# Class10_Halloween_Project

AUTHOR
Zixuan Zeng(A16142927)

## Importing 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
```

```
dim(candy)
```

```
[1] 85 12
```

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

```
[1] 38
```

> Q1. 85 types of candy

> Q2. 38 fruity candy types

```
candy["Reese's Peanut Butter cup", ]$winpercent
```

```
[1] 84.18029
```

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

```
[1] 76.7686
```

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

```
[1] 49.6535
```

# What is your favourite candy

> Q3. Reese's Peanut Butter cup, 84.18%

> Q4. 73.76%

> Q5. 49.65%

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

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 | ▇▁▁▁▁▂ |

| skim_variable | n_missing | complete_rate | mean | sd | p0 | p25 | p50 | p75 | p100 | hist |
|---|---|---|---|---|---|---|---|---|---|---|
| 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. Winpercent seems to be on a different scale. All others are on a scale from 0 to 1 while winpercent is on 0-100.

> Q7. Does this candy contain chocolate or not

> Q8. Plot shown below

```r
library(ggplot2)
ggplot(candy) + aes(winpercent) + geom_histogram()
```

`stat_bin()` using `bins = 30`. Pick better value `binwidth`.

```
win <- candy$winpercent
win.mean <- mean(win)
round(win.mean,2)
```

```
[1] 50.32
```

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

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

> Q9. The distribution seems a little skewed to the left. Most candies have a winpercent between 40 and 80.

> Q10. The center of the distribution seems to be below 50%. Median = 47.83

```
mean(candy$winpercent[as.logical(candy$chocolate)]) > mean(candy$winpercent[as.logical(ca
```

```
[1] TRUE
```

> Q11. On average, the chocolate candy ranked higher than the fruity candy.

```
t.test(candy$winpercent[as.logical(candy$chocolate)],candy$winpercent[as.logical(candy$fr
```

```
	Welch Two Sample t-test

data:  candy$winpercent[as.logical(candy$chocolate)] and
candy$winpercent[as.logical(candy$fruity)]
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
```

> Q12. p-value is less than 0.05, indicating the difference is statistically significant.

## Overall Candy Rankings

