

Data Wrangling in R

Homework

Prof. Bisbee

Due Date: 2023-02-01

Univariate Data Analysis

Univariate is pretty much what it sounds like: one variable. When undertaking univariate data analysis, we need first and foremost to figure what type of variable it is that we're working with. Once we do that, we can choose the appropriate use of the variable, either as an outcome or as a possible predictor.

Motivating Question

Today we'll be working with data from every NBA player who was active during the 2018-19 season.

Here's the data:

```
nba<-readRDS("../data/nba_players_2018.Rds")
```

```
## Error in gzfile(file, "rb"): cannot open the connection
```

This data contains the following variables:

Codebook for NBA Data

Name	Definition
namePlayer	Player name
idPlayer	Unique player id
slugSeason	Season start and end
numberPlayerSeason	Which season for this player
isRookie	Rookie season, true or false
slugTeam	Team short name
idTeam	Unique team id
gp	Games Played
gs	Games Started
fgm	Field goals made
fga	Field goals attempted

Name	Definition
pctFG	Percent of field goals made
fg3m	3 point field goals made
fg3a	3 point field goals attempted
pctFG3	Percent of 3 point field goals made
pctFT	Free Throw percentage
fg2m	2 point field goals made
fg2a	2 point field goals attempted
pctFG2	Percent of 2 point field goals made
agePlayer	Player age
minutes	Minutes played
ftm	Free throws made
fta	Free throws attempted
oreb	Offensive rebounds
dreb	Defensive rebounds
treb	Total rebounds
ast	Assists
blk	Blocks
tov	Turnovers
pf	Personal fouls
pts	Total points
urlNBAAPI	Source url

We're interested in the following questions:

- Do certain colleges produce players that have more field goals? What about free throw percentage above a certain level? Are certain colleges in the east or the west more likely to produce higher scorers? How does this vary as a player has more seasons?

To answer these questions we need to look at the following variables:

- Field goals
- Free throw percentage above .25
- Colleges
- Player seasons
- Region

We're going to go through a pretty standard set of steps for each variable. First, examine some cases. Second, based on our examination, we'll try either a plot or a table. Once we've seen the plot or the table, we'll think a bit about ordering, and then choose an appropriate measure of central tendency, and maybe variation.

Types of Variables

It's really important to understand the types of variables you're working with. Many times analysts are indifferent to this step particularly with larger datasets. This can lead to a great deal of confusion down the road. Below are the variable types we'll be working with this semester and the definition of each.

Continuous Variables

A continuous variable can theoretically be subdivided at any arbitrarily small measure and can still be identified. You may have encountered further subdivision of continuous variables into "interval" or "ratio" data in other classes. We RARELY use these distinctions in practice. The distinction between a continuous and a categorical variable is hugely consequential, but the distinction between interval and ratio is not really all that important in practice.

The mean is the most widely used measure of central tendency for a continuous variable. If the distribution of the variable isn't very symmetric or there are large outliers, then the median is a much better measure of central tendency.

Categorical Variables

A categorical variable divides the sample up into a set of mutually exclusive and exhaustive categories. Mutually exclusive means that each case can only be one, and exhaustive means that the categories cover every possible option. Categorical is sort of the "top" level classification for variables of this type. Within the broad classification of categorical there are multiple types of other variables.

Categorical: ordered

an ordered categorical variable has— you guessed it— some kind of sensible order that can be applied. For instance, the educational attainment of an individual: high school diploma, associates degree, bachelor's degree, graduate degree— is an ordered categorical variable.

Ordered categorical variables should be arranged in the order of the variable, with proportions or percentages associated with each order. The mode, or the category with the highest proportion, is a reasonable measure of central tendency, but with fewer than ten categories the analyst should generally just show the proportion in each category.

Categorical: ordered, binary

An ordered binary variable has just two levels, but can be ordered. For instance, is a bird undertaking its first migration: yes or no? A "no" means that the bird has more than one.

