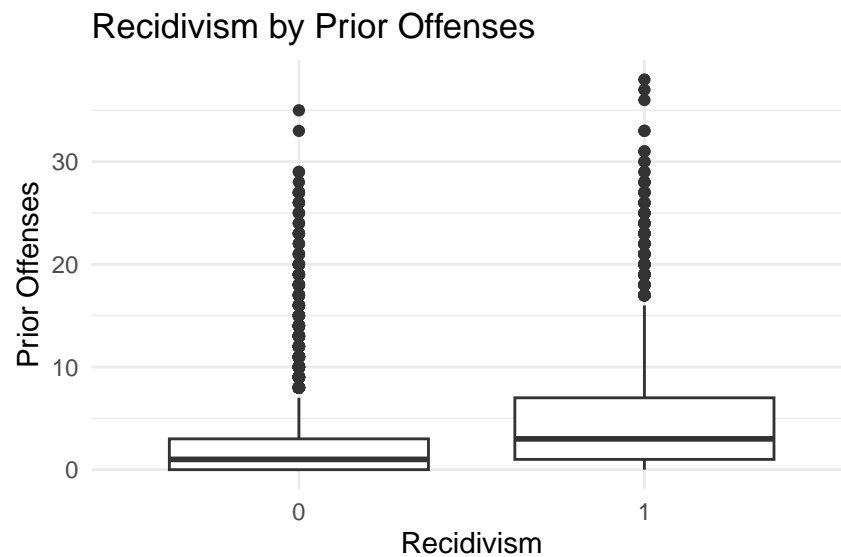
##	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	0	0	0	0	11	0.06355864	0.3966132	6199	0
## 2	1	0	0	0	0	9	0.16844417	0.5768642	3188	0

There is some distributional difference in non-juvenile prior offenses for defendants who recidivate versus those who don't, as is indicated by the different means and medians. Those who recommit a crime within two years tend to have more prior offenses. This will be a useful variable to include in the models.

```
favstats(data = clean_compasdata, priors_count ~ is_recid)
```

```
##   is_recid min Q1 median Q3 max    mean    sd    n missing
## 1         0  0  0      1  3  35 2.157283 3.684641 6199      0
## 2         1  0  1      3  7  38 4.708908 5.589893 3188      0
```

```
ggplot(data = clean_compasdata,
       mapping = aes(x = as.factor(is_recid), y = priors_count)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Recidivism by Prior Offenses",
       x = "Recidivism",
       y = "Prior Offenses")
```



There doesn't appear to be any distributional difference in days between COMPAS screening and arrest for defendants who recidivate versus those who don't. This will not be a useful variable for modeling.

```
favstats(data = clean_compasdata, days_b_screening_arrest ~ is_recid)
```

```
##   is_recid min Q1 median Q3 max    mean    sd    n missing
## 1         0  0  1      1  1  30 2.149218 4.859912 6199      0
## 2         1  0  1      1  1  30 2.122647 4.957777 3188      0
```

While the means differ because of the right-skew nature of the data, there doesn't appear to be much distributional difference in days since COMPAS screening for defendants who recidivate versus those who don't. This will not be a useful variable for modeling.

```
favstats(data = clean_compasdata, c_days_from_compas ~ is_recid)
```

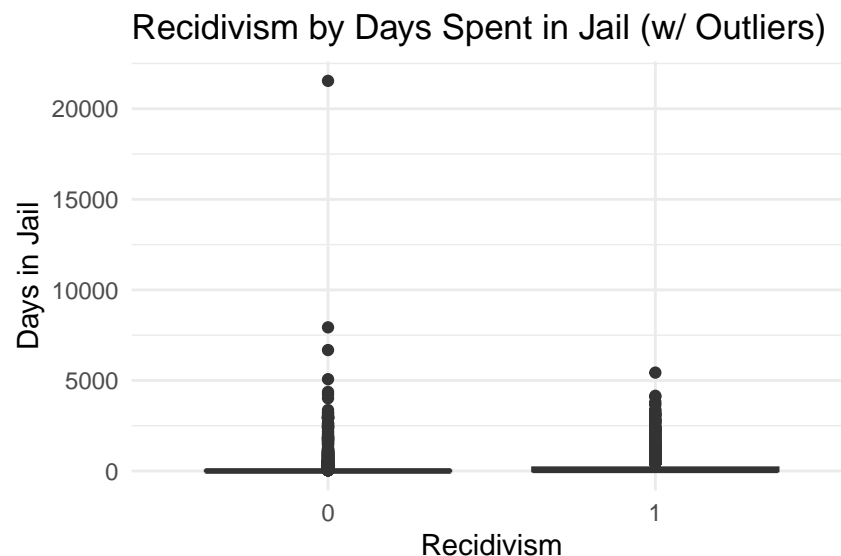
```
##   is_recid min Q1 median Q3 max    mean    sd    n missing
## 1         0  0  1      1  1 9485 31.21052 305.1847 6199      0
## 2         1  0  1      1  1 5450 12.70107 151.5946 3188      0
```

There is an evident difference in the distribution of the number of days spent in jail for participants who recommit a crime within two years versus those who don't. The mean, median, and max value differ significantly, indicating variation that may be useful in modeling. It's hard to visualize the boxplots with all the outliers, so the second boxplot trims the y-axis to better visualize this relationship.

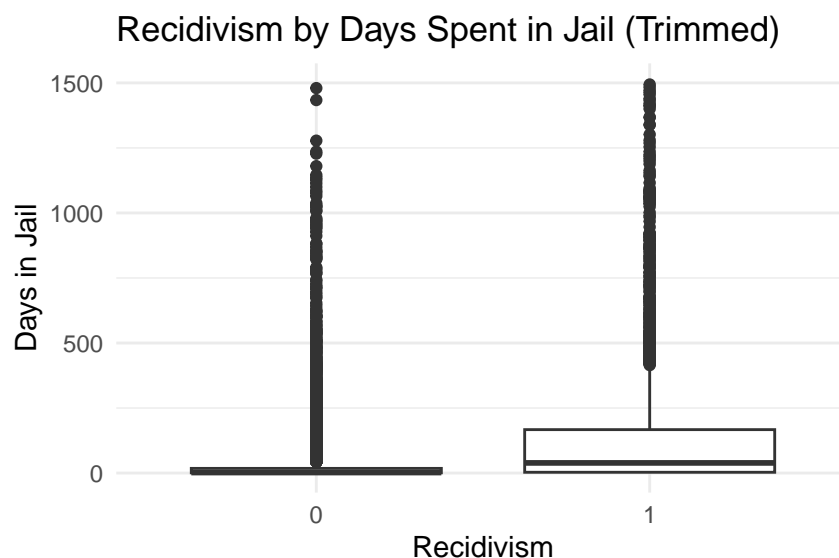
```
favstats(data = clean_compasdata, days_in_jail ~ is_recid)
```

##	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	0	1	2	19.00	21540	65.86562	395.1086	6199	0
## 2	1	0	3	41	178.25	5432	166.87767	380.8262	3188	0

```
ggplot(data = clean_compasdata,
       mapping = aes(x = as.factor(is_recid), y = days_in_jail)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Recidivism by Days Spent in Jail (w/ Outliers)",
       x = "Recidivism",
       y = "Days in Jail")
```



```
ggplot(data = clean_compasdata,
       mapping = aes(x = as.factor(is_recid), y = days_in_jail)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Recidivism by Days Spent in Jail (Trimmed)",
       x = "Recidivism",
       y = "Days in Jail") +
  ylim(0,1500)
```

There doesn't appear to be much distributional difference, rather than the effects of extreme right skews, in days spent in prison for defendants who recidivate versus those who don't. The difference between jail and prison is still not clear, so this will not be a useful variable for modeling.

```
favstats(data = clean_compasdata, days_in_prison ~ is_recid)
```

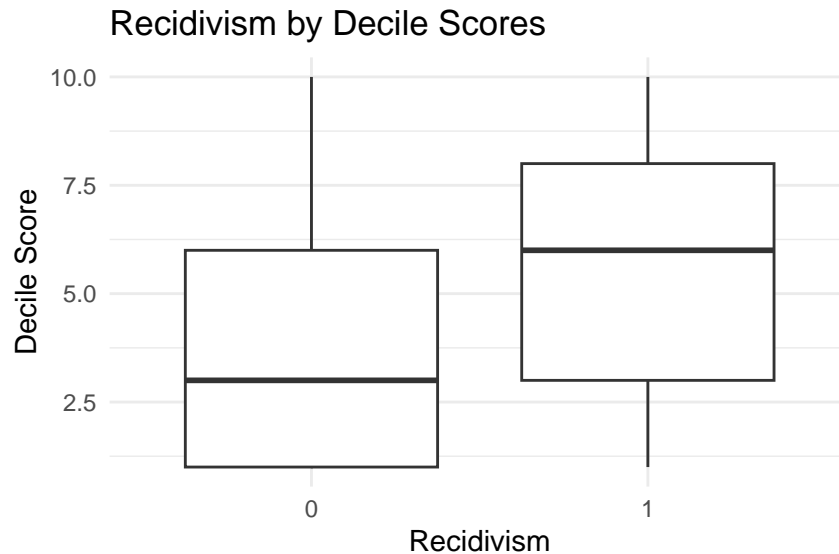
	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	0	0	0	0.00	190739	441.0958	3017.073	6199	0
## 2	1	0	0	0	223.75	67056	1453.1114	4141.345	3188	0