```
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   |

> Q13. The the five least liked candy types are Nik L Nip, Boston Baked Beans, Chiclets, Super Bubble, and Jawbusters.

```
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        |

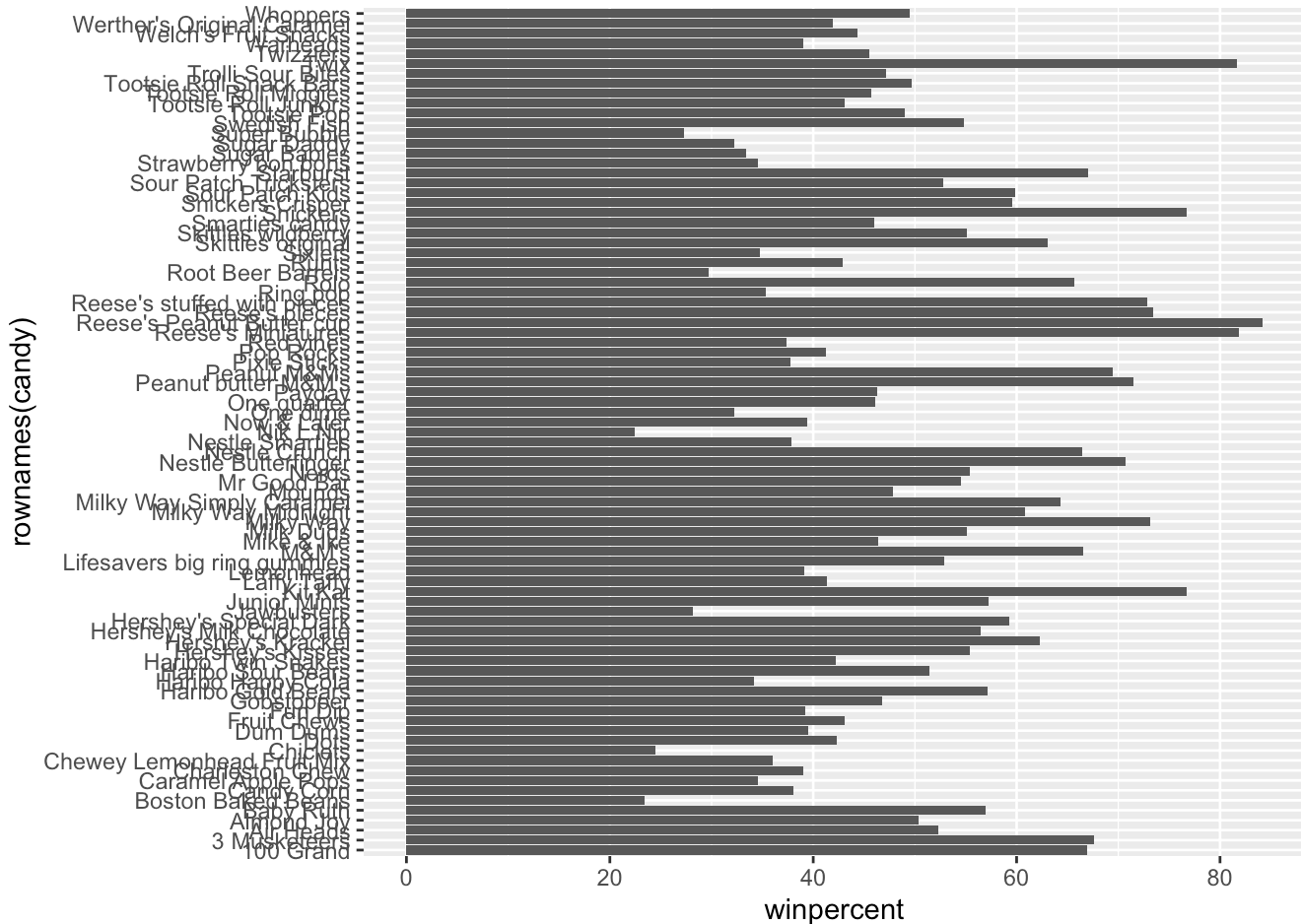|                           | pricepercent | winpercent |
|---------------------------|--------------|------------|
| Snickers                  | 0.651        | 76.67378   |
| Kit Kat                   | 0.511        | 76.76860   |
