

Halloween Mini Project

Lily Huynh (A16929651)

2025-02-04

Table of contents

| | |
|---|----|
| 1.Importing candy data | 1 |
| 2. What is your favorite candy? | 2 |
| 3. Overall Candy Rankings | 8 |
| Time to add some useful color | 13 |
| 4. Taking a look at pricepercent | 15 |
| 5 Exploring the correlation structure | 19 |
| 6. Principal Component Analysis | 20 |

Today we will examine data from 538 on common Halloween candy. In particular we will use ggplot, dplyr, and PCA to make sense of this multivariate dataset

1.Importing candy data

```
candy <- read.csv("candy-data.txt", 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?

Ans. There are 85 different candy types in this dataset.

```
nrow(candy)
```

```
[1] 85
```

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

Ans. There are 38 fruity candy types in this dataset.

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

```
[1] 38
```

How many chocolate candy are there in this dataset?

Ans. There are 37 chocolate candy in this dataset.

```
sum(candy$chocolate)
```

```
[1] 37
```

2. What is your favorite candy?

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

Ans. My favorite candy in the dataset is Peanut M&Ms. The winpercent value is 69.48379.

```
candy["Peanut M&Ms", "winpercent"]
```

```
[1] 69.48379
```

```
candy["Peanut M&Ms",]$winpercent
```

```
[1] 69.48379
```

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

Ans. The winpercent value for Kit Kat is 76.7686.

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

```
[1] 76.7686
```

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

Ans. The winpercent value for Tootsie Roll Snack Bars is 49.6535.

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

```
[1] 49.6535
```

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

Ans. The winpercent variable is on a different scale compared to the majority of the other columns in the dataset. The other variables range from 0 to 1, but the winpercent variable range is much higher than 1.

N.B It looks like the `winpercent` row in the `skim_variable` column is on a different scale than the others (0-100% rather than 0-1). I will need to scale this dataset before analysis like PCA.

```
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 | |

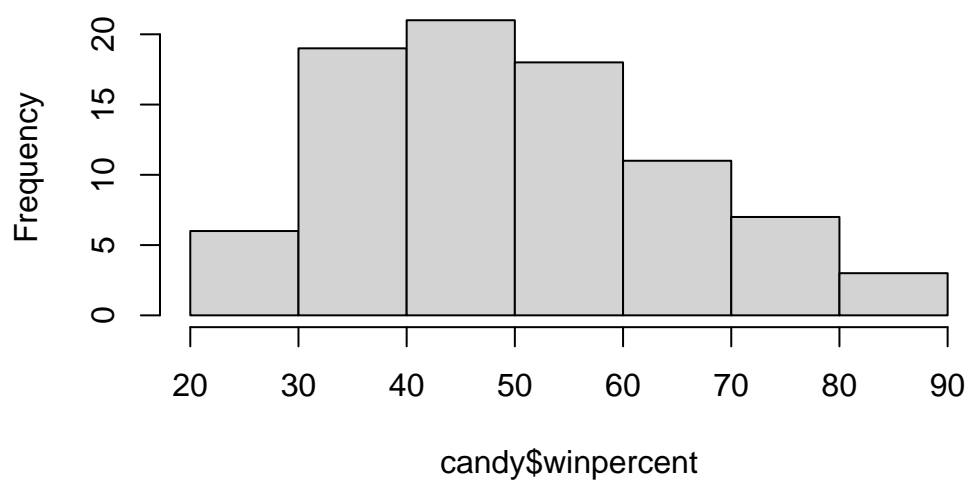
Q7. What do you think a zero and one represent for the candy\$chocolate column?

Ans. The zero for the candy\$chocolate column represents candies that aren't chocolate. The one represents the candies that contain chocolate.

Q8. Plot a histogram of winpercent values.

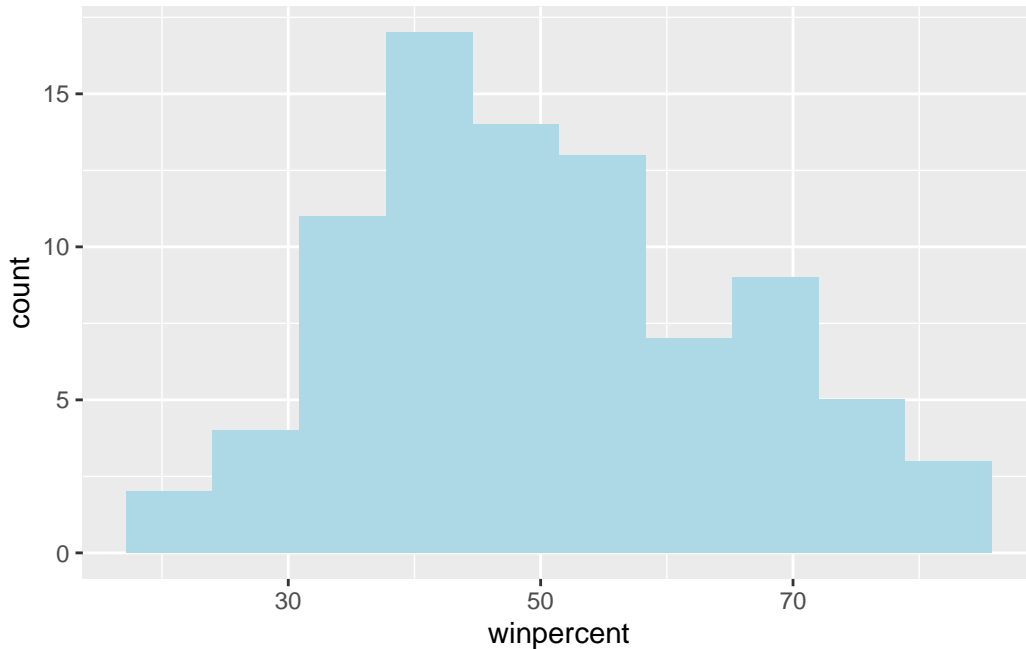
```
hist(candy$winpercent)
```

Histogram of candy\$winpercent



```
library(ggplot2)

ggplot(candy) +
  aes(x=winpercent) +
  geom_histogram(bins=10, fill="lightblue")
```



Q9. Is the distribution of winpercent values symmetrical?

