

lab 9: Halloween candy mini project

Elsa Chen (A16632961)

Table of contents

Importing data	1
What is your favorite candy?	2
Overall candy rankings	5
add some colors	9
Taking a look at pricepercent	11
Exploring the correlation structure	12
Principal Component Analysis	13

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

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

There are 85 candy types in this dataset

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

There are 38 candy types in the dataset

What is your favorite candy?

```
candy["Twix", ]$winpercent
```

```
[1] 81.64291
```

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

```
candy["Warheads", ]$winpercent
```

```
[1] 39.0119
```

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"?

It is 49.653503

```
# install.packages("skimr")
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?

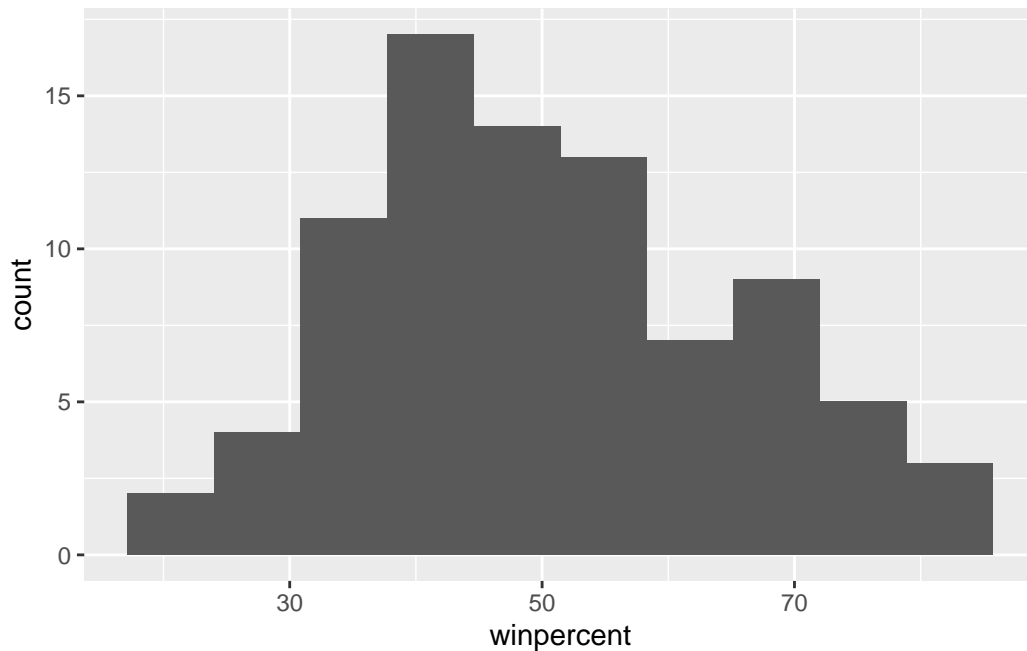
Yes, winpercent

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

0 represents FALSE, 1 represents TRUE

Q8. Plot a histogram of winpercent values

```
library(ggplot2)
ggplot(candy) +
  aes(x = winpercent) +
  geom_histogram(bins = 10)
```



Q9. Is the distribution of winpercent values symmetrical?

No skews to the left a bit

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

below 50%

```
mean(candy$winpercent)
```

```
[1] 50.31676
```

But mean is above 50%

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

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

```
[1] 60.92153
```

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

```
[1] 44.11974
```

chocolate candy is ranked higher than fruity candy

Q12. Is this difference statistically significant?

```
t.test(chocolate, fruity)
```

Welch Two Sample t-test

```
data: chocolate and 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
```

It is

Overall candy rankings

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