| Twix                      | 0.906        | 81.64291   |
| Reese's Miniatures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

> Q14. The five most liked candy types are Reese's Peanut Butter cup, Reese's Miniatures, Twix, Kit Kat, and Snickers.

> I don't have a preference over these two approaches. Base r does not require loading additional packages, while dplyr provides a more streamlined syntax.
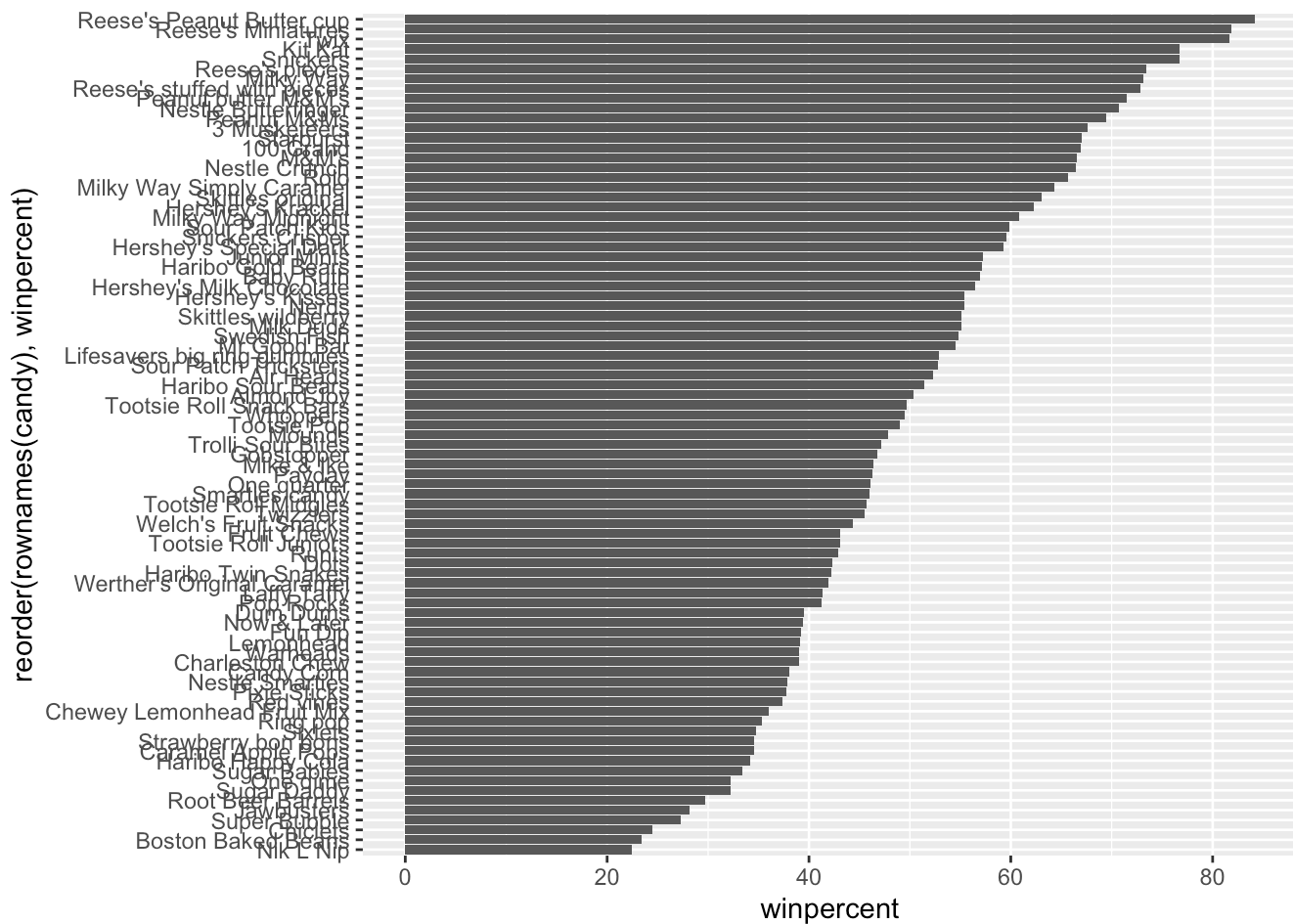
## Q15. Figure shown below