Ans. No, the distribution of winpercent value is skewed to the right.

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

Ans. The center of the distribution is below 50% based on the histogram, with the center being around 45%. We can also see this by using the `summary()` function to figure out the median, 47.83, which is lower than 50.

```
summary(candy$winpercent)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|-------|---------|--------|-------|---------|-------|
| 22.45 | 39.14 | 47.83 | 50.32 | 59.86 | 84.18 |

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

Ans. The chocolate candy is ranked higher than the fruit candy.

- Step 1: Find all “chocolate” candy
- Step 2: Find their “winpercent” values
- Step 3: Summarize these values
- Step 4: Find all “fruit” candy

- Step 5: Find their “winpercent” values
- Step 6: Summarize these values
- Step 7: Compare the two summary values

1. Find all chocolate candy

```
choc.inds <- candy$chocolate == 1
```

2. Find their winpercent values

```
choc.win <- candy[choc.inds,]$winpercent
```

3. Summarize these values

```
choc.mean <- mean(choc.win)
choc.mean
```

```
[1] 60.92153
```

4. Find all fruit candy

```
fruity.inds <- candy$fruity == 1
```

5. Find their winpercent values

```
fruity.win <- candy[fruity.inds,]$winpercent
```

6. Summarize these values

```
fruity.mean <- mean(fruity.win)
fruity.mean
```

```
[1] 44.11974
```

7. Compare the two summary values The chocolate winpercent is higher than the fruity candy.

```
choc.mean
```

```
[1] 60.92153
```

```
fruity.mean
```

```
[1] 44.11974
```

Q12. Is this difference statistically significant?

Ans. This difference is statistically significant, due to the p-value being very low, 2.871e-08.

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

Welch Two Sample t-test

```
data:  choc.win and fruity.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
```

3. Overall Candy Rankings

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

Ans. The 5 least liked candy types in this dataset is Nik L Nip, Boston Baked Beans, Chiclets, Super Bubble, and Jawbusters.

```
# Not that useful - it just sorts the values
sort(candy$winpercent)
```

```
[1] 22.44534 23.41782 24.52499 27.30386 28.12744 29.70369 32.23100 32.26109
[9] 33.43755 34.15896 34.51768 34.57899 34.72200 35.29076 36.01763 37.34852
[17] 37.72234 37.88719 38.01096 38.97504 39.01190 39.14106 39.18550 39.44680
[25] 39.46056 41.26551 41.38956 41.90431 42.17877 42.27208 42.84914 43.06890
[33] 43.08892 44.37552 45.46628 45.73675 45.99583 46.11650 46.29660 46.41172
[41] 46.78335 47.17323 47.82975 48.98265 49.52411 49.65350 50.34755 51.41243
[49] 52.34146 52.82595 52.91139 54.52645 54.86111 55.06407 55.10370 55.35405
```



```
[57] 55.37545 56.49050 56.91455 57.11974 57.21925 59.23612 59.52925 59.86400
[65] 60.80070 62.28448 63.08514 64.35334 65.71629 66.47068 66.57458 66.97173
[73] 67.03763 67.60294 69.48379 70.73564 71.46505 72.88790 73.09956 73.43499
[81] 76.67378 76.76860 81.64291 81.86626 84.18029
```

```
x <- c(10, 1, 100)
order(x)
```

```
[1] 2 1 3
```

```
x[order(x)]
```

```
[1] 1 10 100
```

The `order()` function tells us how it arrange the elements of the input to make them sorted - i.e. how to order them

We can determine the order of winpercent to make them sorted and use that order to arrange the whole dataset.

```
order.inds <- order(candy$winpercent)
head(candy[order.inds, ])
```

| | chocolate | fruity | caramel | peanut | almond | 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 |
| Root Beer Barrels | 0 | 0 | 0 | | 0 | 0 |

| | crisped | rice | wafer | hard | bar | pluribus | sugar | percent | price | percent |
|--------------------|---------|------|-------|------|-----|----------|-------|---------|-------|---------|
| 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 |
| Root Beer Barrels | | 0 | 1 | 0 | | 1 | | 0.732 | | 0.069 |

| | winpercent |
|--------------------|------------|
| Nik L Nip | 22.44534 |
| Boston Baked Beans | 23.41782 |
| Chiclets | 24.52499 |

| | |
|-------------------|----------|
| Super Bubble | 27.30386 |
| Jawbusters | 28.12744 |
| Root Beer Barrels | 29.70369 |

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

Ans. The top 5 all time favorite candy types are Reese's Peanut Butter Cup, Snickers, Kit Kat, Twix, and Reese's Miniatures.