Finally, let's assess the COMPAS decile scores. The median and mean decile scores differ for defendants who recommit a crime within 2 years versus those who don't. The median score for those who don't is 3, which is mapped to low risk. The median score for those who do is 6, which is mapped to medium risk. This indicates that the COMPAS tool has some predictive accuracy. However, the range of scores is the same for both defendants who recidivate versus those who do not, suggesting that the tool is not entirely accurate in its predictions.

```
favstats(data = clean_compasdata, decile_score ~ is_recid)
```

	is_recid	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	0	1	1	3	6	10	3.694467	2.667049	6199	0
## 2	1	1	3	6	8	10	5.494668	2.816137	3188	0

```
ggplot(data = clean_compasdata,
       mapping = aes(x = as.factor(is_recid), y = decile_score)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Recidivism by Decile Scores",
       x = "Recidivism",
       y = "Decile Score")
```



This wraps up our analysis of the bivariate relationships between the continuous variables in the data set and the response variable: `is_recid`.

Multivariate Analysis

Based on the univariate and bivariate analysis, the 8 most informative predictive variables for modeling will be:

- `sex`
- `age`
- `age category`
- `marital status`
- `custody status`
- `prior offenses`
- `charge degree`
- `days in jail`

We will also include `race` in the modeling data set as our demographic variable, though it will not be included in the models themselves. Finally, `is_recid`, the response variable, will also be selected in the data set.

For further analysis, we will also include the COMPAS decile scores to assess which of these variables may have been used to model the COMPAS risk assessment tool.

Now, let's create a new data set with these 11 variables. Below is a glimpse of the data set.

```
compas_final <- clean_compasdata %>%
  dplyr::select(c(race, sex, age, age_cat, marital_status,
                  custody_status, priors_count, c_charge_degree,
                  days_in_jail, decile_score, is_recid))

glimpse(compas_final)

## Rows: 9,387
## Columns: 11
## $ race      <chr> "Other", "African-American", "African-American", "Othe~
## $ sex       <chr> "Male", "Male", "Male", "Male", "Male", "Male", "Femal~
## $ age       <int> 69, 34, 24, 44, 41, 43, 39, 20, 26, 27, 23, 37, 22, 41~
```

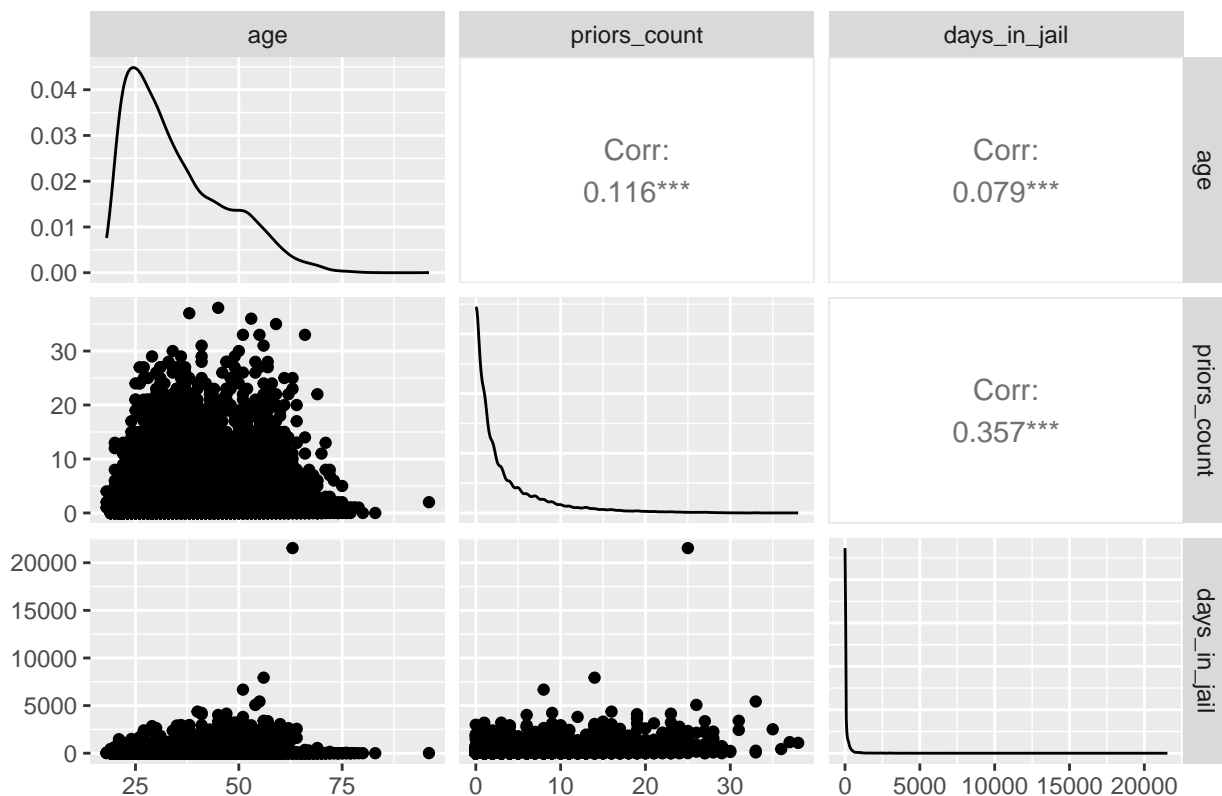
```
## $ age_cat      <chr> "Greater than 45", "25 - 45", "Less than 25", "25 - 45~
## $ marital_status <chr> "Single", "Single", "Single", "Separated", "Single", "~
## $ custody_status <chr> "Jail Inmate", "Jail Inmate", "Jail Inmate", "Jail Inm~
## $ priors_count  <int> 0, 0, 4, 0, 14, 3, 0, 0, 0, 0, 3, 0, 0, 0, 1, 7, 0, 3,~
## $ c_charge_degree <chr> "(F3)", "(F3)", "(F3)", "(M1)", "(F3)", "(F3)", "(M1)"~
## $ days_in_jail  <dbl> 8, 10, 139, 1, 48, 17, 3, 46, 87, 1, 4, 1, 0, 1, 183, ~
## $ decile_score  <int> 1, 3, 4, 1, 6, 4, 1, 10, 5, 4, 6, 1, 3, 4, 1, 3, 1, 10~
## $ is_recid      <int> 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, ~
```

Scatterplot Matrix

First, a scatterplot matrix with just the 3 continuous predictive variables in the final data set, `age`, `prior offenses`, and `days in jail`, will help to elucidate the covariate relationships between the variables. All the variables have moderate to weak correlations, with the strongest correlation of 0.357 being between the number of prior offenses and the number of days spent in jail. There are no concerns for multicollinearity.

```
ggpairs(data = compas_final,
  columns = c("age", "priors_count", "days_in_jail"),
  title = "Scatterplot Matrix of the COMPAS Continuous Variables")
```

Scatterplot Matrix of the COMPAS Continuous Variables



As observed, the variables have a significant right skew and the relationship is non-linear. Let's explore what effect different transformations may have on the covariate relationships. The log transformation resulted in many non-finite values, so we will look at a square root transformation instead.