```
library(ggplot2)

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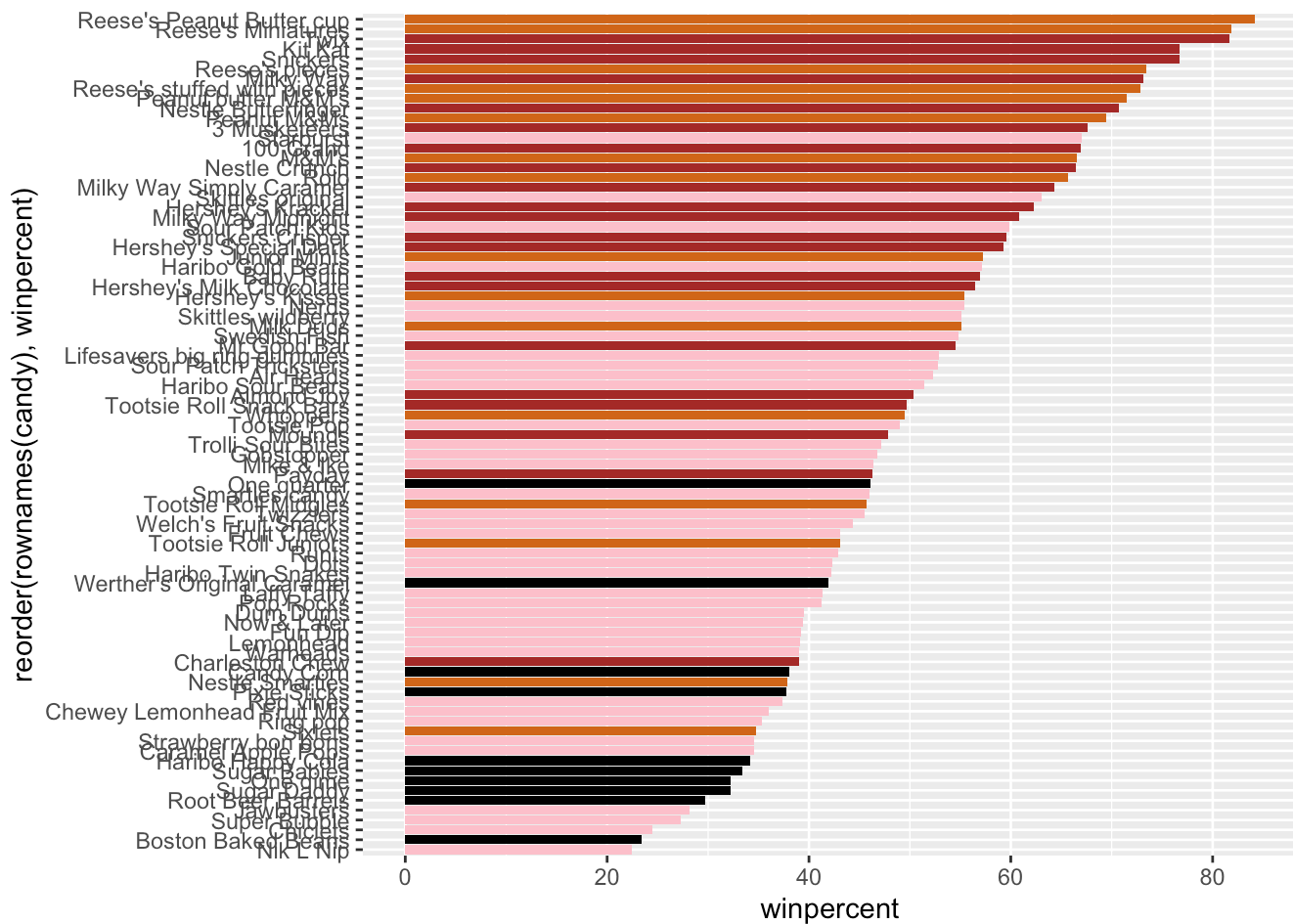


