



Midterm Presentation

Kathryn Doorley

Math 421

Loading and Cleaning Data

- Used haven library to load the SAS dataset
- Filtered for patients discharged in 2018
- Selected only 49 columns
- Removed the following columns with missing values
 - Payfix
 - Preopday
 - Obs_hour
 - Nicu_day



Exploratory Data Analysis

- What month saw the most amount of patients?

```
```{r}
midterm_df_final %>%
 group_by(moa) %>%
 count()
```
```

| moa
<dbl> | n
<int> |
|--------------|------------|
| 1 | 11309 |
| 2 | 10259 |
| 3 | 11073 |
| 4 | 11174 |
| 5 | 11404 |
| 6 | 10690 |
| 7 | 10962 |
| 8 | 11199 |
| 9 | 10666 |
| 10 | 11408 |

| moa
<dbl> | n
<int> |
|--------------|------------|
| 11 | 10636 |
| 12 | 10698 |

EDA Continued

- What sex has the higher mean age?

```
```{r}
midterm_df_final %>%
 filter(sex==1) %>%
 summarise(mean(age))
```
```

| mean(age) |
|-----------|
| <dbl> |
| 51.49705 |

1 row

```
```{r}
midterm_df_final %>%
 filter(sex==2) %>%
 summarise(mean(age))
```
```

