

Thesis Simulation Document for Chapter 4

Dasha Asienaga

2024-03-13

Contents

Data Generation Mechanism #1 (class balance)	1
Reading in the Data	1
Data Subsetting	2
Generating the Parent Simulation Data Set	4
Examining Distributions of the Recidivism in the Parent Data Set	5
Assessing Baseline Predictive Performance of the Parent Data Set	6
Data Generation Mechanism #2 (better predictive performance)	7
Generating the Parent Simulation Data Set	7
Examining Distributions of the Recidivism in the Parent Data Set	7
Assessing Baseline Predictive Performance of the Parent Data Set	9
Data Generation Mechanism #3 (class balance and better predictive performance)	10
Generating the Parent Data Set	10
Examining Distributions of the Recidivism in the Parent Data Set	11
Assessing Baseline Predictive Performance of the Parent Data Set	12
Data Generation Mechanism #4 (class balance, but with some error)	13
Generating the Parent Data Set	13
Examining Distributions of the Recidivism in the Parent Data Set	14
Assessing Baseline Predictive Performance of the Parent Data Set	15

This file is intended to contain all the code and information to set up the simulation study and supplement Chapter 4.

Data Generation Mechanism #1 (class balance)

We're interested in creating a data set that has 50-50 class balance, even across the demographic group, and also has better predictive performance than the COMPAS tool. For this set-up, we will only use 2 variables from the COMPAS data set: 1 continuous variable and 1 categorical variable.

Reading in the Data

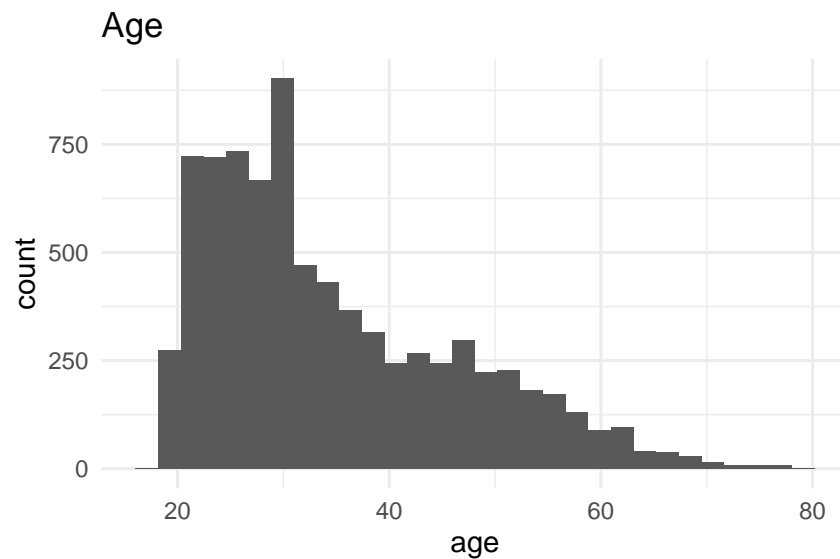
First, let's read in the data.

```
compas_path <- "/home/dasienga24/Statistics-Senior-Honors-Thesis/Data Sets/COMPAS/compas_seldonian_bw.c  
compas_sim <- read.csv(compas_path)
```