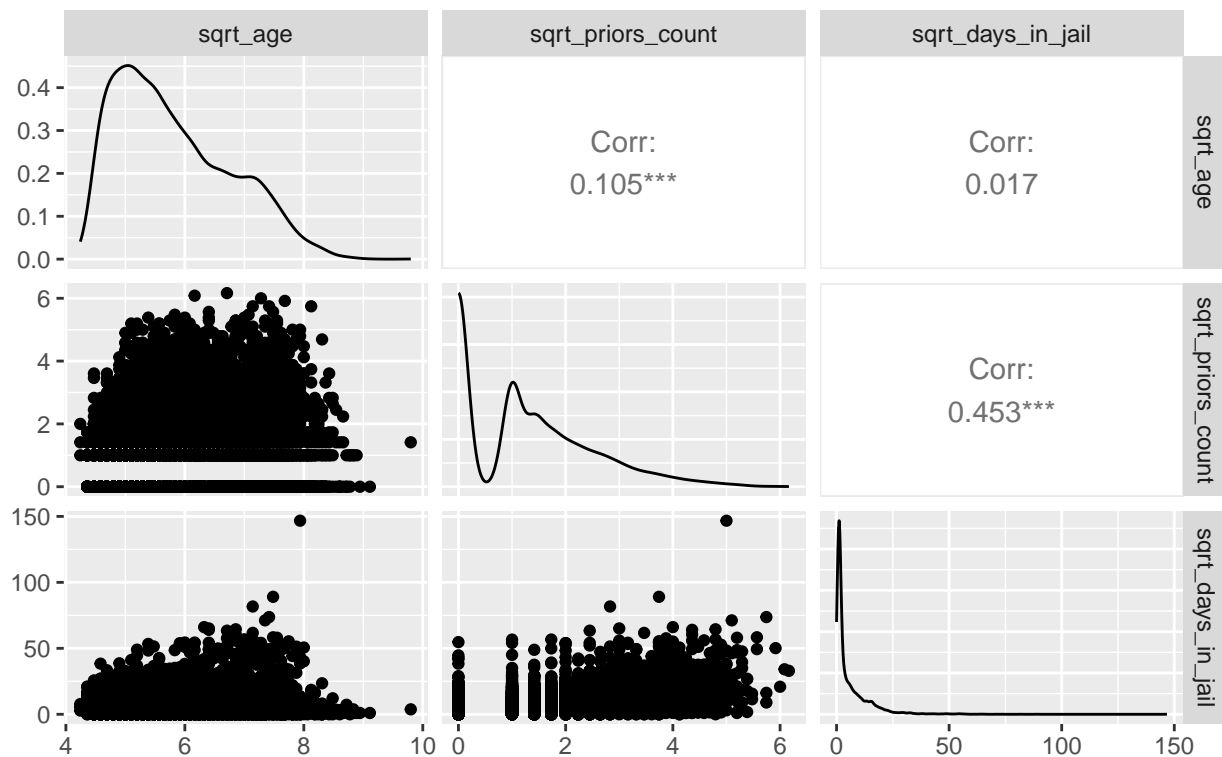
```
compas_final <- compas_final %>%
  mutate(sqrt_age = sqrt(age),
```

```
sqrt_priors_count = sqrt(priors_count),
sqrt_days_in_jail = sqrt(days_in_jail))
```

While this transformation strengthened the relationship between the number of prior offenses and the number of days spent in jail, it weakened the other correlations. Next, we'll assess how much this affects the relationship with the COMPAS decile scores.

```
ggpairs(data = compas_final,
        columns = c("sqrt_age", "sqrt_priors_count", "sqrt_days_in_jail"),
        title = "Scatterplot Matrix of the COMPAS Continuous Variables  
(Square Root)")
```

Scatterplot Matrix of the COMPAS Continuous Variables
(Square Root)



Pearson's Correlation Matrix

Decile scores has a moderate relationship with age, prior offenses, and number of days spent in jail – this suggests that these variables may indeed be useful for modeling recidivism.

```
mycordata1 <- compas_final %>%
  dplyr::rename("Age" = age,
               "Priors" = priors_count,
               "Days in Jail" = days_in_jail,
               "Decile Scores" = decile_score) %>%
  dplyr::select("Age", "Priors", "Days in Jail", "Decile Scores")
```

```
cor(mycordata1) %>%
  kable(digits = 2,
        booktabs = TRUE)
```

	Age	Priors	Days in Jail	Decile Scores
Age	1.00	0.12	0.08	-0.39
Priors	0.12	1.00	0.36	0.45
Days in Jail	0.08	0.36	1.00	0.25
Decile Scores	-0.39	0.45	0.25	1.00

The square root transformation of the predictor variables actually strengthens the relationship with decile scores – this suggests that the square root transformation of these variables may be better for modeling recidivism.

```
mycordata2 <- compas_final %>%
  dplyr::rename("Square Root Age" = sqrt_age,
               "Priors" = sqrt_priors_count,
               "Days in Jail" = sqrt_days_in_jail,
               "Decile Scores" = decile_score) %>%
  dplyr::select("Square Root Age", "Priors", "Days in Jail", "Decile Scores")

cor(mycordata2) %>%
  kable(digits = 2,
        booktabs = TRUE)
```

	Square Root Age	Priors	Days in Jail	Decile Scores
Square Root Age	1.00	0.11	0.02	-0.39
Priors	0.11	1.00	0.45	0.46
Days in Jail	0.02	0.45	1.00	0.45
Decile Scores	-0.39	0.46	0.45	1.00

Spearman's Correlation Matrix

Spearman's Correlation is better at capturing non-linear relationships. Using Spearman correlations reveals stronger correlations between the variables and the COMPAS decile scores.

There is no observable difference when calculating Spearman's correlation with the square-root transformed variables versus the original variables.

```
mycordata3 <- compas_final %>%
  dplyr::rename("Age" = age,
               "Priors" = priors_count,
               "Days in Jail" = days_in_jail,
               "Decile Scores" = decile_score) %>%
  dplyr::select("Age", "Priors", "Days in Jail", "Decile Scores")

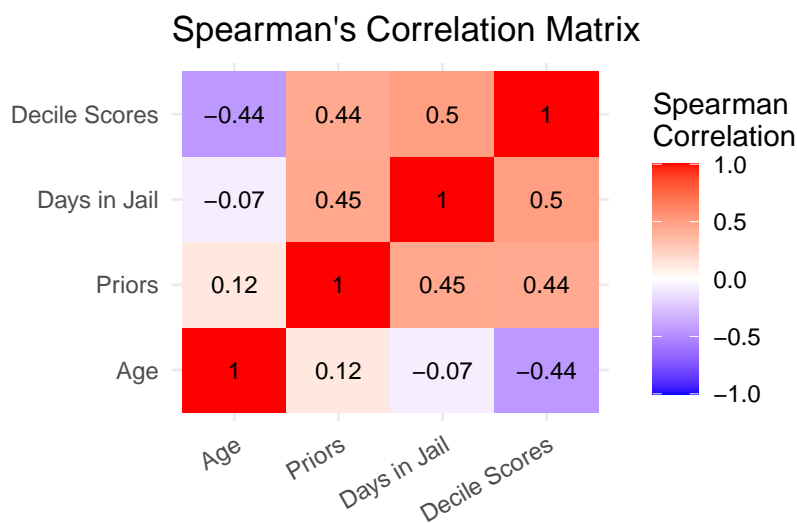
cor(mycordata3, method = "spearman") %>%
  kable(digits = 2,
        booktabs = TRUE)
```

	Age	Priors	Days in Jail	Decile Scores
Age	1.00	0.12	-0.07	-0.44
Priors	0.12	1.00	0.45	0.44
Days in Jail	-0.07	0.45	1.00	0.50
Decile Scores	-0.44	0.44	0.50	1.00

Finally, let's visualize these correlations.

```
mycors <- round(cor(mycordata3, method = "spearman"), 2)
mycorplot <- melt(mycors)

ggplot(data = mycorplot, aes(x = Var1, y = Var2, fill = value)) +
  geom_tile() +
  labs(x = "",
       y = "",
       title = "Spearman's Correlation Matrix") +
  scale_fill_gradient2(
    low = "blue",
    high = "red",
    mid = "white",
    midpoint = 0,
    limit = c(-1, 1),
    space = "Lab",
    name = "Spearman\nCorrelation"
  ) +
  geom_text(aes(Var2, Var1, label = value),
            color = "black",
            size = 3) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 30, vjust = 1, hjust = 1))
```



Now that we have a thorough understanding of the make-up of the data set, we will perform a demographic analysis next to get a better understanding of the racial discrepancies that may be present before, finally,

proceeding with the recidivism risk modeling.