```
candy <- dplyr::arrange(candy, winpercent)
head(candy, 5)
```

	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	
	crisped	rice	wafer	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

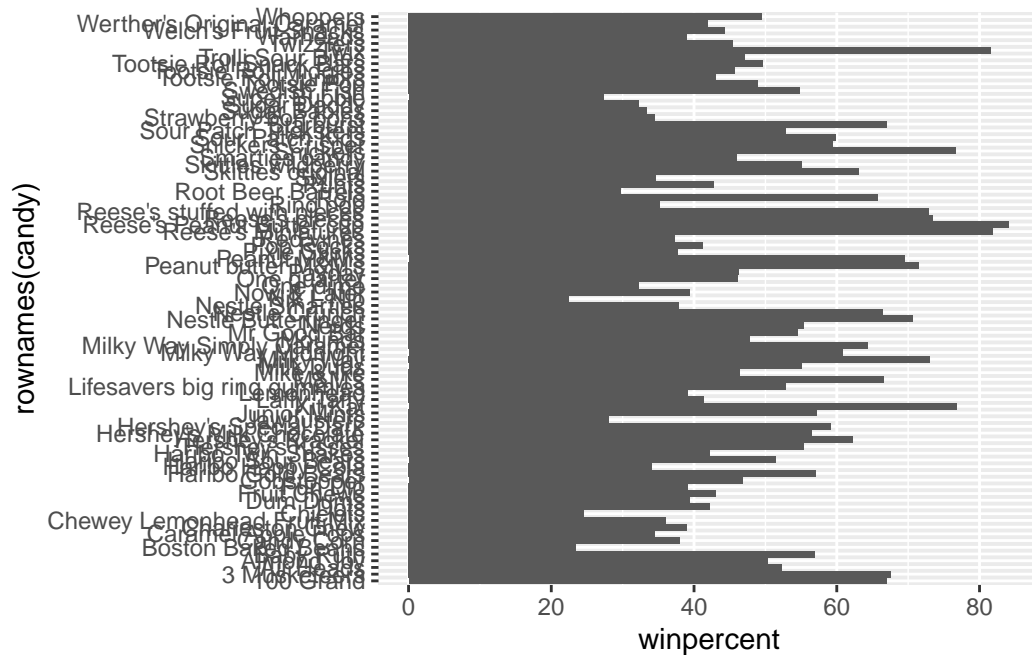
Q14. What are the top 5 all time favorite candy types out of this set?