```
tail(candy[order.inds, ])
```

| | chocolate | fruity | caramel | peanut | almond | nougat |
|---------------------------|-----------|--------|---------|--------|--------|--------|
| Reese's pieces | 1 | 0 | 0 | | 1 | 0 |
| 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 |

| | crisped | rice | wafers | hard | bar | pluribus | sugar | percent |
|---------------------------|---------|------|--------|------|-----|----------|-------|---------|
| Reese's pieces | | | 0 | 0 | 0 | 1 | | 0.406 |
| 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 |

| | price | percent | win | percent |
|---------------------------|-------|---------|--------|---------|
| Reese's pieces | 0.651 | | 73.434 | 99 |
| Snickers | 0.651 | | 76.673 | 78 |
| Kit Kat | 0.511 | | 76.768 | 60 |
| Twix | 0.906 | | 81.642 | 91 |
| Reese's Miniatures | 0.279 | | 81.866 | 26 |
| Reese's Peanut Butter cup | 0.651 | | 84.180 | 29 |

```
order.winpercent.decrease <- order(candy$winpercent, decreasing=TRUE)
head(candy[order.winpercent.decrease, ])
```

| | chocolate | fruity | caramel | peanut | almond | nougat |
|---------------------------|-----------|--------|---------|--------|--------|--------|
| Reese's Peanut Butter cup | 1 | 0 | 0 | | 1 | 0 |
| Reese's 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 |
| Reese's pieces | 1 | 0 | 0 | 1 | 0 |
| | crisped | ricewafer | hard bar | pluribus | sugarpercent |
| Reese's Peanut Butter cup | | 0 | 0 | 0 | 0.720 |
| Reese's Miniatures | | 0 | 0 | 0 | 0.034 |
| Twix | | 1 | 0 | 1 | 0.546 |
| Kit Kat | | 1 | 0 | 1 | 0.313 |
| Snickers | | 0 | 0 | 1 | 0.546 |
| Reese's pieces | | 0 | 0 | 0 | 1 |
| | pricepercent | winpercent | | | |
| Reese's Peanut Butter cup | 0.651 | 84.18029 | | | |
| Reese's Miniatures | 0.279 | 81.86626 | | | |
| Twix | 0.906 | 81.64291 | | | |
| Kit Kat | 0.511 | 76.76860 | | | |
| Snickers | 0.651 | 76.67378 | | | |
| Reese's pieces | 0.651 | 73.43499 | | | |