| mean(age) |
|-----------|
| <dbl> |
| 50.86232 |

EDA Continued

- What provider has the highest cost?

```
{r}  
midterm_df_final$total <- as.numeric(midterm_df_final$total)  
midterm_df_final %>%  
  group_by(provider) %>%  
  summarise(mean(total))  
```
```

R Console

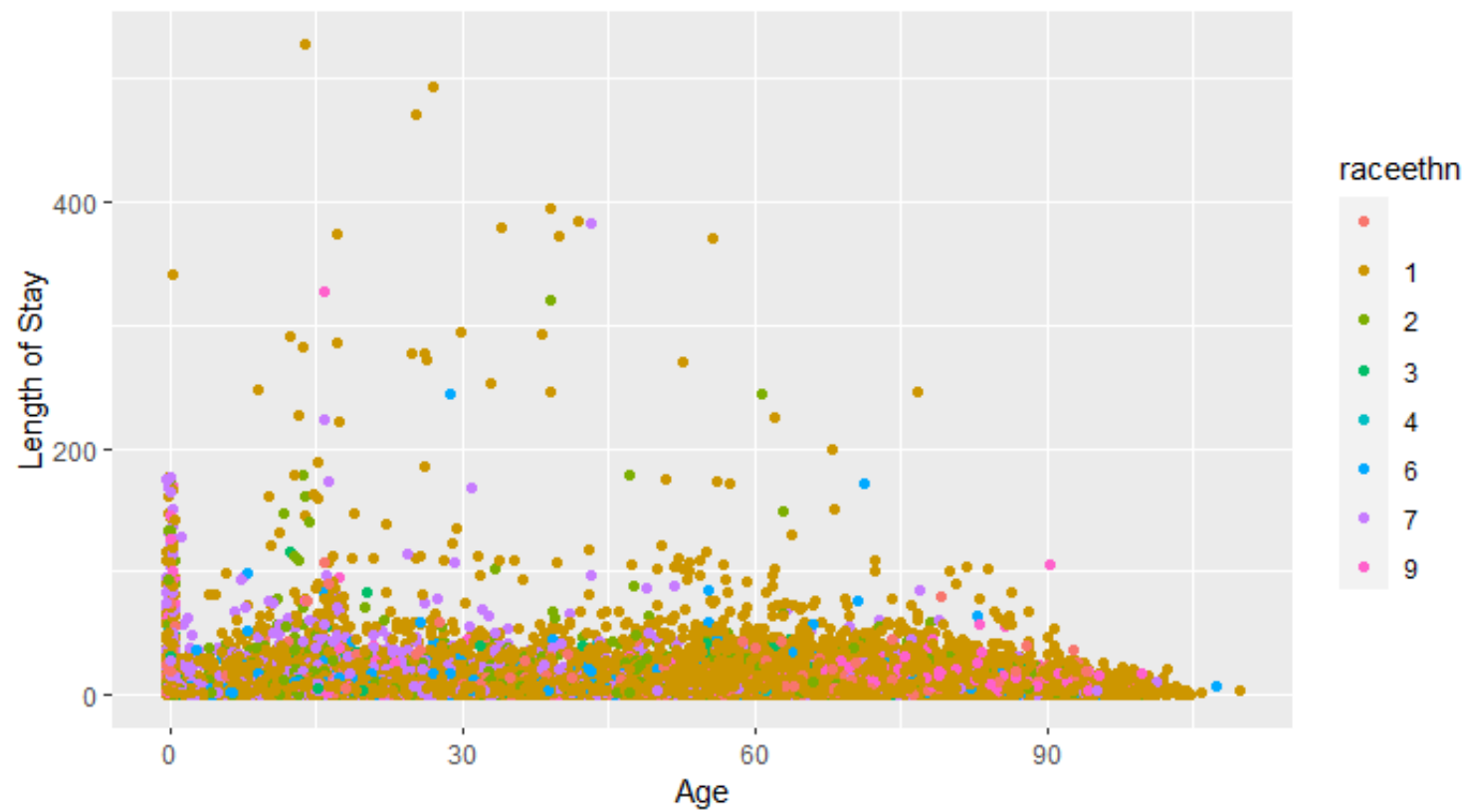
tbl\_df  
12 x 2

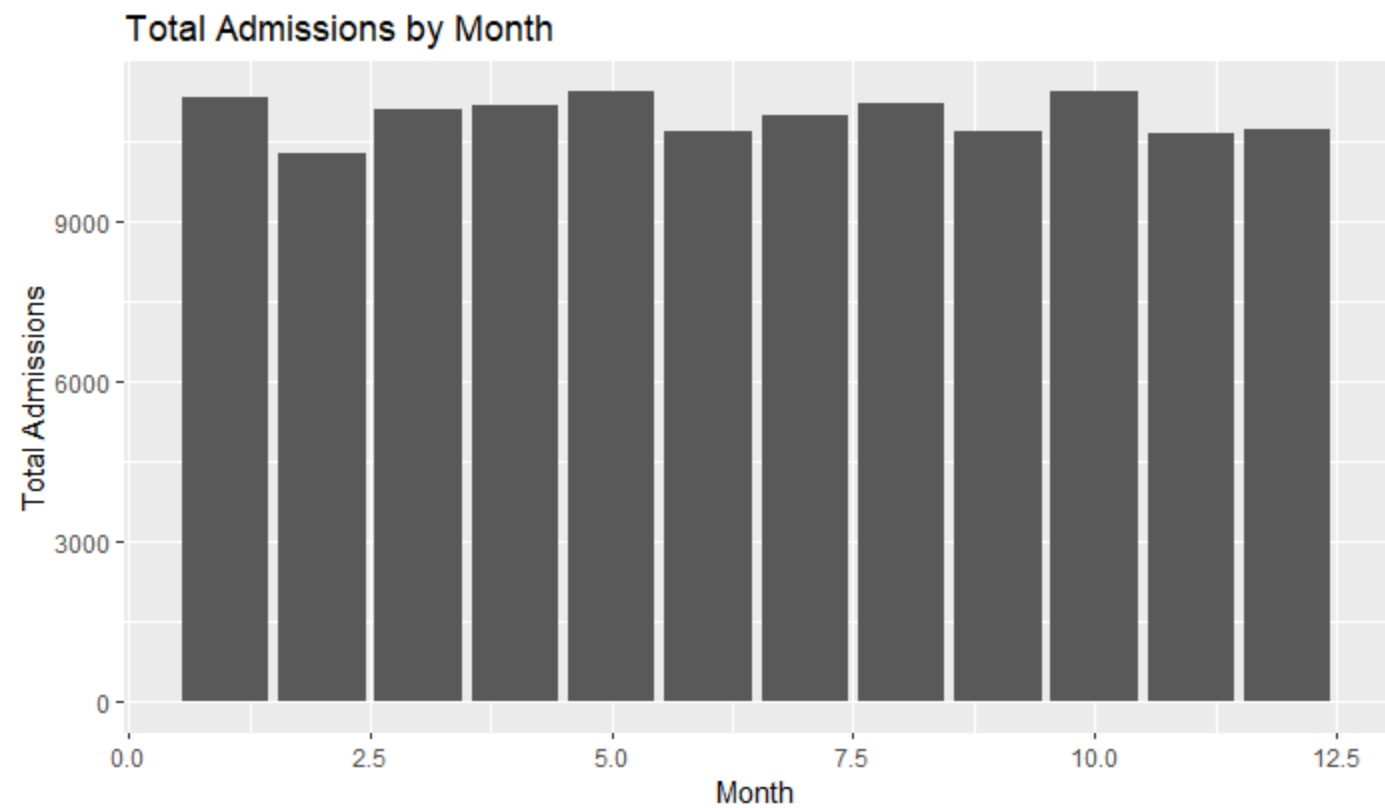
provider <chr>	mean(total) <dbl>	provider <chr>	mean(total) <dbl>
