

Class10

Nicholas Thiphakhinkeo A17686679

Read/Inspect

```
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?

```
num_candy_types <- nrow(candy)
print(num_candy_types)
```

[1] 85

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

```
fruity_candy <- candy[candy$fruity==T,]
num_fruity_candy_types <- nrow(fruity_candy)
print(num_fruity_candy_types)
```

[1] 38

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

```
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

Skim Function

```
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	

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

sugarpercent, pricepercent, and winpercent show percent rather than value.

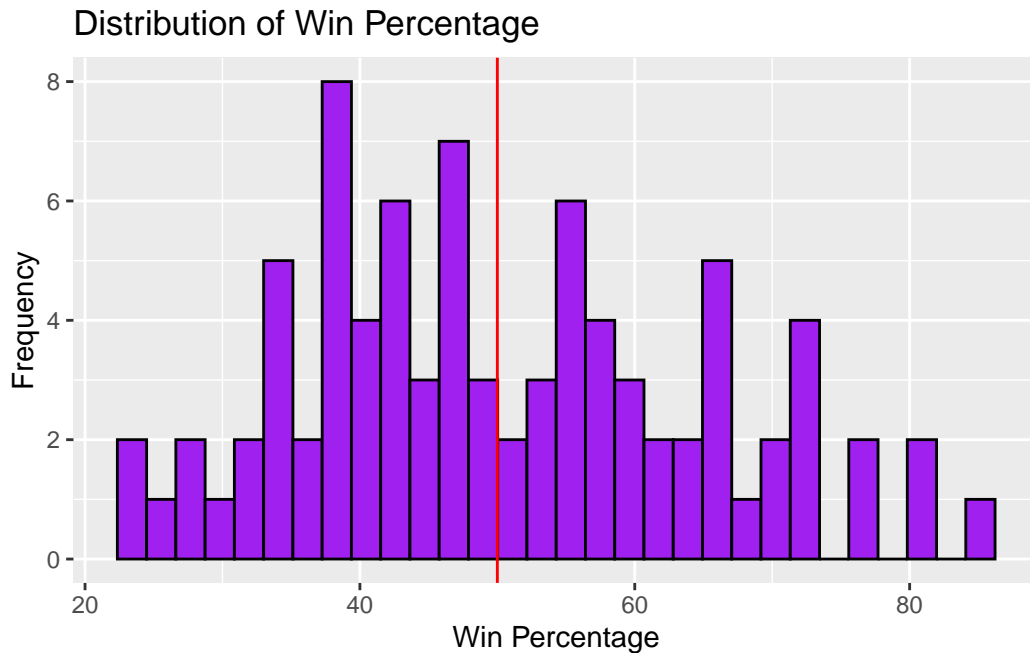
Q7. What do you think a zero and one represent for the candy\$chocolate column?

The 0s are likely for candies not containing chocolate while 1s are for candies that do.

Q8. Plot a histogram of winpercent values

```
library("ggplot2")
ggplot(candy, aes(x = winpercent)) +
  geom_histogram(fill = "purple", color = "black") +
  geom_vline(xintercept = 50, color = "red") +
  labs(title = "Distribution of Win Percentage", x = "Win Percentage", y = "Frequency")
```

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



Q9. Is the distribution of winpercent values symmetrical?

No the distribution is right-skewed

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

Below 50% suggesting most candies have win percentage below 50%

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

```
mean_chocolate_winpercent <- mean(candy$winpercent[as.logical(candy$chocolate)])
mean_fruity_winpercent <- mean(candy$winpercent[as.logical(candy$fruity)])

print(mean_chocolate_winpercent)
```

```
[1] 60.92153
```

```
print(mean_fruity_winpercent)
```

```
[1] 44.11974
```

Chocolate is ranked higher than fruity

Q12. Is this difference statistically significant?