```
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

```
candy %>%
  arrange(winpercent) %>%
  tail(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 |

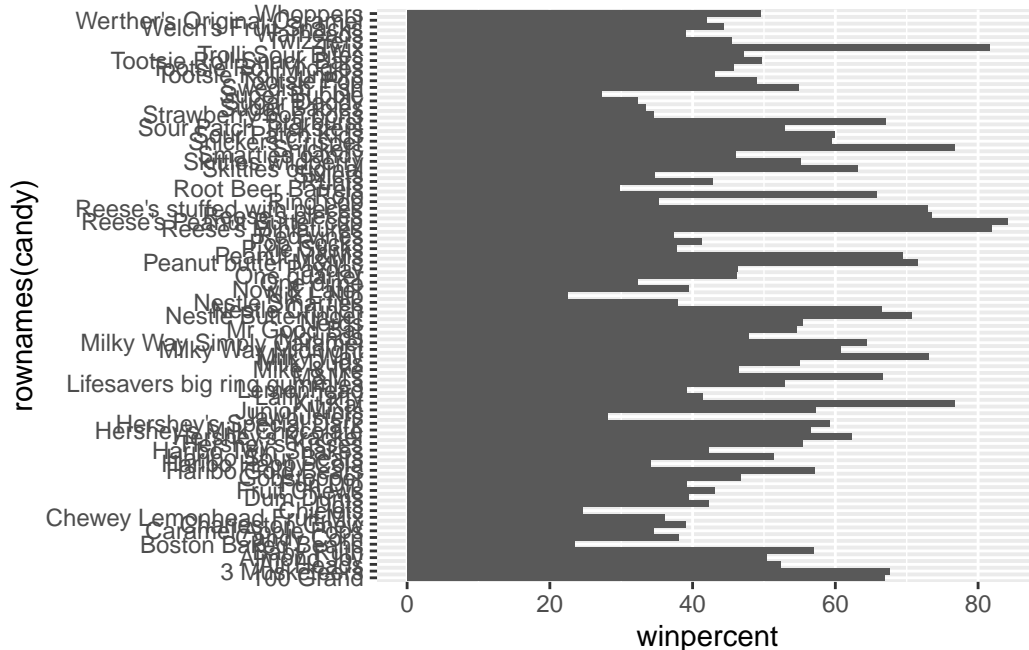
| | crisped | rice | wafer | hard | bar | pluribus | sugar | percent |
|---------------------------|---------|---------|-------|---------|-----|----------|-------|---------|
| Snickers | | | 0 | 0 | 1 | | 0 | 0.546 |
| Kit Kat | | | 1 | 0 | 1 | | 0 | 0.313 |
| Twix | | | 1 | 0 | 1 | | 0 | 0.546 |
| Reese's Miniatures | | | 0 | 0 | 0 | | 0 | 0.034 |
| Reese's Peanut Butter cup | | | 0 | 0 | 0 | | 0 | 0.720 |
| | price | percent | win | percent | | | | |
| Snickers | 0.651 | | 76.67 | 378 | | | | |
| Kit Kat | 0.511 | | 76.76 | 860 | | | | |
| Twix | 0.906 | | 81.64 | 291 | | | | |
| Reese's Miniatures | 0.279 | | 81.86 | 626 | | | | |
| Reese's Peanut Butter cup | 0.651 | | 84.18 | 029 | | | | |

Q Which approach do you prefer and why?

I prefer the dyplr because it requires less parenthesis and commas, which makes it harder for me to make a mistake with a comma.

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

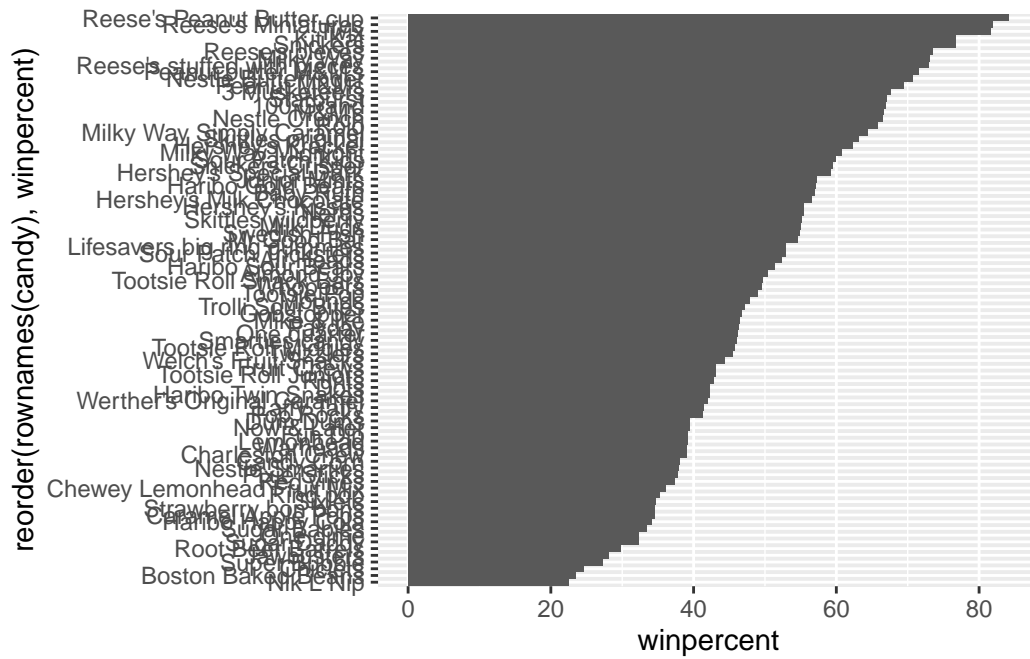
```
ggplot(candy) +  
  aes(x=winpercent, rownames(candy)) +  
  geom_col()
```



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

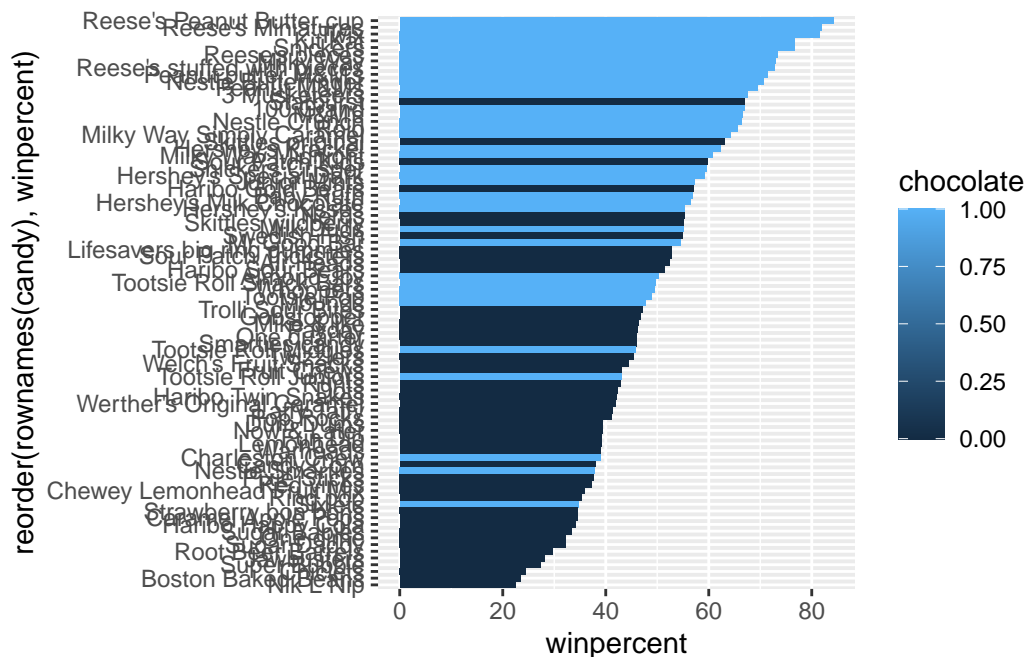
Let's reorder it:

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



Time to add some useful color

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



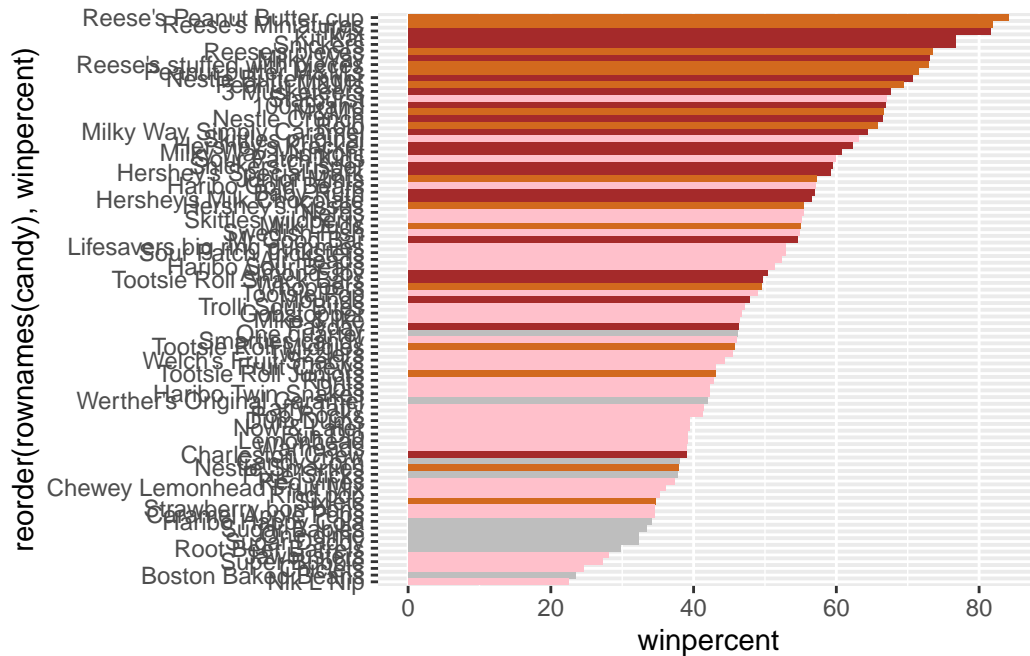
We need to make our own separate color vector where we can spell out what candy is colored a particular color.

