

# Capstone Project - Predicting Patient No-Shows Using Appointment Data

*Derek Samsom*

Missed medical appointments are a major problem in the medical industry, resulting in lost revenue. Medical providers can over-book appointments to try to minimize the lost revenue, but without any way to predict the probability of an appointment being missed, there will still be times where more or fewer patients show up at a given time than expected. The result will be that lost revenue will be reduced but not eliminated, as there will still be times that more appointments are missed than expected. There will also be times more appointments show up than expected, which can overwhelm staff and resources and affect the level of patient care.

This project is a classification problem that will explore the prediction of whether a medical appointment will be missed, and its probability of being kept or missed. The prediction error will result in times where there are too many or too few patients at a given time. The main goal in the prediction will be to minimize the error, as this will reduce instances of having more or fewer patients than desired.

There are countless reasons and circumstances that can lead someone to miss an appointment, such as a last minute work meeting or a family emergency, that aren't directly captured in the data and are impossible to know in advance. Missed appointments can only be predicted based on indirect factors that are known, such as past history and demographics. Because of this, there will be a level of error that cannot be eliminated, however, any reduction in error compared to having no predictive model at all is still beneficial.

Medical providers can use the missed appointment predictions by incorporating them into their booking methods and systems. The methods used in booking will have to consider the implications of the inherent prediction errors and balance the risk the errors represent: too many patients leading to staff/resource shortage, and too few patients leading to lost revenue. The methods of implementing the use of missed appointment predictions into an appointment booking system are client-specific and not included in the scope of this project, which is limited to minimizing the error while predicting the probability that an appointment will be missed or kept.

I will start off by loading the required packages and the data.

```
library(tidyverse)
library(lubridate)
library(caret)
library(randomForest)
library(GGally)
```

Read data and assign to `appointments`

```
appointments <- read_csv("Final_Data.csv")
```

```
## Parsed with column specification:
## cols(
##   kept_status = col_character(),
##   appt_date = col_character(),
##   appt_time = col_time(format = ""),
##   appt_length = col_integer(),
##   date_scheduled = col_character(),
##   patient_age = col_integer(),
##   patient_gender = col_character(),
##   billing_type = col_character(),
##   prior_missed = col_integer(),
```

```
## prior_kept = col_integer(),
## patient_distance = col_integer(),
## office_zip = col_character(),
## provider_specialty = col_character(),
## remind_call_result = col_character()
## )
```

```
appointments_original <- appointments
zipcodes <- read_csv("zipcodes.csv")
```

```
## Parsed with column specification:
## cols(
##   office_zip = col_character(),
##   county_code = col_character()
## )
```

The raw data, which has been named `appointments`, contains information on 342862 past appointments, pre-sorted by the date and time of appointment. The depended variable, `kept_status`, shows whether the appointment was kept or missed.

There is no field that can be used to identify a specific patients in the data set. A patient may have had more than one appointment during the time-period represented in the data, meaning that one individual patient may make up one or multiple observations. If there was a patient ID field, it would allow the data to be grouped by patient and give the option of organizing the data by patient rather than by appointment.

A secondary data set, `zipcodes`, has information about the county the offices are located in. This will be used to see if the location can help predict whether an appointment will be missed. The county names are converted to a 2-letter code for confidentiality.

## Data Summary and Structure

```
summary(appointments)
```

```
## kept_status      appt_date      appt_time      appt_length
## Length:342862    Length:342862    Length:342862    Min.   : 10
## Class :character  Class :character  Class1:hms        1st Qu.: 60
## Mode  :character  Mode  :character  Class2:difftime   Median : 60
##                                     Mode   :numeric   Mean   : 57
##                                     3rd Qu.: 60
##                                     Max.   :600
##
## date_scheduled   patient_age      patient_gender    billing_type
## Length:342862     Min.   : 0.00    Length:342862     Length:342862
## Class :character  1st Qu.: 17.00   Class :character   Class :character
## Mode  :character  Median : 34.00   Mode  :character   Mode  :character
##                                     Mean   : 35.56
##                                     3rd Qu.: 54.00
##                                     Max.   :264.00
##
## prior_missed      prior_kept      patient_distance  office_zip
## Min.   : 0.000     Min.   : 0.00    Min.   : 0.0      Length:342862
## 1st Qu.: 1.000     1st Qu.: 2.00    1st Qu.: 0.0      Class :character
## Median : 2.000     Median : 6.00    Median : 3.0      Mode  :character
## Mean   : 2.451     Mean   : 8.02    Mean   : 10.8
## 3rd Qu.: 3.000     3rd Qu.: 11.00   3rd Qu.: 9.0
```

```
## Max.      :117.000    Max.      :676.00    Max.      :2688.0
##                                     NA's      :974
## provider_specialty remind_call_result
## Length:342862      Length:342862
## Class :character    Class :character
## Mode  :character    Mode  :character
##
##
##
##
```

```
str(appointments, give.attr = FALSE)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 342862 obs. of 14 variables:
## $ kept_status      : chr "Kept" "Kept" "Kept" "Kept" ...
## $ appt_date        : chr "9/1/16" "9/1/16" "9/1/16" "9/1/16" ...
## $ appt_time        :Classes 'hms', 'difftime' atomic [1:342862] 19800 28800 28800 28800 28800 28800 ...
## $ appt_length      : int 90 60 120 60 60 60 60 60 60 90 ...
## $ date_scheduled   : chr "8/1/16" "1/18/16" "2/3/16" "6/8/16" ...
## $ patient_age      : int 7 75 31 45 49 71 49 38 36 13 ...
## $ patient_gender    : chr "Male" "Female" "Male" "Male" ...
## $ billing_type      : chr "DMAP" "Commercial" "DMAP" "DMAP" ...
## $ prior_missed      : int 1 2 1 6 5 6 8 0 2 3 ...
## $ prior_kept        : int 3 5 5 15 6 6 20 0 5 12 ...
## $ patient_distance  : int 41 29 5 5 0 5 0 539 0 4 ...
## $ office_zip        : chr "AP" "BL" "BL" "BL" ...
## $ provider_specialty: chr "A" "A" "A" "B" ...
## $ remind_call_result: chr "Left Message" "Answered - Confirmed" "Left Message" "Answered - No Resp"
```

```
head(appointments[, 1:5])
```

```
## # A tibble: 6 x 5
##   kept_status appt_date appt_time appt_length date_scheduled
##   <chr>      <chr>      <time>      <int> <chr>
## 1 Kept      9/1/16      05:30         90 8/1/16
## 2 Kept      9/1/16      08:00         60 1/18/16
## 3 Kept      9/1/16      08:00        120 2/3/16
## 4 Kept      9/1/16      08:00         60 6/8/16
## 5 Missed    9/1/16      08:00         60 6/28/16
## 6 Kept      9/1/16      08:00         60 7/12/16
```

```
head(appointments[, 6:10])
```

```
## # A tibble: 6 x 5
##   patient_age patient_gender billing_type prior_missed prior_kept
##   <int> <chr>      <chr>      <int>      <int>
## 1      7 Male      DMAP         1         3
## 2     75 Female    Commercial    2         5
## 3     31 Male      DMAP         1         5
## 4     45 Male      DMAP         6        15
## 5     49 Male      Commercial    5         6
## 6     71 Male      DMAP         6         6
```

```
head(appointments[, 11:14])
```

```
## # A tibble: 6 x 4
##   patient_distance office_zip provider_specialty remind_call_result
```

##	<int>	<chr>	<chr>	<chr>
## 1	41	AP	A	Left Message
## 2	29	BL	A	Answered - Confirmed
## 3	5	BL	A	Left Message
## 4	5	BL	B	Answered - No Response
## 5	0	BL	B	Answered - No Response
## 6	5	BL	A	Answered - Confirmed

## Data Dictionary

```
variable_descriptions <- c(
  "Dependent variable: kept or missed",
  "Appointment date",
  "Appointment time",
  "Appointment length in minutes",
  "Date appointment was scheduled",
  "Patient age",
  "Patient gender",
  "Billing type",
  "Number of prior missed appointments",
  "Number of prior kept appointments",
  "Patient distance from office in miles",
  "Office Zip Code - Anonymized",
  "Provider primary specialty code",
  "Reminder Call result")
variable <- colnames(appointments)
variable_type <- unlist(map(appointments, class))
variable_type <- variable_type[-4]
as_data_frame(cbind(c(1:length(variable)), variable, variable_type, variable_descriptions))
```

```
## # A tibble: 14 x 4
##   V1    variable      variable_type variable_descriptions
##   <chr> <chr>          <chr>          <chr>
## 1 1    kept_status    character      Dependent variable: kept or mis~
## 2 2    appt_date      character      Appointment date
## 3 3    appt_time      hms           Appointment time
## 4 4    appt_length    integer        Appointment length in minutes
## 5 5    date_scheduled character      Date appointment was scheduled
## 6 6    patient_age    integer        Patient age
## 7 7    patient_gender  character      Patient gender
## 8 8    billing_type    character      Billing type
## 9 9    prior_missed    integer        Number of prior missed appointm~
## 10 10   prior_kept      integer        Number of prior kept appointmen~
## 11 11   patient_distance integer        Patient distance from office in~
## 12 12   office_zip     character      Office Zip Code - Anonymized
## 13 13   provider_specialty character      Provider primary specialty code
## 14 14   remind_call_result character      Reminder Call result
```

The appt\_date and appt\_time variables can be combined into one variable, appt\_datetime.

```
appointments <- appointments %>%
  mutate(appt_datetime = lubridate::mdy_hms(paste(appt_date, appt_time)))

appointments$date_scheduled <- lubridate::as_date(
```

```
appointments$date_scheduled, format = "%m/%d/%y", tz = "UTC")
```

## Data Exploration

First I want to calculate the percent of missed appointments overall by creating a logical variable `missed`, where 1 represents a missed appointment and 0 represents a kept appointment. This will determine the degree of class imbalance.

```
appointments <- appointments %>%  
  mutate(missed = ifelse(appointments$kept_status == "Missed", 1, 0))  
missed_rate <- mean(appointments$missed)  
missed_rate
```

```
## [1] 0.1592944
```

15.93 % of the total appointments are missed. This is an imbalanced classification, which will have implications in the modeling. For example, the model could predict all of the appointments will be kept and be correct 84.07 % of the time. This results in a high accuracy without providing any useful prediction of which appointments will be missed.

Next I want to check the data to see if there are any missing values that could indicate reduced data integrity or adversely affect the modelling.

```
map_dbl(appointments, ~sum(is.na(.)))
```

```
##      kept_status      appt_date      appt_time  
##           0           0           0  
##      appt_length      date_scheduled      patient_age  
##           0           0           0  
##      patient_gender      billing_type      prior_missed  
##           0           0           0  
##           prior_kept      patient_distance      office_zip  
##           0           974           0  
## provider_specialty remind_call_result      appt_datetime  
##           0           0           0  
##           missed  
##           0
```

One variable, `patient_distance` has 974 missing value. This is fairly minor and will be evaluated later on when exploring the variable further.

### patient\_age

I expected missed appointments to have to vary across age ranges. Perhaps older patients have fewer commitments with kids or work, and make their appointments more regularly, or perhaps younger adults might skip more appointments because they aren't as critical? I will break the data into age groups to make the plot simpler to evaluate.