Demographic Group Analysis

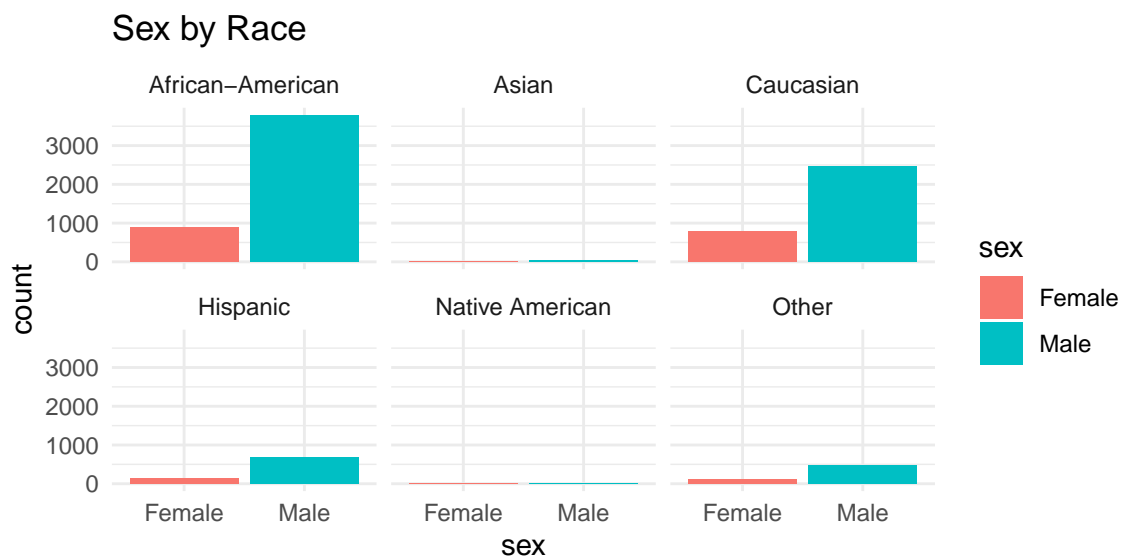
Bivariate Analysis by Race

The demographic variable of interest for this analysis is **race**, and we will employ the Seldonian algorithm in a hope to achieve fairer recidivism risk predictions. With that goal in mind, it's important to perform a demographic group analysis along the defendants' races to better understand any underlying or proxy relationships with the variables.

First, let's analyze the bivariate relationships of some of the most important variables with race.

For all the races, with the exception of native Americans for who there are not many data points available, there are more male defendants than female defendants.

```
ggplot(data = compas_final, mapping = aes(x = sex, fill = sex)) +  
  geom_bar() +  
  theme_minimal() +  
  facet_wrap(~race) +  
  labs(title = "Sex by Race")
```

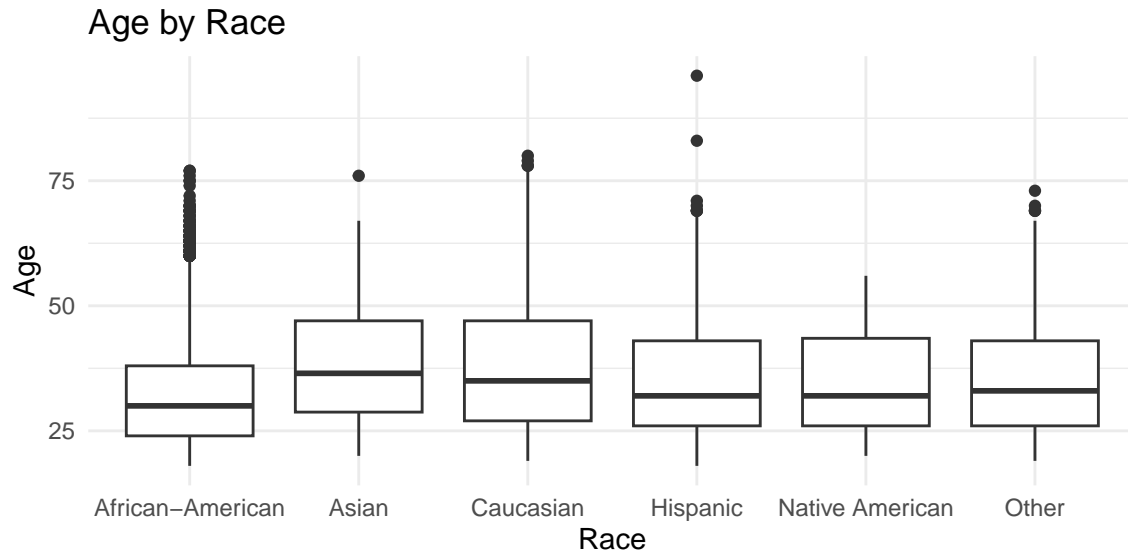


African-American defendants tend to be, on average, the youngest compared to all the other races. Asian defendants, followed by Caucasian defendants, tend to be the oldest. However, there is considerable overlap among all the races, and the relationships is visualized in the boxplot below.

```
favstats(data = compas_final, age ~ race)
```

##	race	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	African-American	18	24.00	30.0	38.0	77	32.60312	10.77213	4674	0
## 2	Asian	20	28.75	36.5	47.0	76	38.20833	12.21607	48	0
## 3	Caucasian	19	27.00	35.0	47.0	80	37.51969	12.60537	3250	0
## 4	Hispanic	18	26.00	32.0	43.0	96	35.22494	11.86236	818	0
## 5	Native American	20	26.00	32.0	43.5	56	34.29630	10.57870	27	0
## 6	Other	19	26.00	33.0	43.0	73	35.67895	11.68420	570	0

```
ggplot(data = compas_final,
       mapping = aes(x = as.factor(race), y = age)) +
  geom_boxplot() +
  theme_minimal() +
  labs(title = "Age by Race",
       x = "Race",
       y = "Age")
```

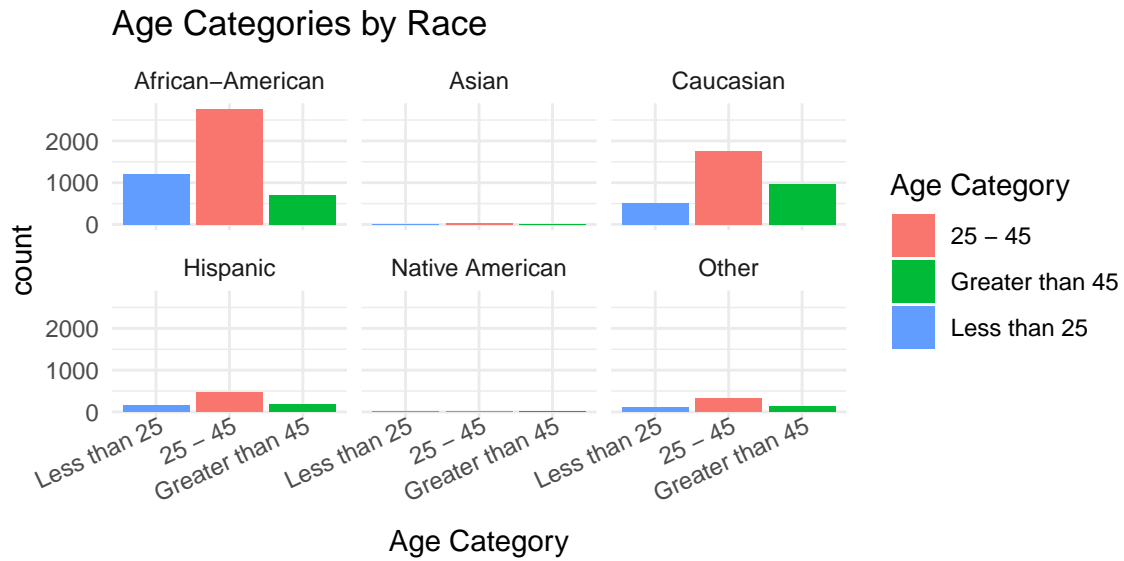


The bar plot below illustrates that for all races, most defendants fall between the ages of 25 and 45. Notice, however, that for African-American defendants, there are more defendants that are less than 25 than those that are greater than 45. The converse is true for Caucasians, with more defendants that are greater than 45 in comparison to those less than 25.