```
mycols <- rep("gray", nrow(candy))
mycols[candy$chocolate == 1] <- "chocolate"
mycols[candy$bar == 1] <- "brown"
mycols[candy$fruity == 1] <- "pink"
mycols
```

```
[1] "brown"    "brown"    "gray"      "gray"      "pink"      "brown"
[7] "brown"    "gray"      "gray"      "pink"      "brown"      "pink"
[13] "pink"     "pink"     "pink"     "pink"     "pink"     "pink"
[19] "pink"     "gray"     "pink"     "pink"     "chocolate" "brown"
[25] "brown"    "brown"    "pink"     "chocolate" "brown"     "pink"
[31] "pink"     "pink"     "chocolate" "chocolate" "pink"     "chocolate"
[37] "brown"    "brown"    "brown"    "brown"    "brown"     "pink"
[43] "brown"    "brown"    "pink"     "pink"     "brown"     "chocolate"
[49] "gray"     "pink"     "pink"     "chocolate" "chocolate" "chocolate"
[55] "chocolate" "pink"     "chocolate" "gray"     "pink"     "chocolate"
[61] "pink"     "pink"     "chocolate" "pink"     "brown"     "brown"
[67] "pink"     "pink"     "pink"     "pink"     "gray"     "gray"
[73] "pink"     "pink"     "pink"     "chocolate" "chocolate" "brown"
[79] "pink"     "brown"    "pink"     "pink"     "pink"     "gray"
```

```
[85] "chocolate"
```

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



```
as.logical(c(1,0,1))
```

```
[1] TRUE FALSE TRUE
```

Q17. What is the worst ranked chocolate candy?

Ans. The worst ranked chocolate candy is Sixlets.

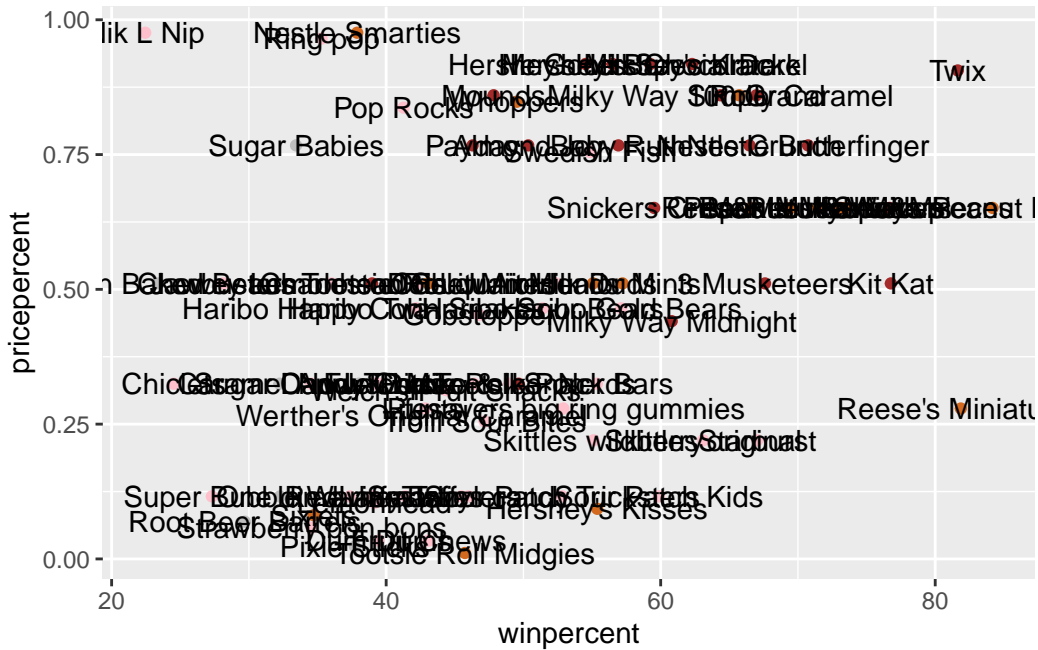
Q18. What is the best ranked fruity candy?

Ans. The best ranked fruity candy is Starburst.

4. Taking a look at pricepercent

Make a plot of winpercent (x-axis) vs pricepercent (y-axis)