```
candy <- dplyr::arrange(candy, desc(winpercent))
head(candy, 5)
```

	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	
	crisped	rice	wafer	hard	bar	pluribus	sugar
Reese's Peanut Butter cup		0	0	0		0	0.720
Reese's 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
	price	percent	win	percent			
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					

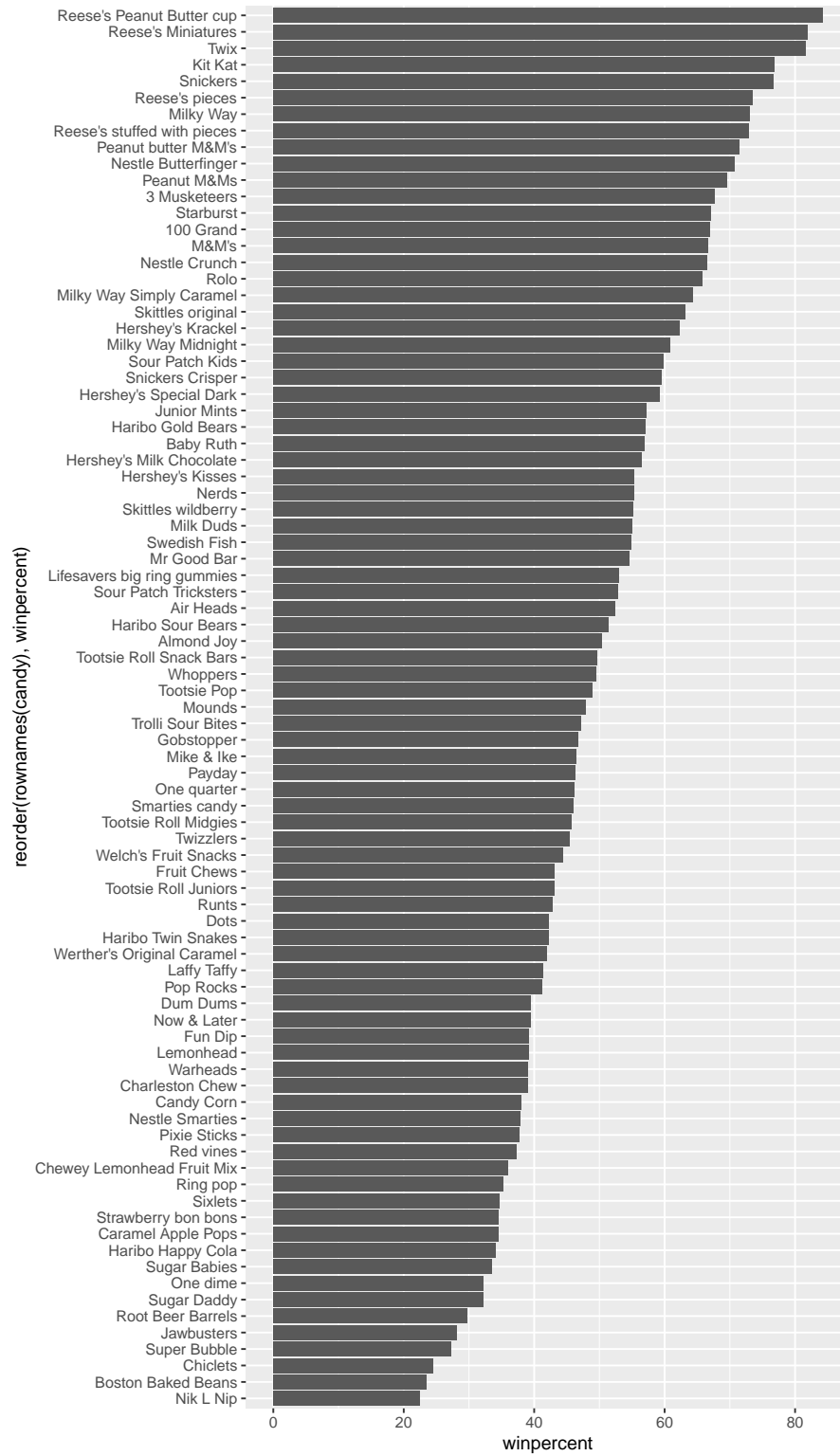
Q15. Make a first barplot of candy ranking based on winpercent values.

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



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

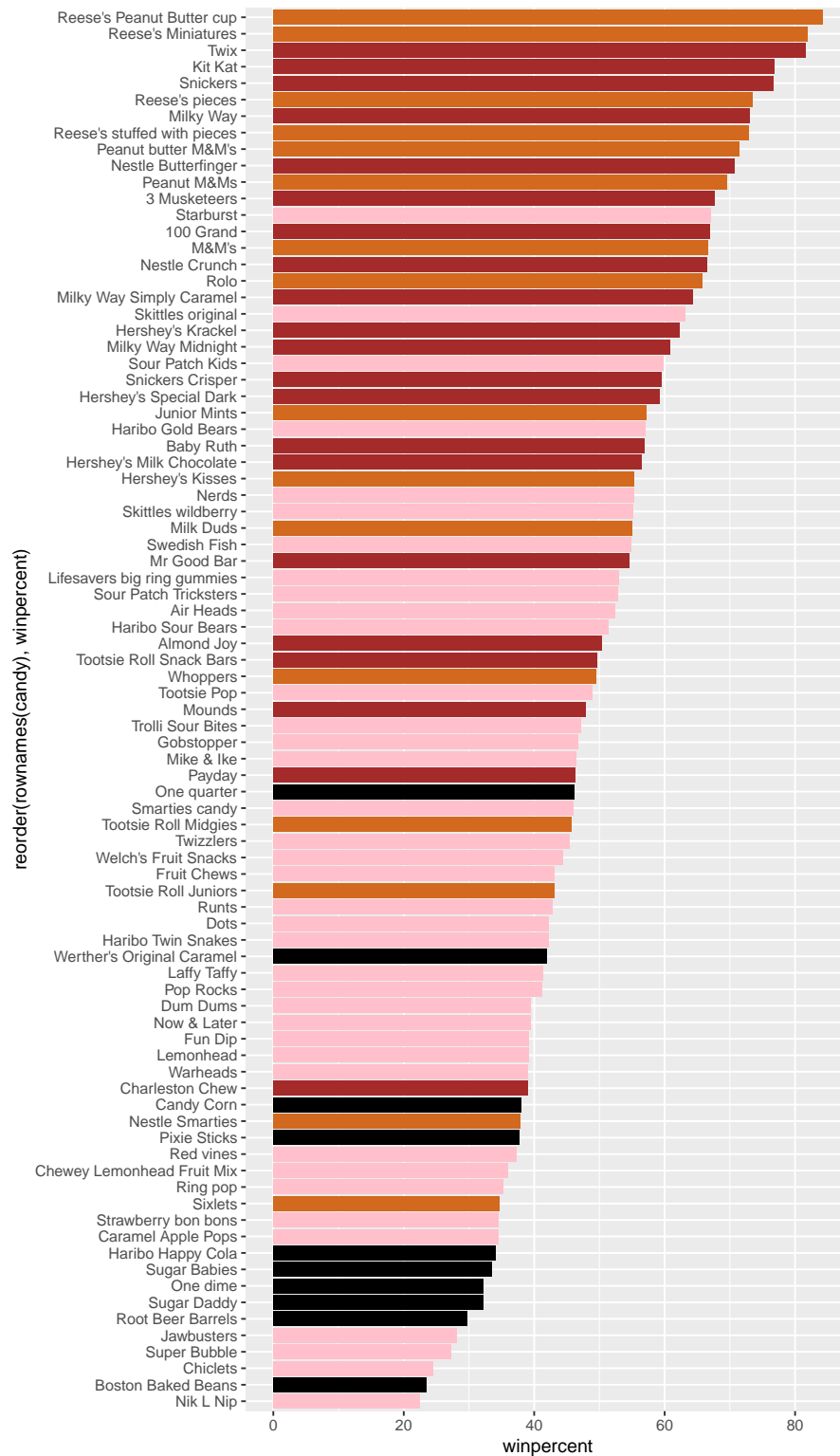
```
ggplot(candy) +
  aes(x = winpercent, y = reorder(rownames(candy), winpercent)) +
  geom_col()
```