```
t_test1 <- t.test(candy$winpercent[as.logical(candy$chocolate)], candy$winpercent[as.logical(candy$fruit)], var.equal = FALSE)
print(t_test1)
```

Welch Two Sample t-test

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

p value <-0.05 so likely significant

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

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

	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?

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

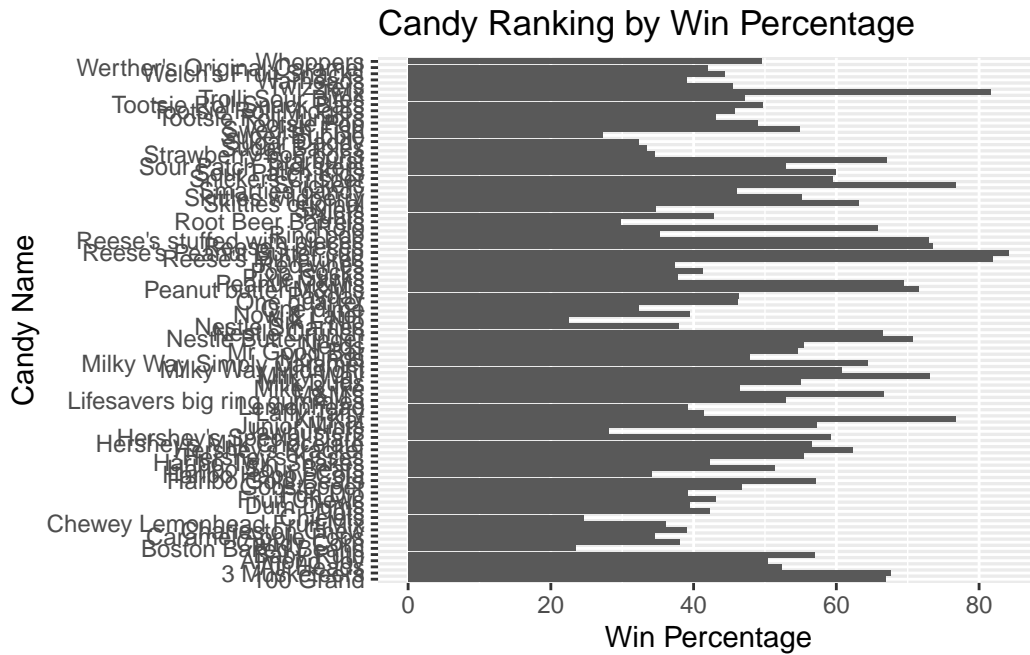
	chocolate	fruity	caramel	peanut	almond	nougat
Reese's Peanut Butter cup	1	0	0		1	0
Reese's Miniatures	1	0	0		1	0
Twix	1	0	1		0	0
Kit Kat	1	0	0		0	0
Snickers	1	0	1		1	1

	crisped	rice	wafer	hard	bar	pluribus	sugar	percent
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	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

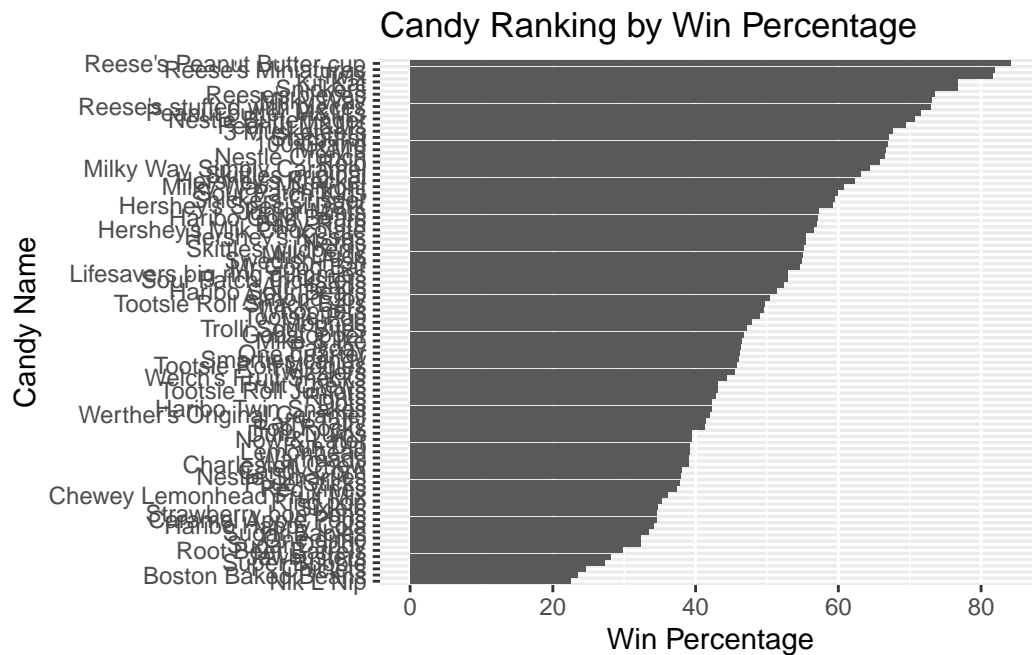
Q15. Make a first barplot of candy ranking based on winpercent values.