## Q16. Figure shown below

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

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

Q17. Sixlets

Q18. Starburst

# Taking a look at pricepercent

```
library(ggrepel)

# How about a plot of price vs win
my_cols[candy$fruity==1] = "red"
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)
```

Warning: ggrepel: 50 unlabeled data points (too many overlaps). Consider
increasing max.overlaps

Q19. Reese's Miniatures

Q20. Nik L Nip, Nestle Smarties, Ring pop, Hershey's Krackel, Hershey's Milk Chocolate

```
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   |

Q21. Figure shown below

```
# Make a lollipop chart of pricepercent
ggplot(candy) +
  aes(pricepercent, reorder(rownames(candy), pricepercent)) +
  geom_segment(aes(yend = reorder(rownames(candy), pricepercent),
```

```
                        xend = 0), col="gray40") +
        geom_point()
```



# Exploring the correlation structure

```
library(corrplot)
```

```
corrplot 0.95 loaded
```

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

|  | chocolate | fruity | caramel | peanutyalmondy | nougat |
|---|---|---|---|---|---|
| chocolate | 1.0000000 | -0.74172106 | 0.24987535 | 0.37782357 | 0.25489183 |
| fruity | -0.7417211 | 1.00000000 | -0.33548538 | -0.39928014 | -0.26936712 |
| caramel | 0.2498753 | -0.33548538 | 1.00000000 | 0.05935614 | 0.32849280 |
| peanutyalmondy | 0.3778236 | -0.39928014 | 0.05935614 | 1.00000000 | 0.21311310 |
| nougat | 0.2548918 | -0.26936712 | 0.32849280 | 0.21311310 | 1.00000000 |
| crispedricewafer | 0.3412098 | -0.26936712 | 0.21311310 | -0.01764631 | -0.08974359 |
| hard | -0.3441769 | 0.39067750 | -0.12235513 | -0.20555661 | -0.13867505 |
| bar | 0.5974211 | -0.51506558 | 0.33396002 | 0.26041960 | 0.52297636 |

| | | | | | |
|---|---|---|---|---|---|
| pluribus | −0.3396752 | 0.29972522 | −0.26958501 | −0.20610932 | −0.31033884 |
| sugarpercent | 0.1041691 | −0.03439296 | 0.22193335 | 0.08788927 | 0.12308135 |
| pricepercent | 0.5046754 | −0.43096853 | 0.25432709 | 0.30915323 | 0.15319643 |
| winpercent | 0.6365167 | −0.38093814 | 0.21341630 | 0.40619220 | 0.19937530 |

| | crispedricewafer | hard | bar | pluribus |
|---|---|---|---|---|
| chocolate | 0.34120978 | −0.34417691 | 0.59742114 | −0.33967519 |
| fruity | −0.26936712 | 0.39067750 | −0.51506558 | 0.29972522 |
| caramel | 0.21311310 | −0.12235513 | 0.33396002 | −0.26958501 |
| peanutyalmondy | −0.01764631 | −0.20555661 | 0.26041960 | −0.20610932 |
| nougat | −0.08974359 | −0.13867505 | 0.52297636 | −0.31033884 |
| crispedricewafer | 1.00000000 | −0.13867505 | 0.42375093 | −0.22469338 |
| hard | −0.13867505 | 1.00000000 | −0.26516504 | 0.01453172 |
| bar | 0.42375093 | −0.26516504 | 1.00000000 | −0.59340892 |
| pluribus | −0.22469338 | 0.01453172 | −0.59340892 | 1.00000000 |
| sugarpercent | 0.06994969 | 0.09180975 | 0.09998516 | 0.04552282 |
| pricepercent | 0.32826539 | −0.24436534 | 0.51840654 | −0.22079363 |
| winpercent | 0.32467965 | −0.31038158 | 0.42992933 | −0.24744787 |

| | sugarpercent | pricepercent | winpercent |
|---|---|---|---|
| chocolate | 0.10416906 | 0.5046754 | 0.6365167 |
| fruity | −0.03439296 | −0.4309685 | −0.3809381 |
| caramel | 0.22193335 | 0.2543271 | 0.2134163 |
| peanutyalmondy | 0.08788927 | 0.3091532 | 0.4061922 |
| nougat | 0.12308135 | 0.1531964 | 0.1993753 |
| crispedricewafer | 0.06994969 | 0.3282654 | 0.3246797 |
| hard | 0.09180975 | −0.2443653 | −0.3103816 |
| bar | 0.09998516 | 0.5184065 | 0.4299293 |
| pluribus | 0.04552282 | −0.2207936 | −0.2474479 |
| sugarpercent | 1.00000000 | 0.3297064 | 0.2291507 |
| pricepercent | 0.32970639 | 1.0000000 | 0.3453254 |
| winpercent | 0.22915066 | 0.3453254 | 1.0000000 |

