# class10miniproect

## James Woolley A16440072

```r
candy.df <- read.csv("candy-data.csv", row.names= 1)
candy <- candy.df
head(candy.df)
```

|               | chocolate | fruity | caramel | peanutyalmondy | nougat | crispedricewafer |
|---------------|-----------|--------|---------|----------------|--------|------------------|
| 100 Grand     | 1         | 0      | 1       | 0              | 0      | 1                |
| 3 Musketeers  | 1         | 0      | 0       | 0              | 1      | 0                |
| One dime      | 0         | 0      | 0       | 0              | 0      | 0                |
| One quarter   | 0         | 0      | 0       | 0              | 0      | 0                |
| Air Heads     | 0         | 1      | 0       | 0              | 0      | 0                |
| Almond Joy    | 1         | 0      | 0       | 1              | 0      | 0                |

|               | hard | bar | pluribus | sugarpercent | pricepercent | winpercent |
|---------------|------|-----|----------|--------------|--------------|------------|
| 100 Grand     | 0    | 1   | 0        | 0.732        | 0.860        | 66.97173   |
| 3 Musketeers  | 0    | 1   | 0        | 0.604        | 0.511        | 67.60294   |
| One dime      | 0    | 0   | 0        | 0.011        | 0.116        | 32.26109   |
| One quarter   | 0    | 0   | 0        | 0.011        | 0.511        | 46.11650   |
| Air Heads     | 0    | 0   | 0        | 0.906        | 0.511        | 52.34146   |
| Almond Joy    | 0    | 1   | 0        | 0.465        | 0.767        | 50.34755   |

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

```r
nrow(candy)
```

```
[1] 85
```

We can use the `ncol()` function to find that there are 85 different types of candy being compared.

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

```r
sum(candy$fruity)
```

[1] 38

We can use the sum function with a row argument to find that there are 38 fruity types.

```r
candy[as.logical(candy$chocolate),]
```

|                          | chocolate | fruity | caramel | peanutyalmondy | nougat |
|--------------------------|-----------|--------|---------|----------------|--------|
| 100 Grand                | 1         | 0      | 1       | 0              | 0      |
| 3 Musketeers             | 1         | 0      | 0       | 0              | 1      |
| Almond Joy               | 1         | 0      | 0       | 1              | 0      |
| Baby Ruth                | 1         | 0      | 1       | 1              | 1      |
| Charleston Chew          | 1         | 0      | 0       | 0              | 1      |
| Hershey's Kisses         | 1         | 0      | 0       | 0              | 0      |
| Hershey's Krackel        | 1         | 0      | 0       | 0              | 0      |
| Hershey's Milk Chocolate | 1         | 0      | 0       | 0              | 0      |
| Hershey's Special Dark   | 1         | 0      | 0       | 0              | 0      |
| Junior Mints             | 1         | 0      | 0       | 0              | 0      |
| Kit Kat                  | 1         | 0      | 0       | 0              | 0      |
| Peanut butter M&M's      | 1         | 0      | 0       | 1              | 0      |
| M&M's                    | 1         | 0      | 0       | 0              | 0      |
| Milk Duds                | 1         | 0      | 1       | 0              | 0      |
| Milky Way                | 1         | 0      | 1       | 0              | 1      |
| Milky Way Midnight       | 1         | 0      | 1       | 0              | 1      |
| Milky Way Simply Caramel | 1         | 0      | 1       | 0              | 0      |
| Mounds                   | 1         | 0      | 0       | 0              | 0      |
| Mr Good Bar              | 1         | 0      | 0       | 1              | 0      |
| Nestle Butterfinger      | 1         | 0      | 0       | 1              | 0      |
| Nestle Crunch            | 1         | 0      | 0       | 0              | 0      |
| Peanut M&Ms              | 1         | 0      | 0       | 1              | 0      |
| Reese's Miniatures       | 1         | 0      | 0       | 1              | 0      |
| Reese's Peanut Butter cup| 1         | 0      | 0       | 1              | 0      |
| Reese's pieces           | 1         | 0      | 0       | 1              | 0      |
| Reese's stuffed with pieces | 1      | 0      | 0       | 1              | 0      |
| Rolo                     | 1         | 0      | 1       | 0              | 0      |
| Sixlets                  | 1         | 0      | 0       | 0              | 0      |
| Nestle Smarties          | 1         | 0      | 0       | 0              | 0      |
| Snickers                 | 1         | 0      | 1       | 1              | 1      |
| Snickers Crisper         | 1         | 0      | 1       | 1              | 0      |