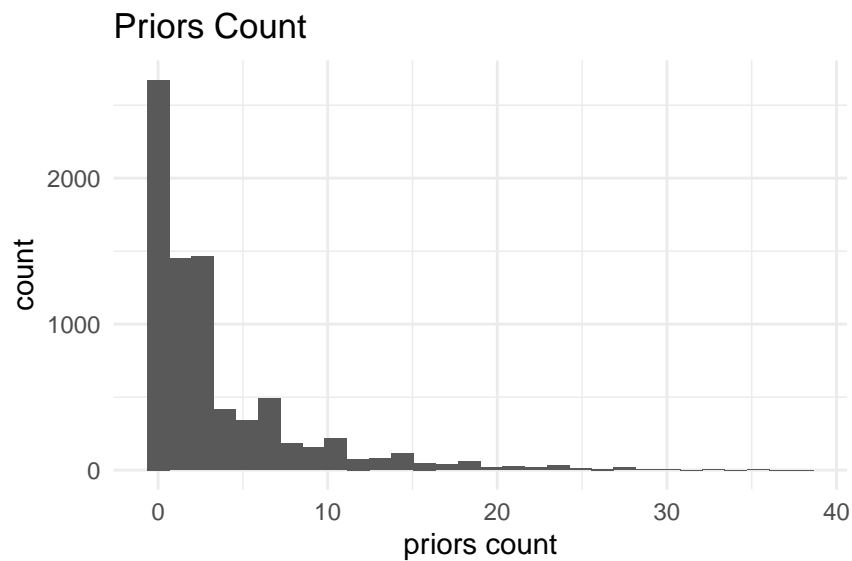
Data Subsetting

Next, let's plot the distributions of the continuous variables to choose which one we'll proceed with.

```
compas_sim %>%  
  ggplot(mapping = aes(x = age)) +  
  geom_histogram() +  
  theme_minimal() +  
  labs(title = "Age")
```



```
compas_sim %>%  
  ggplot(mapping = aes(x = priors_count)) +  
  geom_histogram() +  
  theme_minimal() +  
  labs(title = "Priors Count",  
       x = "priors count")
```



Because age has more variation, we'll use it as our continuous variable. We'll convert `priors_count` into a categorical variable.

```
compas_sim <- compas_sim %>%  
  mutate(prior_offense = ifelse(priors_count > 0, 1, 0)) %>%  
  dplyr::select(c(race, prior_offense, age, is_recid))
```

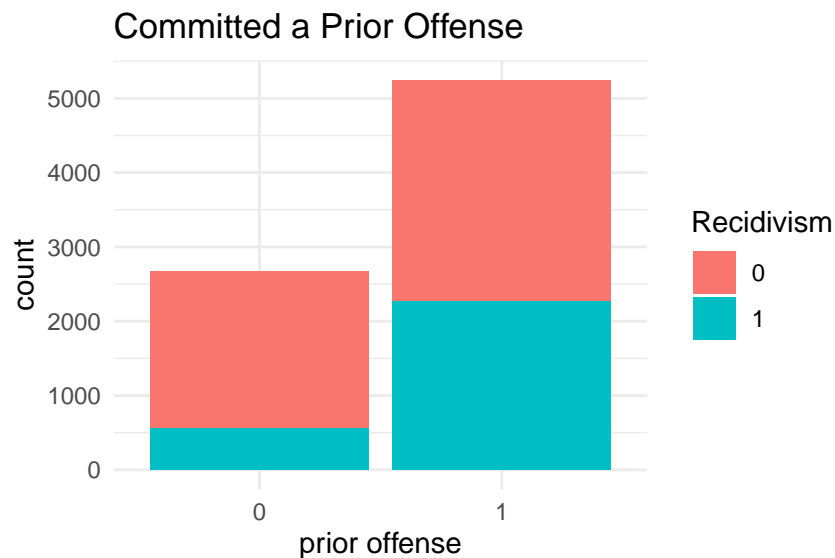
age seems to be a useful predictor for recidivism.

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

```
##   is_recid min Q1 median Q3 max   mean    sd  n missing  
## 1      0  18 26   33 44  80 35.92591 12.24397 5102      0  
## 2      1  18 24   29 37  78 32.25797 10.57842 2822      0
```

Whether a defendant has committed a prior offense or not appears to be a useful predictor for recidivism as well.

```
compas_sim %>%  
  ggplot(mapping = aes(x = as.factor(prior_offense), fill = as.factor(is_recid))) +  
  geom_bar() +  
  theme_minimal() +  
  labs(title = "Committed a Prior Offense",  
       fill = "Recidivism",  
       x = "prior offense")
```



We'll proceed with these 2 variables – age and prior_offense for the simulation study. A glimpse of the data is shown below.

```
compas_sim <- compas_sim %>%  
  mutate(prior_offense = as.factor(prior_offense),  
         is_recid = as.factor(is_recid))
```

```
glimpse(compas_sim)

## Rows: 7,924
## Columns: 4
## $ race      <chr> "African-American", "African-American", "Caucasian", "Ca~
## $ prior_offense <fct> 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0,~
## $ age        <int> 34, 24, 41, 39, 20, 26, 27, 23, 37, 22, 41, 47, 31, 25, ~
## $ is_recid    <fct> 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1,~
head(compas_sim)
```

```
##           race prior_offense age is_recid
## 1 African-American          0  34         1
## 2 African-American          1  24         1
## 3      Caucasian            1  41         1
## 4      Caucasian            0  39         0
## 5      Caucasian            0  20         0
## 6      Caucasian            0  26         0
```

Generating the Parent Simulation Data Set

We want a setting with 50-50 class balance for each combination of race and recidivism status. To achieve that, we'll perform sample observations with replacement. Let's create a data set with 1250 observations in each of these 4 groups, hence, 5000 observations total.