There are a small number of observations where the age is higher than plausible. Therefore, the observations greater than age 110 will be removed from the data.

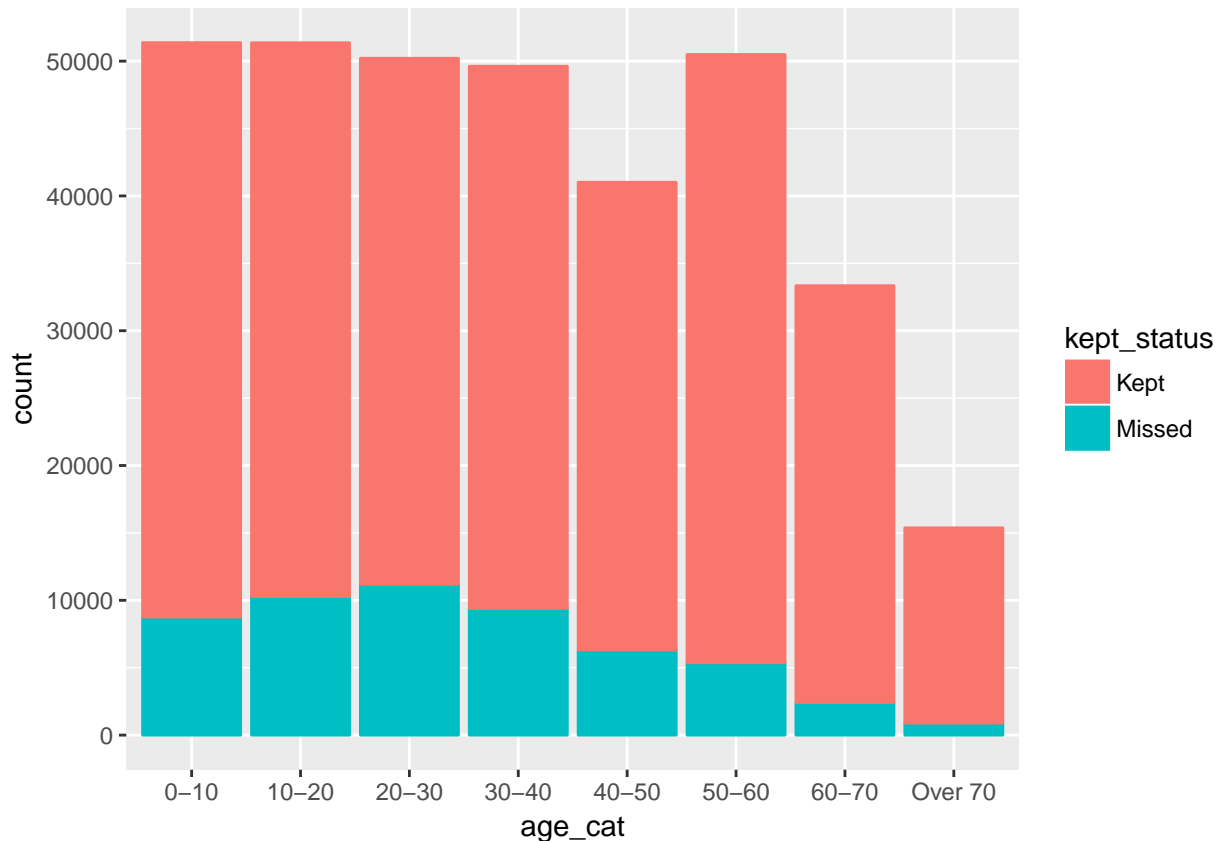
```
age_labels <- c("0-10", "10-20", "20-30", "30-40", "40-50", "50-60", "60-70",  
              "Over 70")  
age_breaks <- c(-1, 10, 20, 30, 40, 50, 60, 70, 111)  
appointments <- appointments %>%
```

```

filter(patient_age <= 110) %>%
mutate(
  age_cat = cut(patient_age, breaks = age_breaks, labels = age_labels))

ggplot(
  appointments,
  aes(x = age_cat, color = kept_status, fill = kept_status)
) +
  stat_count()

```



Missed appointments are highest with young adults, and decrease with older and younger patients.

### billing\_type

```
table(appointments$billing_type)
```

```
##
##   Commercial      DMAP To Be Assigned
##      78282      264500           1
```

There is only one observation of “To Be Assigned”, therefore it will be removed from the data.

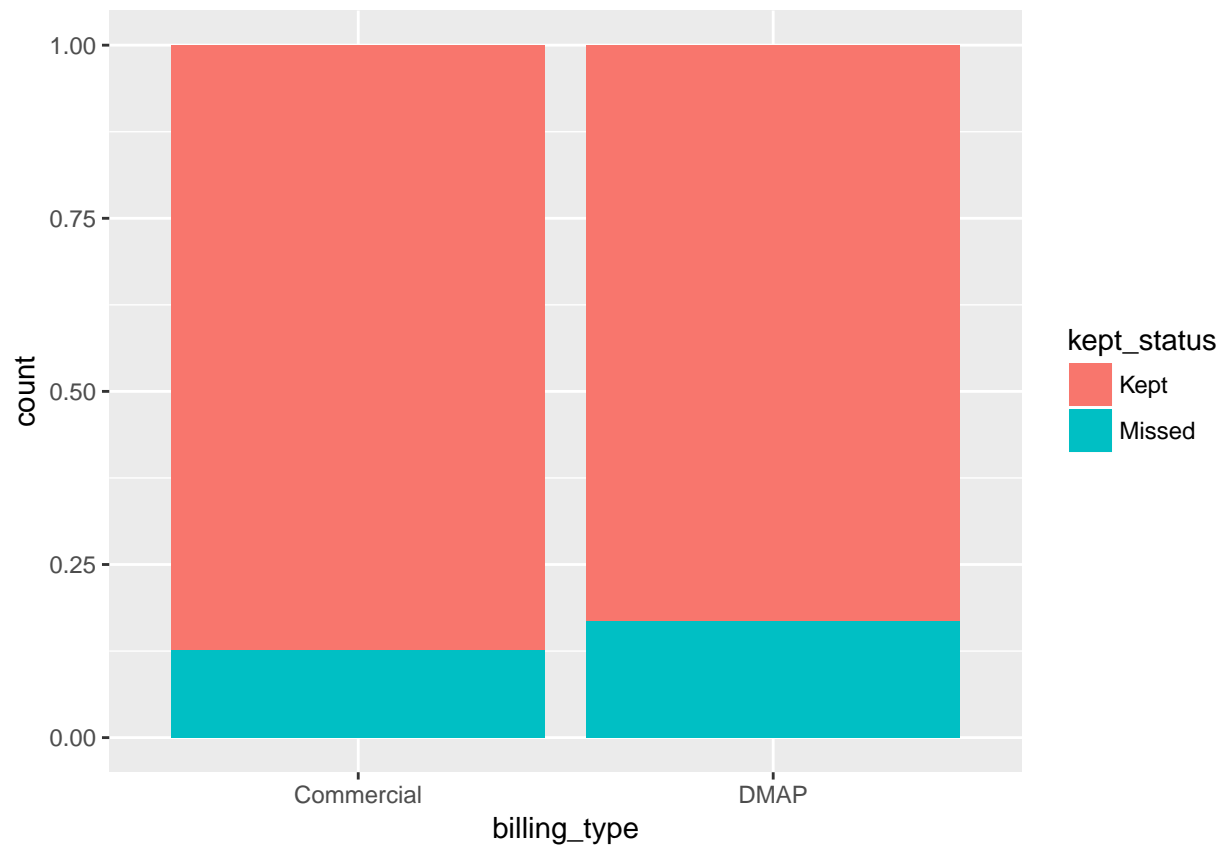
```
appointments <- subset(appointments, billing_type != "To Be Assigned")
```

```
ggplot(
  appointments,
```

```

aes(x = billing_type, fill = kept_status)
) +
geom_bar(position = "fill")

```



There is a minor difference between billing types. DMAP has a higher proportion of missed appointments than commercial.

## appt\_datetime

For the variable `appt_datetime`, I will create an `hour` variable to see the variation in missed appointments by hour of day.

```

appointments <- appointments %>%
  mutate(hour = lubridate::hour(appointments$appt_datetime))
table(appointments$hour)

```

```

##
##      0      5      6      7      8      9     10     11     12     13     14     15
##      7     24     25    98 42816 42133 45326 38033    321 48307 43787 44033
##     16     17     18     19     20     21
## 35449 2180  205   33    3    2

```

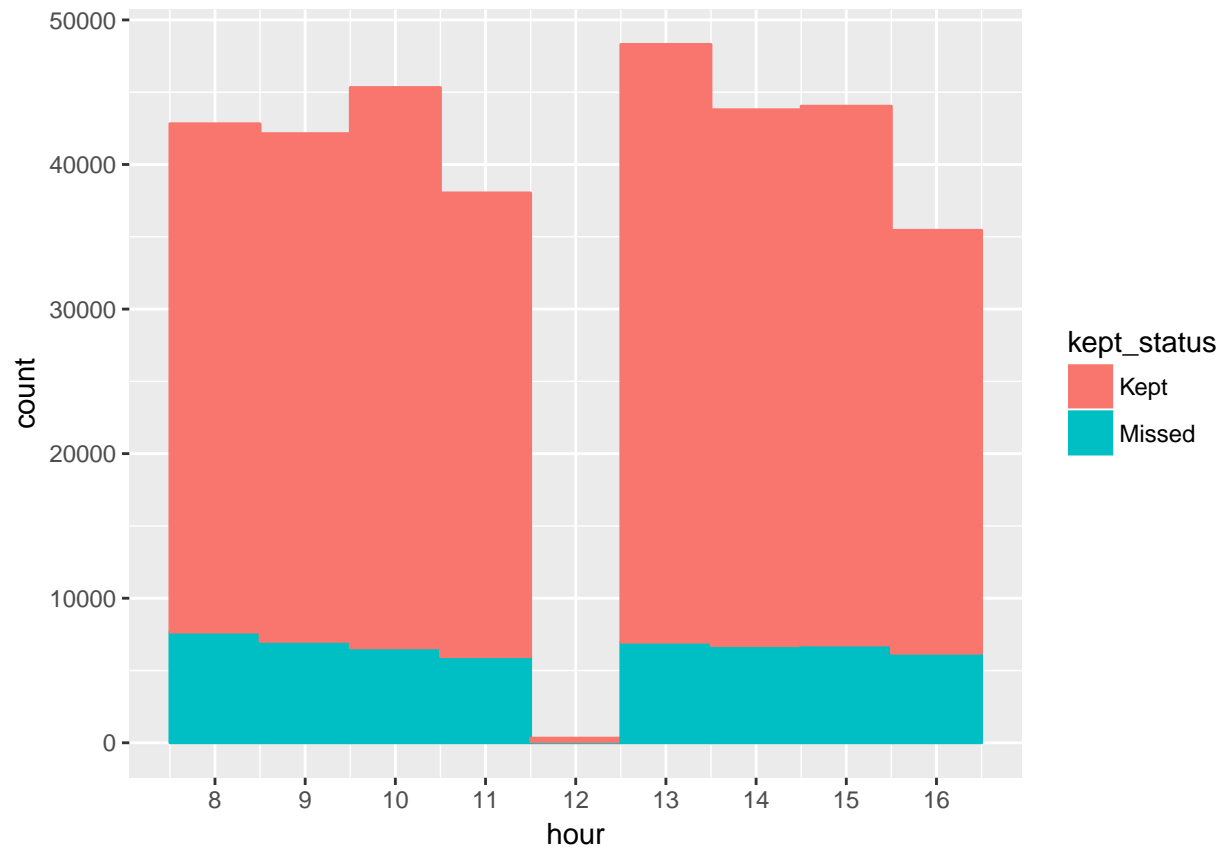
Most appointments are scheduled between 8:00 AM and 5:00 PM, with a one hour gap starting at 12:00.

```

appointments_hour <- appointments %>%
  select(kept_status, hour) %>%
  filter(hour >= 8 & hour <= 16)