```
# ggsave("barplot.png", height = 12) # if want to save graph as figure  
 if want to insert an image
```

add some colors

```
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_cols)
```



Q17. What is the worst ranked chocolate candy?

Sixlets > Q18. What is the best ranked fruity candy?

Starburst

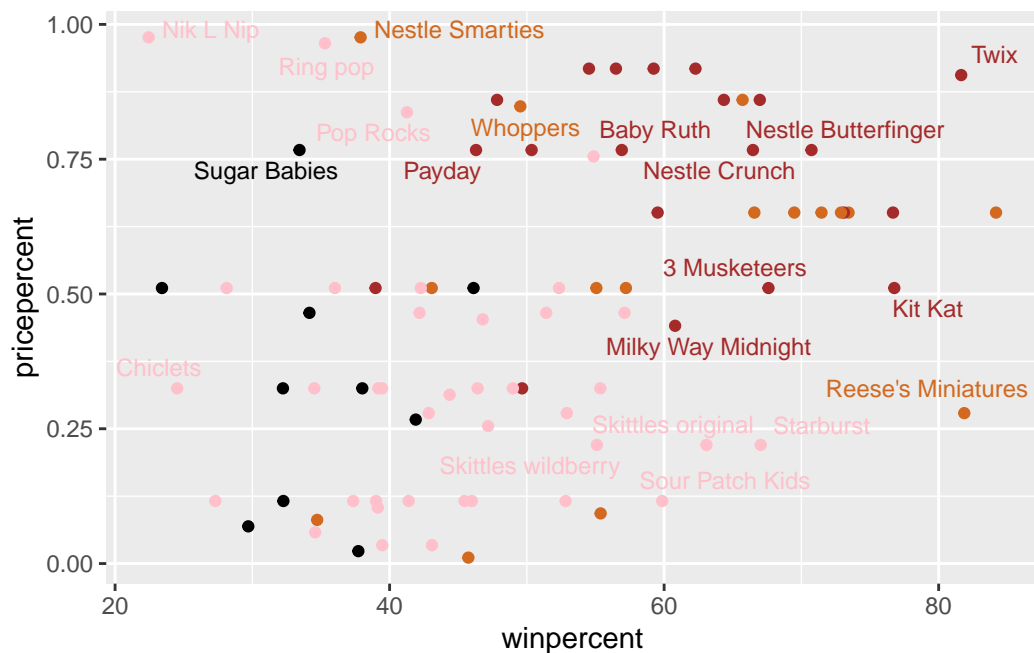
Taking a look at pricepercent

```
library(ggrepel)
```

Warning: package 'ggrepel' was built under R version 4.3.3

```
# 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=3.3, max.overlaps = 5)
```

Warning: ggrepel: 65 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?