First, let's subset these 4 groups.

```
compas_b_y <- compas_sim %>%
  filter(race == "African-American" & is_recid == 1)

compas_b_n <- compas_sim %>%
  filter(race == "African-American" & is_recid == 0)

compas_w_y <- compas_sim %>%
  filter(race == "Caucasian" & is_recid == 1)

compas_w_n <- compas_sim %>%
  filter(race == "Caucasian" & is_recid == 0)
```

Next, let's randomly sample 1250 observations from each of these groups.

```
compas_b_y_balanced <- compas_b_y[sample(nrow(compas_b_y), 1250, replace = TRUE),]
compas_b_n_balanced <- compas_b_n[sample(nrow(compas_b_n), 1250, replace = TRUE),]
compas_w_y_balanced <- compas_w_y[sample(nrow(compas_w_y), 1250, replace = TRUE),]
compas_w_n_balanced <- compas_w_n[sample(nrow(compas_w_n), 1250, replace = TRUE),]
```

Finally, let's union all these together into a single data set.

```
compas_sim_balanced <- rbind(compas_b_y_balanced,
                             compas_b_n_balanced,
                             compas_w_y_balanced,
                             compas_w_n_balanced)
```

Let's also shuffle the data set row orderings to aid the machine learning algorithms later.

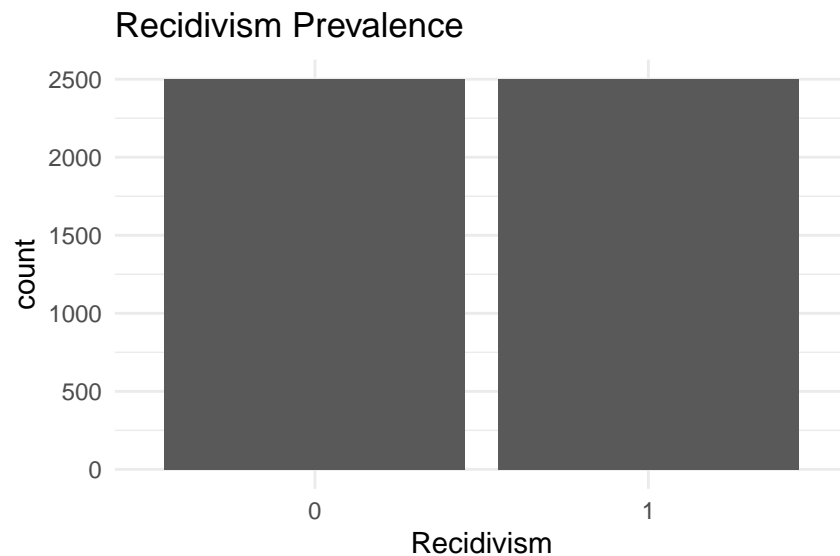
```
compas_sim_balanced <- compas_sim_balanced[sample(nrow(compas_sim_balanced),
                                                  5000,
                                                  replace = FALSE),]
```

The parent data set is now ready.

Examining Distributions of the Recidivism in the Parent Data Set

The bar plot below shows that we've achieve perfect class balance.

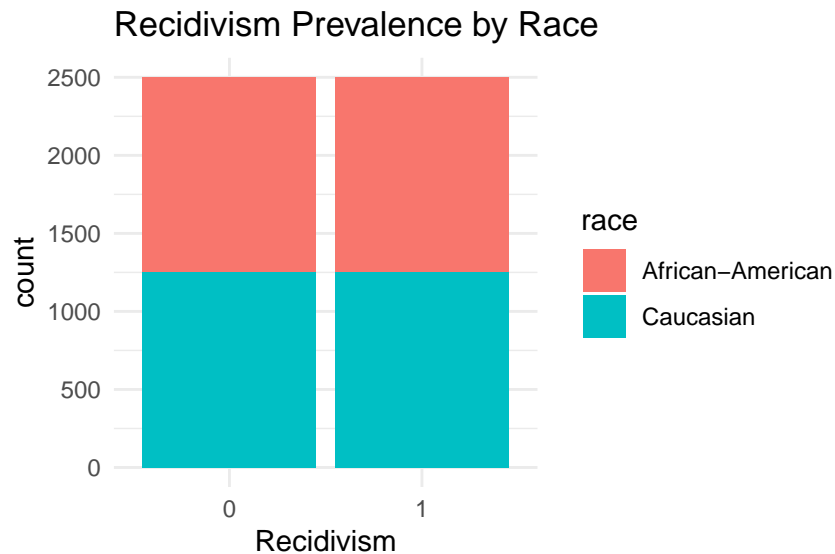
```
compas_sim_balanced %>%
  ggplot(mapping = aes(x = is_recid)) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Recidivism",
       title = "Recidivism Prevalence")
```