This illustrates that there is some relationship between age and race in this data set, particularly for African-Americans versus Caucasians.

```
order <- c("Less than 25", "25 - 45", "Greater than 45")

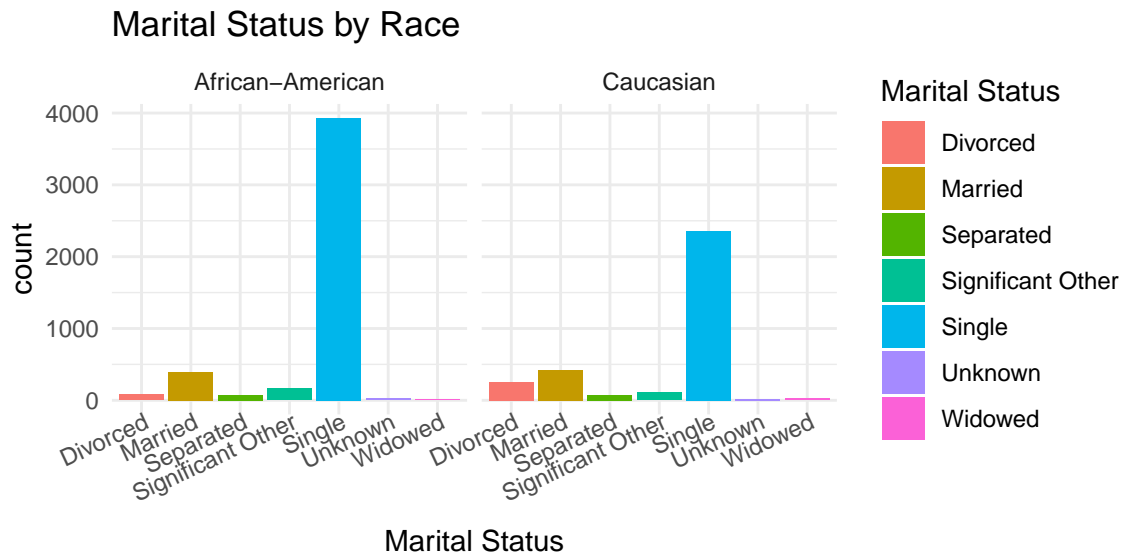
ggplot(data = compas_final,
       mapping = aes(x = age_cat, fill = age_cat)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~race) +
  labs(title = "Age Categories by Race",
       x = "Age Category",
       fill = "Age Category") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1)) +
  scale_x_discrete(limits = order)
```

For simpler visualization, let's assess the marital status just for Black and White defendants. For both races, most defendants are single. There are no distinct distributional differences.

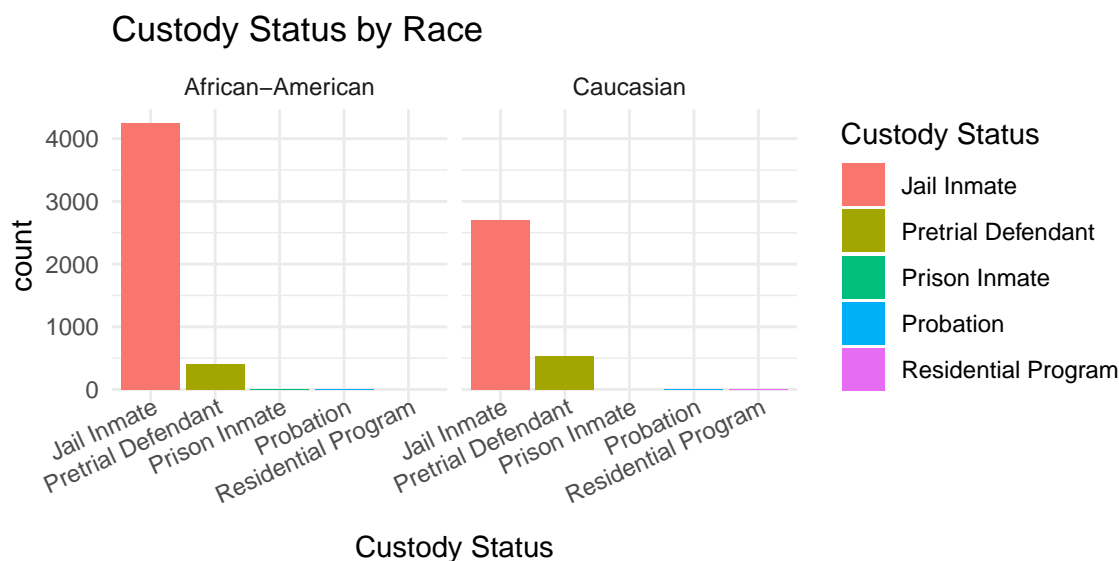
```
compas_final_bw <- compas_final %>%
  filter(race %in% c("African-American", "Caucasian"))

ggplot(data = compas_final_bw,
  mapping = aes(x = marital_status, fill = marital_status)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~race) +
  labs(title = "Marital Status by Race",
    x = "Marital Status",
    fill = "Marital Status") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



Similarly, there is no observable distribution difference in custody status by race. However, as much as there are less Caucasian defendants overall in comparison to African-American defendants, there are slightly more Caucasian prison defendants.

```
ggplot(data = compas_final_bw,
       mapping = aes(x = custody_status, fill = custody_status)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~race) +
  labs(title = "Custody Status by Race",
       x = "Custody Status",
       fill = "Custody Status") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



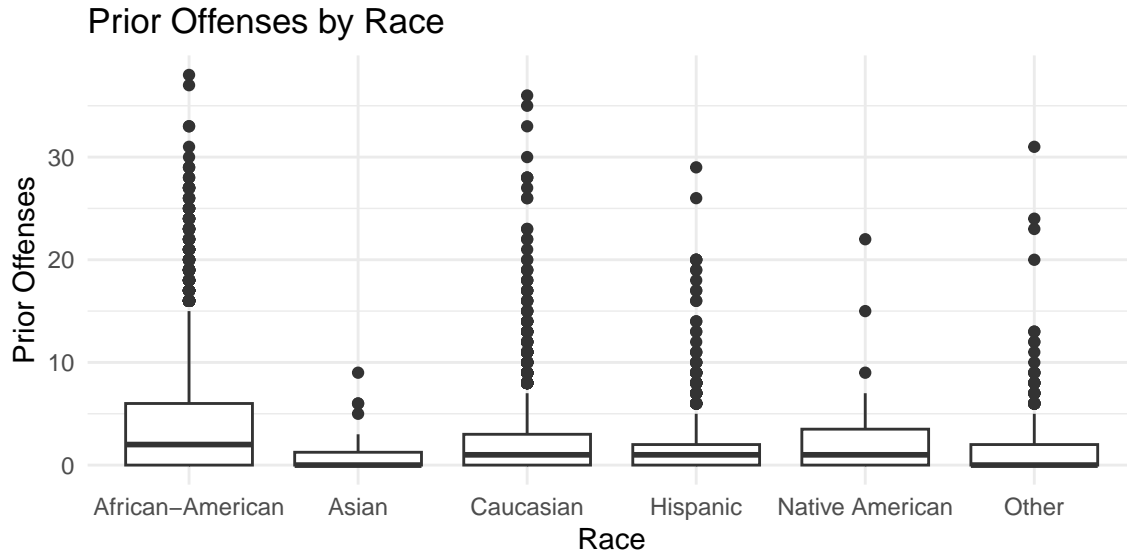
Notably, African-Americans have the most prior offenses, on average, followed by Native Americans. When comparing with Caucasian defendants, African-Americans have almost twice as many prior offenses, suggesting a strong proxy relationship between race and prior offenses. Asian defendants have the least prior offenses. This is an important result and illustrates how a system that pre-disposes certain races to prison can perpetuate that discriminatory trend by using those same variables to predict risk of recommitting another crime. The boxplot helps to visualize this relationship more clearly.

```
favstats(data = compas_final, priors_count ~ race)
```

##	race	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	African-American	0	0	2	6.00	38	4.042576	5.345310	4674	0
## 2	Asian	0	0	0	1.25	9	1.083333	1.888750	48	0
## 3	Caucasian	0	0	1	3.00	36	2.146462	3.472545	3250	0
## 4	Hispanic	0	0	1	2.00	29	1.806846	3.321139	818	0
## 5	Native American	0	0	1	3.50	22	3.185185	5.076621	27	0
## 6	Other	0	0	0	2.00	31	1.575439	2.949861	570	0

```
ggplot(data = compas_final,
       mapping = aes(x = as.factor(race), y = priors_count)) +
```

```
geom_boxplot() +
theme_minimal() +
labs(title = "Prior Offenses by Race",
x = "Race",
y = "Prior Offenses")
```

