# Class10 - CANDY

Madeleine Duquette

10/28/22

## Background

In this mini-project we will examine 538 Halloween Candy data. What is your favorite candy? What is nougat anyway?

First step is to read the data...

```
candy_file <- "candy-data.csv"

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

| | 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?

```
nrow(candy)
```

[1] 85

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

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

[1] 38

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

```
#rownames(candy)
#Sour Patch kids are the best...
candy["Sour Patch Kids", ]$winpercent
```

[1] 59.864

## Q4. What is the winpercent value for "Kit Kat"?

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

[1] 76.7686

# Q5. What is the winpercent value for "Tootsie Roll Snack Bars"?

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

[1] 49.6535

```
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 | |

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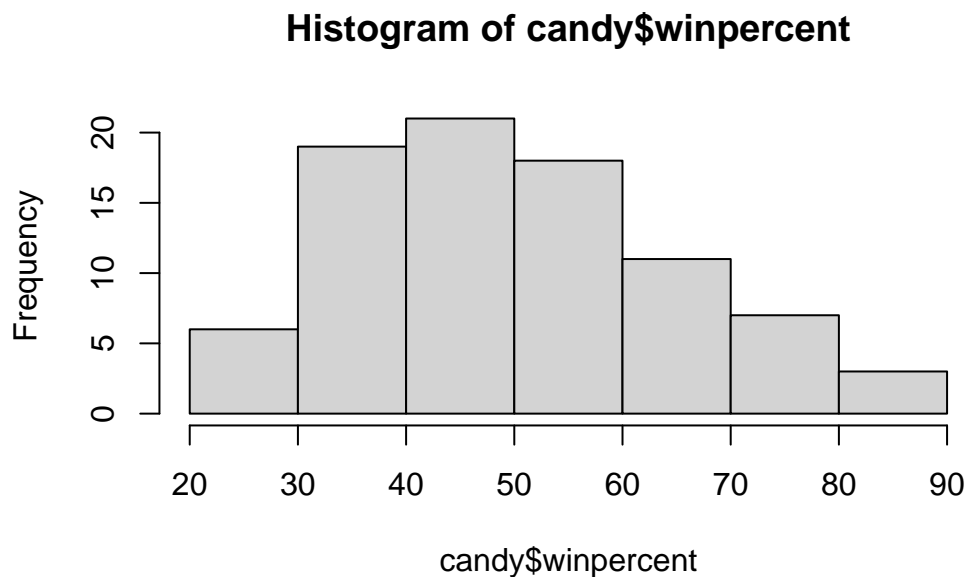
winpercent is on a different scale

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

Whether or not the candy is a chocolate
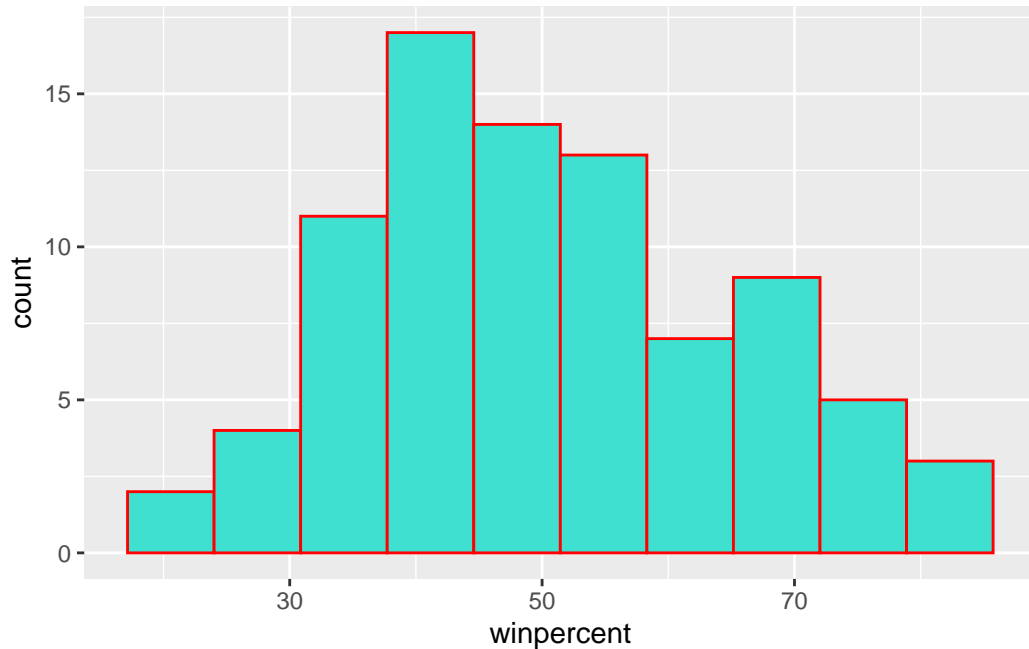
### Q8. Plot a histogram of winpercent values

```
hist(candy$winpercent)
```