```
corrplot(cij)
```

Q22. Chocolate and fruity

Q23. Chocolate and bar

# Principal Component Analysis

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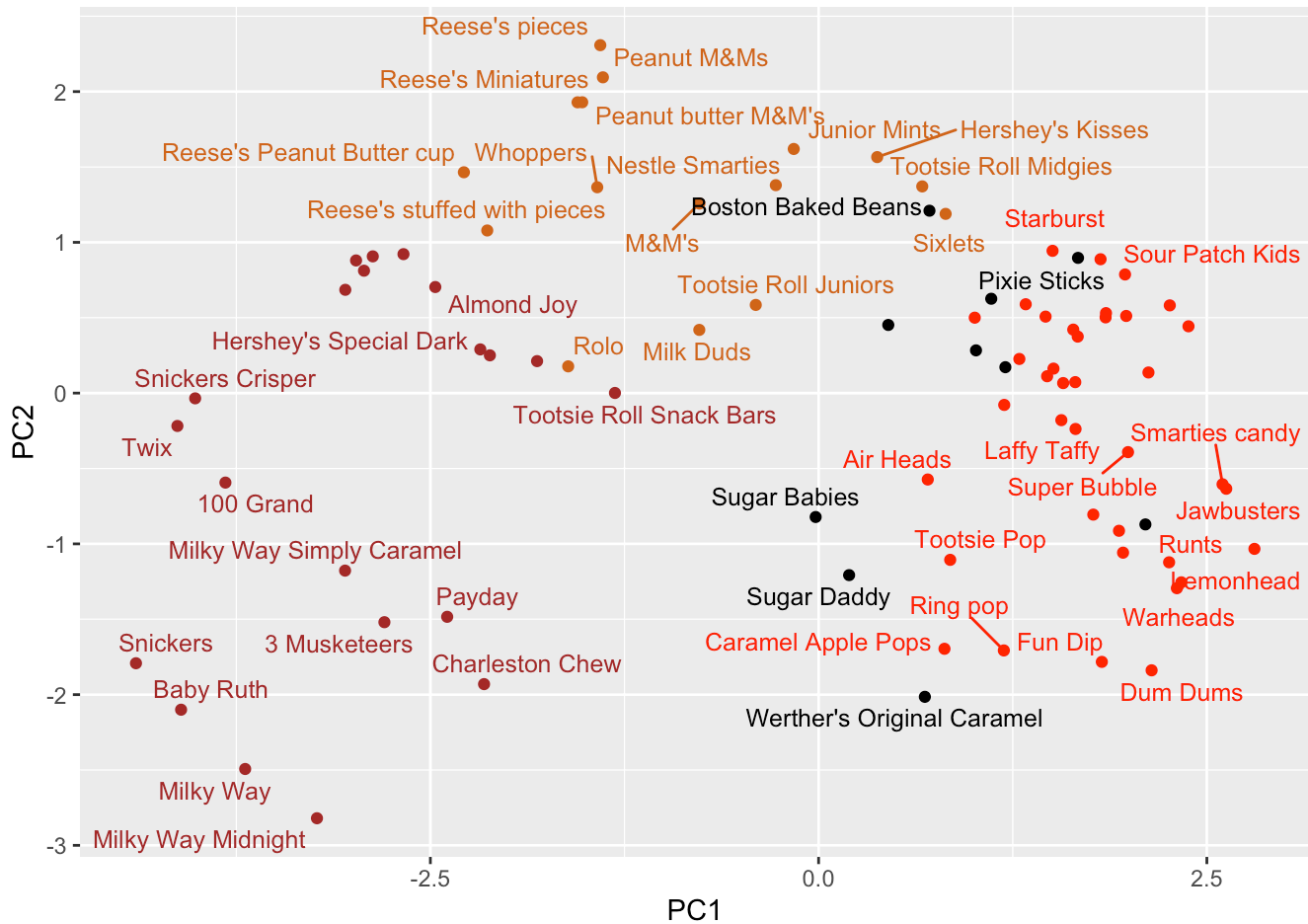```

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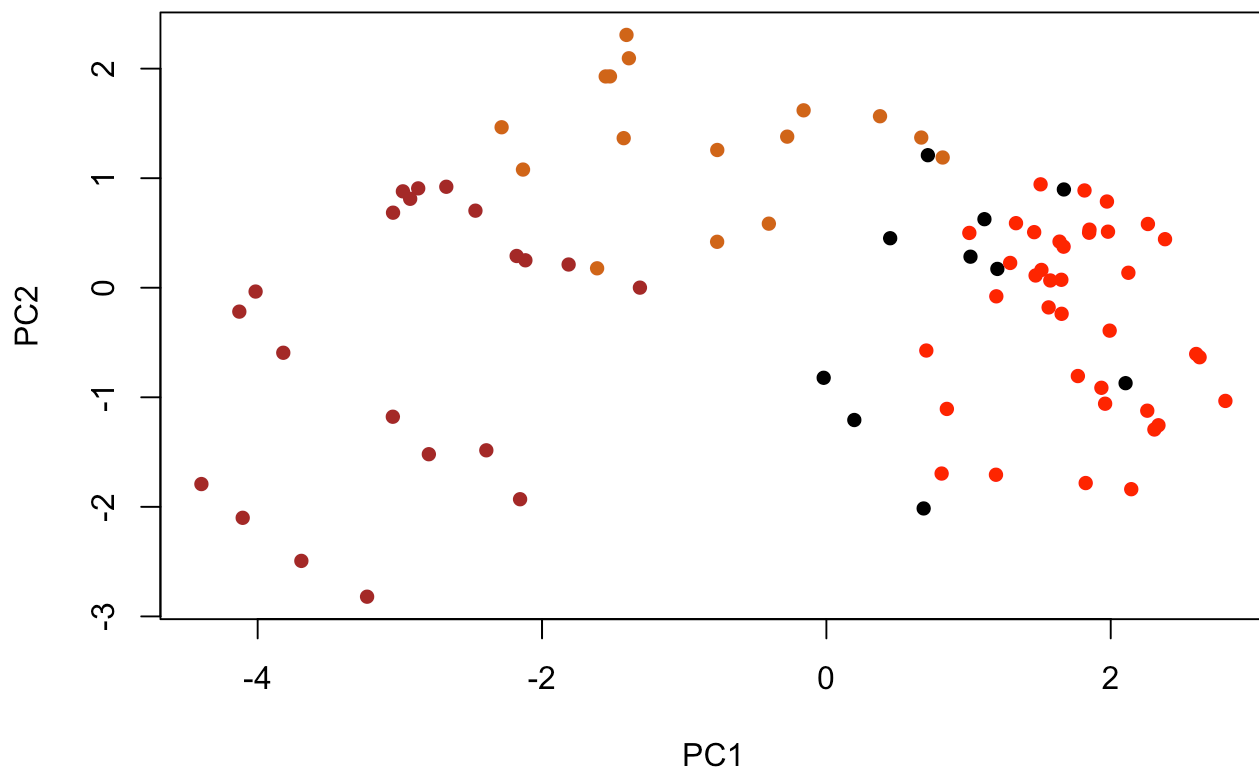
```
ggplot(pca$x) + aes(PC1,PC2, label=rownames(pca$x)) + geom_point(col=my_cols) + geom_text
```

