

```

---
title: "Class09miniprojectinclass"
author: 'San Luc (PID: A59010657)'
date: "10/29/2021"
output:
  html_document: default
  pdf_document: default
---

```

Importing candy data

```

```{r}
candy_file <- "candy-data.csv"

candy = read.csv(candy_file, row.names=1)
head(candy)
```

```

Q1. How many different candy types are in this dataset?

```

```{r}
dim(candy)
nrow(candy)
```

```

Q2. How many fruity candy types are in the dataset?

```

```{r}
table(candy$fruity)
```

```

The functions `dim()`, `nrow()`, `table()` and `sum()` may be useful for answering the first 2 questions.

We can find the `winpercent` value for Twix by using its name to access the corresponding row of the dataset. This is because the dataset has each candy name as rownames (recall that we set this when we imported the original CSV file). For example the code for Twix is:

```

# What is your favorite candy?
```{r}
candy["Twix",]$winpercent
```

```

Q3. What is your favorite candy in the dataset and what is its `winpercent` value?

```

```{r}
candy["Swedish Fish",]$winpercent
```

```

Q4. What is the `winpercent` value for “Kit Kat”?

```

```{r}
candy["Kit Kat",]$winpercent
```

```

Q5. What is the `winpercent` value for “Tootsie Roll Snack Bars”?

```

```{r}
candy["Tootsie Roll Snack Bars",]$winpercent
```

```

Side-note: the `skimr::skim()` function

There is a useful `skim()` function in the `skimr` package that can help give you a quick overview of a given dataset. Let’s install this package and try it on our candy data.

```

```{r}

```

```
library("skimr")
skim(candy)
```

```

Q6. Is there any variable/column that looks to be on a different scale to the majority of the other columns in the dataset?

Q7. What do you think a zero and one represent for the candy\$chocolate column?

Hint: look at the “Variable type” print out from the skim() function. Most variables (i.e. columns) are on the zero to one scale but not all. Some columns such as chocolate are exclusively either zero or one values.

A good place to start any exploratory analysis is with a histogram. You can do this most easily with the base R function hist(). Alternatively, you can use ggplot() with geom_hist(). Either works well in this case and (as always) its your choice.

Q8. Plot a histogram of winpercent values

```
```{r}
hist(candy$winpercent)
```

```

Q9. Is the distribution of winpercent values symmetrical?
non-symmetrical

Q10. Is the center of the distribution above or below 50%?

```
```{r}
summary(candy$winpercent)
```
below 50%
```

Q11. On average is chocolate candy higher or lower ranked than fruity candy?

```
```{r}
candy_chocolate <- candy$winpercent[as.logical(candy$chocolate)]
mean(candy_chocolate)
```
```{r}
candy_fruity <- candy$winpercent[as.logical(candy$fruity)]
mean(candy_fruity)
```

```

Cho

Q12. Is this difference statistically significant?

Hint: The chocolate, fruity, nougat etc. columns indicate if a given candy has this feature (i.e. one if it has nougat, zero if it does not etc.). We can turn these into logical (a.k.a. TRUE/FALSE) values with the as.logical() function. We can then use this logical vector to access the corresponding candy rows (those with TRUE values). For example to get the winpercent values for all nougat containing candy we can use the code: candy\$winpercent[as.logical(candy\$nougat)]. In addition the functions mean() and t.test() should help you answer the last two questions here.

```
```{r}
t.test(candy_chocolate, candy_fruity)
```

```

Overall Candy Rankings

Q13. What are the five least liked candy types in this set?

```
```{r}
head(candy[order(candy$winpercent),], n=5)
```

```

Q14. What are the top 5 all time favorite candy types out of this set?

```

```{r}
head(candy[order(candy$winpercent, decreasing = TRUE),], n=5)
```

```

To examine more of the dataset in this vain we can make a barplot to visualize the overall rankings. We will use an iterative approach to building a useful visulization by getting a rough starting plot and then refining and adding useful details in a stepwise process.