The mean of a binary variable is exactly the same thing as the proportion of the sample with that characteristic. So, the mean of a binary variable for "first migration" where 1="yes" will give the proportion of birds migrating for the first time.

An ordered binary variable coded as 0 or 1 can be summarized using the mean which is the same thing as the proportion of the sample with that characteristic.

Categorical: unordered

An unordered categorical variable has no sensible ordering that can be applied. Think about something like college major. There's no "number" we might apply to philosophy that has any meaningful distance from a number we might apply to chemical engineering.

Unlike an ordered variable, an unordered categorical variable should be ordered in terms of the proportions falling into each of the categories. As with an unordered variable, it's best just to show the proportions in each category for variables with less than ten levels. The mode is a reasonable single variable summary of an unordered categorical variable.

Categorical: unordered, binary

This kind of variable has no particular order, but can be just binary. A "1" means that the case has that characteristics, a "0" means the case does not have that characteristic. For instance, whether a tree is deciduous or not.

An unordered binary variable coded as 0 or 1 can also be summarized by the mean, which is the same thing as the proportion of the sample with that characteristic.

Formats for categorical variables

In R, categorical variables CAN be stored as text, numbers or even logicals. Don't count on the data to help you out— you as the analyst need to figure this out.

Factors

We probably need to talk about factors. In R, a factor is a way of storing categorical variables. The factor provides additional information, including an ordering of the variable and a number assigned to each "level" of the factor. A categorical variable is a general term that's understood across statistics. A factor variable is a specific R term. Most of the time it's best not to have a categorical variable structured as a factor unless you know you want it to be a factor. More on this later . . .

The Process: #TrustTheProcess

I'm going to walk you through how an analyst might typically decide what type of variables they're working with. It generally works like this:

1. Take a look at a few observations and form a guess as to what type of variable it is.
2. Based on that guess, create an appropriate plot or table.
3. If the plot or table looks as expected, calculate some summary measures. If not, go back to 1.

"Glimpse" to start: what's in here anyway?

The first thing we're going to do with any dataset is just to take a quick look. We can call the data itself, but that will just show the first few cases and the first few variables. Far better is the `glimpse` command, which shows us all variables and the first few observations for all of the variables. Here's a link to the codebook for this dataset:

The six variables we're going to think about are field goals, free throw percentage, seasons played, rookie season, college attended, and conference played in.

```
glimpse(nba)
```

```
## Error in glimpse(nba): could not find function "glimpse"
```

Continuous

Let's start by taking a look at field goals. It seems pretty likely that this is a continuous variable. Let's take a look at the top 50 spots.

```
nba%>% ## Start with the dataset
  select(namePlayer, slugTeam, fgm)%>% ## and then select a few variables
  arrange(-fgm)%>% ## arrange in reverse order of field goals
  print(n=50) ## print out the top 50
```

```
## Error in nba %>% select(namePlayer, slugTeam, fgm) %>% arrange(-fgm) %>% : could not
find function "%>%"
```

So what I'm seeing here is that field goals aren't "clumped" at certain levels. Let's confirm that by looking at a kernel density plot.

```
nba%>%
  ggplot(aes(x=fgm)) +
  geom_density()
```

```
## Error in nba %>% ggplot(aes(x = fgm)): could not find function "%>%"
```

We can also use a histogram to figure out much the same thing.

```
nba%>%
  ggplot(aes(x=fgm)) +
  geom_histogram()
```

```
## Error in nba %>% ggplot(aes(x = fgm)): could not find function "%>%"
```

Now, technically field goals don't meet the definition I set out above as being a continuous variable because they aren't divisible below a certain amount. Usually in practice though we just ignore this— this variable is "as good as" continuous, given that it varies smoothly over the range and isn't confined to a relatively small set of possible values.