7201	22775.33	7215	69945.55
7202	35504.44	7216	17781.83
7204	35276.94		
7205	48738.82		
7206	31017.81		
7209	24538.79		
7210	27690.88		
7211	24088.58		
7213	38200.23		
7214	22362.14		

# Plots

The bottom of the slide features two horizontal blue bars. The first bar is a solid medium blue and spans the entire width of the slide. The second bar is a slightly lighter shade of blue, positioned to the right of the first bar, creating a layered or 3D effect.

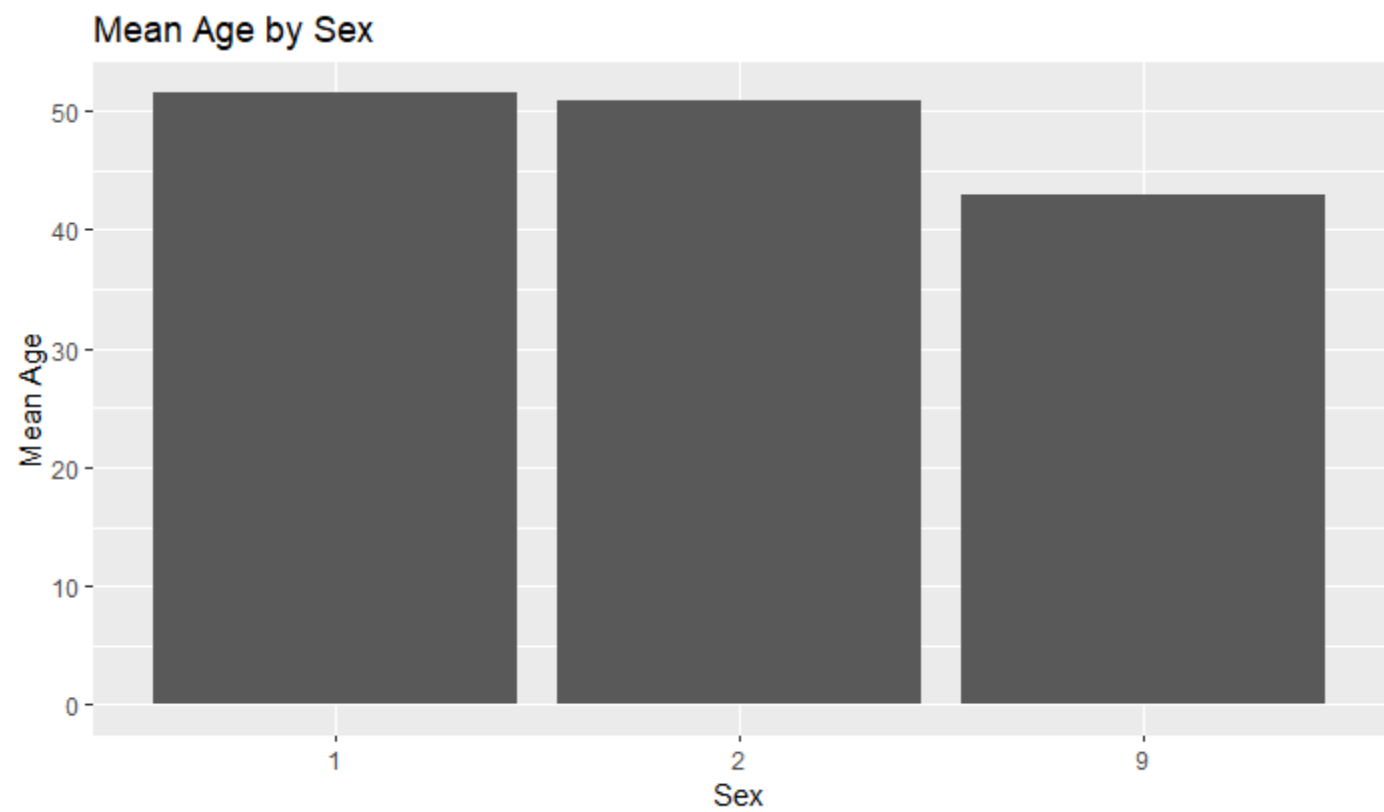
Age vs Length of Stay by Race



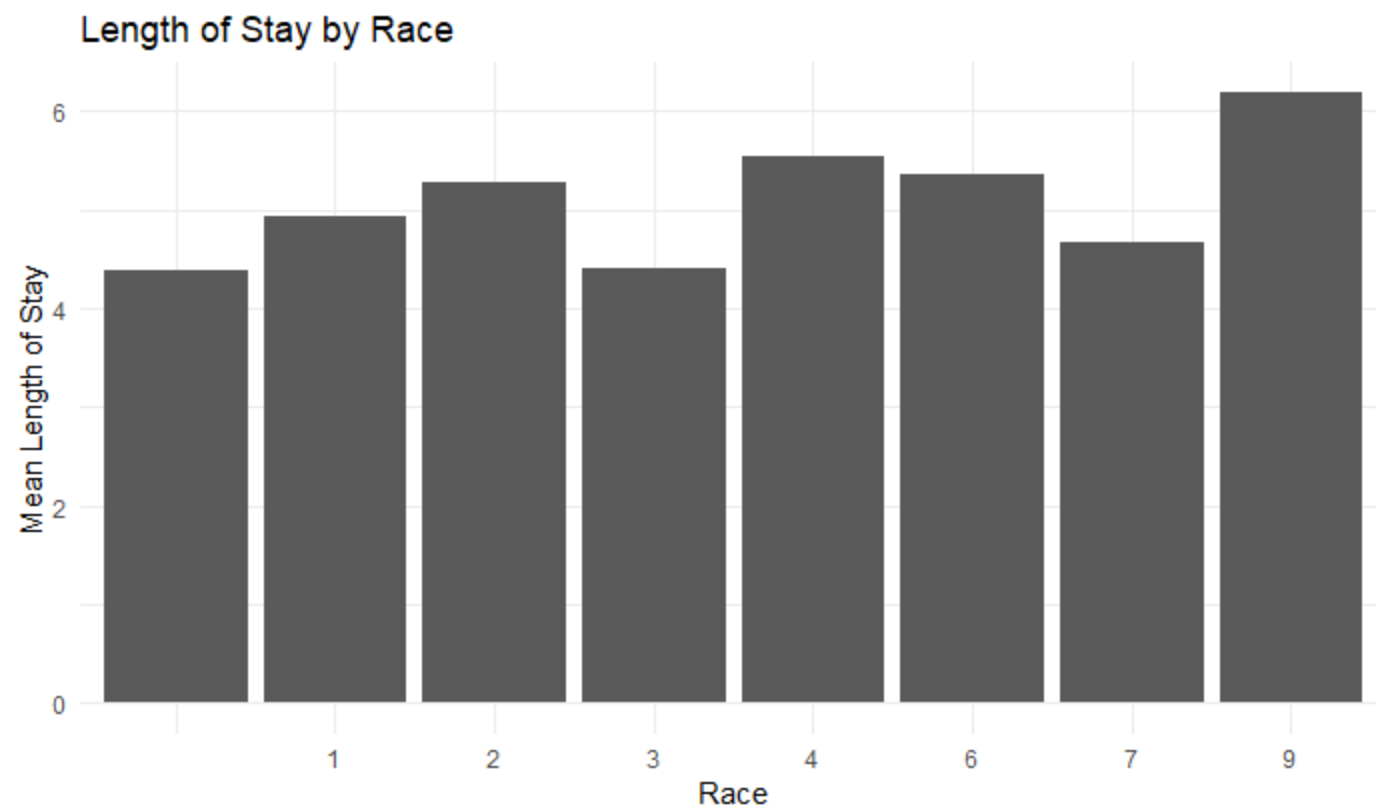


A plot of total admissions by month, January has the highest total where the next month of February has the lowest.

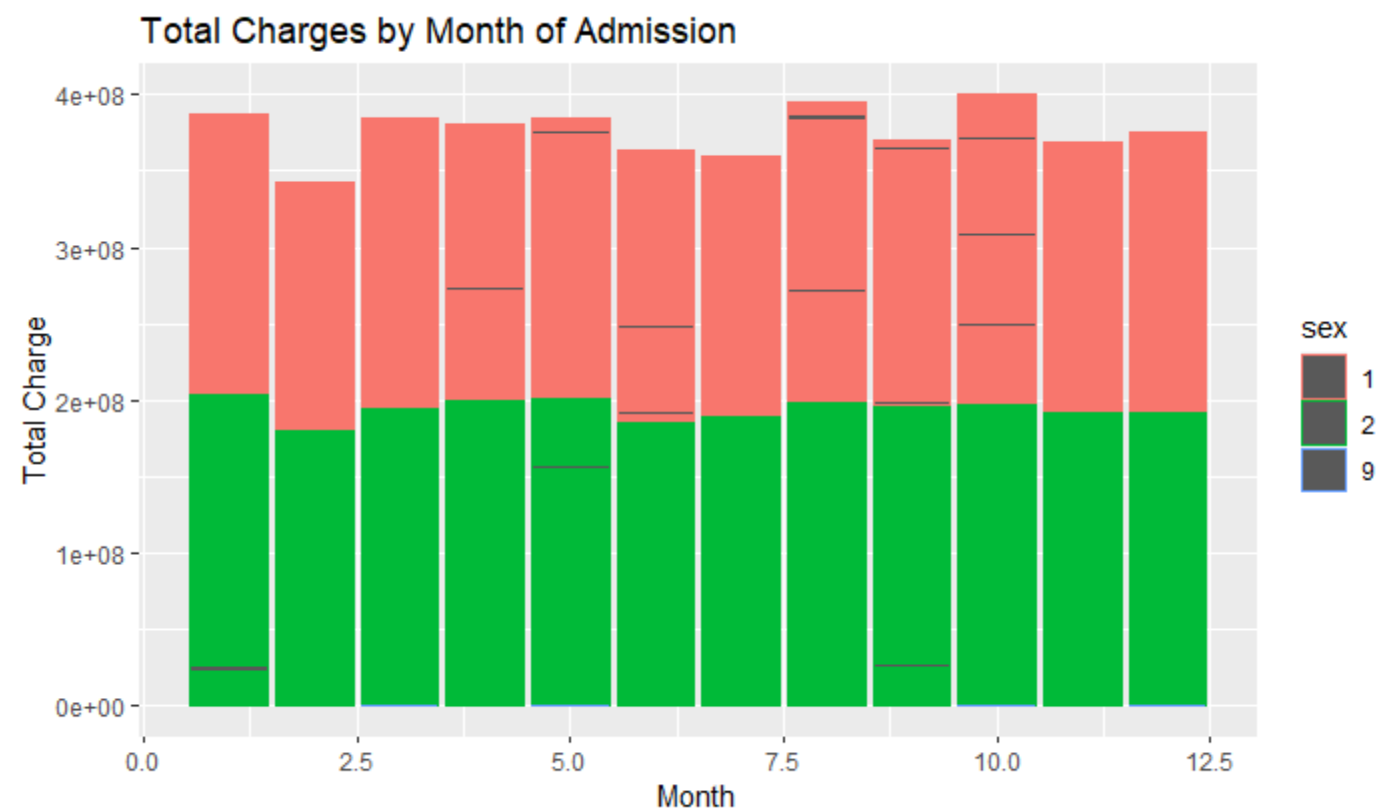




A plot of average age by sex, the female admitted patients have a lower age than the male patients