```
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=mycols) +
  geom_text()
```

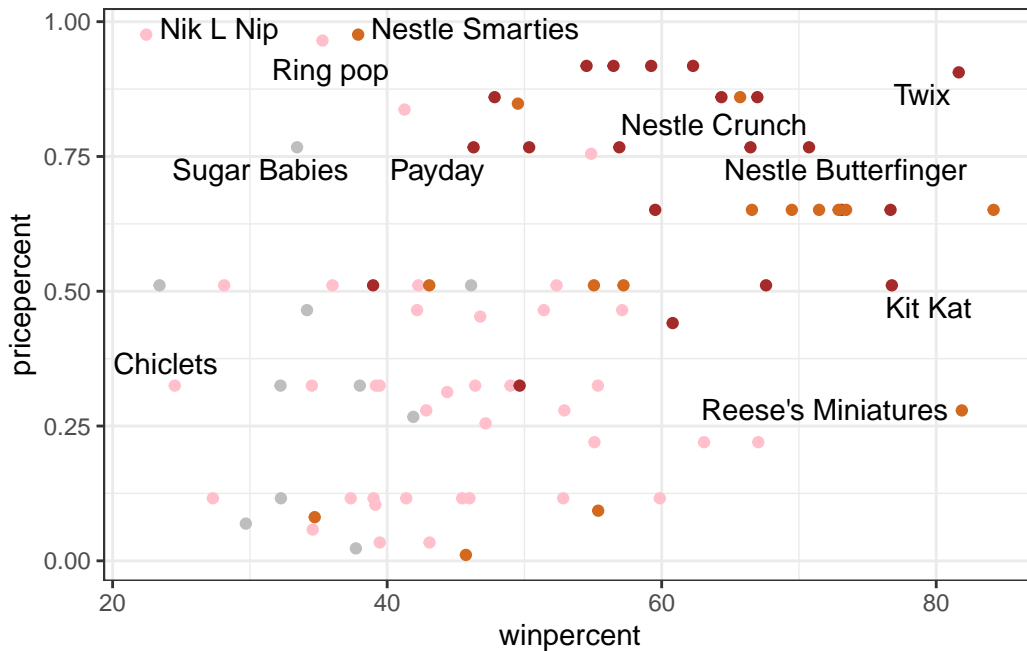


To avoid the overplotting of the text labels, we can use the add on package **ggrepel**

```
library(ggrepel)

ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=mycols) +
  geom_text_repel(max.overlaps = 5) +
  theme_bw()
```

Warning: ggrepel: 74 unlabeled data points (too many overlaps). Consider increasing max.overlaps



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?

Ans. Reeses Miniatures has one of the highest winpercent and one of the lowest pricepercent.

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

| | pricepercent | winpercent |
|--------------------------|--------------|------------|
| Nik L Nip | 0.976 | 22.44534 |
| Nestle Smarties | 0.976 | 37.88719 |
| Ring pop | 0.965 | 35.29076 |
| Hershey's Krackel | 0.918 | 62.28448 |
| Hershey's Milk Chocolate | 0.918 | 56.49050 |

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

| | pricepercent | winpercent |
|---------------------------|--------------|------------|
| Reese's Peanut Butter cup | 0.651 | 84.18029 |
| Reese's Miniatures | 0.279 | 81.86626 |

| | | |
|----------|-------|----------|
| Twix | 0.906 | 81.64291 |
| Kit Kat | 0.511 | 76.76860 |
| Snickers | 0.651 | 76.67378 |

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

Ans. The top 5 most expensive candy types in the dataset are Hershey's Special Dark, Mr Good Bar, Ring pop, Nik L Nip, and Nestle Smarties. The least popular of these 5 is the Nik L Nip.