Quick Exercise: Do the same thing for field goal percentage and think about what kind of variable it is.

```
# INSERT CODE HERE
```

Measures for Continuous Variables

The mean is used most of the time for continuous variables, but it's VERY sensitive to outliers. The median (50th percentile) is usually better, but it can be difficult to explain to general audiences.

```
nba%>%
  summarize(mean_fgm=mean(fgm))
```

```
## Error in nba %>% summarize(mean_fgm = mean(fgm)): could not find function "%>%"
```

```
nba%>%
  summarize(median_fgm=median(fgm))
```

```
## Error in nba %>% summarize(median_fgm = median(fgm)): could not find function "%>%"
```

In this case I'd really prefer the mean as a single measure of field goal production, but depending on the audience I still might just go ahead and use the median.

Quick Exercise What measure would you prefer for field goal percentage? Calculate that measure.

```
# INSERT CODE HERE
```

Categorical: ordered

Let's take a look at player seasons.

```
nba%>%
  select(namePlayer,numberPlayerSeason)%>%
  arrange(-numberPlayerSeason)%>%
  print(n=50)
```

```
## Error in nba %>% select(namePlayer, numberPlayerSeason) %>% arrange(-numberPlayerSeason) %>% : could not find function "%>%"
```

Looks like it might be continuous? Let's plot it:

```
nba%>%
  ggplot(aes(x=numberPlayerSeason)) +
  geom_histogram(binwidth = 1)
```

```
## Error in nba %>% ggplot(aes(x = numberPlayerSeason)): could not find function "%>%"
```

Nope. See how it falls into a small set of possible categories? This is an ordered categorical variable. That means we should calculate the proportions in each category

```
nba%>%
  group_by(numberPlayerSeason)%>%
  count(name="total_in_group")%>%
  ungroup()%>%
  mutate(proportion=total_in_group/sum(total_in_group))
```

```
## Error in nba %>% group_by(numberPlayerSeason) %>% count(name = "total_in_group") %>%  
: could not find function "%>%"
```

What does this tell us?

Quick Exercise Create a histogram for player age. What does that tell us about the NBA?

```
# INSERT CODE HERE
```

Categorical: ordered, binary

Let's take a look at the variable for Rookie season.

```
nba%>%select(namePlayer,isRookie)
```

```
## Error in nba %>% select(namePlayer, isRookie): could not find function "%>%"
```

Okay, so that's set to a logical. In R, TRUE or FALSE are special values that indicate the result of a logical question. In this it's whether or not the player is a rookie.

Usually we want a binary variable to have at least one version that's structured so that 1= TRUE and 2=FALSE. This makes data analysis much easier. Let's do that with this variable.

This code uses `ifelse` to create a new variable called `isRookiebin` that's set to 1 if the `isRookie` variable is true, and 0 otherwise.

```
nba<-nba%>%  
  mutate(isRookie_bin=ifelse(isRookie==TRUE,1,0))
```

```
## Error in nba %>% mutate(isRookie_bin = ifelse(isRookie == TRUE, 1, 0)): could not find  
d function "%>%"
```

Now that it's coded 0,1 we can calculate the mean, which is the same thing as the proportion of the players that are rookies.

```
nba%>%summarize(mean=mean(isRookie_bin))
```

```
## Error in nba %>% summarize(mean = mean(isRookie_bin)): could not find function "%>%"
```

Categorical: unordered

Let's take a look at which college a player attended, which is a good example of an unordered categorical variable. The `org` variable tells us which organization the player was in before playing in the NBA.

```
nba%>%  
  select(org)%>%  
  glimpse()
```

```
## Error in nba %>% select(org) %>% glimpse(): could not find function "%>%"
```

This looks like team or college names, so this would be a categorical variable. Let's take a look at the counts of players from different organizations:

```
nba%>%  
  group_by(org)%>%  
  count()%>%  
  arrange(-n)%>%  
  print(n=50)
```

```
## Error in nba %>% group_by(org) %>% count() %>% arrange(-n) %>% print(n = 50): could not find function "%>%"
```

Here we have a problem. If we're interested just in colleges, we're going to need to structure this a bit more. The code below filters out three categories that we don't want: missing data, anything classified as others, and sports teams from other countries. The last is incomplete— I probably missed some! If I were doing this for real, I would use a list of colleges and only include those names.