```
candy <- dplyr::arrange(candy, desc(pricepercent))
head( candy[,c(11,12)], n=5 )
```

	pricepercent	winpercent
Nestle Smarties	0.976	37.88719
Nik L Nip	0.976	22.44534
Ring pop	0.965	35.29076
Hershey's Krackel	0.918	62.28448
Hershey's Special Dark	0.918	59.23612

Hershey's Krackel

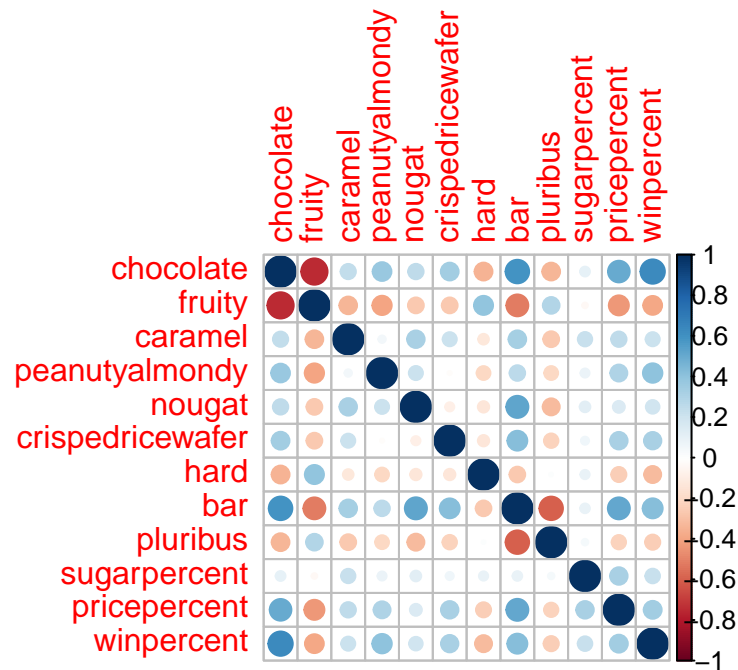
Exploring the correlation structure