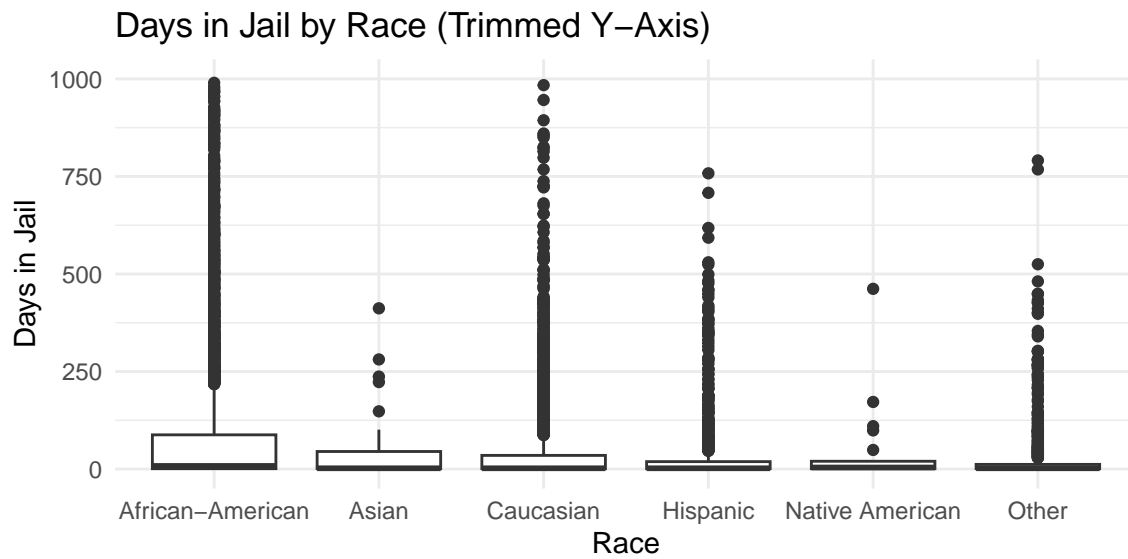


The same trend is observed when observing the difference in the days spent in jail. On average, African-Americans and Native Americans spend the most time in jail. In fact, their average jail time is more than double that of Caucasian defendants. Asian defendants spend the least time on average. When looking at the medians, 50% of the African-American defendants spent 9 or more days in jail, as compared to only 2 or more days for 50% of the Caucasian defendants. The boxplot below helps to visualize this better.

```
favstats(data = compas_final, days_in_jail ~ race)
```

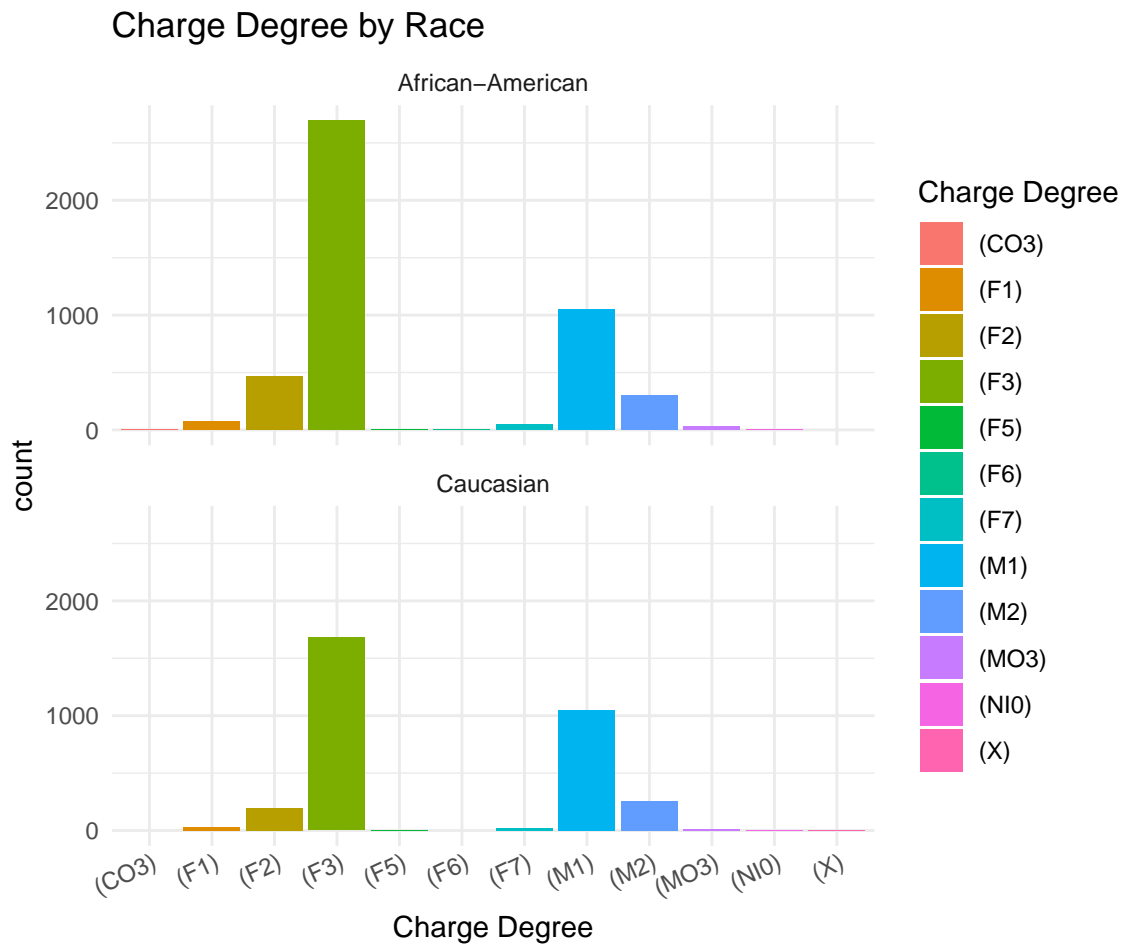
##	race	min	Q1	median	Q3	max	mean	sd	n	missing
## 1	African-American	0	1	9	106.00	21540	142.00941	509.0963	4674	0
## 2	Asian	0	1	2	45.00	412	41.10417	84.9779	48	0
## 3	Caucasian	0	1	2	37.00	5432	62.97723	229.2571	3250	0
## 4	Hispanic	0	1	2	20.00	2916	51.43032	208.9485	818	0
## 5	Native American	0	1	5	36.00	1536	142.59259	381.0188	27	0
## 6	Other	0	1	2	12.75	2440	42.08070	150.0490	570	0

```
ggplot(data = compas_final,
mapping = aes(x = as.factor(race), y = days_in_jail)) +
geom_boxplot() +
theme_minimal() +
labs(title = "Days in Jail by Race (Trimmed Y-Axis)",
x = "Race",
y = "Days in Jail") +
ylim(0,1000)
```



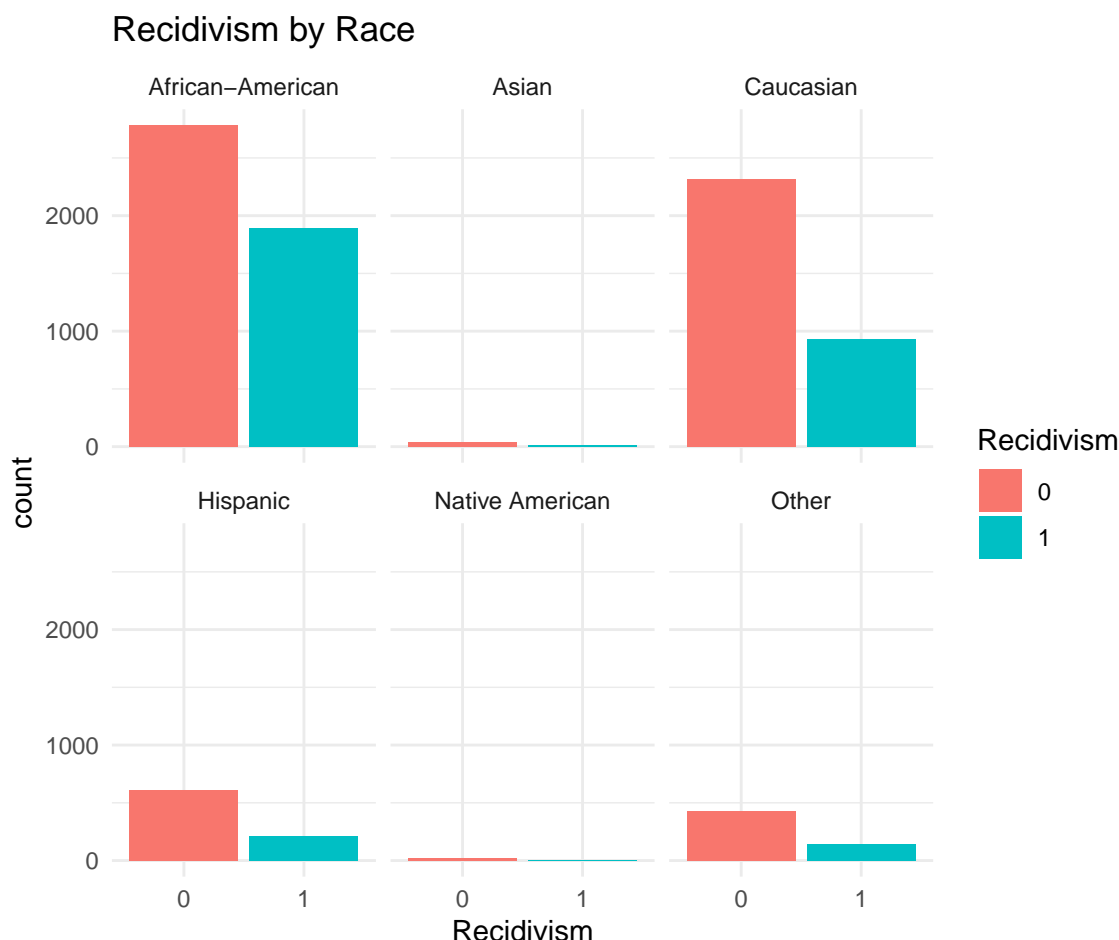
Next, let's look at the distribution of the charge degrees just for Black and White defendants. There are no notable distributional differences.

```
ggplot(data = compas_final_bw,
       mapping = aes(x = c_charge_degree, fill = c_charge_degree)) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~race, ncol = 1) +
  labs(title = "Charge Degree by Race",
       x = "Charge Degree",
       fill = "Charge Degree") +
  theme(axis.text.x = element_text(angle = 25, vjust = 1.2, hjust=1))
```