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| Tootsie Pop | 1 | 1 | 0 | 0 | 0 |
| Tootsie Roll Juniors | 1 | 0 | 0 | 0 | 0 |
| Tootsie Roll Midgies | 1 | 0 | 0 | 0 | 0 |
| Tootsie Roll Snack Bars | 1 | 0 | 0 | 0 | 0 |
| Twix | 1 | 0 | 1 | 0 | 0 |
| Whoppers | 1 | 0 | 0 | 0 | 0 |

|  | crispedricewafer | hard | bar | pluribus | sugarpercent |
|---|---|---|---|---|---|
| 100 Grand | 1 | 0 | 1 | 0 | 0.732 |
| 3 Musketeers | 0 | 0 | 1 | 0 | 0.604 |
| Almond Joy | 0 | 0 | 1 | 0 | 0.465 |
| Baby Ruth | 0 | 0 | 1 | 0 | 0.604 |
| Charleston Chew | 0 | 0 | 1 | 0 | 0.604 |
| Hershey's Kisses | 0 | 0 | 0 | 1 | 0.127 |
| Hershey's Krackel | 1 | 0 | 1 | 0 | 0.430 |
| Hershey's Milk Chocolate | 0 | 0 | 1 | 0 | 0.430 |
| Hershey's Special Dark | 0 | 0 | 1 | 0 | 0.430 |
| Junior Mints | 0 | 0 | 0 | 1 | 0.197 |
| Kit Kat | 1 | 0 | 1 | 0 | 0.313 |
| Peanut butter M&M's | 0 | 0 | 0 | 1 | 0.825 |
| M&M's | 0 | 0 | 0 | 1 | 0.825 |
| Milk Duds | 0 | 0 | 0 | 1 | 0.302 |
| Milky Way | 0 | 0 | 1 | 0 | 0.604 |
| Milky Way Midnight | 0 | 0 | 1 | 0 | 0.732 |
| Milky Way Simply Caramel | 0 | 0 | 1 | 0 | 0.965 |
| Mounds | 0 | 0 | 1 | 0 | 0.313 |
| Mr Good Bar | 0 | 0 | 1 | 0 | 0.313 |
| Nestle Butterfinger | 0 | 0 | 1 | 0 | 0.604 |
| Nestle Crunch | 1 | 0 | 1 | 0 | 0.313 |
| Peanut M&Ms | 0 | 0 | 0 | 1 | 0.593 |
| Reese's Miniatures | 0 | 0 | 0 | 0 | 0.034 |
| Reese's Peanut Butter cup | 0 | 0 | 0 | 0 | 0.720 |
| Reese's pieces | 0 | 0 | 0 | 1 | 0.406 |
| Reese's stuffed with pieces | 0 | 0 | 0 | 0 | 0.988 |
| Rolo | 0 | 0 | 0 | 1 | 0.860 |
| Sixlets | 0 | 0 | 0 | 1 | 0.220 |
| Nestle Smarties | 0 | 0 | 0 | 1 | 0.267 |
| Snickers | 0 | 0 | 1 | 0 | 0.546 |
| Snickers Crisper | 1 | 0 | 1 | 0 | 0.604 |
| Tootsie Pop | 0 | 1 | 0 | 0 | 0.604 |
| Tootsie Roll Juniors | 0 | 0 | 0 | 0 | 0.313 |
| Tootsie Roll Midgies | 0 | 0 | 0 | 1 | 0.174 |
| Tootsie Roll Snack Bars | 0 | 0 | 1 | 0 | 0.465 |
| Twix | 1 | 0 | 1 | 0 | 0.546 |

```
Whoppers                                1    0    0         1         0.872
                              pricepercent winpercent
100 Grand                            0.860   66.97173
3 Musketeers                         0.511   67.60294
Almond Joy                           0.767   50.34755
Baby Ruth                            0.767   56.91455
Charleston Chew                      0.511   38.97504
Hershey's Kisses                     0.093   55.37545
Hershey's Krackel                    0.918   62.28448
Hershey's Milk Chocolate             0.918   56.49050
Hershey's Special Dark               0.918   59.23612
Junior Mints                         0.511   57.21925
Kit Kat                              0.511   76.76860
Peanut butter M&M's                  0.651   71.46505
M&M's                                0.651   66.57458
Milk Duds                            0.511   55.06407
Milky Way                            0.651   73.09956
Milky Way Midnight                   0.441   60.80070
Milky Way Simply Caramel             0.860   64.35334
Mounds                               0.860   47.82975
Mr Good Bar                          0.918   54.52645
Nestle Butterfinger                  0.767   70.73564
Nestle Crunch                        0.767   66.47068
Peanut M&Ms                          0.651   69.48379
Reese's Miniatures                   0.279   81.86626
Reese's Peanut Butter cup            0.651   84.18029
Reese's pieces                       0.651   73.43499
Reese's stuffed with pieces          0.651   72.88790
Rolo                                 0.860   65.71629
Sixlets                              0.081   34.72200
Nestle Smarties                      0.976   37.88719
Snickers                             0.651   76.67378
Snickers Crisper                     0.651   59.52925
Tootsie Pop                          0.325   48.98265
Tootsie Roll Juniors                 0.511   43.06890
Tootsie Roll Midgies                 0.011   45.73675
Tootsie Roll Snack Bars              0.325   49.65350
Twix                                 0.906   81.64291
Whoppers                             0.848   49.52411
```