```
library(corrplot)
```

Warning: package 'corrplot' was built under R version 4.3.3

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 bar

Principal Component Analysis

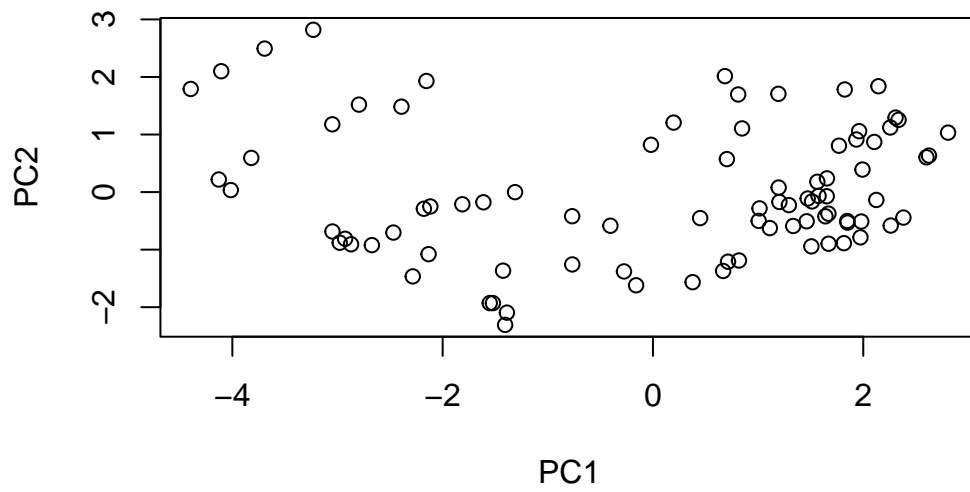
```
pca <- prcomp(candy, scale = TRUE)
summary(pca)
```

Importance of components:

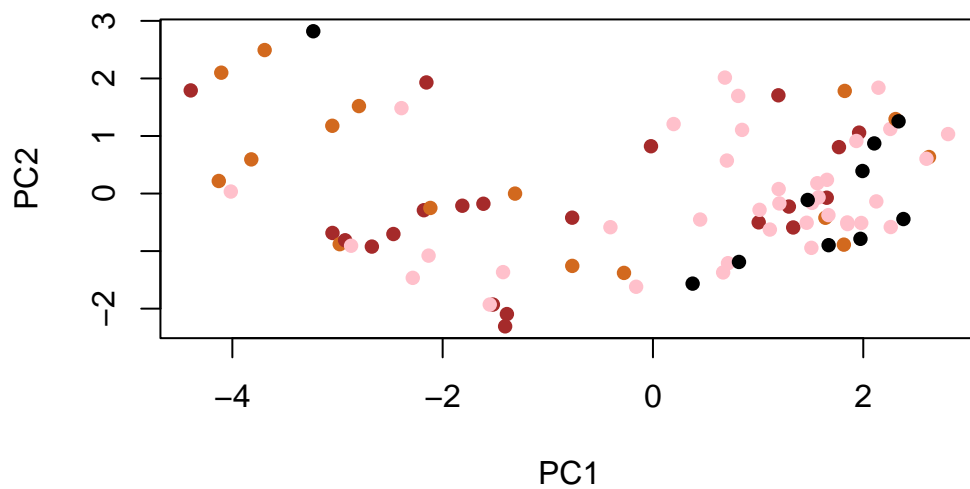
	PC1	PC2	PC3	PC4	PC5	PC6	PC7
Standard deviation	2.0788	1.1378	1.1092	1.07533	0.9518	0.81923	0.81530
Proportion of Variance	0.3601	0.1079	0.1025	0.09636	0.0755	0.05593	0.05539
Cumulative Proportion	0.3601	0.4680	0.5705	0.66688	0.7424	0.79830	0.85369

	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

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



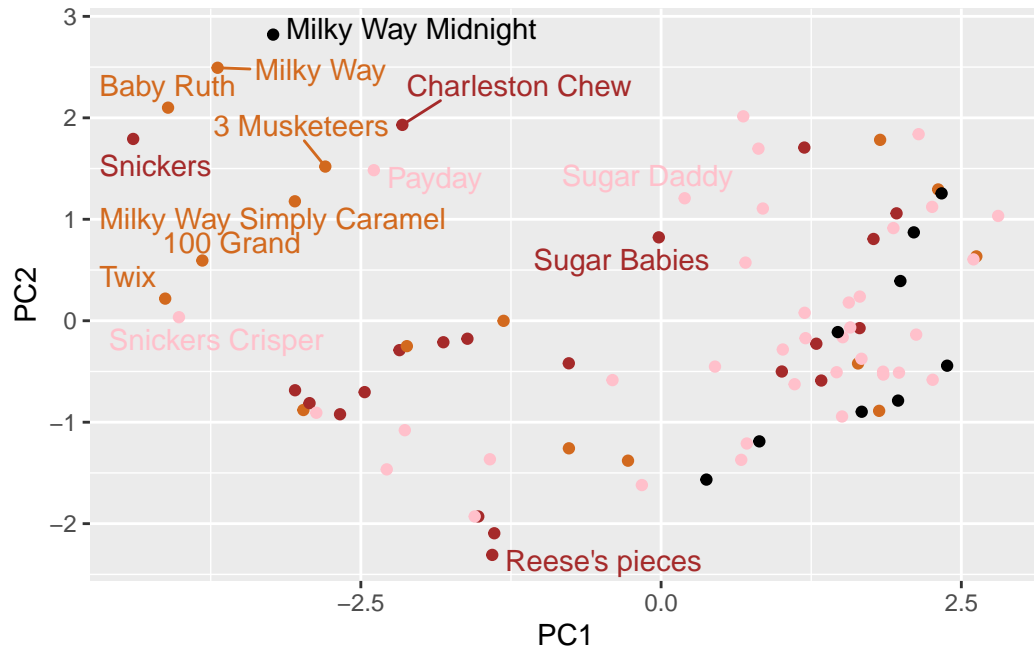
```
plot(pca$x[,1:2], col=my_cols, pch=16)
```



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