Finally, let's look at the response variable: `is_recid`. Most defendants do not re-commit a crime, although the ratio of those who don't against those who do appears to be larger for African-American defendants than for Caucasian defendants. Nevertheless, of important note is that the class imbalance is in the same direction.

```
ggplot(data = compas_final,
       mapping = aes(x = as.factor(is_recid), fill = as.factor(is_recid))) +
  geom_bar() +
  theme_minimal() +
  facet_wrap(~race) +
  labs(title = "Recidivism by Race",
       x = "Recidivism",
       fill = "Recidivism")
```



In conclusion, our demographic variable, **race** (A), has proxy relationships with some of the predictor variables (X). Even a group-blind classifier will not be entirely blind to race because of the correlations present and the information that can be gained from the proxy variables.

Additionally, we've learned that the distribution of those who recommit a crime within 2 years versus those who don't is not drastically different for Black v White defendants – the class imbalance is in the same direction, so ideally, the model should not misclassify defendants in different directions on the basis of race.

Next, let's see how the COMPAS tool actually performs with regard to race and whether it is in line with these expectations.

COMPAS Analysis

Before we can analyze the COMPAS performance, recall that that the decile scores are mapped to 3 different risk levels – low, medium, or high – yet the response variable only has 2 levels – 0 or 1. Let's recreate this mapping such that there are only 2 risk levels to line up with the 2 levels of the response variable. Instead, decile scores of 1 to 5 will be associated with lower risk of recidivism and decile scores of 6 to 10 will be associated with higher risk as displayed in the table below.

```
compas <- compas_final %>%
  mutate(risk = ifelse(decile_score %in% c(1, 2, 3, 4, 5), 'Lower', 'Higher'))
```

```

compas %>%
  dplyr::select(decile_score, risk) %>%
  rename("Risk" = risk) %>%
  group_by(Risk) %>%
  summarise("Min" = min(decile_score),
            "Max" = max(decile_score)) %>%
  arrange(Min) %>%
  kable(booktabs = TRUE)

```

Risk	Min	Max
Lower	1	5
Higher	6	10

Now, let's recreate the table from Chapter 1 using this data set to assess the false positive and false negative rates for Black v White defendants.

While the numbers are different, which could likely be due to a host of reasons such as different subsets of the data set, the same trends are evident. 16.34% of White defendants who did not re-offend are labelled as higher risk, compared to more than twice as many Black defendants (37.71%). Similarly, 62.26% of White defendants who do re-offend are labelled as lower risk, compared to 39.01% of Black defendants. This is in line with ProPublica's findings, and is alarming given that race was not included in the model.

```

compas_table <- compas %>%
  filter(race %in% c("African-American", "Caucasian")) %>%
  dplyr::select(race, risk, is_recid) %>%
  rename("Risk" = risk,
         "Race" = race) %>%
  group_by(Race, is_recid) %>%
  mutate(Total = n()) %>%
  group_by(Risk, Race, Total) %>%
  summarise("Reoffended" = count(is_recid == 1),
            "Did Not Reoffend" = count(is_recid == 0)) %>%
  pivot_longer(cols = c("Reoffended", "Did Not Reoffend"),
               names_to = "Recidivism") %>%
  pivot_wider(
    id_cols = c("Risk", "Recidivism", "Total"),
    names_from = "Race",
    values_from = value
  ) %>%
  rename("Black" = `African-American`,
         "White" = `Caucasian`) %>%
  mutate(Black = round(100 * Black / Total, 2),
         White = round(100 * White / Total, 2)) %>%
  dplyr::select(-Total) %>%
  group_by(Risk, Recidivism) %>%
  summarize(Black = max(Black, na.rm = TRUE),
            White = max(White, na.rm = TRUE)) %>%
  filter((Risk == "Higher" & Recidivism == "Did Not Reoffend") |
         (Risk == "Lower" & Recidivism == "Reoffended"))
)

compas_table %>%
  kable(booktabs = TRUE)

```

Risk	Recidivism	Black	White
Higher	Did Not Reoffend	37.71	16.34
Lower	Reoffended	39.01	62.26

The results in this section also show that the model is more wrong in predicting whether defendants will re-offend versus predicting defendants who do not re-offend. This is expected because of the class imbalance we observed when performing exploratory data analysis – there are more defendants that don’t re-offend, so the model maximizes performance for those defendants.

However, it also raises a question of what type of prediction is more important: the risk of recidivism or the risk of non-recidivism. Is wrongly attributing a defendant as higher risk may cause or wrongly attributing a defendant as lower risk worse?

The table below shows the confusion matrix of the COMPAS model as a whole, including all races. As observed, a larger proportion (49.25%) of defendants who re-offended are incorrectly labelled as lower risk in comparison to the proportion (25.23%) that did not re-offend but are incorrectly labelled as higher risk.

Observe that while the overall FPR is 25.23%, it’s much higher for Black defendants (37.71%) and much lower for White defendants (16.34%). Similarly, while the overall FNR is 49.25%, it’s much higher for White defendants (62.26%) and much lower for Black defendants (39.01%).

```
compas %>%
  dplyr::select(risk, is_recid) %>%
  rename("Risk" = risk) %>%
  group_by(is_recid) %>%
  mutate(Total = n()) %>%
  group_by(Risk, Total) %>%
  summarise("Reoffended" = count(is_recid == 1),
            "Did Not Reoffend" = count(is_recid == 0)) %>%
  mutate(Reoffended = round(100 * Reoffended / Total, 2),
         `Did Not Reoffend` = round(100 * `Did Not Reoffend` / Total, 2)) %>%
  dplyr::select(-Total) %>%
  group_by(Risk) %>%
  summarize(Reoffended = max(Reoffended, na.rm = TRUE),
            `Did Not Reoffend` = max(`Did Not Reoffend`, na.rm = TRUE)) %>%
  kable()
```

Risk	Reoffended	Did Not Reoffend
Higher	50.75	25.23
Lower	49.25	74.77

The table below replicates the above table with raw numbers instead of proportions. This reveals that the model has an overall accuracy of 66.61%, but the analysis above reveals the discrepancies (and unfairness) in the model results for Black v White defendants.

However, if we classified every observation in the majority class (no recidivism), we’d expect an accuracy of 66.04%. The COMPAS model, thus, is not useful in a predictive sense, and on the contrary, introduces more bias into the judicial system.

```
compas %>%
  dplyr::select(risk, is_recid) %>%
  rename("Risk" = risk) %>%
```



```
group_by(Risk) %>%
  summarise("Reoffended" = count(is_recid == 1),
            "Did Not Reoffend" = count(is_recid == 0)) %>%
  kable()
```

Risk	Reoffended	Did Not Reoffend
Higher	1618	1564
Lower	1570	4635

```
(1618 + 4635) / 9387
```

```
## [1] 0.666134
```

```
count(compas$is_recid == 0) / 9387
```

```
##      n_TRUE
```

```
## 0.6603814
```

Now that we've recreated the table from Chapter 1, let's examine the false negative and false positive rates for all the races in the data set.