```
ggplot(candy) +
  aes(x = rownames(candy), y = winpercent) +
  geom_col() +
  coord_flip() +
  labs(title = "Candy Ranking by Win Percentage",
       x = "Candy Name",
       y = "Win Percentage")
```



Q16. use the `reorder()` function to get the bars sorted by winpercent?

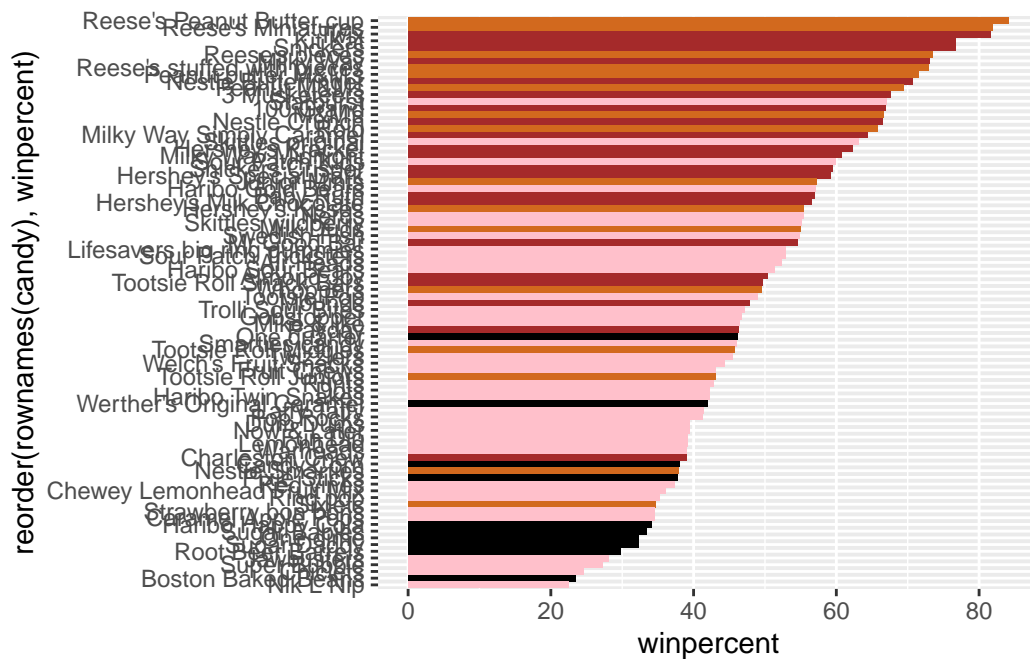
```
ggplot(candy) +
  aes(x = reorder(rownames(candy), winpercent), y=winpercent) +
  geom_col() +
  coord_flip() +
  labs(title = "Candy Ranking by Win Percentage",
       x = "Candy Name",
       y = "Win Percentage")
```



Adding Color

```
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?