Q3. What is your favorite candy in the dataset and what is it's winpercent value?

```
candy["Nestle Butterfinger", ]$winpercent
```

`[1] 70.73564`

My favourite candy is butterfingers, which is very popular with a `winpercent` of 70.74%. >Q4. What is the winpercent value for "Kit Kat"?

```
candy["Kit Kat", ]$winpercent
```

`[1] 76.7686`

Kit Kats have a `winpercent` of 76.77. Not sure why given that they're vile bars of cardboard. >Q5. What is the winpercent value for "Tootsie Roll Snack Bars"?

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

`[1] 49.6535`

This candy has a 49.65% `winpercent`, and isn't very popular.

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

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

Table 1: Data summary

| Name | candy |
|---|---|
| Number of rows | 85 |
| Number of columns | 12 |
| | |
| Column type frequency: | |
| numeric | 12 |
| | |
| Group variables | None |

**Variable type: numeric**

| skim_variable | n_missing | complete_rate | mean | sd | p0 | p25 | p50 | p75 | p100 | hist |
|---|---|---|---|---|---|---|---|---|---|---|
| chocolate | 0 | 1 | 0.44 | 0.50 | 0.00 | 0.00 | 0.00 | 1.00 | 1.00 | |
| fruity | 0 | 1 | 0.45 | 0.50 | 0.00 | 0.00 | 0.00 | 1.00 | 1.00 | |
| caramel | 0 | 1 | 0.16 | 0.37 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| peanutyalmondy | 0 | 1 | 0.16 | 0.37 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| nougat | 0 | 1 | 0.08 | 0.28 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| crispedricewafer | 0 | 1 | 0.08 | 0.28 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| hard | 0 | 1 | 0.18 | 0.38 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| bar | 0 | 1 | 0.25 | 0.43 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | |
| pluribus | 0 | 1 | 0.52 | 0.50 | 0.00 | 0.00 | 1.00 | 1.00 | 1.00 | |
| sugarpercent | 0 | 1 | 0.48 | 0.28 | 0.01 | 0.22 | 0.47 | 0.73 | 0.99 | |
| pricepercent | 0 | 1 | 0.47 | 0.29 | 0.01 | 0.26 | 0.47 | 0.65 | 0.98 | |
| winpercent | 0 | 1 | 50.32 | 14.71 | 22.45 | 39.14 | 47.83 | 59.86 | 84.18 | |

Yes. The `winpercent` column is on a very different scale.

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

The 1s and 0s show whether or not the candy include chocolate.

Q8. Plot a histogram of winpercent values.

```
library(ggplot2)

ggplot(candy) +
  aes(winpercent) +
geom_histogram(binwidth=10)
```

6

Q9. Is the distribution of winpercent values symmetrical?

No. The distribution of `winpercent` values is not symmetrical >Q10. Is the center of the distribution above or below 50%?

The center of the distribution is above 50%. >Q11. On average is chocolate candy higher or lower ranked than fruit candy?

```
choc.inds <- as.logical(candy$chocolate)
choc.win <- candy[choc.inds,]$winpercent
mean(choc.win)
```

```
[1] 60.92153
```

We can see that chocolate wins about 61% of the time.

```
fruit.inds <- as.logical(candy$fruit)
fruit.win <- candy[fruit.inds,]$winpercent
mean(fruit.win)
```

```
[1] 44.11974
```

and that fruit wins about 44% of the time.

Q12. Is this difference statistically significant?

```
t.test(choc.win, fruit.win)
```

```
	Welch Two Sample t-test

data:  choc.win and fruit.win
t = 6.2582, df = 68.882, p-value = 2.871e-08
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 11.44563 22.15795
sample estimates:
mean of x mean of y
 60.92153  44.11974
```

We can use the `ttest()` function to show that this is a statistically significant difference.