```

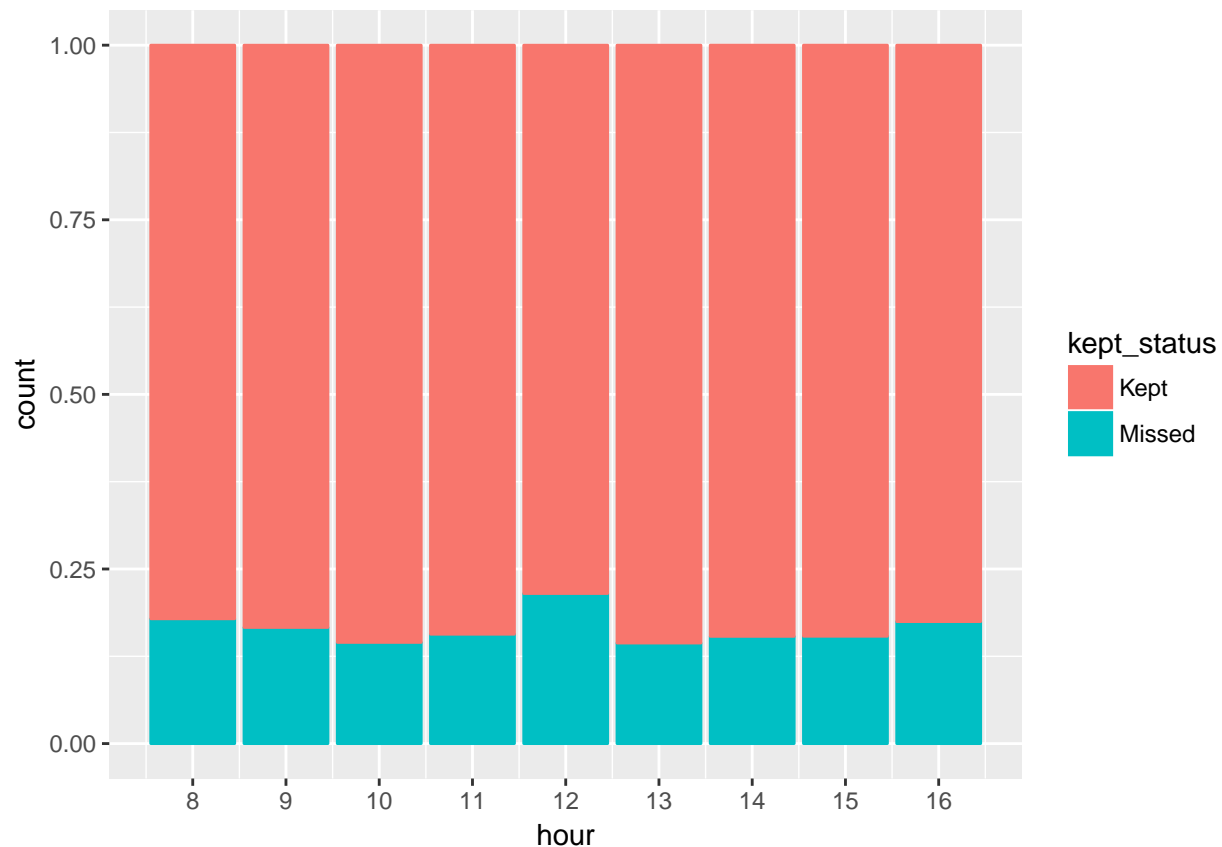
```
ggplot(
  appointments_hour,
  aes(x = hour, col = kept_status, fill = kept_status)
) +
  geom_histogram(binwidth = 1) +
  scale_x_continuous(breaks = seq(8, 17, 1))
```



There is a decline in the total number of missed appointments as both the morning and afternoon period progress, however, there are fewer appointments towards the end of the two periods.

```
ggplot(
  appointments_hour,
  aes(x = hour, col = kept_status, fill = kept_status)
) +
  geom_bar(position = "fill") +
  scale_x_continuous(breaks = seq(8, 17, 1))
```





Proportionally more appointments are missed at the beginning and end of the typical scheduling hours, and during the few noon appointments.

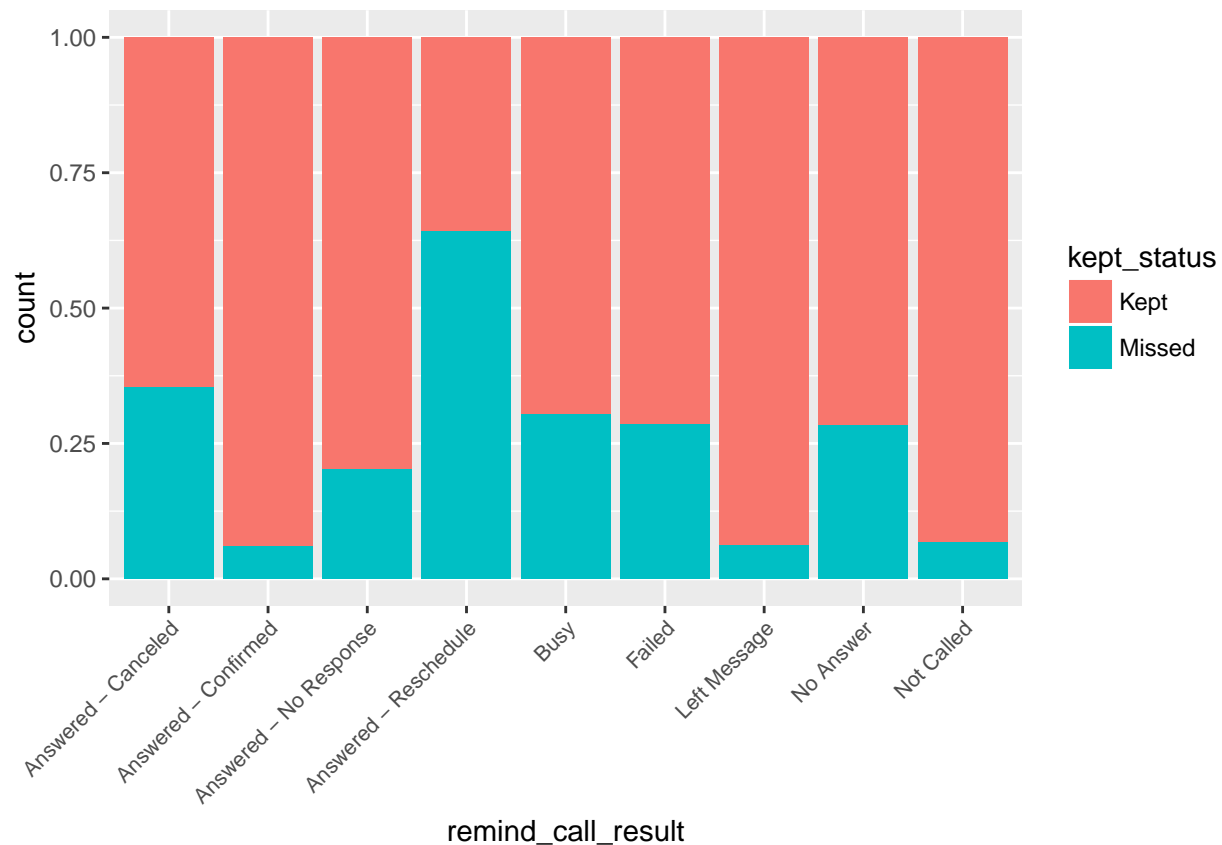
#### remind\_call\_result

```
table(appointments$remind_call_result)
```

```
##
##      Answered - Canceled      Answered - Confirmed Answered - No Response
##              152              49108              180869
## Answered - Reschedule              Busy              Failed
##              1369              1104              27944
##           Left Message              No Answer              Not Called
##              18430              377              63429
```

Low counts of “Answered - Cancelled”, “Answered - Reschedule”, “Busy”, and “No Answer”

```
ggplot(
  appointments,
  aes(x = remind_call_result, fill = kept_status)
) +
  geom_bar(position = "fill") +
  theme(axis.text.x = element_text(size = 8, angle = 45, hjust = 1, vjust = 1))
```



~65% of appointments with “Answered - Cancelled” and ~35% with “Answered-Reschedule” still kept their appointments, however, very few observations in these categories.

### provider\_specialty

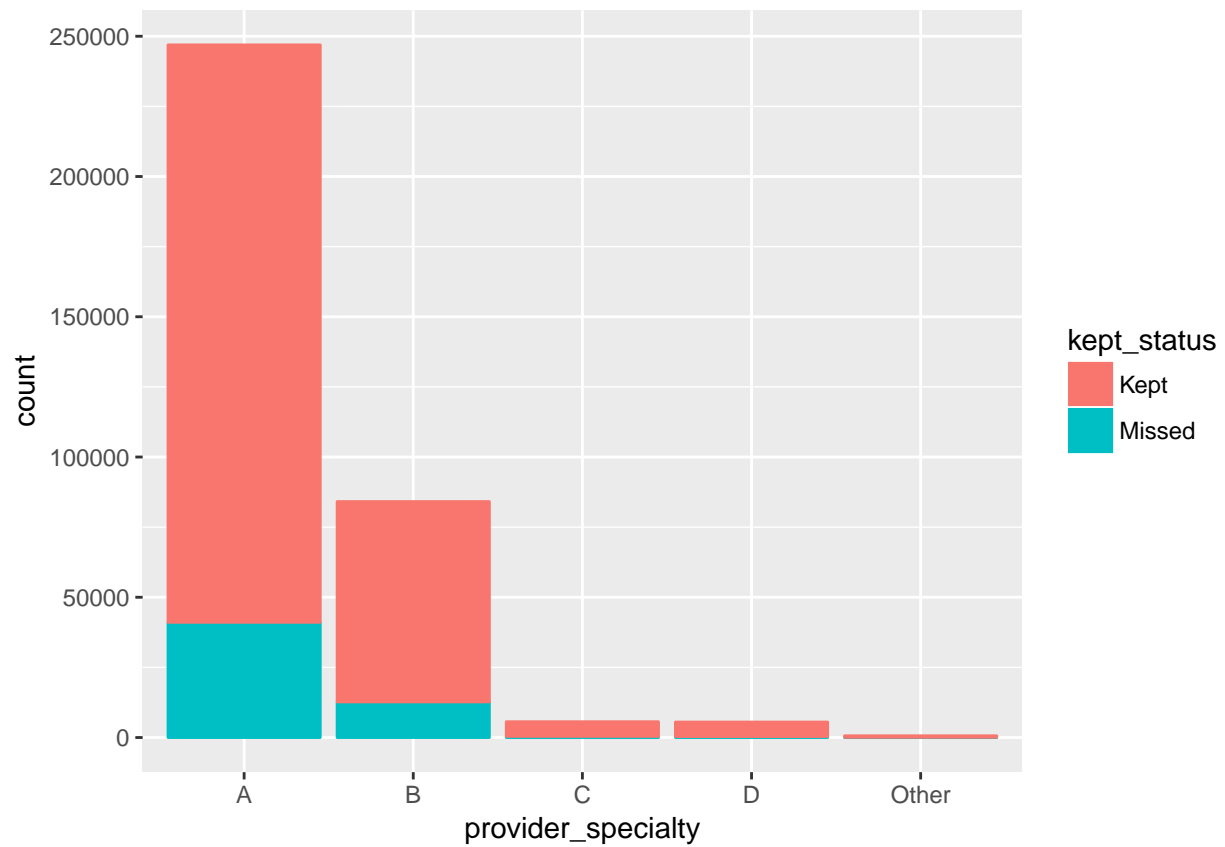
```
table(appointments$provider_specialty)
```

```
##
##      A      B      C      D      E      F      G
## 246917 84115  5623  5512   42   525   48
```