**Histogram of candy$winpercent**



```
library(ggplot2)
ggplot(candy, aes(x=winpercent)) + geom_histogram(bins=10, col="red", fill="turquoise")
```

4

## Q9. Is the distribution of winpercent values symmetrical?

No it does not appear to be symmetrical

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

Below

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

```
chocolate <- candy$chocolate == 1
choco_score <- mean(candy[chocolate,]$winpercent)

fruity <- candy$fruity == 1
fruity_score <- mean(candy[fruity,]$winpercent)
```

```
  choco_score > fruity_score
```

[1] TRUE

## Q12. Is this difference statistically significant?

```
  t.test(candy[chocolate,]$winpercent, candy[fruity,]$winpercent)
```

```
    Welch Two Sample t-test

data:  candy[chocolate, ]$winpercent and candy[fruity, ]$winpercent
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
```

Yes!

## The base R `sort()` and `order()` functions are very useful!

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

```
  library(dplyr)
```

```
Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

    filter, lag
```

```
The following objects are masked from 'package:base':

    intersect, setdiff, setequal, union
```

```r
candy %>%
  arrange(winpercent) %>%
  head(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
```

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

```r
library(dplyr)
candy %>%
  arrange(desc(winpercent)) %>%
  head(5)
```

```
                        chocolate fruity caramel peanutyalmondy nougat
ReeseÕs Peanut Butter cup        1      0       0              1      0
ReeseÕs Miniatures               1      0       0              1      0
Twix                             1      0       1              0      0
```

```
Kit Kat                             1      0      0            0      0
Snickers                            1      0      1            1      1
                       crispedricewafer hard bar pluribus sugarpercent
ReeseÕs Peanut Butter cup              0    0   0        0         0.720
ReeseÕs Miniatures                     0    0   0        0         0.034
Twix                                   1    0   1        0         0.546
Kit Kat                                1    0   1        0         0.313
Snickers                               0    0   1        0         0.546
                       pricepercent winpercent
ReeseÕs Peanut Butter cup     0.651    84.18029
ReeseÕs Miniatures            0.279    81.86626
Twix                          0.906    81.64291
Kit Kat                       0.511    76.76860
Snickers                      0.651    76.67378
```

## Q15. Make a first barplot of candy ranking based on winpercent values.
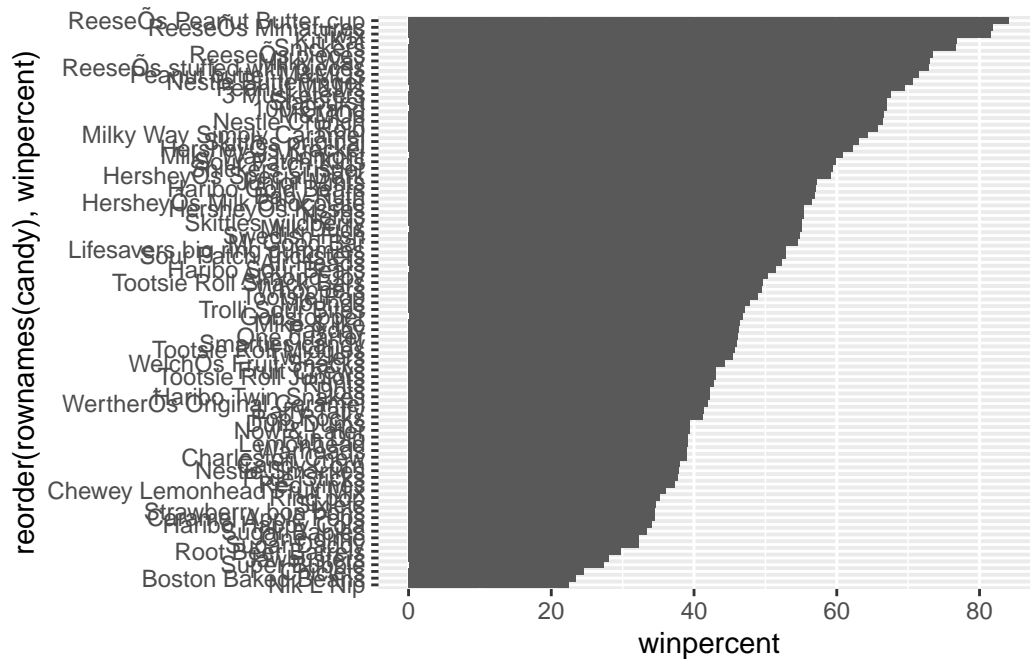
```r
library(ggplot2)

ggplot(candy) +
  aes(winpercent, rownames(candy)) +
  geom_col()
```