```
Warning: ggrepel: 35 unlabeled data points (too many overlaps). Consider
increasing max.overlaps
```



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

```
my_data <- cbind(candy, pca$x[,1:3])
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
```
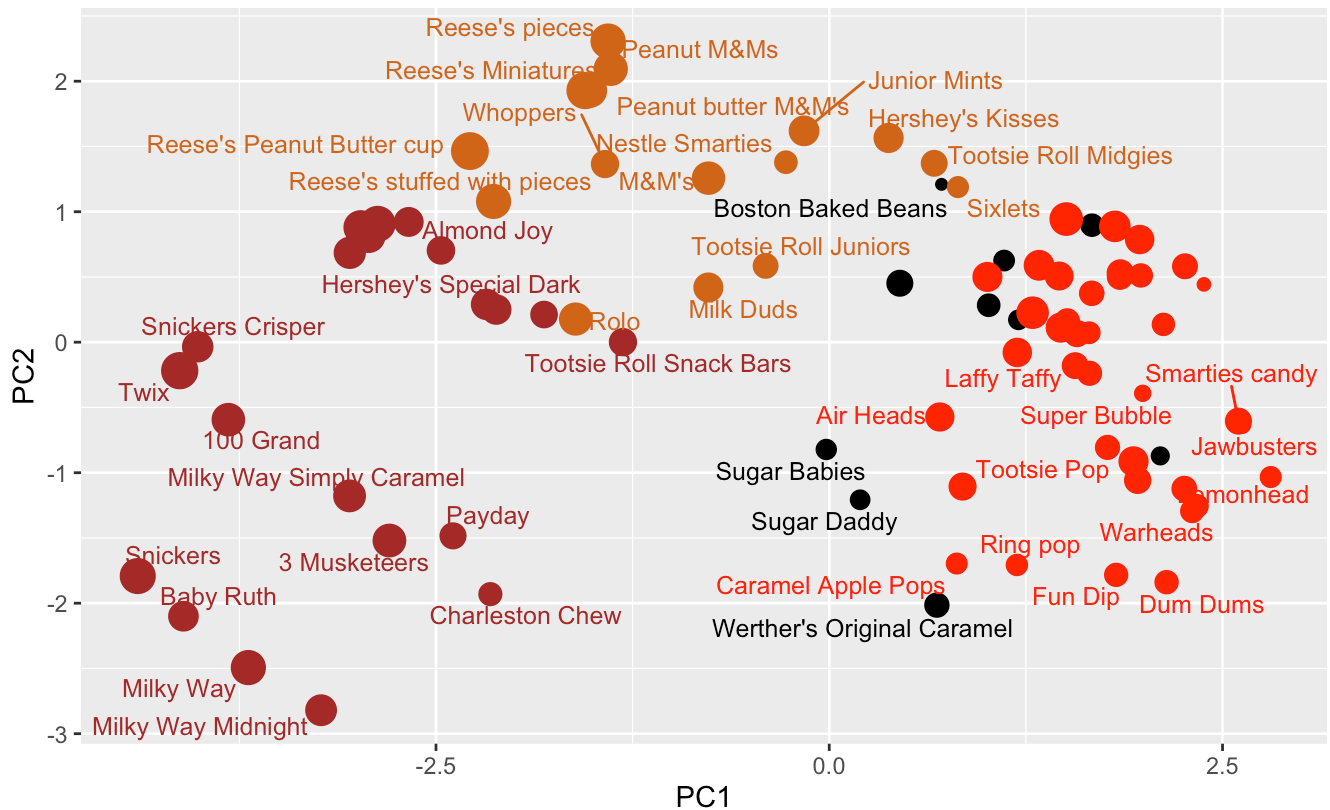
```
p + geom_text_repel(size=3.3, col=my_cols, max.overlaps = 7)  +
  theme(legend.position = "none") +
  labs(title="Halloween Candy PCA Space",
       subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brow
       caption="Data from 538")
```

Warning: ggrepel: 39 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), fruity (red), other (blac



Data from 538