The bar plot below reveals that the balance is preserved by race as well.

```
compas_sim_balanced %>%
  ggplot(mapping = aes(x = is_recid, fill = race)) +
  geom_bar() +
  theme_minimal() +
```

```
labs(x = "Recidivism",
     title = "Recidivism Prevalence by Race")
```



Assessing Baseline Predictive Performance of the Parent Data Set

We want to make sure that our data set also has good predictive performance. We'll fit a logistic regression and assess baseline accuracy.

```
glm1 <- glm(is_recid ~ age + prior_offense,
            data = compas_sim_balanced,
            family = binomial(logit))
```

```
msummary(glm1)
```

```
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.175836   0.100091   1.757   0.079 .
## age         -0.025471   0.002551  -9.983  <2e-16 ***
## prior_offense1 1.022128   0.064427  15.865  <2e-16 ***
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 6931.5  on 4999  degrees of freedom
## Residual deviance: 6573.7  on 4997  degrees of freedom
## AIC: 6579.7
##
## Number of Fisher Scoring iterations: 4
```

```
glm1augment <- glm1 %>%
  broom::augment(type.predict = "response")
glm1augment <- mutate(glm1augment, binprediction = round(.fitted, 0))
with(glm1augment, table(is_recid, binprediction))
```

```
##      binprediction
## is_recid    0    1
##          0 1349 1151
##          1  790 1710
```

```
(1334 + 1713)/5000
```

```
## [1] 0.6094
```

This model has 61% accuracy.

Data Generation Mechanism #2 (better predictive performance)

Using the balanced data set created above, we'll simulate Y values that have a stronger correlation with the 2 predictors. This is in an aim to hopefully get better predictive performance, in addition to the class balance.

Generating the Parent Simulation Data Set

```
# define a linear combination of predictors as desired
linear_combination = - 15 - 24 * compas_sim_balanced$age + ifelse(compas_sim_balanced$prior_offense == 1, 10, 0)

# pass through an inverse-logit function
probs = exp(linear_combination) / (1 + exp(linear_combination))

# generate 5000 Bernoulli RVs for y
is_recid_sim = rbinom(5000, 1, probs)

# join to original data frame
compas_sim_balanced_2 <- cbind(compas_sim_balanced, is_recid_sim)
```

There are 2487 defendants who did not recidivate in this data set.

```
count(compas_sim_balanced_2$is_recid_sim == 0)
```

```
## n_TRUE
##    2469
```

There are 2513 defendants who recidivated in this data set.

```
count(compas_sim_balanced_2$is_recid_sim == 1)
```

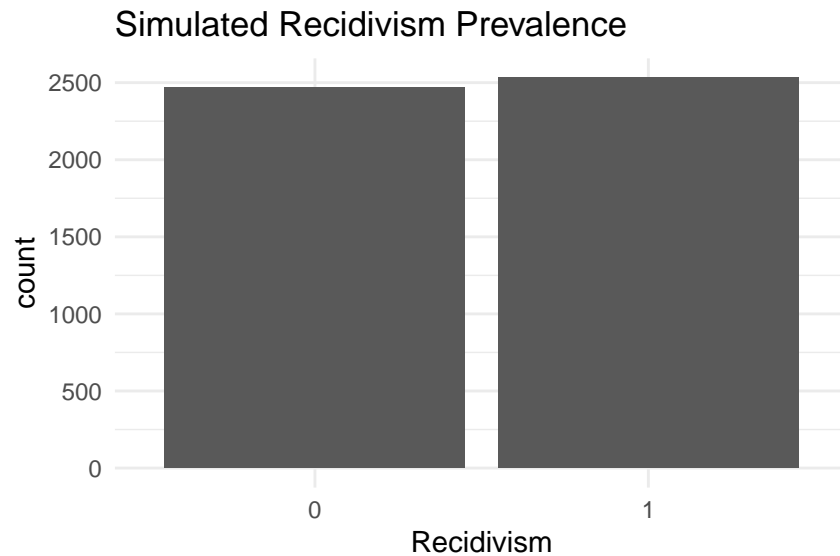
```
## n_TRUE
##    2531
```

The data set seems pretty balanced. Let's look at this distribution in more detail, though.