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

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

```
ggsave("mybarplot.png")
```
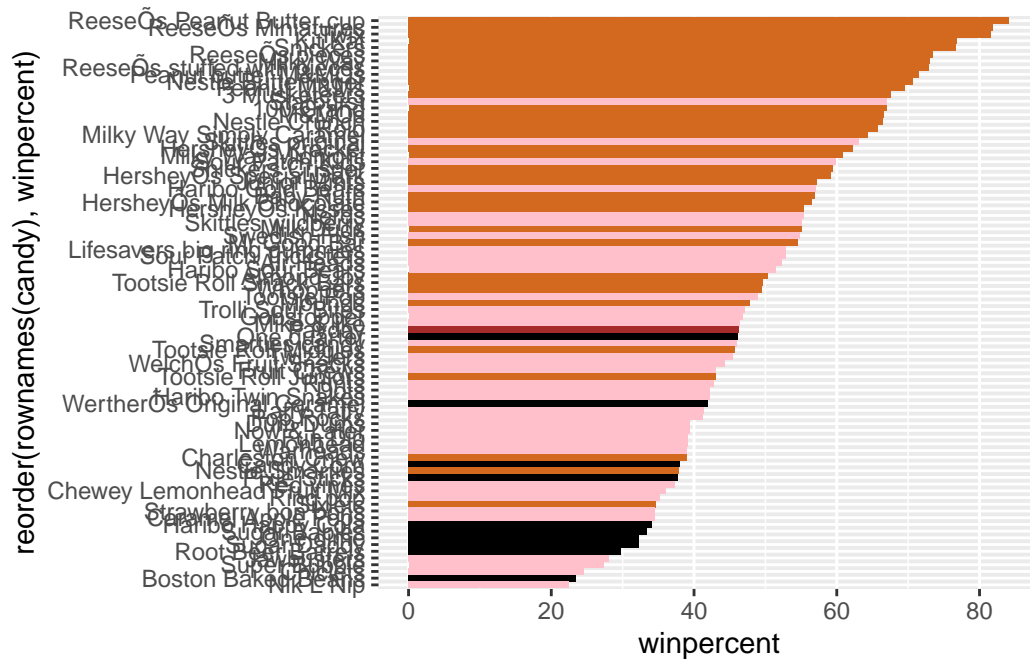
Saving 5.5 x 3.5 in image

```
my_cols=rep("black", nrow(candy))
my_cols[as.logical(candy$bar)] = "brown"
my_cols[as.logical(candy$chocolate)] = "chocolate"
# my_cols[as.logical(candy$bar & candy$chocolate)] = "blue"
my_cols[as.logical(candy$fruity)] = "pink"
```

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

## Q17. What is the worst ranked chocolate candy?

Sixlets

## Q18. What is the best ranked fruity candy?

Starburst

```
library(ggrepel)

# How about a plot of price vs win
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(col=my_cols, size=2, max.overlaps = 20)
```

**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?**

Reese's Mini

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

```
candy %>%
  arrange(desc(pricepercent)) %>%
  head(5)
```