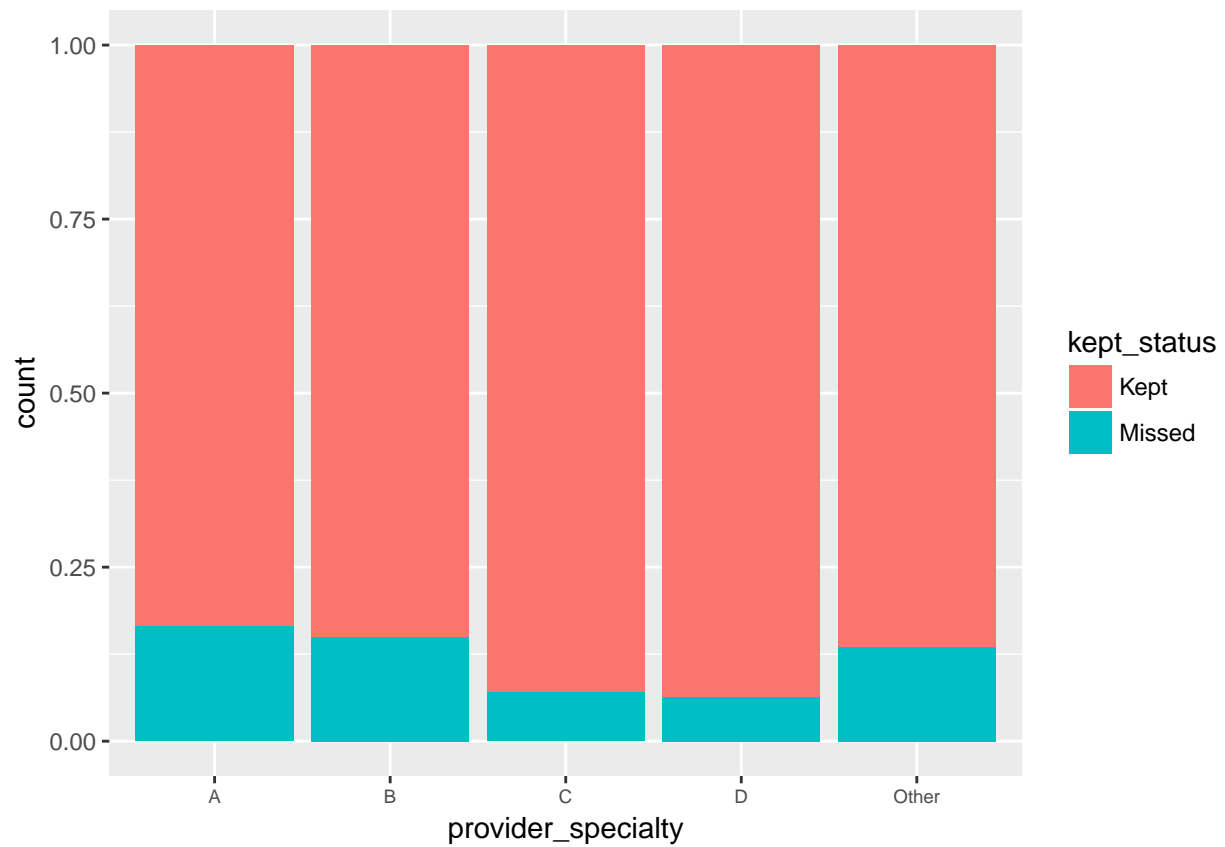
Most observations are specialty A and B. Specialties E, F, and G have very few observations and will be grouped as “Other”.

```
appointments$provider_specialty <- appointments$provider_specialty %>%
  fct_collapse(Other = c("E", "F", "G"))
```

```
ggplot(
  appointments,
  aes(x = provider_specialty, col = kept_status, fill = kept_status)
) +
  stat_count()
```

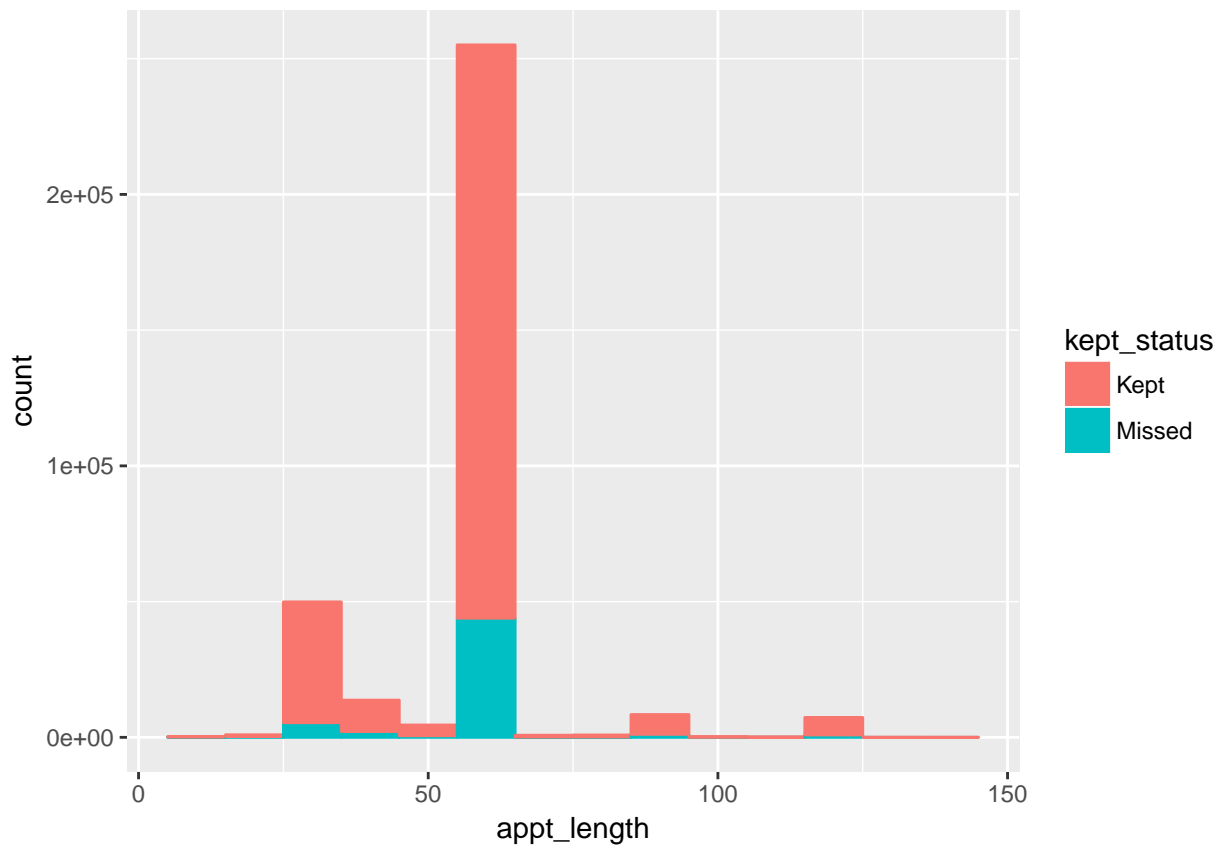


```
ggplot(  
  appointments,  
  aes(x = provider_specialty, fill = kept_status)  
) +  
  geom_bar(position = "fill") +  
  theme(axis.text.x = element_text(size = 7))
```



appt\_length

```
appointments %>%  
  filter(appt_length < 150) %>%  
  ggplot(  
    aes(x = appt_length, color = kept_status, fill = kept_status)  
  ) +  
    geom_histogram(binwidth = 10)
```



Most Appointments are 60 minutes long. 30-minute appointments are next most common.