Examining Distributions of the Recidivism in the Parent Data Set

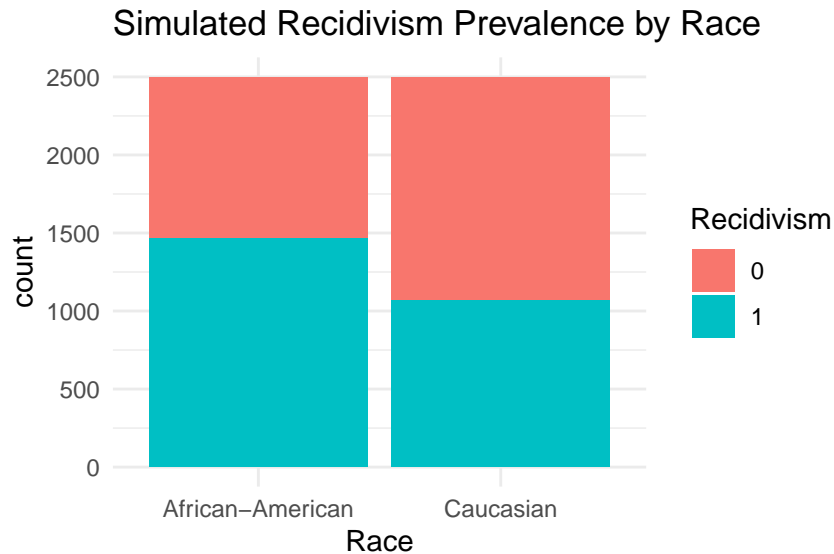
The bar plot below confirms the balanced nature of this new data set.

```
compas_sim_balanced_2 %>%
  ggplot(mapping = aes(x = as.factor(is_recid_sim))) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Recidivism",
       title = "Simulated Recidivism Prevalence")
```



However, the balance is not perfectly achieved by race. Some of the proxy relationships present in the predictor variables have influenced the distribution of the Y variable by race.

```
compas_sim_balanced_2 %>%
  ggplot(mapping = aes(x = race, fill = as.factor(is_recid_sim))) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Race",
       title = "Simulated Recidivism Prevalence by Race",
       fill = "Recidivism")
```

Assessing Baseline Predictive Performance of the Parent Data Set

We want to make sure that our data set also has good predictive performance. We'll fit a logistic regression and assess baseline accuracy.

```
glm2 <- glm(is_recid_sim ~ age + prior_offense,
            data = compas_sim_balanced_2,
            family = binomial(logit))
```

```
msummary(glm2)
```

```
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    336.97   12527.07   0.027   0.979
## age           -18.81    611.46  -0.031   0.975
## prior_offense1  434.93   14244.70   0.031   0.976
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 6930.703  on 4999  degrees of freedom
## Residual deviance:   66.508  on 4997  degrees of freedom
## AIC: 72.508
##
## Number of Fisher Scoring iterations: 25
```

```
glm2augment <- glm2 %>%
  broom::augment(type.predict = "response")
glm2augment <- mutate(glm2augment, binprediction = round(.fitted, 0))
with(glm2augment, table(is_recid_sim, binprediction))
```

```
##           binprediction
## is_recid_sim    0     1
##              0 2452   17
##              1    0 2531
```

```
(2475 + 2513) / 5000
```

```
## [1] 0.9976
```

This data set has 99.76% accuracy, although the distribution of race is affected. There is barely any error, so this may not be too useful for our fairness definition.

Data Generation Mechanism #3 (class balance and better predictive performance)

Let's introduce some balance across racial lines.

Generating the Parent Data Set

First, let's subset these 4 groups.

```
compas_b_y <- compas_sim_balanced_2 %>%  
  filter(race == "African-American" & is_recid_sim == 1)  
  
compas_b_n <- compas_sim_balanced_2 %>%  
  filter(race == "African-American" & is_recid_sim == 0)  
  
compas_w_y <- compas_sim_balanced_2 %>%  
  filter(race == "Caucasian" & is_recid_sim == 1)  
  
compas_w_n <- compas_sim_balanced_2 %>%  
  filter(race == "Caucasian" & is_recid_sim == 0)
```

Next, let's randomly sample 1250 observations with replacement from each of the 4 groups in the new data set.

```
compas_b_y_balanced <- compas_b_y[sample(nrow(compas_b_y), 1250, replace = TRUE),]  
compas_b_n_balanced <- compas_b_n[sample(nrow(compas_b_n), 1250, replace = TRUE),]  
compas_w_y_balanced <- compas_w_y[sample(nrow(compas_w_y), 1250, replace = TRUE),]  
compas_w_n_balanced <- compas_w_n[sample(nrow(compas_w_n), 1250, replace = TRUE),]
```

Finally, let's union all these together into a single data set.

```
compas_sim_balanced_3 <- rbind(compas_b_y_balanced,  
                               compas_b_n_balanced,  
                               compas_w_y_balanced,  
                               compas_w_n_balanced)
```

Let's also shuffle the data set row orderings to aid the machine learning algorithms later.