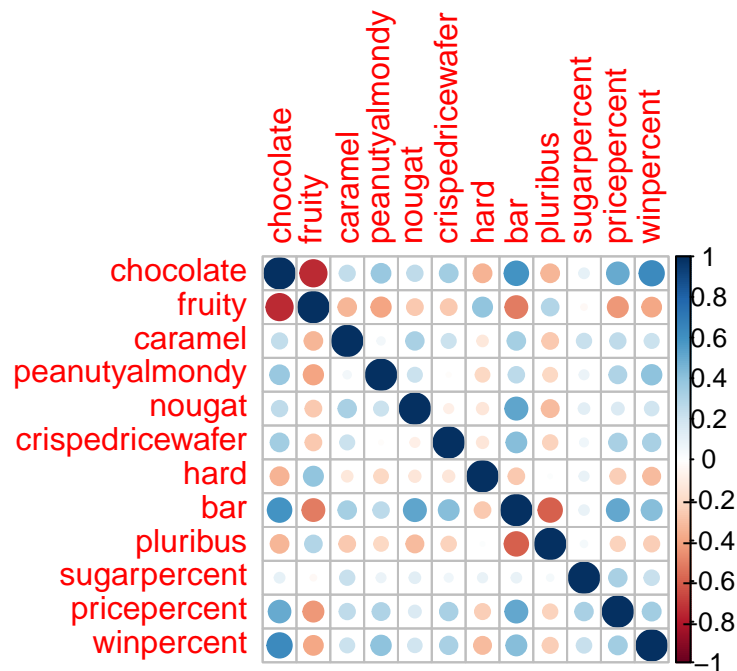|                  | chocolate | fruity | caramel | peanutyalmondy | nougat |
|------------------|-----------|--------|---------|----------------|--------|
| Nik L Nip        | 0         | 1      | 0       | 0              | 0      |
| Nestle Smarties  | 1         | 0      | 0       | 0              | 0      |
| Ring pop         | 0         | 1      | 0       | 0              | 0      |
| HersheyÕs Krackel| 1         | 0      | 0       | 0              | 0      |

```
HersheyÕs Milk Chocolate               1      0         0                0       0
                           crispedricewafer hard bar pluribus sugarpercent
Nik L Nip                                 0    0   0        1       0.197
Nestle Smarties                           0    0   0        1       0.267
Ring pop                                  0    1   0        0       0.732
HersheyÕs Krackel                         1    0   1        0       0.430
HersheyÕs Milk Chocolate                  0    0   1        0       0.430
                           pricepercent winpercent
Nik L Nip                         0.976   22.44534
Nestle Smarties                   0.976   37.88719
Ring pop                          0.965   35.29076
HersheyÕs Krackel                 0.918   62.28448
HersheyÕs Milk Chocolate          0.918   56.49050
```

Nik L Nip is the least popular

```
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)?

Chocolate and fruity

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

Chocolate and winpercent

## PCA

The main function that is always there for us is `prcomp()`. It has an important argument that is set to `scale=FALSE` but needs to be set to `TRUE` in this case because our data is not on the same scale.