```
order.pricepercent <- order(candy$pricepercent)
tail(candy[order.pricepercent, ],n=5)
```

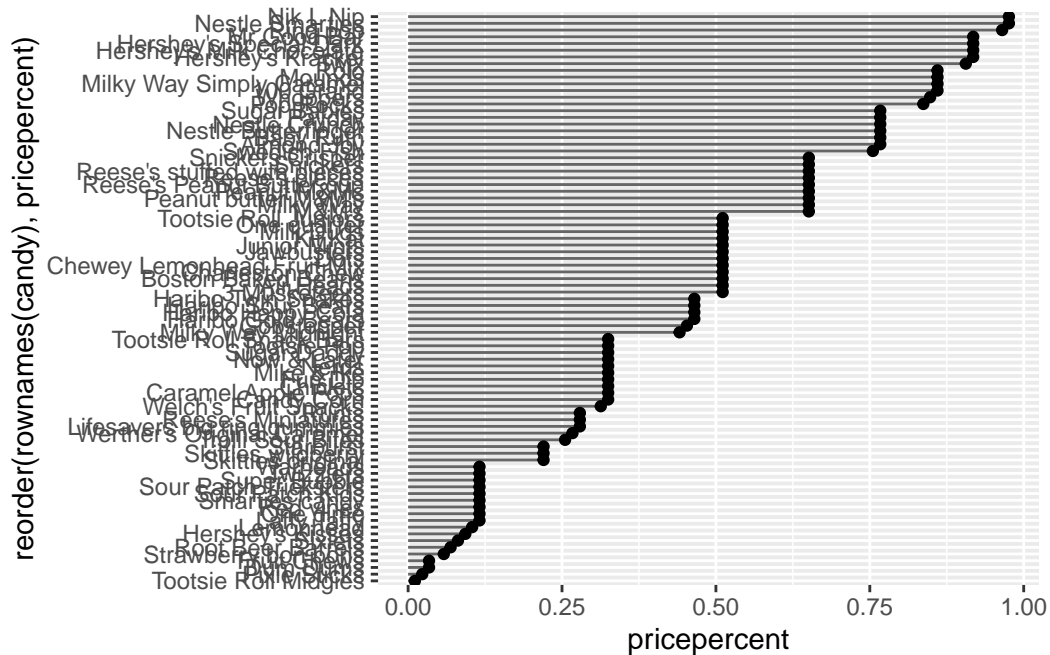
| | chocolate | fruity | caramel | peanut | almond | nougat |
|------------------------|-----------|--------|---------|--------|--------|--------|
| Hershey's Special Dark | 1 | 0 | 0 | | 0 | 0 |
| Mr Good Bar | 1 | 0 | 0 | | 1 | 0 |
| Ring pop | 0 | 1 | 0 | | 0 | 0 |
| Nik L Nip | 0 | 1 | 0 | | 0 | 0 |
| Nestle Smarties | 1 | 0 | 0 | | 0 | 0 |

| | crisped | rice | wafer | hard | bar | pluribus | sugar | percent |
|------------------------|---------|------|-------|------|-----|----------|-------|---------|
| Hershey's Special Dark | | 0 | 0 | 1 | | 0 | | 0.430 |
| Mr Good Bar | | 0 | 0 | 1 | | 0 | | 0.313 |
| Ring pop | | 0 | 1 | 0 | | 0 | | 0.732 |
| Nik L Nip | | 0 | 0 | 0 | | 1 | | 0.197 |
| Nestle Smarties | | 0 | 0 | 0 | | 1 | | 0.267 |

| | pricepercent | winpercent |
|------------------------|--------------|------------|
| Hershey's Special Dark | 0.918 | 59.23612 |
| Mr Good Bar | 0.918 | 54.52645 |
| Ring pop | 0.965 | 35.29076 |
| Nik L Nip | 0.976 | 22.44534 |
| Nestle Smarties | 0.976 | 37.88719 |

Q21. Make a barplot again with `geom_col()` this time using `pricepercent` and then improve this step by step, first ordering the x-axis by value and finally making a so called “dot chat” or “lollipop” chart by swapping `geom_col()` for `geom_point()` + `geom_segment()`.

```
ggplot(candy) +
  aes(pricepercent, reorder(rownames(candy), pricepercent)) +
  geom_segment(aes(yend = reorder(rownames(candy), pricepercent),
                  xend = 0), col="gray40") +
  geom_point()
```



5 Exploring the correlation structure

Now that we have explored the data set a little, we will see how the variables interact with one another.

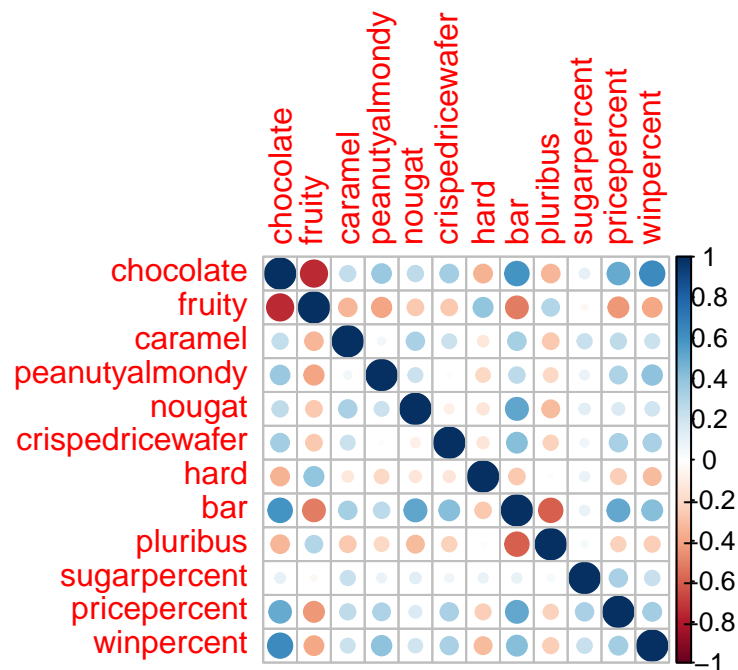
First we will use correlation and view the results with the **corrplot** package to plot a correlation matrix.

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

```
library(corrplot)
```

```
corrplot 0.95 loaded
```

```
corrplot(cij)
```



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

Ans. Two variables that are anti-correlated are chocolate and fruity candy.

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

Ans. Two variables that are most positively correlated are chocolate and winpercent.

6. Principal Component Analysis

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

```
pca <- prcomp(candy, scale=TRUE)
```

```
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 | | |

```
attributes(pca)
```