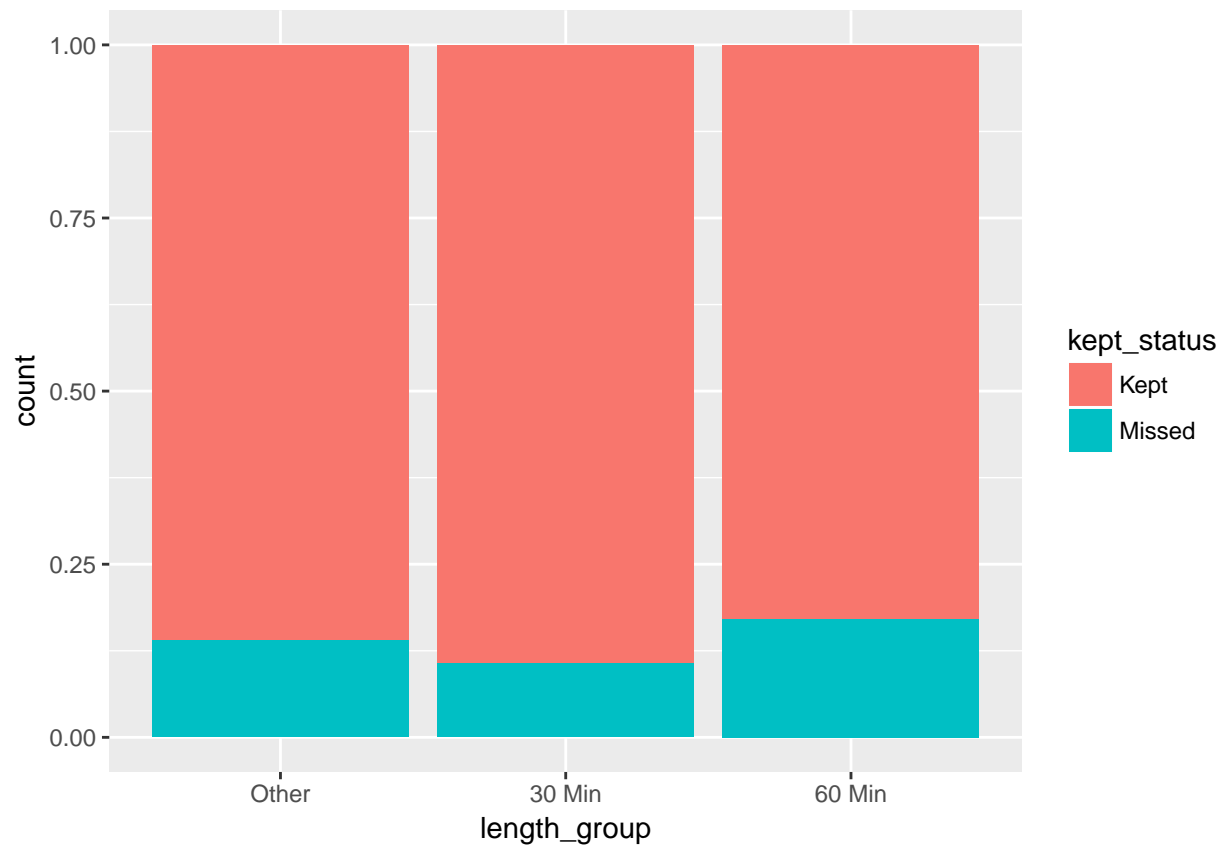
```
length_breaks <- c(-1, 29, 30, 59, 60, 1000)

length_labels <- c("Other1", "30 Min", "Other2", "60 Min", "Other3")

appointments <- appointments %>%
  mutate(
    length_group = cut(
      appt_length, breaks = length_breaks, labels = length_labels)
  )

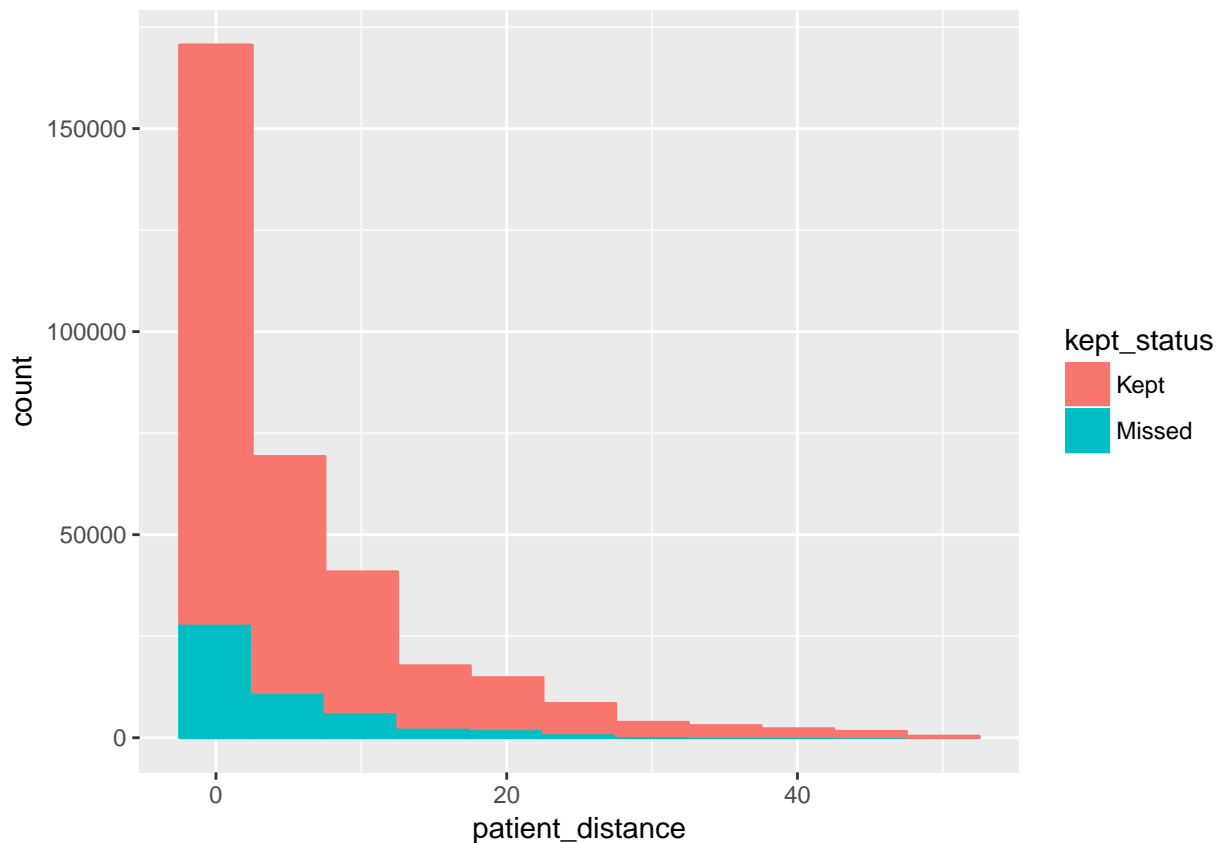
appointments$length_group <- appointments$length_group %>%
  fct_collapse(Other = c("Other1", "Other2", "Other3"))

ggplot(
  data = appointments,
  mapping = aes(x = length_group, fill = kept_status)
) +
  geom_bar(position = "fill")
```



patient\_distance

```
appointments %>%  
  filter(patient_distance < 50) %>%  
  ggplot(  
    aes(x = patient_distance, color = kept_status, fill = kept_status)  
  ) +  
    geom_histogram(binwidth = 5)
```



patient\_distance is very right-skewed, therefore NA values will be replaced with median rather than mean.

```
appointments$patient_distance <- appointments$patient_distance %>%
  replace_na(median(appointments$patient_distance, na.rm = TRUE))
```

## New Variables

In addition to the original variables, there are several additional variables that can be calculated based on the originals.

The `percent_missed` variable is the percentage of prior appointments missed, calculated by dividing the prior missed appointments by the total number of prior appointments. For new patients, this calculation will result in an error because it will be attempting to divide by zero.

The variable `is_new_patient` will specify whether a patient is new, represented by a 1, or existing, represented by 0. My hypothesis is that new patients are more likely to keep their appointments, since I think it is human nature to try to give a good first impression. This is calculated by searching for appointments where `prior_missed` and `prior_kept` are both 0.

The variable `appt_lead_time` will calculate how far in advance an appointment was booked. This is calculated by taking the difference between `date_scheduled` and `appt_date`. If people are more likely to forget appointments booked farther in advance, or they are more likely to be for less urgent preventative care than last minute appointments, this will pick that up.

The variable `appt_weekday` is the day of the week the appointments occurs, and `weekday_scheduled` is the date the appointment was booked.

```
appointments <- appointments %>%
  mutate(percent_missed = prior_missed / (prior_missed + prior_kept)) %>%
```

```

mutate(
  is_new_patient = ifelse(prior_missed == 0 & prior_kept == 0, 1, 0)) %>%
mutate(appt_lead_time = date(appt_datetime) - date(date_scheduled)) %>%
mutate(appt_weekday = strftime(appt_datetime, "%A")) %>%
mutate(weekday_scheduled = strftime(date_scheduled, "%A"))

appointments$percent_missed <- as.integer(appointments$percent_missed * 100)
appointments$percent_missed <- appointments$percent_missed %>%
  tidyr::replace_na(0)

```

Add county\_code from zipcode data.

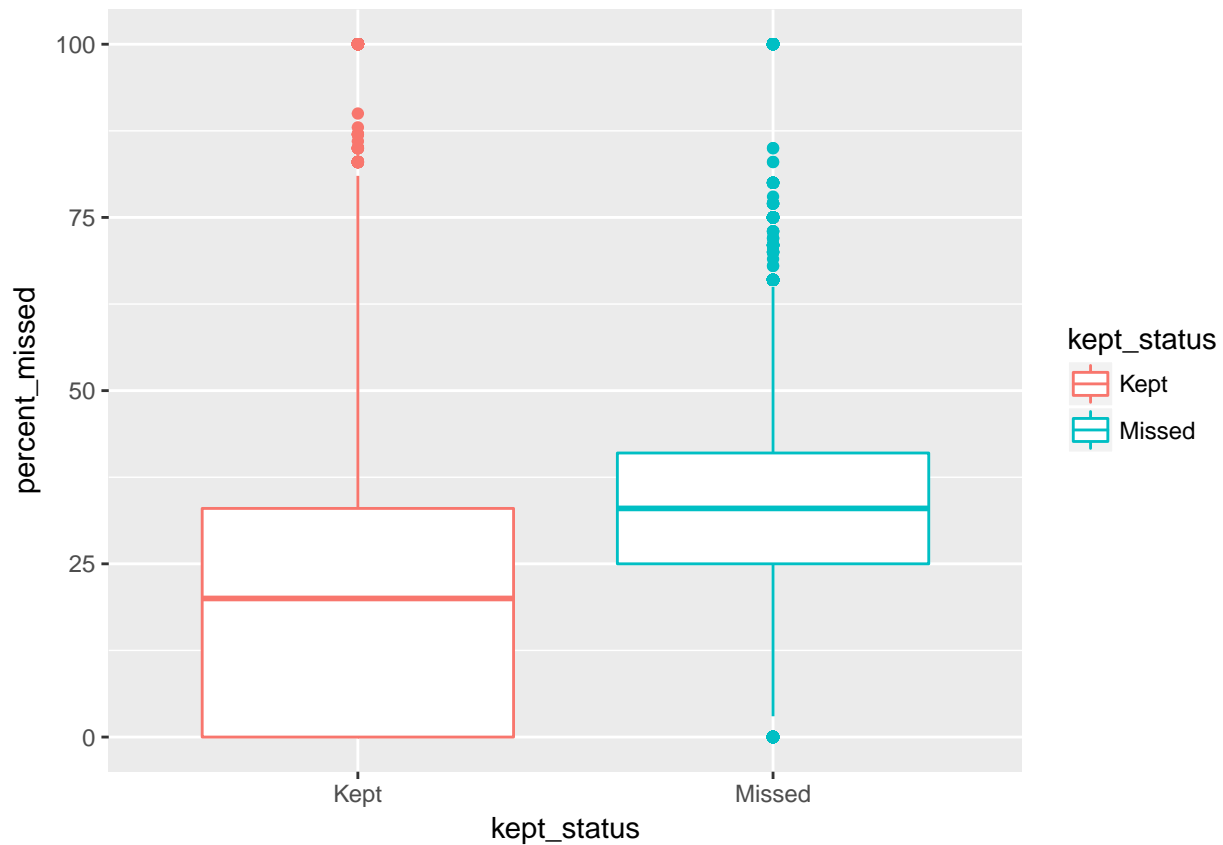
```
appointments <- dplyr::left_join(appointments, zipcodes, by = "office_zip")
```

percent\_missed

```

ggplot(
  data = appointments,
  aes(x = kept_status, y = percent_missed, col = kept_status)
) +
  geom_boxplot()