```
pca <- prcomp(candy, scale=T)
summary(pca)
```
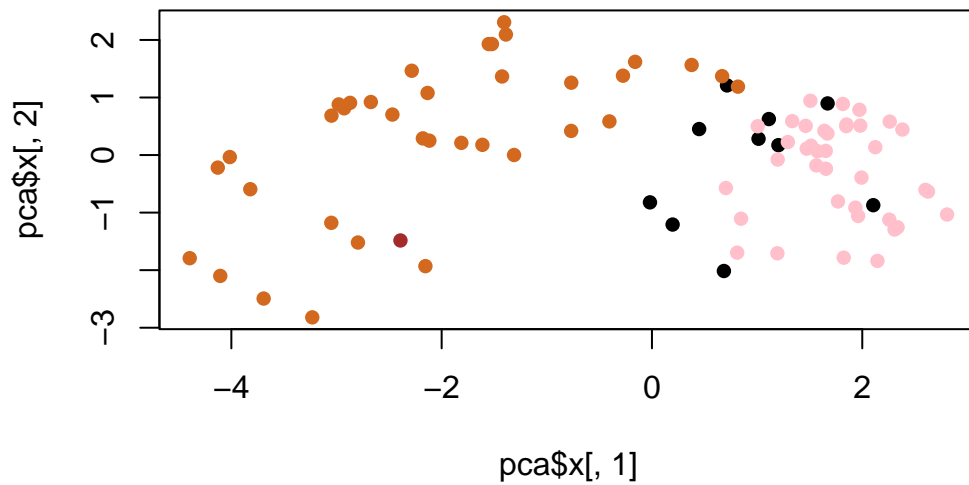
```
Importance of components:
                          PC1    PC2    PC3     PC4    PC5     PC6     PC7
Standard deviation     2.0788 1.1378 1.1092 1.07533 0.9518 0.81923 0.81530
Proportion of Variance 0.3601 0.1079 0.1025 0.09636 0.0755 0.05593 0.05539
Cumulative Proportion  0.3601 0.4680 0.5705 0.66688 0.7424 0.79830 0.85369
                          PC8     PC9    PC10    PC11    PC12
Standard deviation     0.74530 0.67824 0.62349 0.43974 0.39760
Proportion of Variance 0.04629 0.03833 0.03239 0.01611 0.01317
Cumulative Proportion  0.89998 0.93832 0.97071 0.98683 1.00000
```

My PCA plot (aka) PC1 vs PC2 score plot.

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

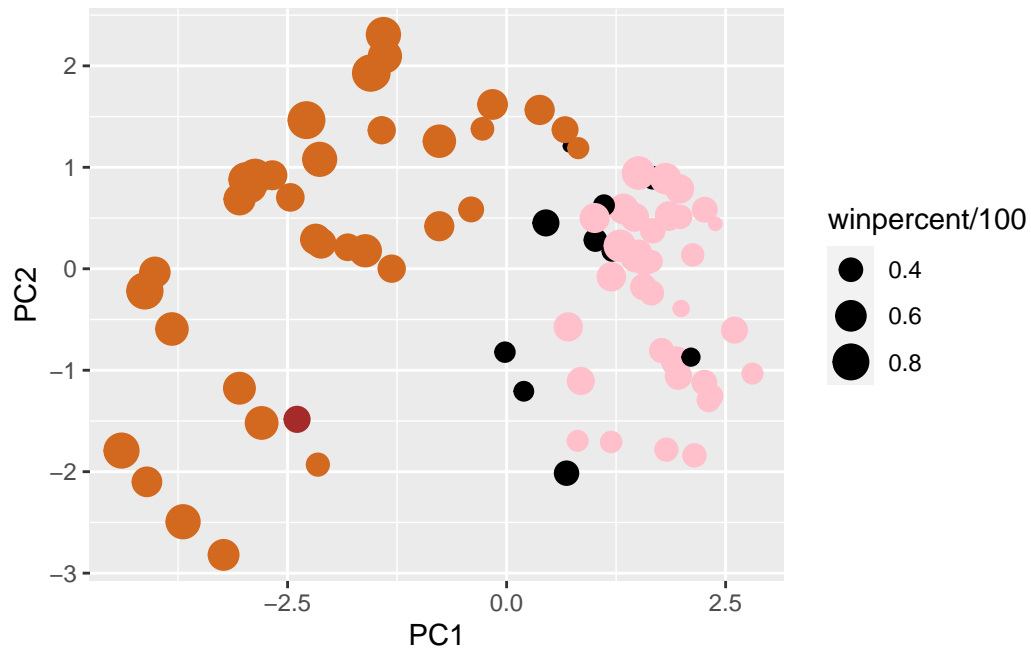I will make a "nicer" plot with ggplot ggplot only works with data.frames as input so I need to make one...

```
# Make a new data-frame with our PCA results and candy data
my_data <- cbind(candy, pca$x[,1:3])
```

```
library(ggrepel)
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
```
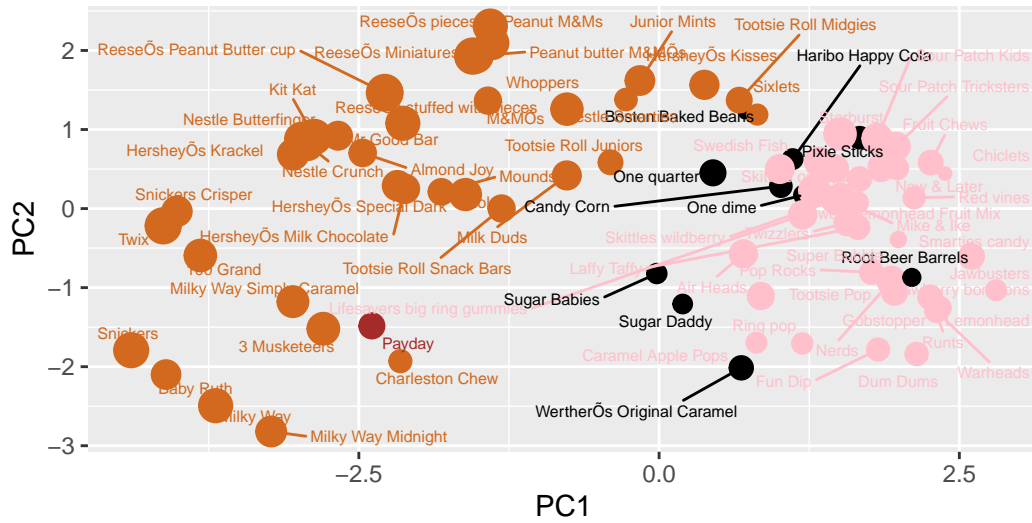
```
library(ggrepel)