Asian defendants who did not re-offend were the least likely to be labelled as higher risk – Black defendants were the most likely. Conversely, excluding the “Other” group, White defendants who re-offended were the most likely to be labelled as lower risk – Native Americans were the least likely. This, and previous analysis, suggest disparities with favorable outcomes for white and Asian defendants, and unfavorable outcomes for Black and Native American defendants.

```
compas %>%
  dplyr::select(race, risk, is_recid) %>%
  rename("Risk" = risk,
         "Race" = race) %>%
  group_by(Race, is_recid) %>%
  mutate(Total = n()) %>%
  group_by(Risk, Race, Total) %>%
  summarise("Reoffended" = count(is_recid == 1),
            "Did Not Reoffend" = count(is_recid == 0)) %>%
  pivot_longer(cols = c("Reoffended", "Did Not Reoffend"),
               names_to = "Recidivism") %>%
  pivot_wider(
    id_cols = c("Risk", "Recidivism", "Total"),
    names_from = "Race",
    values_from = value
  ) %>%
  rename("Black" = `African-American`,
         "White" = `Caucasian`) %>%
  mutate(Black = round(100 * Black / Total, 2),
         White = round(100 * White / Total, 2),
         Asian = round(100 * Asian / Total, 2),
         `Native American` = round(100 * `Native American` / Total, 2),
         Other = round(100 * Other / Total, 2)) %>%
  dplyr::select(-Total) %>%
  group_by(Risk, Recidivism) %>%
  summarize(Black = max(Black, na.rm = TRUE),
            White = max(White, na.rm = TRUE),
```

```

Asian = max(Asian, na.rm = TRUE),
`Native American` = max(`Native American`, na.rm = TRUE),
Other = max(Other, na.rm = TRUE)) %>%
filter((Risk == "Higher" & Recidivism == "Did Not Reoffend") |
(Risk == "Lower" & Recidivism == "Reoffended"))
) %>%
kable(booktabs = TRUE)

```

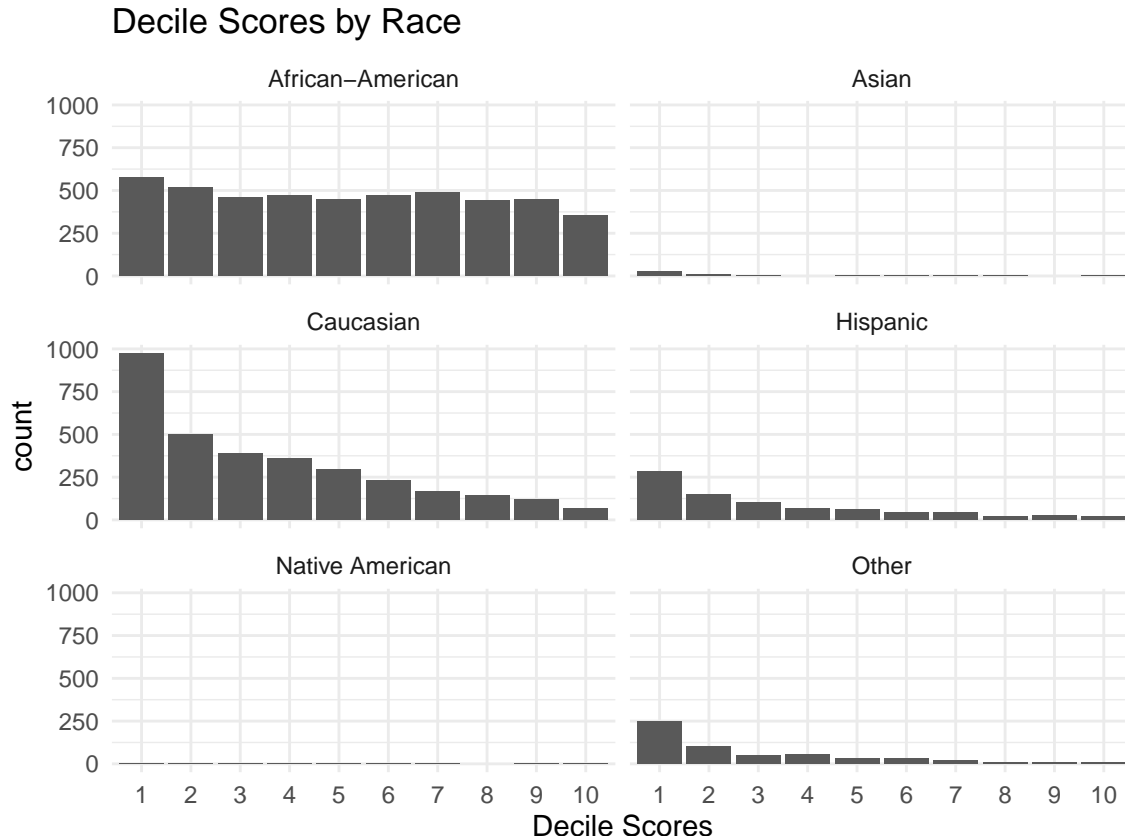
Risk	Recidivism	Black	White	Asian	Native American	Other
Higher	Did Not Reoffend	37.71	16.34	7.89	21.05	9.77
Lower	Reoffended	39.01	62.26	50.00	25.00	74.29

Finally, to visualize these results better, let's plot the distribution of the decile scores by race. Observe that while we observed a similar distribution of recidivism for all races, the African-American decile scores are distributed differently from the Caucasian decile scores. The decile scores for the Caucasian defendants is quite right-skewed in comparison to that of the African-Americans which appears to be more evenly distributed, further emphasizing the racial disparity in the risk scores.

```

ggplot(data = compas,
mapping = aes(x = as.factor(decile_score))) +
geom_bar() +
theme_minimal() +
facet_wrap(~race, ncol = 2) +
labs(title = "Decile Scores by Race",
x = "Decile Scores")

```



We have a thorough understanding of the data set now and the results we would expect from the COMPAS tool.

We're ready to proceed with modeling now. First, we will fit a logistic regression, which uses the standard ML approach discussed in Chapter 2. We expect it to behave similarly to the COMPAS tool. Then, we will define a fairness metric and fit a Seldonian classification algorithm on this data set. We expect to observe less disparities along racial lines, but perhaps with a slight trade-off in overall accuracy.

Logistic Regression

Logistic regression is a statistical generalized linear model (GLM) that is particularly designed for predicting dichotomous/ binary outcomes such as ours. We will use logistic regression to model the probability of a defendant re-committing a crime within two years in Broward County, Florida. We'll then set cutoffs to divide the probabilities into 2 bins: 0 for 'no' and 1 for 'yes' based on the probability predictions. This will also allow us to analyze which features may be most important in predicting recidivism.

Recall that if we classified every observation in the majority class, we'd expect an accuracy 66%. This serves as a benchmark for the race-blind logistic regression model we will implement. Note that logistic regression follows the standard ML process [will describe in more detail in the thesis body] and is one of the most widely used classification algorithms – this will allow us to assess how we might expect state-of-the-art traditional algorithms that do not take fairness guarantees into account to perform.

Check for Missing Data

Linear models do not handle missing observations well. Let's ensure that there are no missing observations. There were no missing observations!

```
comopas <- tidyr::drop_na(compas)
count(compas)
```

```
##      n
## 1 9387
```

Train and Test Split

Before we proceed to fitting the models, we need to perform a train/ test split. The reason for this is to be able to test how our model would perform on unseen data and to avoid over-fitting. Therefore, we will train our models using the **train** data, and then assess performance on “new” data, that is, the **test** set. Let's partition 70% of the data into the train set and 30% into the test set. We will use randomization without replacement for this.

```
set.seed(123)
n <- nrow(compas)
train_index <- sample(1:n, 0.70 * n)
test_index <- setdiff(1:n, train_index)
train <- compas[train_index, ]
test <- compas[test_index, ]
```

There are now 6570 observations in the train set and 2817 observations in the test set. Let's ensure that the distribution of the response variable is preserved in each set. There split appears to be stratified so there is no concern.

```
tally(train$is_recid)
```

```
## X
##    0    1
## 4333 2237
```

```
tally(test$is_recid)
```

```
## X
##    0    1
## 1866  951
```

Finally, let's also make sure that there are enough observations in each race category for both the train and test splits. There is concern for the Native American and Asian race categories, which have very few observations. We may fit a separate model with just the Caucasian and African-American offenders if this poses a challenge.

```
tally(train$race)
```

```
## X
## African-American      Asian      Caucasian      Hispanic
##           3270           34           2280           573
## Native American      Other
##           18           395
```

```
tally(test$race)
```

```
## X
## African-American      Asian      Caucasian      Hispanic
##           1404           14           970           245
## Native American      Other
##           9           175
```

We're ready for model fitting!

Modeling

- run logistic regression and analyze results.
- refer to compas analysis (racial bias in compas).

Seldonian Classification

run seldonian framework on the data set and analyze results.

Results

Cox Proportional Hazards model for the length of time to recidivate? probably not necessary since the focus is on classification.

Northeastpoint: <https://www.equivant.com/wp-content/uploads/Practitioners-Guide-to-COMPAS-Core-040419.pdf>