```



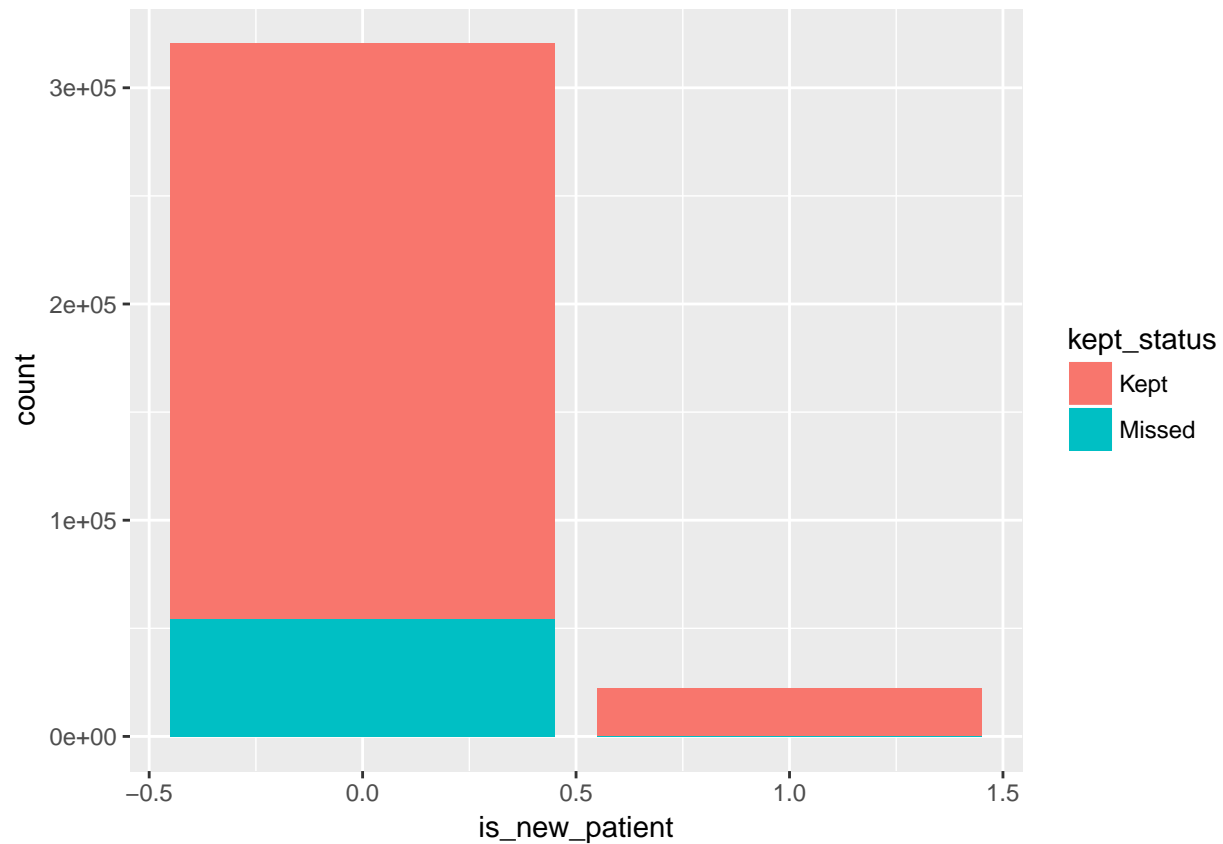


is\_new\_patient

```
table(appointments$is_new_patient)
```

```
##  
##      0      1  
## 320442 22340
```

```
ggplot(  
  appointments,  
  aes(x = is_new_patient, fill = kept_status)  
) +  
  geom_bar()
```



New patients have a very high percentage of kept appointments.

appt\_lead\_time

```
length(which(appointments$appt_lead_time < 0))
```

```
## [1] 82
```

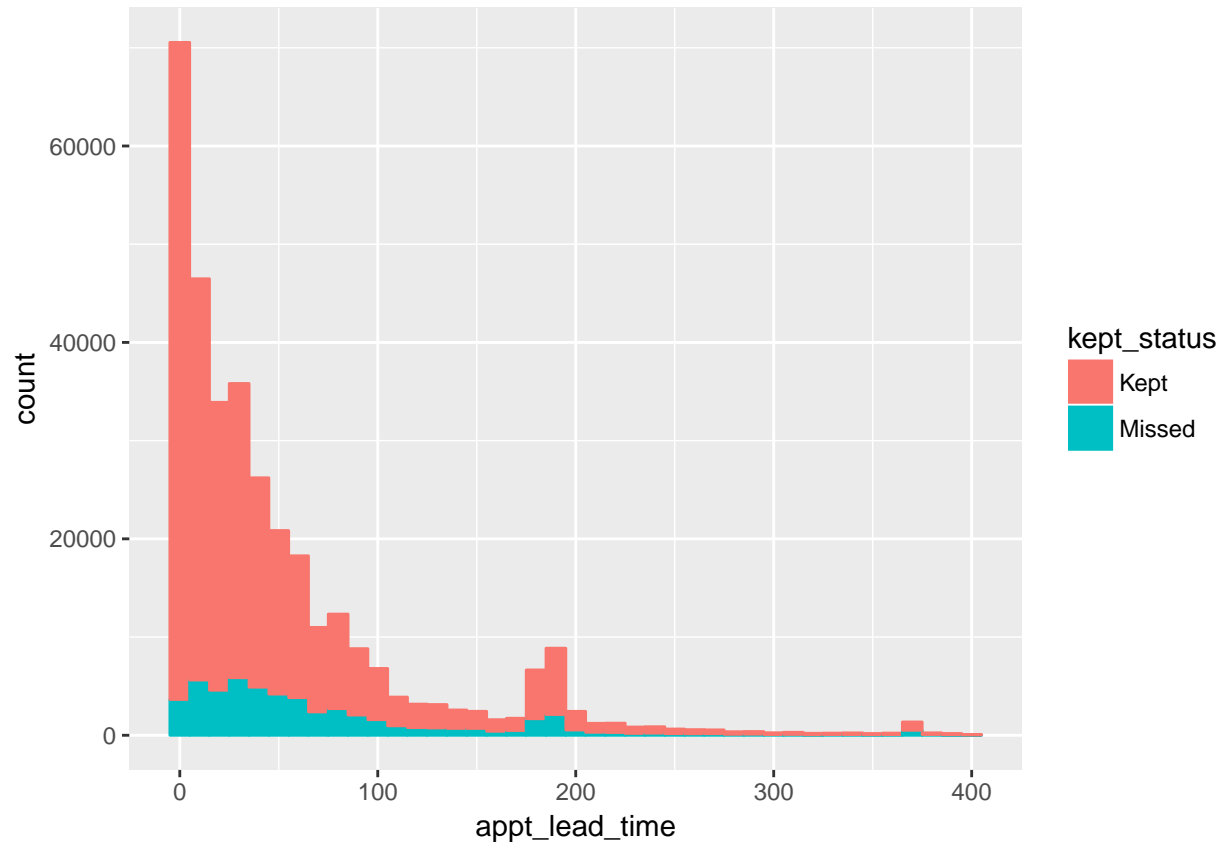
There are 82 negative values, indicating observations that suggest the appointment was booked after it occurred. This is an impossible value and must represent entry errors, but the number is small. These will be replaced with 0.

```

appointments$appt_lead_time <- ifelse(appointments$appt_lead_time < 0, 0, appointments$appt_lead_time)

appointments %>%
  filter(appt_lead_time >= 0 & appt_lead_time < 400) %>%
  ggplot(
    aes(x = appt_lead_time, color = kept_status, fill = kept_status)
  ) +
    geom_histogram(binwidth = 10)

```

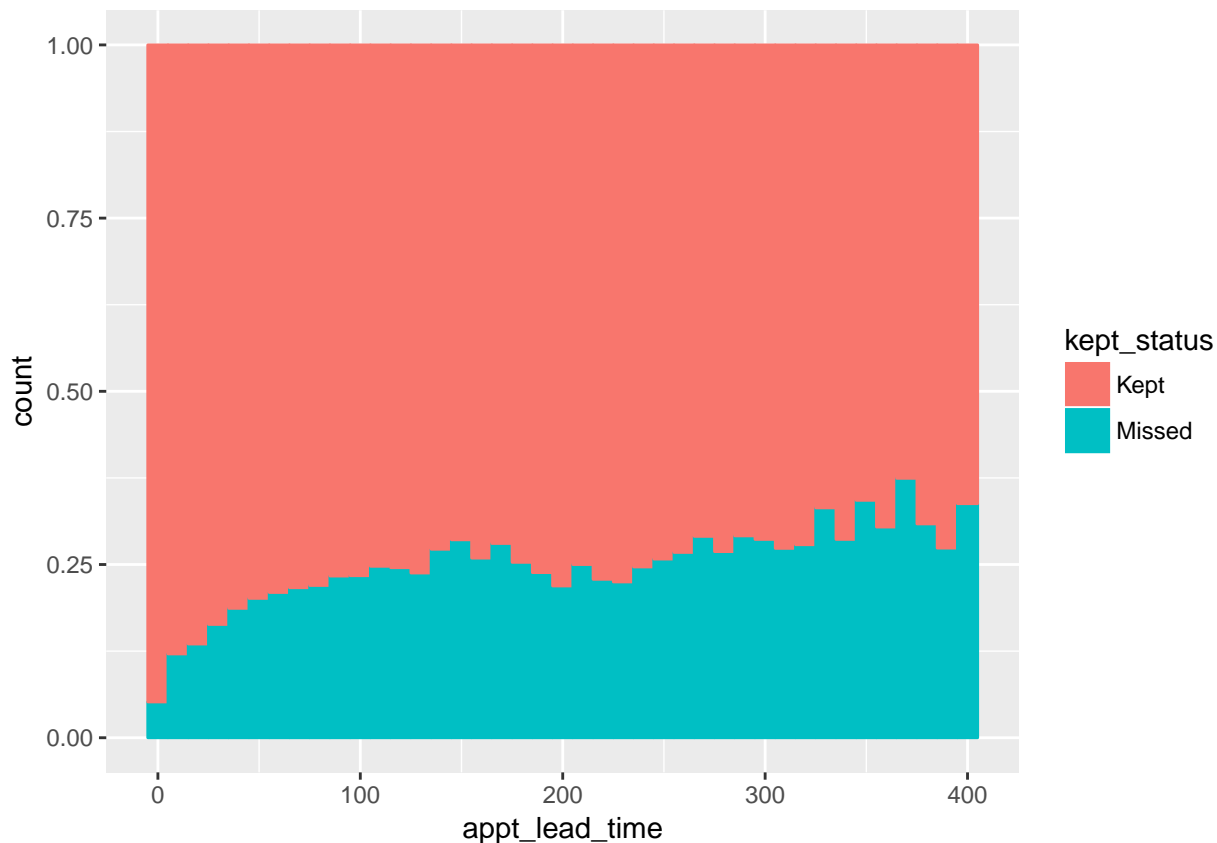


```

appointments %>%
  filter(appt_lead_time >= 0 & appt_lead_time < 400) %>%
  ggplot(
    aes(x = appt_lead_time, color = kept_status, fill = kept_status)
  ) +
    geom_bar(binwidth = 10, position = "fill")