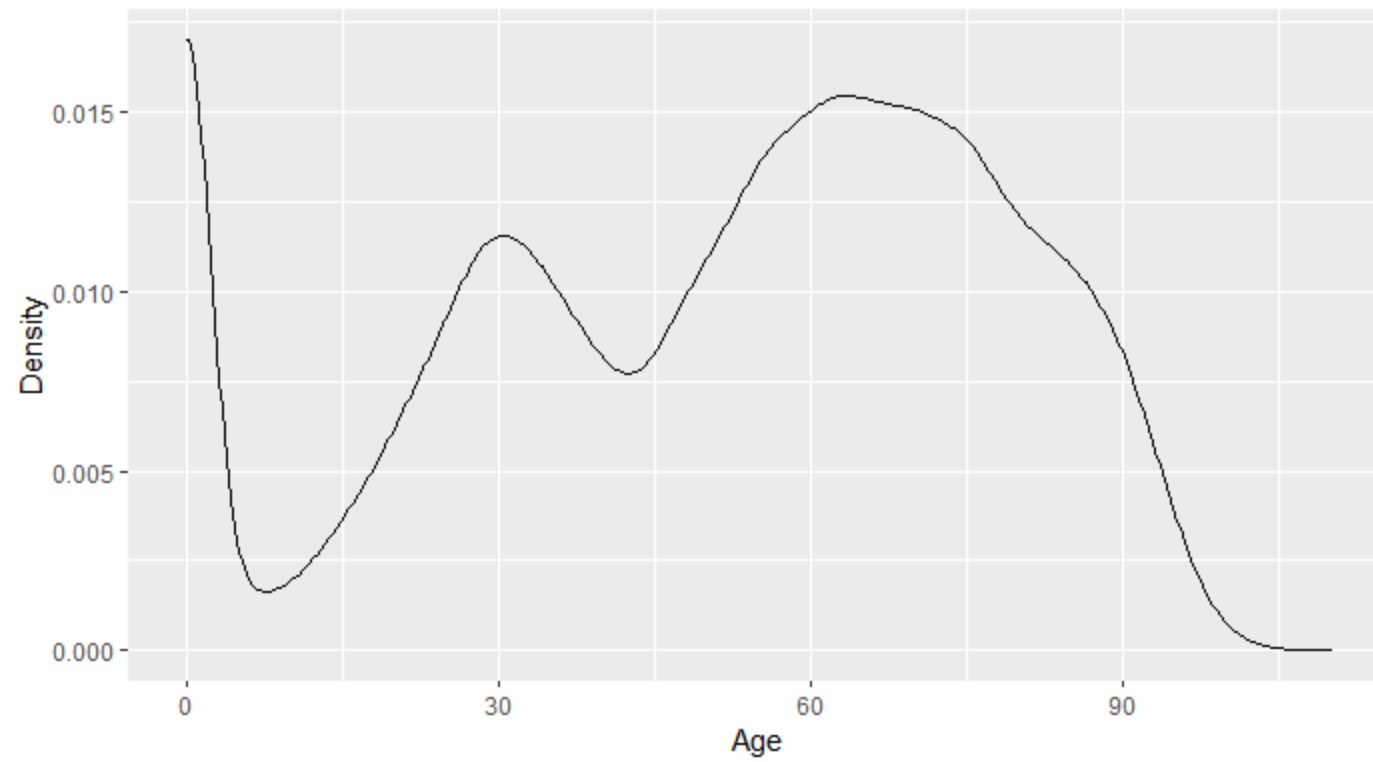


A plot of average length of stay by race. From the plot you can see Unknown race (9) has the highest average length of stay.

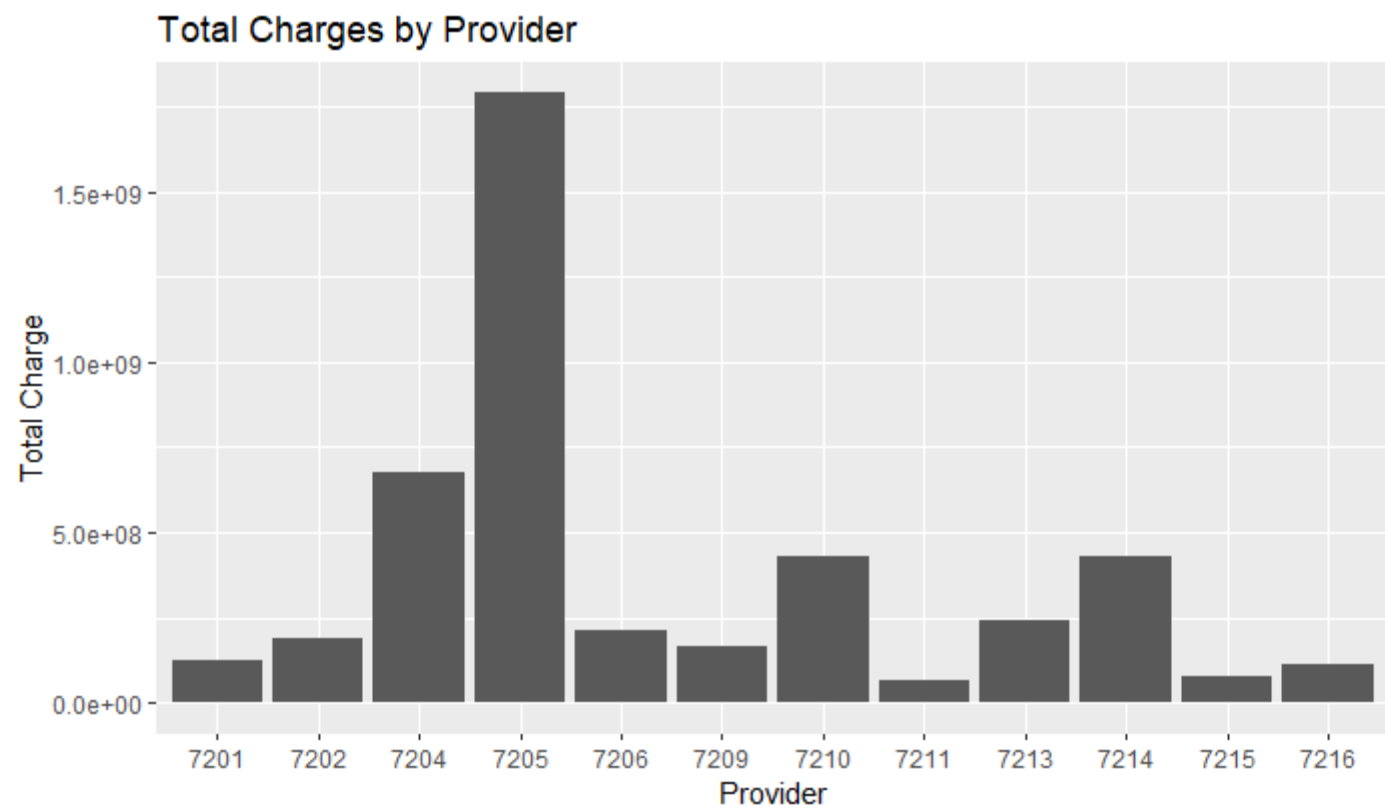


A plot of total charges by month and colored by sex

Density plot of age



Density plot of the age variable. There are three distinct peaks, one for babies, 30-35 year olds, and elderly people