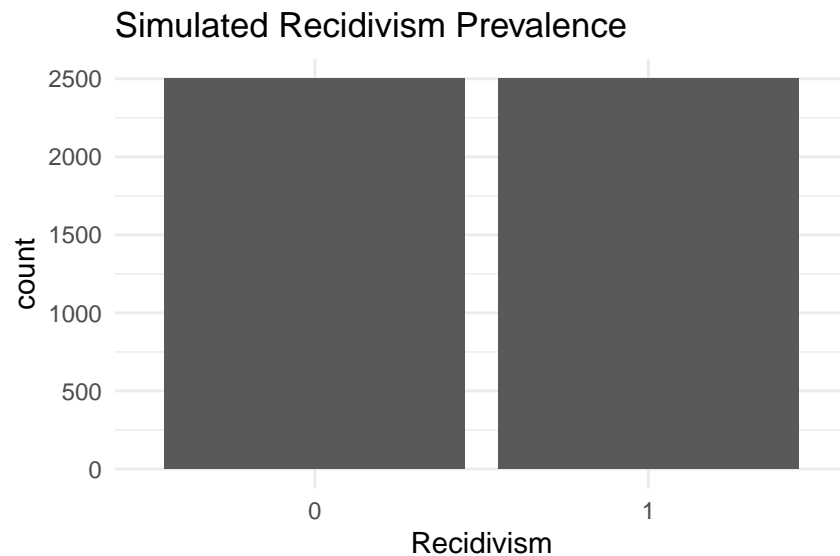
```
compas_sim_balanced_3 <- compas_sim_balanced_3[sample(nrow(compas_sim_balanced_3),  
                                                      5000,  
                                                      replace = FALSE),]
```

The data set is now ready.

Examining Distributions of the Recidivism in the Parent Data Set

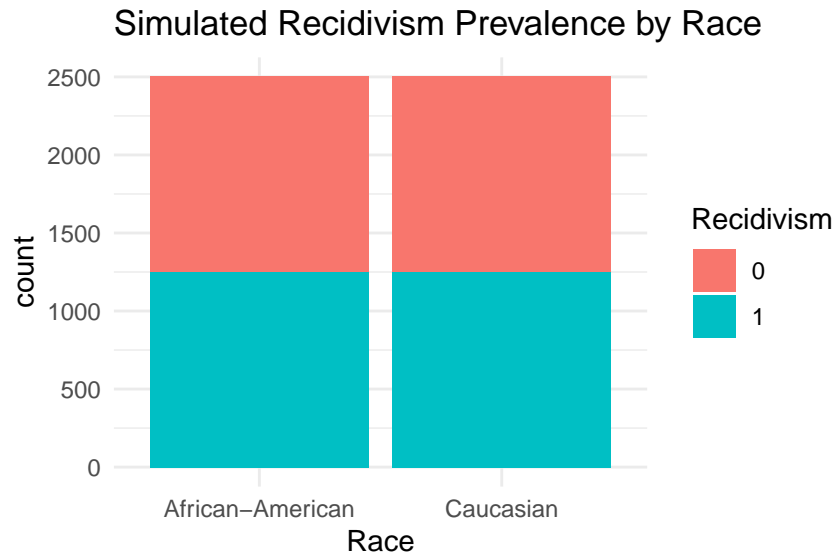
The bar plot below confirms the balanced nature of this new data set.

```
compas_sim_balanced_3 %>%  
  ggplot(mapping = aes(x = as.factor(is_recid_sim))) +  
  geom_bar() +  
  theme_minimal() +  
  labs(x = "Recidivism",  
        title = "Simulated Recidivism Prevalence")
```



The balance is also preserved by race.

```
compas_sim_balanced_3 %>%  
  ggplot(mapping = aes(x = race, fill = as.factor(is_recid_sim))) +  
  geom_bar() +  
  theme_minimal() +  
  labs(x = "Race",  
        title = "Simulated Recidivism Prevalence by Race",  
        fill = "Recidivism")
```