p + geom_text_repel(size=2, col=my_cols, max.overlaps = 25)  +
   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")
```

Warning: ggrepel: 6 unlabeled data points (too many overlaps). Consider
increasing max.overlaps
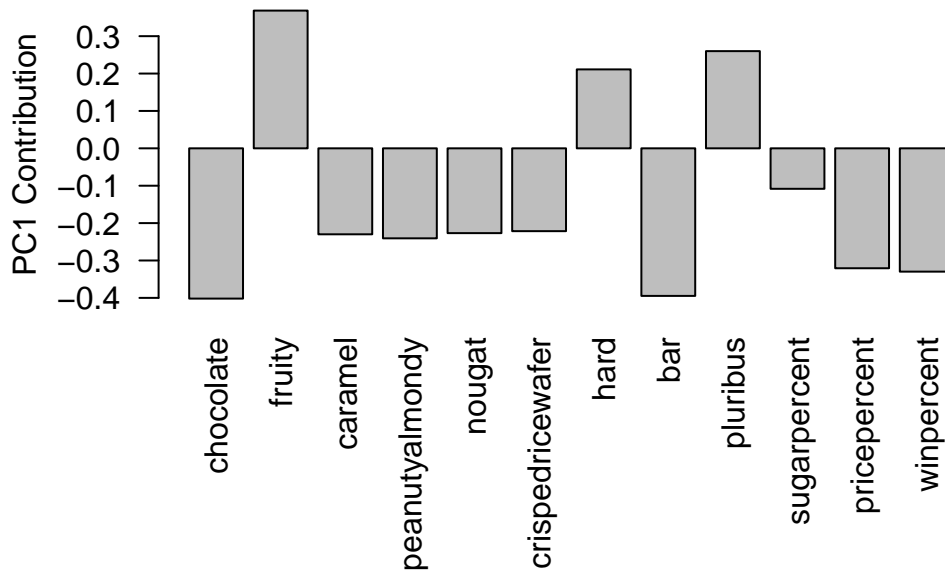
## Halloween Candy PCA Space
### Colored by type: chocolate bar (dark brown), chocolate other (light brown),



ReeseÕs Peanut Butter cup · ReeseÕs Miniatures — ReeseÕs pieces — Peanut M&Ms — Junior Mints — Tootsie Roll Midgies — Haribo Happy Cola — Sour Patch Kids — Peanut butter M&MÕs — HersheyÕs Kisses — Sour Patch Tricksters — Sixlets — Whoppers — Kit Kat — Reese stuffed with pieces — Starburst — Fruit Chews — Nestle Butterfinger — M&MÕs — Nestle Boston Baked Beans — Chiclets — HersheyÕs Krackel — Mr Good Bar — Tootsie Roll Juniors — Swedish Fish — Pixie Sticks — Snickers Crisper — Nestle Crunch — Almond Joy — Mounds — New & Later — Red vines — HersheyÕs Special Dark — Candy Corn — One dime — Mike & Ike — Twix — HersheyÕs Milk Chocolate — Milk Duds — One quarter — Skittles wildberry — Twizzlers — Smarties candy — 100 Grand — Tootsie Roll Snack Bars — Laffy Taffy — Super Bubble — Root Beer Barrels — Milky Way Simply Caramel — Sugar Babies — Air Heads — Pop Rocks — Tootsie Pop — Jawbusters — fruity bon bons — Lifesavers big ring gummies — Sugar Daddy — Ring pop — Gobstopper — Lemonhead — Snickers — 3 Musketeers — Payday — Caramel Apple Pops — Nerds — Runts — Warheads — Baby Ruth — Charleston Chew — Fun Dip — Dum Dums — Milky Way — WertherÕs Original Caramel — Milky Way Midnight

Data from 538

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



17

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

Fruity, hard, and pluribus. It makes sense these are all grouping together in one direction because that's the typical form for fruity candies... Fruit candies thrive on variety! And these characteristics are not common for chocolates.