```
ggplot(my_data) +  
  aes(x=PC1, y=PC2,  
      label=rownames(my_data)) +  
  geom_point(col=my_cols) +  
  geom_text_repel(col = my_cols, max.overlaps = 5)
```

Warning: ggrepel: 71 unlabeled data points (too many overlaps). Consider increasing max.overlaps

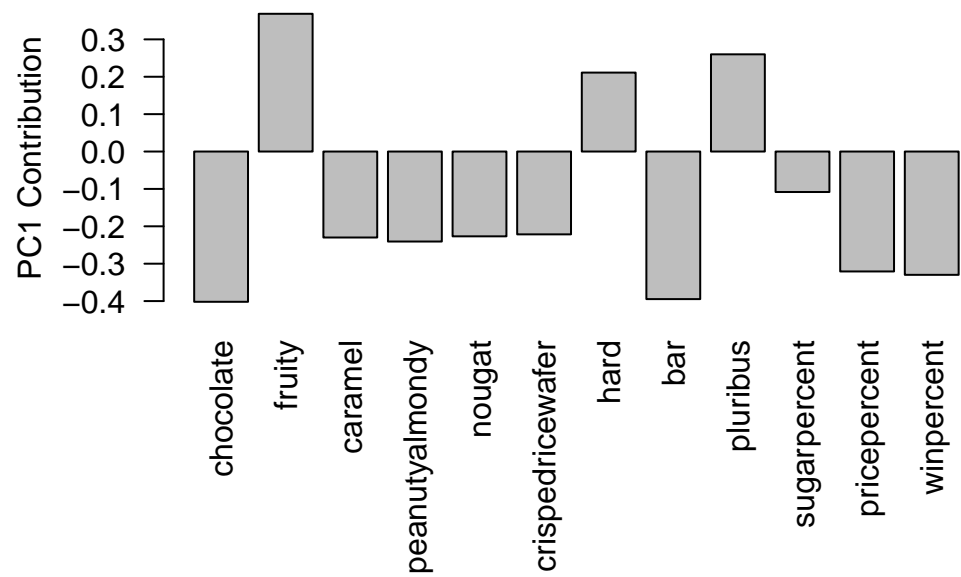


For interactive plot:

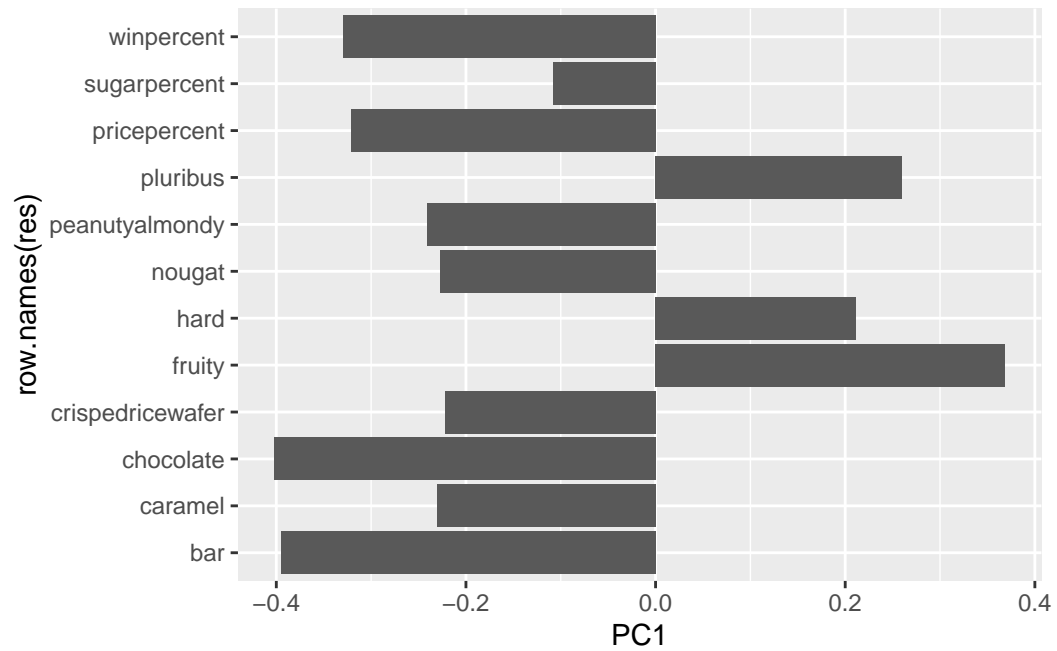
```
#library(plotly)
```

```
#p <- ggplot(my_data) +
#   aes(x=PC1, y=PC2,
#       size=winpercent/100,
#       text=rownames(my_data),
#       label=rownames(my_data)) +
#ggplotly(p)
```

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

```
res <- as.data.frame(pca$rotation)
ggplot(res) +
  aes(PC1, row.names(res)) +
  geom_col()
```



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

pluribus, hard, and fruity