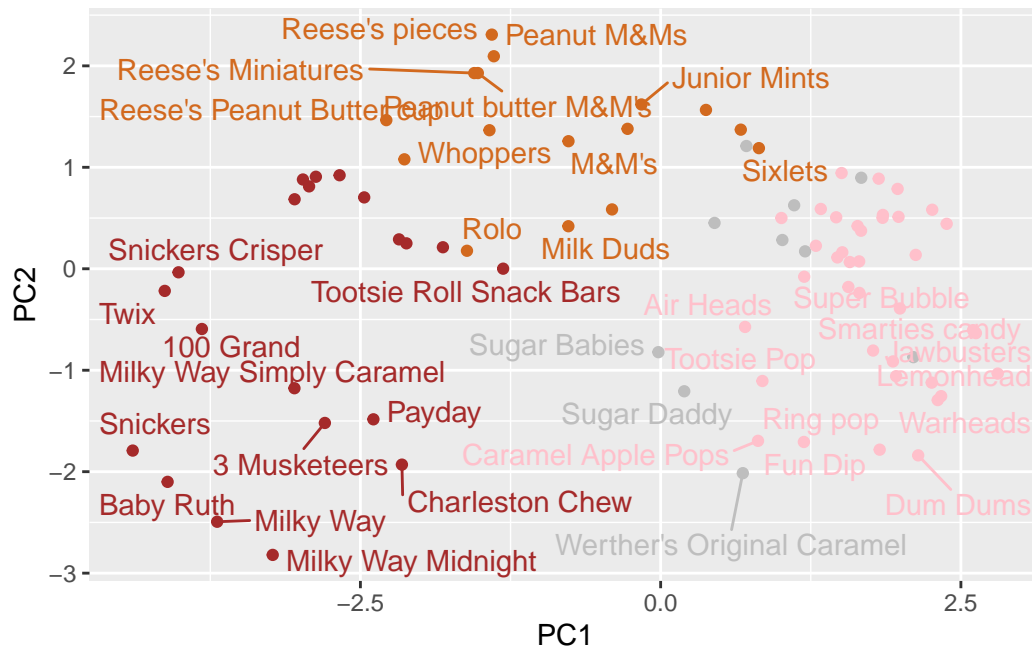
```
$names
[1] "sdev"      "rotation" "center"    "scale"     "x"

$class
[1] "prcomp"
```

Let's plot our main results as our PCA "score plot"

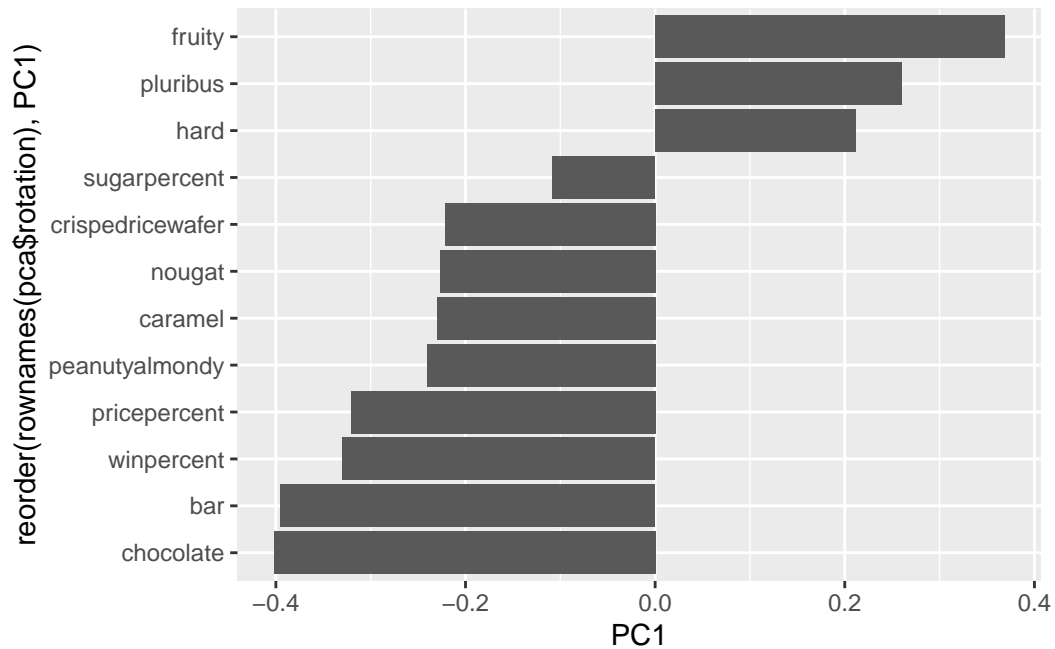
```
ggplot(pca$x) +
  aes(PC1, PC2, label=rownames(pca$x)) +
  geom_point(col=mycols) +
  geom_text_repel(col=mycols)
```

Warning: ggrepel: 48 unlabeled data points (too many overlaps). Consider increasing max.overlaps



Finally, let's look at how the original variables contribute to the PCs, start with PC1

```
ggplot(pca$rotation) +
  aes(PC1, reorder(rownames(pca$rotation), PC1)) +
  geom_col()
```



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

Ans. The original variables that are picked up strongly by PC1 in the positive direction are fruity, pluribus, and hard. This makes sense to me because the graphs shown above have all the fruity candy on the left side of the graph with PC1 on the x-axis.