Assessing Baseline Predictive Performance of the Parent Data Set

We want to make sure that our data set also has good, but not perfect, predictive performance. We'll fit a logistic regression and assess baseline accuracy.

```
glm3 <- glm(is_recid_sim ~ age + prior_offense,
            data = compas_sim_balanced_3,
            family = binomial(logit))
```

```
msummary(glm3)
```

```
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    335.47   12636.01   0.027   0.979
## age           -18.75    616.03  -0.030   0.976
## prior_offense1  433.84   14364.60   0.030   0.976
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 6931.472  on 4999  degrees of freedom
## Residual deviance:   73.365  on 4997  degrees of freedom
## AIC: 79.365
##
## Number of Fisher Scoring iterations: 25
```

```
glm3augment <- glm3 %>%
  broom::augment(type.predict = "response")
glm3augment <- mutate(glm3augment, binprediction = round(.fitted, 0))
with(glm3augment, table(is_recid_sim, binprediction))
```

```
##              binprediction
## is_recid_sim    0      1
##              0 2481   19
##              1    0 2500
```

The accuracy is a little bit too good. For the final mechanism, we will follow this framework, but weaken the strengths of the correlations.

Data Generation Mechanism #4 (class balance, but with some error)

Generating the Parent Data Set

```
# define a linear combination of predictors as desired
linear_combination = - 0.34 - 16 * compas_sim_balanced$age + ifelse(compas_sim_balanced$prior_offense == 1, 10, 0)

# pass through an inverse-logit function
probs = exp(linear_combination) / (1 + exp(linear_combination))

# generate 5000 Bernoulli RVs for y
is_recid_sim = rbinom(5000, 1, probs)

# join to original data frame
compas_sim_balanced_4 <- cbind(compas_sim_balanced, is_recid_sim)
```

There are 2739 defendants who did not recidivate in this data set.

```
count(compas_sim_balanced_4$is_recid_sim == 0)
```

```
## n_TRUE
##    2723
```

There are 2261 defendants who recidivated in this data set.

```
count(compas_sim_balanced_2$is_recid_sim == 1)
```

```
## n_TRUE
##    2531
```

Next, let's induce balance into this data set. First, let's subset these 4 groups.

```
compas_b_y <- compas_sim_balanced_4 %>%
  filter(race == "African-American" & is_recid_sim == 1)

compas_b_n <- compas_sim_balanced_4 %>%
  filter(race == "African-American" & is_recid_sim == 0)

compas_w_y <- compas_sim_balanced_4 %>%
  filter(race == "Caucasian" & is_recid_sim == 1)

compas_w_n <- compas_sim_balanced_4 %>%
  filter(race == "Caucasian" & is_recid_sim == 0)
```

Next, let's randomly sample 1250 observations with replacement from each of the 4 groups in the new data set.

```
compas_b_y_balanced <- compas_b_y[sample(nrow(compas_b_y), 1250, replace = TRUE),]  
compas_b_n_balanced <- compas_b_n[sample(nrow(compas_b_n), 1250, replace = TRUE),]  
compas_w_y_balanced <- compas_w_y[sample(nrow(compas_w_y), 1250, replace = TRUE),]  
compas_w_n_balanced <- compas_w_n[sample(nrow(compas_w_n), 1250, replace = TRUE),]
```

Finally, let's union all these together into a single data set.

```
compas_sim_balanced_4 <- rbind(compas_b_y_balanced,  
                               compas_b_n_balanced,  
                               compas_w_y_balanced,  
                               compas_w_n_balanced)
```

Let's also shuffle the data set row orderings to aid the machine learning algorithms later.