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

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

|  | chocolate | fruity | caramel | peanutyalmondy | nougat |
|---|---|---|---|---|---|
| Nik L Nip | 0 | 1 | 0 | 0 | 0 |
| Boston Baked Beans | 0 | 0 | 0 | 1 | 0 |
| Chiclets | 0 | 1 | 0 | 0 | 0 |
| Super Bubble | 0 | 1 | 0 | 0 | 0 |
| Jawbusters | 0 | 1 | 0 | 0 | 0 |

|  | crispedricewafer | hard | bar | pluribus | sugarpercent | pricepercent |
|---|---|---|---|---|---|---|
| Nik L Nip | 0 | 0 | 0 | 1 | 0.197 | 0.976 |
| Boston Baked Beans | 0 | 0 | 0 | 1 | 0.313 | 0.511 |
| Chiclets | 0 | 0 | 0 | 1 | 0.046 | 0.325 |
| Super Bubble | 0 | 0 | 0 | 0 | 0.162 | 0.116 |
| Jawbusters | 0 | 1 | 0 | 1 | 0.093 | 0.511 |

|  | winpercent |
|---|---|
| Nik L Nip | 22.44534 |
| Boston Baked Beans | 23.41782 |
| Chiclets | 24.52499 |
| Super Bubble | 27.30386 |
| Jawbusters | 28.12744 |

We can list the set by `winpercent` and then list the bottom five to see the least liked candies are Nik L Nips, Boston Baked Beans, Chiclets, Super Bubbles, and Jawbusters. >Q14.What are the top 5 all time favorite candy types out of this set?

```
tail(candy[order(candy$winpercent),], n=5)
```

|  | chocolate | fruity | caramel | peanutyalmondy | nougat |
|---|---|---|---|---|---|
| Snickers | 1 | 0 | 1 | 1 | 1 |
| Kit Kat | 1 | 0 | 0 | 0 | 0 |
| Twix | 1 | 0 | 1 | 0 | 0 |
| Reese's Miniatures | 1 | 0 | 0 | 1 | 0 |
| Reese's Peanut Butter cup | 1 | 0 | 0 | 1 | 0 |

|  | crispedricewafer | hard | bar | pluribus | sugarpercent |
|---|---|---|---|---|---|
| Snickers | 0 | 0 | 1 | 0 | 0.546 |
| Kit Kat | 1 | 0 | 1 | 0 | 0.313 |
| Twix | 1 | 0 | 1 | 0 | 0.546 |
| Reese's Miniatures | 0 | 0 | 0 | 0 | 0.034 |
| Reese's Peanut Butter cup | 0 | 0 | 0 | 0 | 0.720 |

|  | pricepercent | winpercent |
|---|---|---|
| Snickers | 0.651 | 76.67378 |
| Kit Kat | 0.511 | 76.76860 |
| Twix | 0.906 | 81.64291 |
| Reese's Miniatures | 0.279 | 81.86626 |
| Reese's Peanut Butter cup | 0.651 | 84.18029 |

You can use the tail function to see that Snickers, Kit Kats (idk why), Twix bars, Reese's minis, and Reese's PB cups are the most popular.

Q15. Make a first barplot of candy ranking based on winpercent values.

```
library(ggplot2)

ggplot(candy) +
  aes(winpercent, reorder(rownames(candy), winpercent)) +
  geom_col(fill="blue")
```

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

```
library(ggplot2)

ggplot(candy) +
  aes(winpercent, reorder(rownames(candy), winpercent)) +
  geom_col(fill="blue")
```

```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"

library(ggplot2)

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

Q17. What is the worst ranked chocolate candy?

We can see from the helpfully colour-coded graph that Sixlets are the lowest ranking chocolate candy.

Q18. What is the best ranked fruity candy?

We can see from the helpfully colour-coded graph that Starburst are the highest ranking fruity candy.

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
library(ggrepel)

ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(col=my_cols, size=1.8, max.overlaps = 50)
```

We can generate a plot that compares `pricepercent` and `winpercent` and see that Reese's miniatures and Starburst are the most popular for the price point. >Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?

The most expensive candies are Nik L Dips, Ring Pops, Nestle Smarties, Hershey's Krackel, and Hershey's Milk Chocolates. The least popular are Nik L Dips.

```
library(corrplot)
```

```
corrplot 0.92 loaded
```

```
cij <- cor(candy)
corrplot(cij)
```

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

Fruity chocolates are highly anti-correlated.