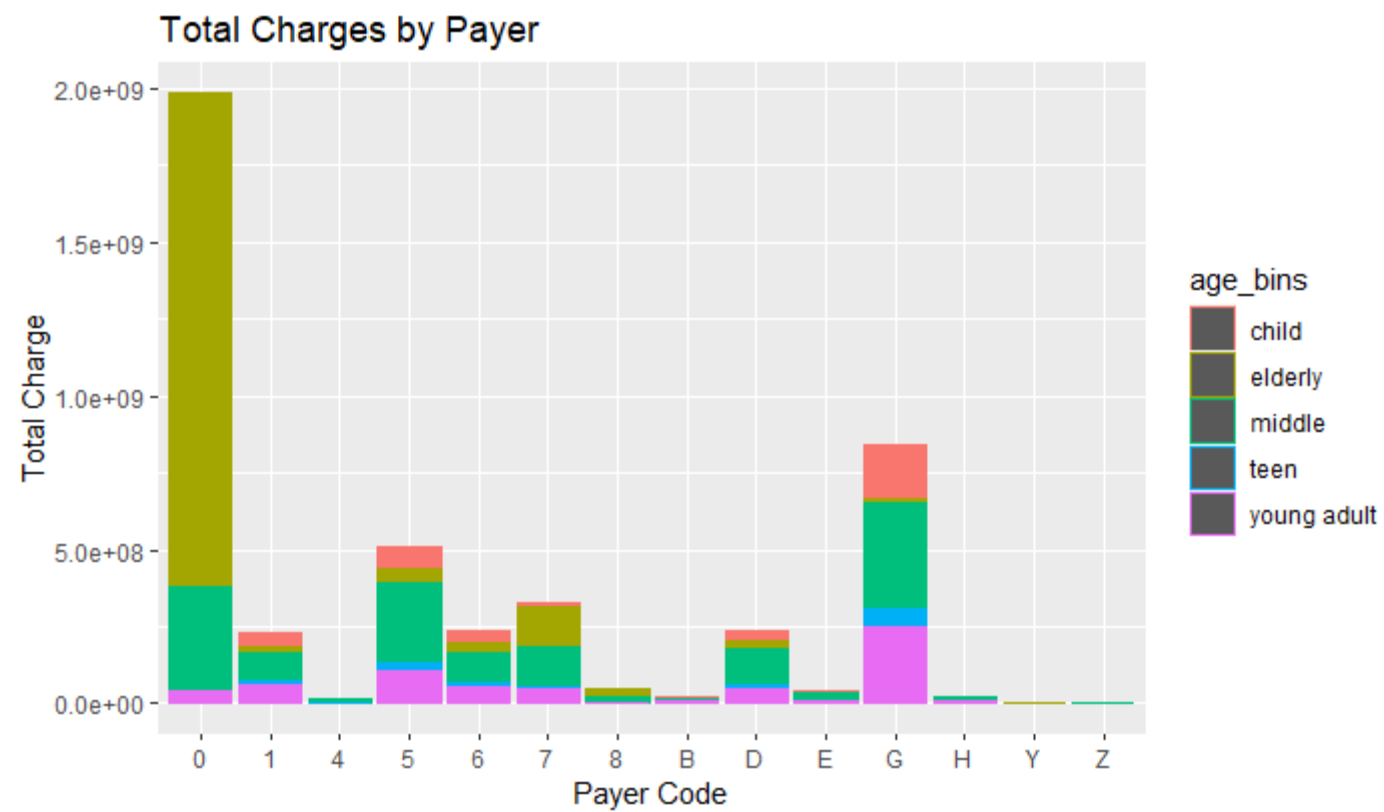


id Hospital, this makes sense because they are the largest hospital in Rhode Island and they do more complicated procedures

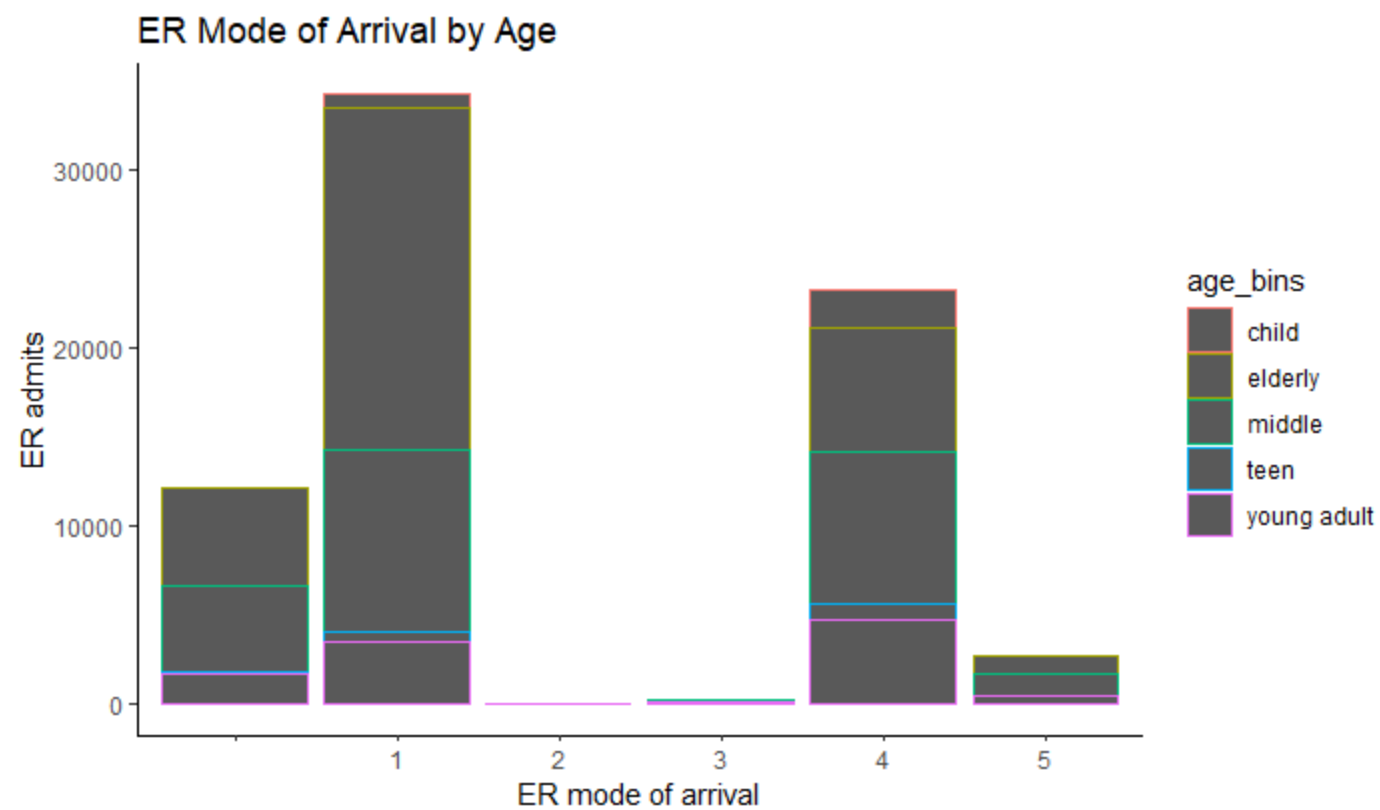
Total Charges by Baby's Birthweight



The lower the baby's birthweight, the higher the charges



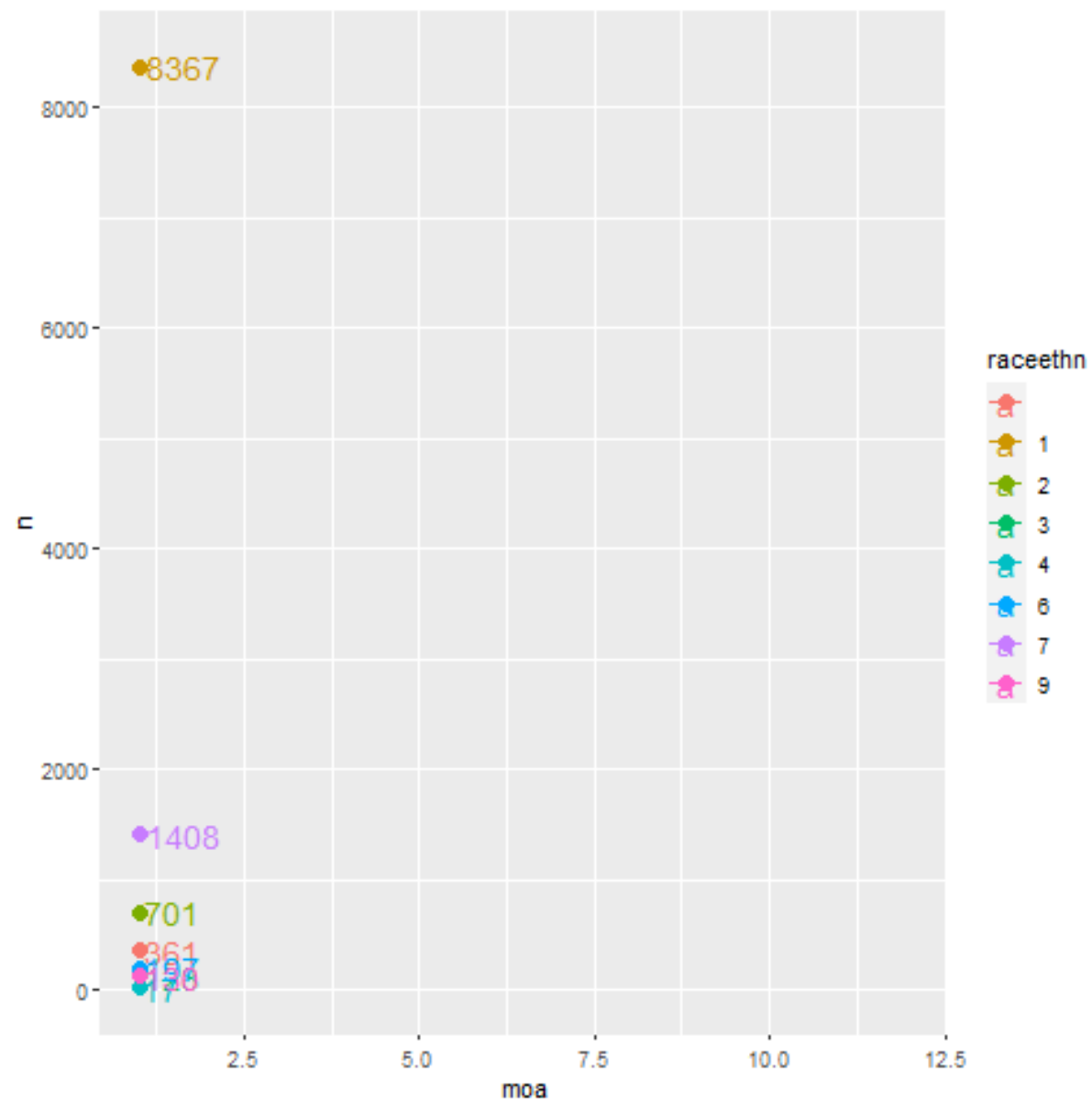
Medicare is the highest payment, especially in the elderly age range.