Q15. Make a first barplot of candy ranking based on winpercent values.
 HINT: Use the aes(winpercent, rownames(candy)) for your first ggplot like so:

```

```{r}
library("ggplot2")

```{r}
ggplot(candy) +
  aes(winpercent, rownames(candy)) +
  geom_col()
```

```

Q16. This is quite ugly, use the reorder() function to get the bars sorted by winpercent?

```

```{r}
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy),winpercent)) +
  geom_col()
```

```

Time to add some useful color

Let's setup a color vector (that signifies candy type) that we can then use for some future plots. We start by making a vector of all black values (one for each candy). Then we overwrite chocolate (for chocolate candy), brown (for candy bars) and red (for fruity candy) values.

```

```{r}
my_cols=rep("black", nrow(candy))
my_cols[as.logical(candy$chocolate)] = "chocolate"
my_cols[as.logical(candy$bar)] = "brown"
my_cols[as.logical(candy$fruity)] = "pink"
```

```

Now let's try our barplot with these colors. Note that we use fill=my\_cols for geom\_col(). Experement to see what happens if you use col=mycols.

```

```{r}
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy),winpercent)) +
  geom_col(fill=my_cols)
```

```

Now, for the first time, using this plot we can answer questions like:

- Q17. What is the worst ranked chocolate candy?  
Sixlets
- Q18. What is the best ranked fruity candy?  
Starburst

# Taking a look at pricepercent

```

```{r}
library(ggrepel)
```

```

# How about a plot of price vs win

```

```{r}
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(col=my_cols, size=3.3, max.overlaps = 5)
```

```

Q19. Which candy type is the highest ranked in terms of winpercent for the least money - i.e. offers the most bang for your buck?

```

```{r}
ord <- order(candy$pricepercent, decreasing = FALSE)
head(candy[ord,c(11,12)], n=5 )
```

```

Tootsie Roll Midgies is the best budgeted candy.

Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?

```

```{r}
ord <- order(candy$pricepercent, decreasing = TRUE)
head(candy[ord,c(11,12)], n=5 )
```

```

The most expensive ones are Nik L Nip, and it has one of the worst rating.

# Exploring the correlation structure

```

```{r}
library(corrplot)
## corrplot 0.90 loaded

```{r}
cij <- cor(candy)
corrplot(cij)
```

```

Q22. Examining this plot what two variables are anti-correlated (i.e. have minus values)?
Chocolate and fruity

Q23. Similarly, what two variables are most positively correlated?
Chocolate and winpercent

Let's apply PCA using the prcomp() function to our candy dataset remembering to set the scale=TRUE argument.

```

```{r}
pca <- prcomp(candy, scale=TRUE)
summary(pca)
```

```

Now we can plot our main PCA score plot of PC1 vs PC2.

```

```{r}
plot(pca$x[,1:2])
```

```

We can change the plotting character and add some color:

```

```{r}
plot(pca$x[,1:2], col=my_cols, pch=16)
```

```

Make a new data-frame with our PCA results and candy data

```

```{r}
my_data <- cbind(candy, pca$x[,1:3])
```

```

```

```{r}
p <- ggplot(my_data) +
 aes(x=PC1, y=PC2,
 size=winpercent/100,
 text=rownames(my_data),
 label=rownames(my_data)) +
 geom_point(col=my_cols)

p,

```{r}
library(ggrepel)

p + geom_text_repel(size=3.3, col=my_cols, max.overlaps = 7) +
  theme(legend.position = "none") +
  labs(title="Halloween Candy PCA Space",
       subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brown), fruity (red), other
(black)",
       caption="Data from 538")

```

```{r}
library(plotly)

```{r}
ggplotly(p)
```

```

Let's finish by taking a quick look at PCA our loadings. Do these make sense to you? Notice the opposite effects of chocolate and fruity and the similar effects of chocolate and bar (i.e. we already know they are correlated).

```

```{r}
par(mar=c(8,4,2,2))
barplot(pca$rotation[,1], las=2, ylab="PC1 Contribution")
```

```

Q24. What original variables are picked up strongly by PC1 in the positive direction? Do these make sense to you? Fruity and pluribus.

HINT. pluribus means the candy comes in a bag or box of multiple candies.