```

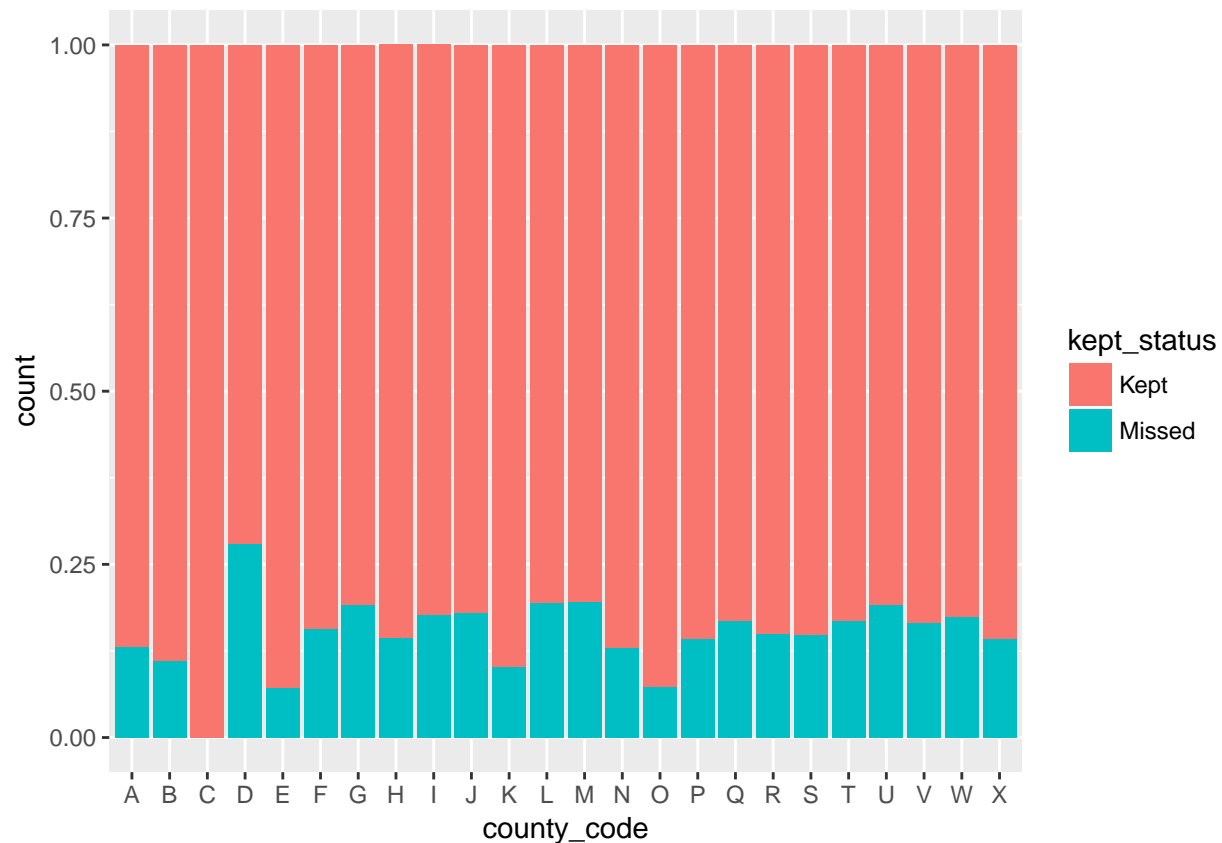
## Warning: `geom\_bar()` no longer has a `binwidth` parameter. Please use  
## `geom\_histogram()` instead.



There is a lower proportion of missed appointments among those with the shortest lead times. This backs up my theory that shorter lead times could indicate more urgent health matters which I wouldn't expect a patient to be as likely to skip. Also, there are small bumps in the histogram around 180 and 360 days, which is probably indicative of regular 6-month and 12-month checkups.

**county\_code**

```
ggplot(  
  appointments,  
  aes(x = county_code, fill = kept_status)  
) +  
  geom_bar(position = "fill")
```



appt\_weekday

```
table(appointments$appt_weekday)
```

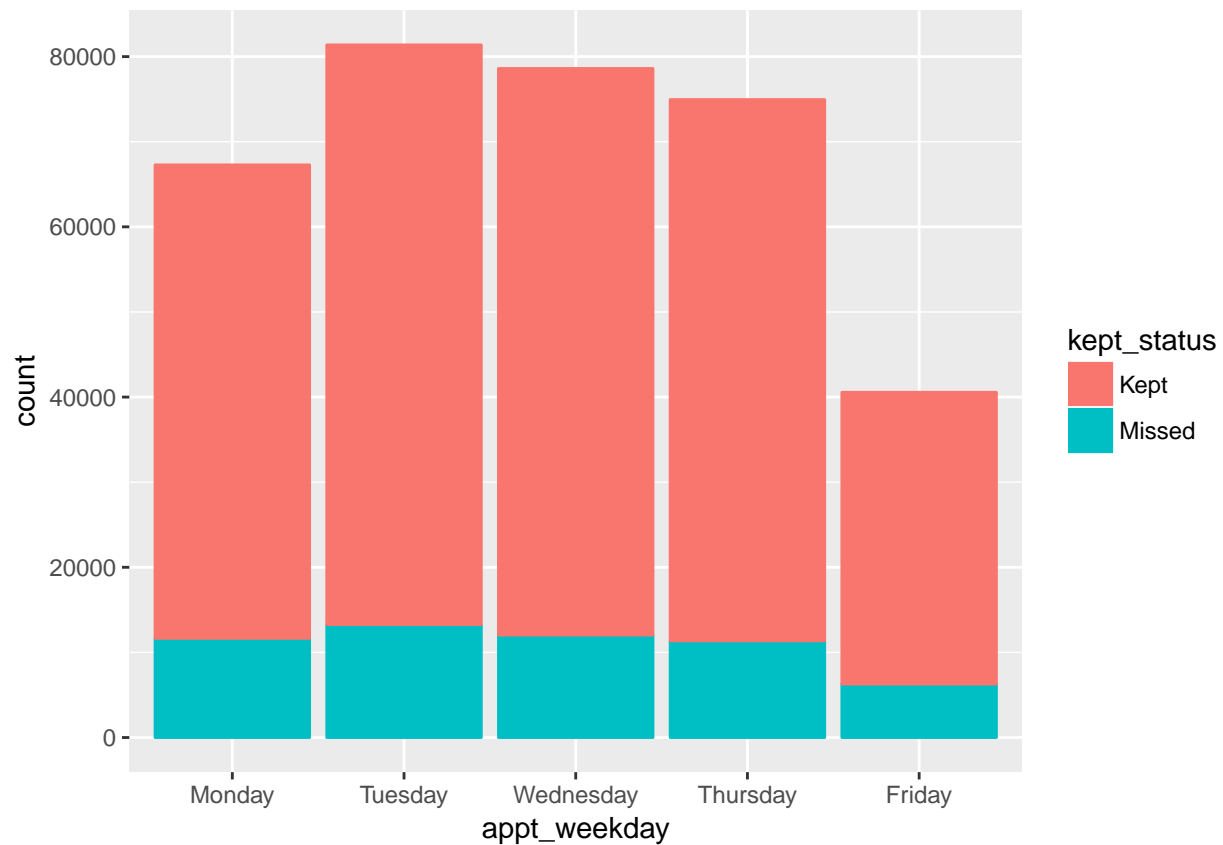
```
##
##    Friday    Monday    Sunday  Thursday  Tuesday Wednesday
##    40558    67280      12    74952    81376    78604
```

There are only 12 Sunday appointments, so I will remove the observations. I will also convert the character variable to an ordered factor to see the days of the week in the correct order.

```
appointments <- appointments %>%
  filter(appt_weekday != "Sunday")
```

```
appointments$appt_weekday <- factor(appointments$appt_weekday, levels = c("Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday"))
```

```
ggplot(
  appointments,
  aes(x = appt_weekday, color = kept_status, fill = kept_status)
) +
  geom_bar()
```



weekday\_scheduled

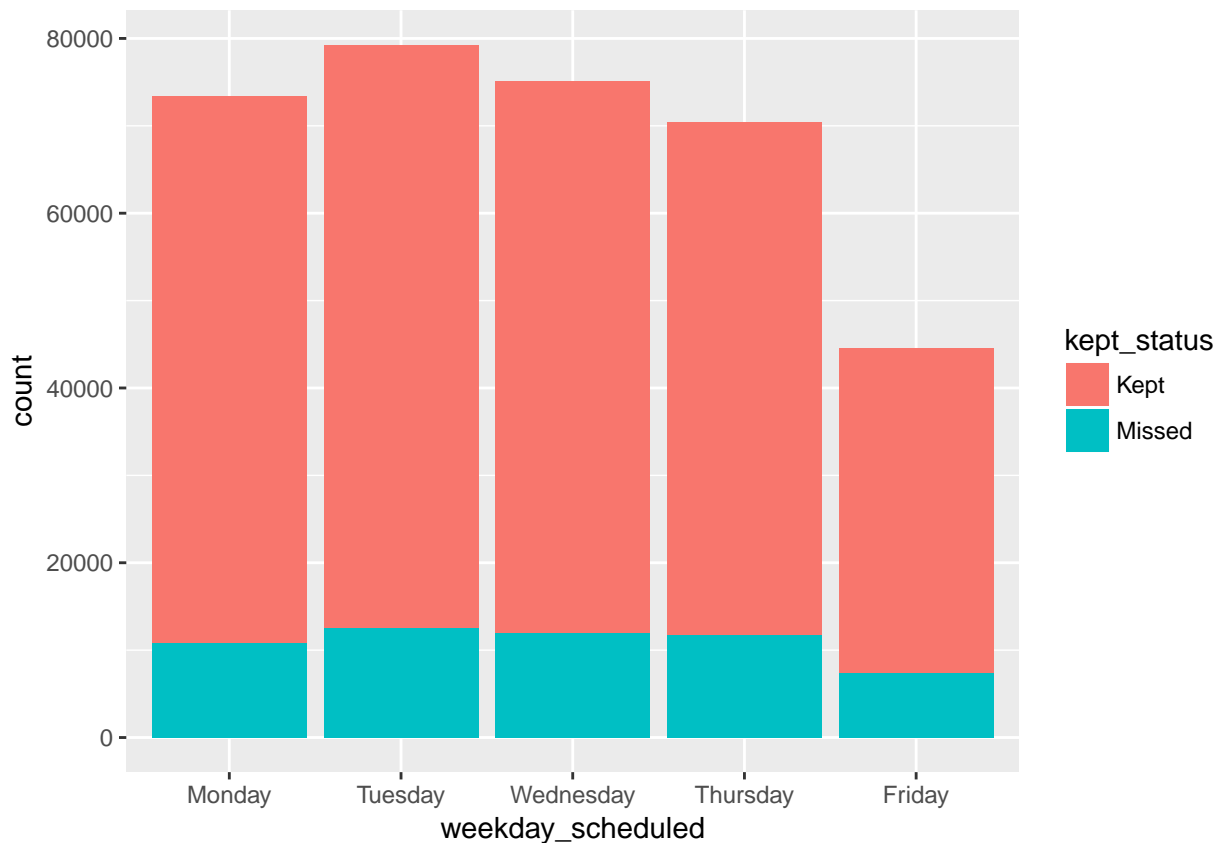
```
table(appointments$weekday_scheduled)
```

```
##
##    Friday    Monday  Saturday    Sunday  Thursday    Tuesday  Wednesday
##    44553    73413      43         7    70441    79261    75052
```

```
appointments <- appointments %>%
  filter(weekday_scheduled != "Sunday") %>%
  filter(weekday_scheduled != "Saturday")
```

```
appointments$weekday_scheduled <- factor(appointments$weekday_scheduled, levels = c("Monday", "Tuesday"
```

```
ggplot(
  appointments,
  aes(x = weekday_scheduled, fill = kept_status)
) +
  geom_bar()
```