What I do below is to negate the `str_detect` variable by placing the `!` in front of it. This means I want all of the cases that don't match the pattern I supplied. The pattern makes heavy use of the OR operator `|`. I'm saying I don't want to include players whose organization included the letters `CB` `r` `KK` and so on (these are common prefixes for sports organizations in other countries, I definitely did not look that up on Wikipedia. Ok, I did.).

```
nba%>%  
  filter(!is.na(org))%>%  
  filter(!org=="Other")%>%  
  filter(!str_detect(org,"CB|KK|rytas|FC|B.C.|S.K.|Madrid"))%>%  
  group_by(org)%>%  
  count()%>%  
  arrange(-n)%>%  
  print(n=50)
```

```
## Error in nba %>% filter(!is.na(org)) %>% filter(!org == "Other") %>% filter(!str_detect(org, : could not find function "%>%"
```

That looks better. Which are the most common colleges and universities that send players to the NBA?

Quick Exercise Arrange the number of players by team in descending order.

```
# INSERT CODE HERE
```

Categorical: unordered, binary

There are two conferences in the NBA, eastern and western. Let's take a look at the variable that indicates which conference the player played in that season.


```
nba%>%select(idConference)%>%
  glimpse()
```

```
## Error in nba %>% select(idConference) %>% glimpse(): could not find function "%>%"
```

It looks like conference is structured as numeric, but a “1” or a “2”. Because it’s best to have binary variables structured as “has the characteristic” or “doesn’t have the characteristic” we’re going to create a variable for western conference that’s set to 1 if the player was playing in the western conference and 0 if the player was not (this is the same as playing in the eastern conference).

```
nba<-nba%>%
  mutate(west_conference=ifelse(idConference==1,1,0))
```

```
## Error in nba %>% mutate(west_conference = ifelse(idConference == 1, 1, : could not fi
nd function "%>%"
```

Once we’ve done that, we can see how many players played in each conference.

```
nba%>%
  summarize(mean(west_conference))
```

```
## Error in nba %>% summarize(mean(west_conference)): could not find function "%>%"
```

Makes sense!

Quick Exercise:* create a variable for whether or not the player is from the USA. Calculate the proportion of players from the USA in the NBA. The coding on country is ... decidedly US-centric, so you’ll need to think about this one a bit.

```
# INSERT CODE HERE
```

Analysis

Ok, now that we know how this works, we can do some summary analysis. First of all, what does the total number of field goals made look like by college?

We know that field goals are continuous (sort of) so let’s summarize them via the mean. We know that college is a categorical variable, so we’ll use that to group the data. This is one of our first examples of a conditional mean, which we’ll use a lot.

Top 50 Colleges by Total FG

```
nba%>%
  filter(!is.na(org))%>%
  filter(!org=="Other")%>%
  filter(!str_detect(org,"CB|KK|rytas|FC|B.C.|S.K.|Madrid"))%>%
  group_by(org)%>%
  summarize(mean_fg=sum(fgm))%>%
  arrange(-mean_fg)%>%
  print(n=50)
```

```
## Error in nba %>% filter(!is.na(org)) %>% filter(!org == "Other") %>% filter(!str_dete
ct(org, : could not find function "%>%"
```

Next, what about field goal percentage?

Top 50 Colleges by Average Field Goal Percent

```
nba%>%
  filter(!is.na(org))%>%
  filter(!org=="Other")%>%
  filter(!str_detect(org,"CB|KK|rytas|FC|B.C.|S.K.|Madrid"))%>%
  group_by(org)%>%
  summarize(mean_ftp=mean(pctFT))%>%
  arrange(-mean_ftp)%>%
  print(n=50)
```

```
## Error in nba %>% filter(!is.na(org)) %>% filter(!org == "Other") %>% filter(!str_dete
ct(org, : could not find function "%>%"
```

Quick Exercise Calculate field goals made by player season.

```
# INSERT CODE HERE
```

Quick Exercise Calculate free throw percent made by player season.

```
# INSERT CODE HERE
```