Filtered out for NA and Unknown, the highest was through an ambulance or personal transportation



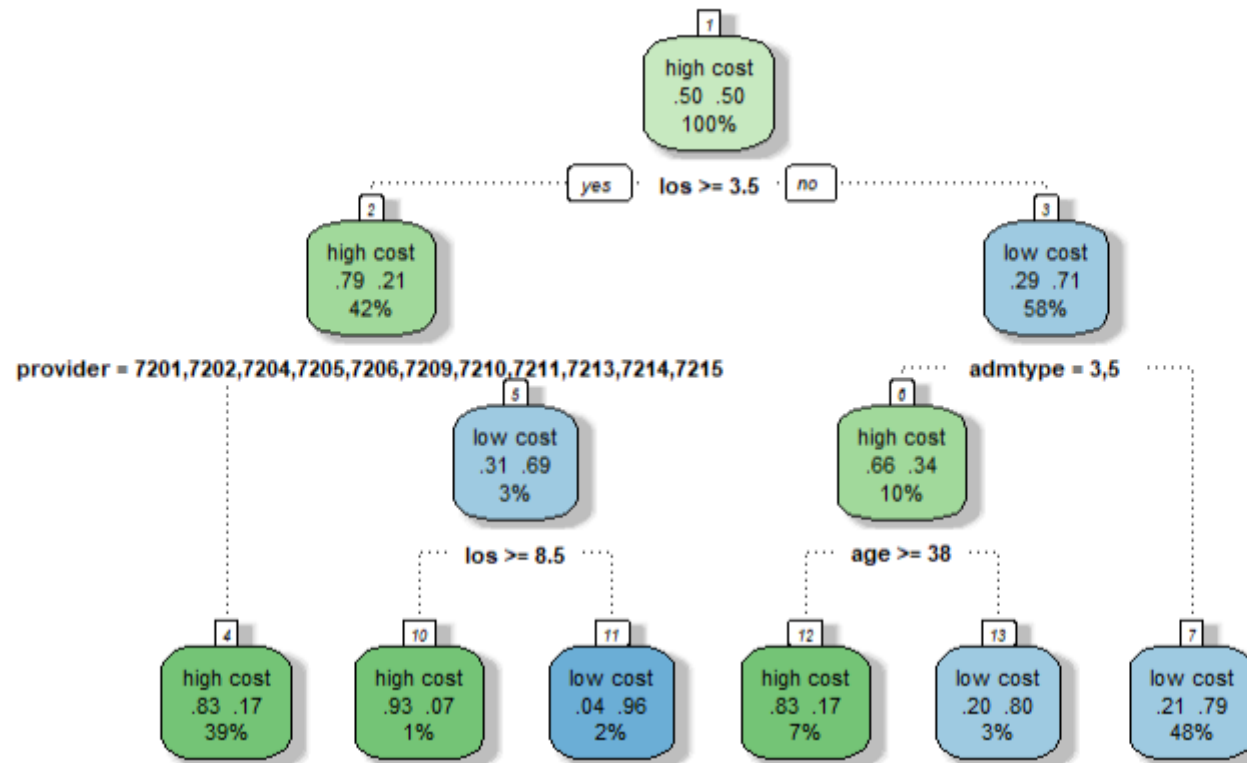
Hospital Admissions by Month Colored  
by Race



# Predictive Modeling



# Tree Model



# Model Choice 1- High/Low Cost

- Random Forest method had highest accuracy:

5. what is your final selection for the model? Test the accuracy of your final model on the test data.

```
```{r}
pred <- predict(forest_cv, df_test)
cm <- confusionMatrix(data = pred, reference = df_test$target, positive = "high cost")
cm$overall[1]
```
```

```
Accuracy
0.8275247
```

# Model Choice 2- Old/Young

- Parallel Random Forest had highest accuracy

```
```{r}
pred <- predict(parRF_cv2, df_test2)
cm <- confusionMatrix(data = pred, reference = df_test2$target, positive = "old")
cm$overall[1]
```
```

```
Accuracy
0.867074
```