## Modeling

### Create Modeling Data

```
model_data <- appointments %>%
  select(
    kept_status, appt_length, patient_age, patient_gender, billing_type,
    patient_distance, provider_specialty, remind_call_result, hour,
    percent_missed, is_new_patient, appt_lead_time, appt_weekday,
    weekday_scheduled, county_code)

factor_columns <- c(
  "kept_status", "patient_gender", "billing_type",
  "provider_specialty", "remind_call_result", "hour", "is_new_patient",
  "appt_weekday", "weekday_scheduled",
  "county_code")
model_data[factor_columns] <- map(model_data[factor_columns], factor)

dummy_vars <- caret::dummyVars(~ ., data = model_data)

model_data_dummy <- data.frame(predict(dummy_vars, newdata = model_data))

linear_combos <- caret::findLinearCombos(model_data_dummy)
linear_combos$remove
```

```
## [1] 8 10 16 25 43 46 52 57 81
lin_combo_index <- linear_combos$remove

model_data_dummy <- model_data_dummy[,-linear_combos$remove]

cor_matrix <- cor(model_data_dummy)
high_cor <- as.data.frame(which(abs(cor_matrix) > 0.90, arr.ind = TRUE))
index <- high_cor %>% filter(row != col)

cor_matrix[index[, 1], index[, 2]]

##                kept_status.Kept kept_status.Missed
## kept_status.Missed      -1.000000000      1.000000000
## kept_status.Kept        1.000000000     -1.000000000
## patient_gender.Male     -0.003350511      0.003350511
## patient_gender.Female    0.003832990     -0.003832990
## provider_specialty.B      0.013063134     -0.013063134
## provider_specialty.A     -0.030816952      0.030816952
##
##                patient_gender.Female patient_gender.Male
## kept_status.Missed     -0.003832990      0.003350511
## kept_status.Kept        0.003832990     -0.003350511
## patient_gender.Male     -0.997632962      1.000000000
## patient_gender.Female    1.000000000     -0.997632962
## provider_specialty.B     -0.004807503      0.004841524
## provider_specialty.A      0.008441594     -0.008524185
##
##                provider_specialty.A provider_specialty.B
## kept_status.Missed      0.030816952     -0.013063134
## kept_status.Kept       -0.030816952      0.013063134
## patient_gender.Male     -0.008524185      0.004841524
## patient_gender.Female    0.008441594     -0.004807503
## provider_specialty.B     -0.915215490      1.000000000
## provider_specialty.A      1.000000000     -0.915215490

colnames(model_data_dummy[index[,1], index[, 2]])

## [1] "kept_status.Kept"      "kept_status.Missed"      "patient_gender.Female"
## [4] "patient_gender.Male"      "provider_specialty.A"      "provider_specialty.B"

model_data_dummy <- model_data_dummy[,-c(2, 6, 11)]

model_data_dummy$kept_status.Kept <- as.factor(model_data_dummy$kept_status.Kept)

near_zero_var <- caret::nearZeroVar(model_data_dummy)
model_data_dummy <- model_data_dummy[,-near_zero_var]
```

## Divide model\_data into train, validate, and test sets

The data will be divided into three sets: train, validate, and test. I will use 60% of the data for training, and 20% each for validation and test.

Because the data is sorted by date, I want to use the most recent data for the test data, the next oldest for validation, and the oldest for training. Since I am trying to predict future appointments, testing on the most

recent data will result in the best measure of the model's performance.

```
# check out caret::downSample
train <- model_data_dummy[1:205660,]
validate <- model_data_dummy[205661:274200,]
test <- model_data_dummy[274201:nrow(model_data_dummy),]
table(train$kept_status.Kept)

##
##      0      1
## 31050 174610

train_balance_subset <- train[168750:205660,]
table(train_balance_subset$kept_status.Kept)

##
##      0      1
##  5861 31050

train_kept <- train_balance_subset[train_balance_subset$kept_status.Kept == 1,]
train_missed <- train[train$kept_status.Kept == 0,]
train_balanced <- rbind(train_kept, train_missed)
table(train_balanced$kept_status.Kept)

##
##      0      1
## 31050 31050
```

need to fix negative appt lead time

### Logistic Regression Model

```
glm_control <- caret::trainControl(method = "none")

glm_model <- caret::train(
  kept_status.Kept ~ .,
  data = train_balanced,
  method = "glm",
  trControl = glm_control
)

summary(glm_model)

##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.1007  -0.9289   0.0396   0.9165   4.5309
##
## Coefficients:
##                                Estimate Std. Error z value
## (Intercept)                3.5598325  0.1590940  22.376
## appt_length               -0.0034995  0.0005922  -5.909
```



## patient_age	0.0106901	0.0005065	21.106
## patient_gender.Female	0.0733094	0.0186968	3.921
## billing_type.Commercial	0.0243503	0.0238422	1.021
## patient_distance	0.0001211	0.0001376	0.880
## provider_specialty.A	-0.4248489	0.0216625	-19.612
## remind_call_result.Answered...Confirmed	0.1104729	0.0408281	2.706
## remind_call_result.Answered...No.Response	-0.9741566	0.0298509	-32.634
## remind_call_result.Failed	-1.5388743	0.0394032	-39.055
## remind_call_result.Left.Message	0.0639698	0.0574839	1.113
## hour.8	0.8119520	0.0875127	9.278
## hour.9	0.8574780	0.0876229	9.786
## hour.10	0.9704986	0.0874921	11.092
## hour.11	0.8118597	0.0881394	9.211
## hour.13	0.9322486	0.0873771	10.669
## hour.14	0.9380374	0.0876698	10.700
## hour.15	0.9978371	0.0875776	11.394
## hour.16	0.9177606	0.0882145	10.404
## percent_missed	-0.0466504	0.0006655	-70.102
## is_new_patient.0	-2.1378325	0.1270566	-16.826
## appt_lead_time	-0.0039104	0.0001155	-33.867
## appt_weekday.Monday	-0.0303453	0.0349868	-0.867
## appt_weekday.Tuesday	0.2074436	0.0333514	6.220
## appt_weekday.Wednesday	0.3158265	0.0338260	9.337
## appt_weekday.Thursday	0.1630567	0.0344419	4.734
## weekday_scheduled.Monday	0.0373300	0.0328507	1.136
## weekday_scheduled.Tuesday	-0.0858016	0.0321183	-2.671
## weekday_scheduled.Wednesday	-0.0584037	0.0325169	-1.796
## weekday_scheduled.Thursday	-0.1095184	0.0326975	-3.349
## county_code.F	0.3861678	0.0364003	10.609
## county_code.I	-0.2882517	0.0303012	-9.513
## county_code.J	-0.0198197	0.0347708	-0.570
## county_code.P	0.0097547	0.0254057	0.384
## county_code.U	-0.3424635	0.0364727	-9.390
##	Pr(> z )		
## (Intercept)	< 2e-16	***	
## appt_length	3.44e-09	***	
## patient_age	< 2e-16	***	
## patient_gender.Female	8.82e-05	***	
## billing_type.Commercial	0.30711		
## patient_distance	0.37876		
## provider_specialty.A	< 2e-16	***	
## remind_call_result.Answered...Confirmed	0.00681	**	
## remind_call_result.Answered...No.Response	< 2e-16	***	
## remind_call_result.Failed	< 2e-16	***	
## remind_call_result.Left.Message	0.26578		
## hour.8	< 2e-16	***	
## hour.9	< 2e-16	***	
## hour.10	< 2e-16	***	
## hour.11	< 2e-16	***	
## hour.13	< 2e-16	***	
## hour.14	< 2e-16	***	
## hour.15	< 2e-16	***	
## hour.16	< 2e-16	***	
## percent_missed	< 2e-16	***	

```
## is_new_patient.0 < 2e-16 ***
## appt_lead_time < 2e-16 ***
## appt_weekday.Monday 0.38576
## appt_weekday.Tuesday 4.97e-10 ***
## appt_weekday.Wednesday < 2e-16 ***
## appt_weekday.Thursday 2.20e-06 ***
## weekday_scheduled.Monday 0.25581
## weekday_scheduled.Tuesday 0.00755 **
## weekday_scheduled.Wednesday 0.07248 .
## weekday_scheduled.Thursday 0.00081 ***
## county_code.F < 2e-16 ***
## county_code.I < 2e-16 ***
## county_code.J 0.56867
## county_code.P 0.70101
## county_code.U < 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: 86089 on 62099 degrees of freedom
## Residual deviance: 68770 on 62065 degrees of freedom
## AIC: 68840
##
## Number of Fisher Scoring iterations: 6
```

## Random Forest Model

Using caret Package

```
rf_control <- caret::trainControl(method = "cv", number = 2, classProbs = FALSE)
seed <- 7
metric <- "Accuracy"
set.seed(seed)
mtry <- 3
tuneGrid <- expand.grid(.mtry = mtry)

rf_model <- caret::train(
  kept_status.Kept ~ ., data = train_balanced, method = "rf", metric = metric,
  tuneGrid = tuneGrid, trControl = rf_control)
```

## Model Comparison

```
pred_glm <- predict(glm_model, validate)

conf_mat_glm <- caret::confusionMatrix(
  pred_glm, validate$kept_status, positive = "0")

conf_mat_glm

## Confusion Matrix and Statistics
##
## Reference
```

```

## Prediction      0      1
##              0 8113 15451
##              1 3343 41633
##
##              Accuracy : 0.7258
##              95% CI : (0.7224, 0.7291)
##      No Information Rate : 0.8329
##      P-Value [Acc > NIR] : 1
##
##              Kappa : 0.3076
##      McNemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.7082
##              Specificity : 0.7293
##      Pos Pred Value : 0.3443
##      Neg Pred Value : 0.9257
##              Prevalence : 0.1671
##      Detection Rate : 0.1184
##      Detection Prevalence : 0.3438
##      Balanced Accuracy : 0.7188
##
##      'Positive' Class : 0
##
conf_mat_glm$byClass["F1"]

##      F1
## 0.4633352

pred_rf <- predict(rf_model, validate)

caret::confusionMatrix(rf_model)

## Cross-Validated (2 fold) Confusion Matrix
##
## (entries are percentual average cell counts across resamples)
##
##      Reference
## Prediction    0      1
##      0 41.4 15.4
##      1  8.6 34.6
##
##      Accuracy (average) : 0.7606

conf_mat_rf <- caret::confusionMatrix(
  pred_rf, validate$kept_status, positive = "0")

conf_mat_rf

## Confusion Matrix and Statistics
##
##      Reference
## Prediction    0      1
##      0 9044 16968
##      1 2412 40116
##

```

```

##           Accuracy : 0.7172
##           95% CI : (0.7139, 0.7206)
##      No Information Rate : 0.8329
##      P-Value [Acc > NIR] : 1
##
##           Kappa : 0.3264
##  McNemar's Test P-Value : <2e-16
##
##      Sensitivity : 0.7895
##      Specificity : 0.7028
##      Pos Pred Value : 0.3477
##      Neg Pred Value : 0.9433
##      Prevalence : 0.1671
##      Detection Rate : 0.1320
##      Detection Prevalence : 0.3795
##      Balanced Accuracy : 0.7461
##
##      'Positive' Class : 0
##

```

```

conf_mat_rf$byClass["F1"]

```

```

##      F1
## 0.4827586

```