```r
library(plotly)
```

```
Attaching package: 'plotly'

The following object is masked from 'package:ggplot2':

    last_plot

The following object is masked from 'package:stats':

    filter

The following object is masked from 'package:graphics':

    layout
```
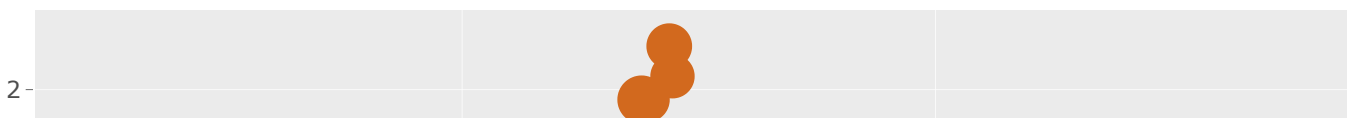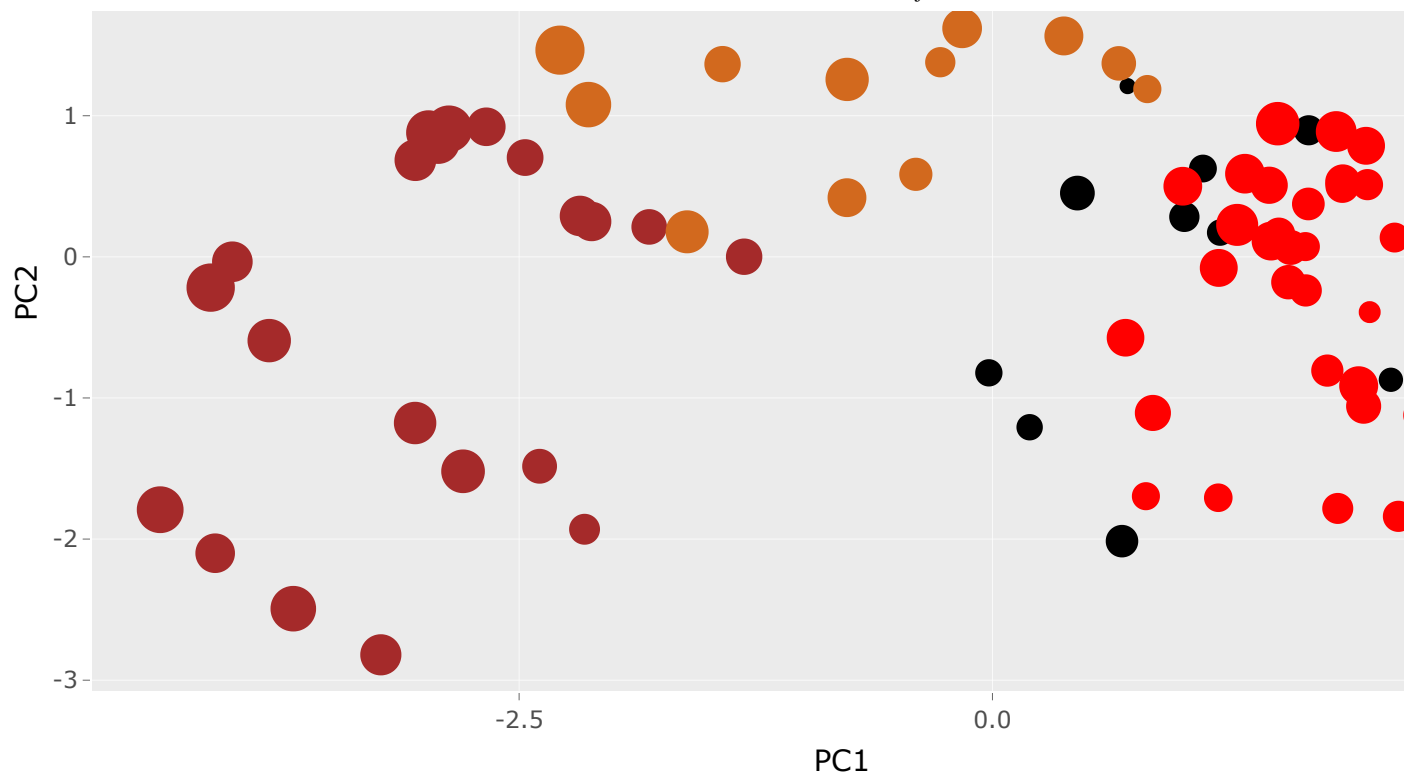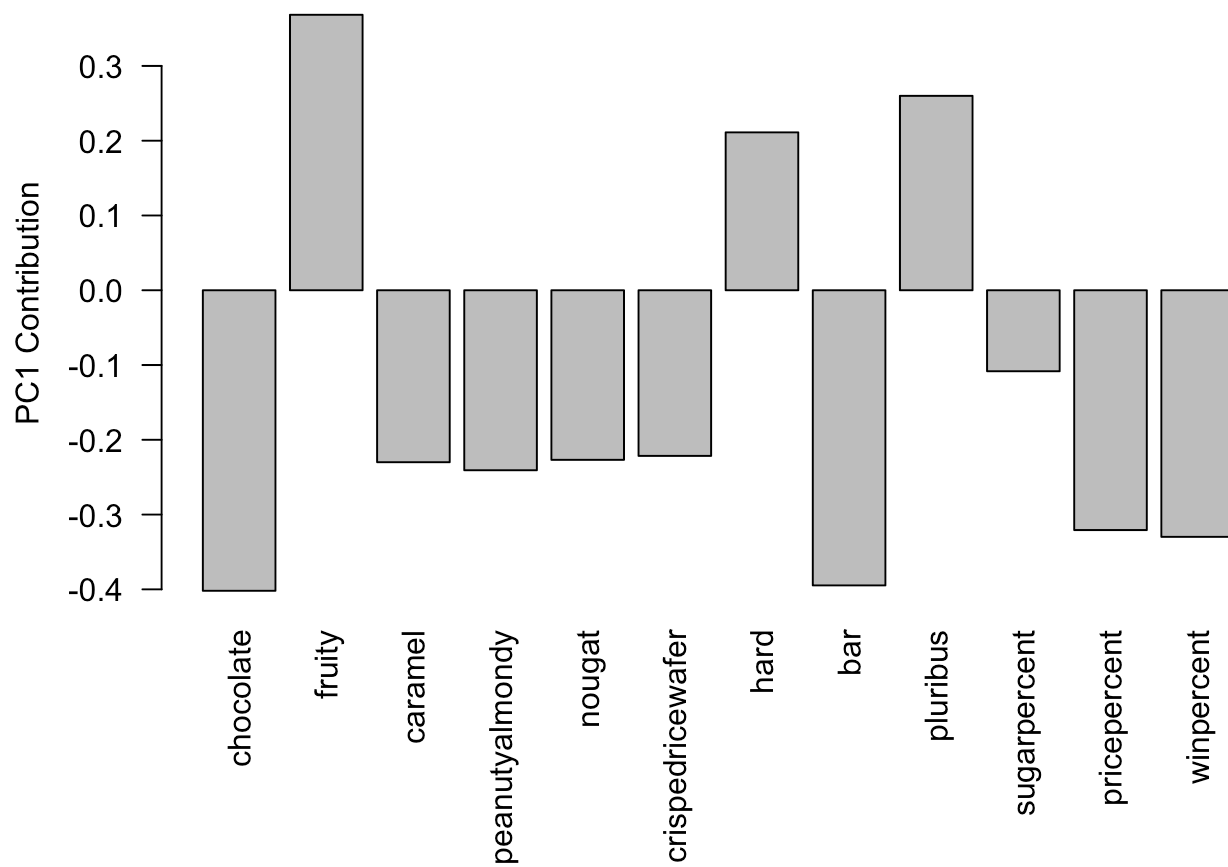
```r
ggplotly(p)
```

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

Q24. Fruity, hard, and pluribus are picked up strongly by PC1 in the positive direction. It makes sense as the first dimension separates fruity candies from chocolate candies. In the correlationn plot, we can also see that fruity candies only correlate with hard and pluribus.