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

The most positively correlated are the same variables, but also chocolate and winpercent.

```r
pca <- prcomp(candy, scale=TRUE)#First we run a PCA of the data
plot(pca$x[,1:2], col=my_cols, pch=16) #then we give it some colors and make the dot size
```

```r
my_data <- cbind(candy, pca$x[,1:3])
p <- ggplot(my_data) + #plot the data nicely in ggplot
        aes(x=PC1, y=PC2,
            size=winpercent/100,
            text=rownames(my_data),
            label=rownames(my_data)) +
        geom_point(col=my_cols)


p
```
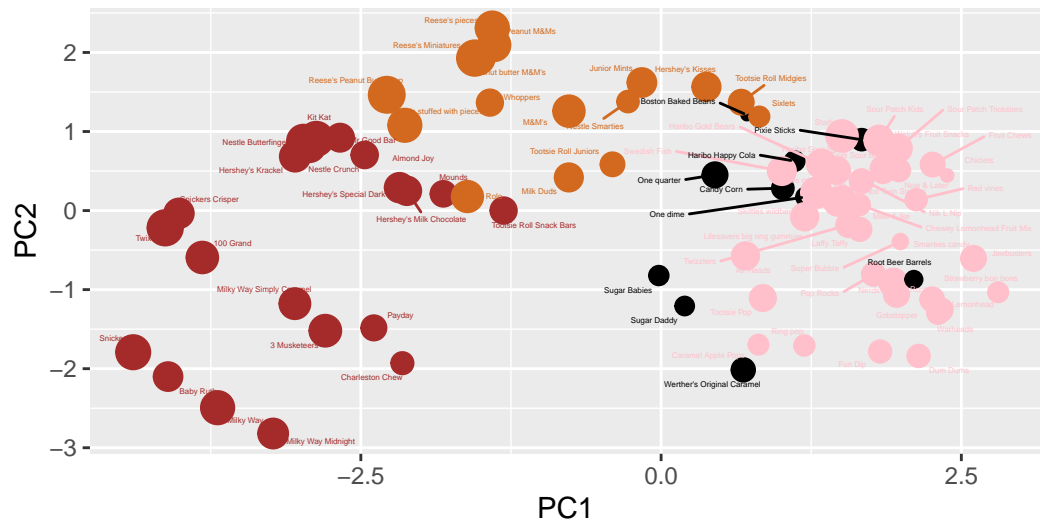
```r
library(ggrepel)

p + geom_text_repel(size=1, col=my_cols, max.overlaps = 50)  + #remember to change the siz
    theme(legend.position = "none") +
    labs(title="Halloween Candy PCA Space",
        subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brown
        caption="Data from 538")
```
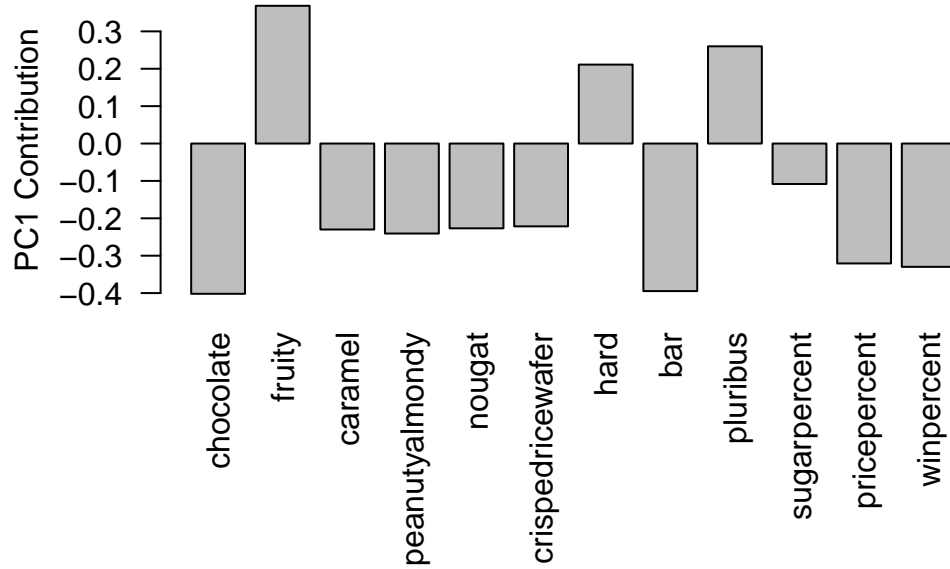
## Halloween Candy PCA Space

Colored by type: chocolate bar (dark brown), chocolate other (light brown),



Data from 538

```
#library(plotly) this would load up the ability to mouse over the data to see what it is,
#ggplotly(p)
```

```
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, Hard, and Pluribus are picked up in the positive direction because they're highly associated, which changes the way that it's plotted in PCA. This makes sense because many fruity candies come in an assorted package.