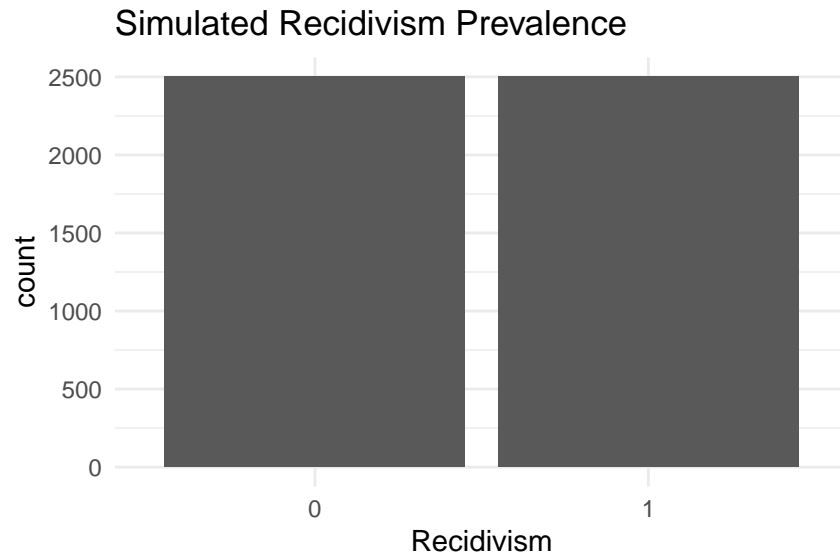
```
compas_sim_balanced_4 <- compas_sim_balanced_4[sample(nrow(compas_sim_balanced_4),  
                                                      5000,  
                                                      replace = FALSE),]
```

The data set is now ready.

Examining Distributions of the Recidivism in the Parent Data Set

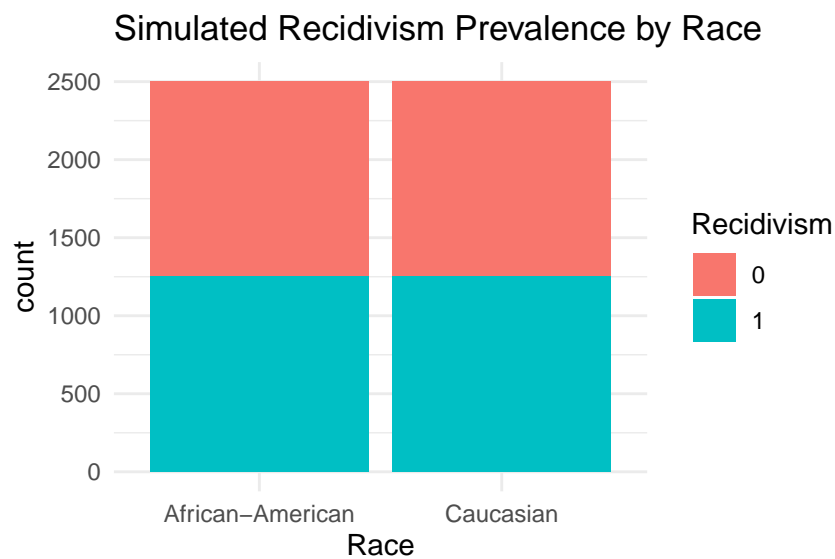
The bar plot below confirms the balanced nature of this new data set.

```
compas_sim_balanced_4 %>%  
  ggplot(mapping = aes(x = as.factor(is_recid_sim))) +  
  geom_bar() +  
  theme_minimal() +  
  labs(x = "Recidivism",  
       title = "Simulated Recidivism Prevalence")
```



The balance is also preserved by race.

```
compas_sim_balanced_4 %>%
  ggplot(mapping = aes(x = race, fill = as.factor(is_recid_sim))) +
  geom_bar() +
  theme_minimal() +
  labs(x = "Race",
       title = "Simulated Recidivism Prevalence by Race",
       fill = "Recidivism")
```



Assessing Baseline Predictive Performance of the Parent Data Set

We want to make sure that our data set also has good, but not perfect, predictive performance. We'll fit a logistic regression and assess baseline accuracy.

```
glm4 <- glm(is_recid_sim ~ age + prior_offense,
            data = compas_sim_balanced_4,
```

```

family = binomial(logit))

msummary(glm4)

## Coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)    674.52   15936.72   0.042   0.966
## age           -36.55    809.77  -0.045   0.964
## prior_offense1  696.28   15539.46   0.045   0.964
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 6.9315e+03  on 4999  degrees of freedom
## Residual deviance: 4.5350e-06  on 4997  degrees of freedom
## AIC: 6
##
## Number of Fisher Scoring iterations: 25

glm4augment <- glm4 %>%
  broom::augment(type.predict = "response")
glm4augment <- mutate(glm4augment, binprediction = round(.fitted, 0))
with(glm4augment, table(is_recid_sim, binprediction))

##           binprediction
## is_recid_sim    0     1
##           0 2500     0
##           1     0 2500

```

Hmm, maybe we need to introduce an error term.

```

mean(linear_combination)

## [1] -145.3272

sd(linear_combination)

## [1] 335.1698

```

Maybe it's okay to not have class balance? Think about this some more (goal of finalizing it today).