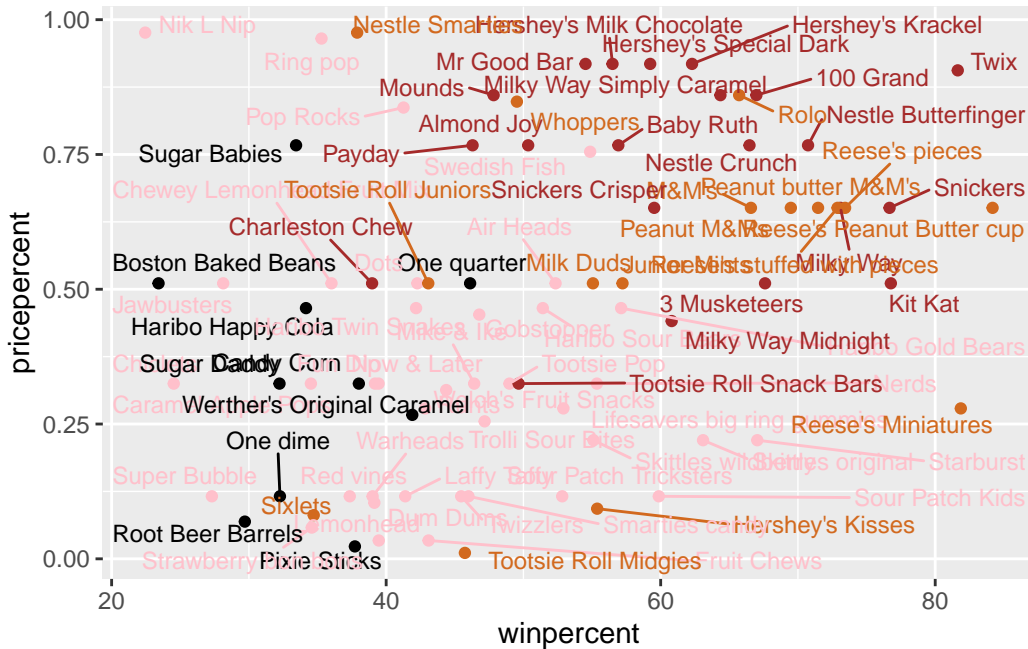
The worst ranked chocolate candies are Sixlets

Q18. What is the best ranked fruity candy?

The best ranked fruity candy is Starburst

Taking a look at pricepercent

```
library("ggrepel")
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(col=my_cols, size=3.3, max.overlaps = 78)
```



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 Bars

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

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

	chocolate	fruity	caramel	peanut	almond	nougat
Nik L Nip	0	1	0		0	0
Nestle Smarties	1	0	0		0	0
Ring pop	0	1	0		0	0
Hershey's Krackel	1	0	0		0	0
Hershey's Milk Chocolate	1	0	0		0	0

	crispedrice	wafer	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
Hershey's Krackel		1	0	1	0	0.430
Hershey's 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
Hershey's Krackel	0.918	62.28448
Hershey's Milk Chocolate	0.918	56.49050

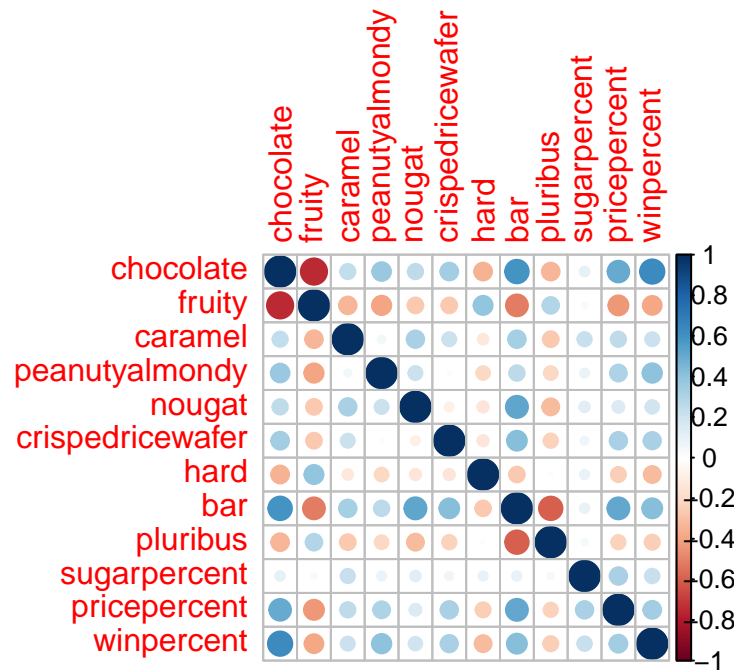
Nik L Nips are the most expensive and least popular

Corrplot

```
library(corrplot)
```

```
corrplot 0.95 loaded
```

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



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

Fruity winpercentage is least correlated

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

Chocolate winpercentage is most positively correlated

PCA

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

Importance of components:

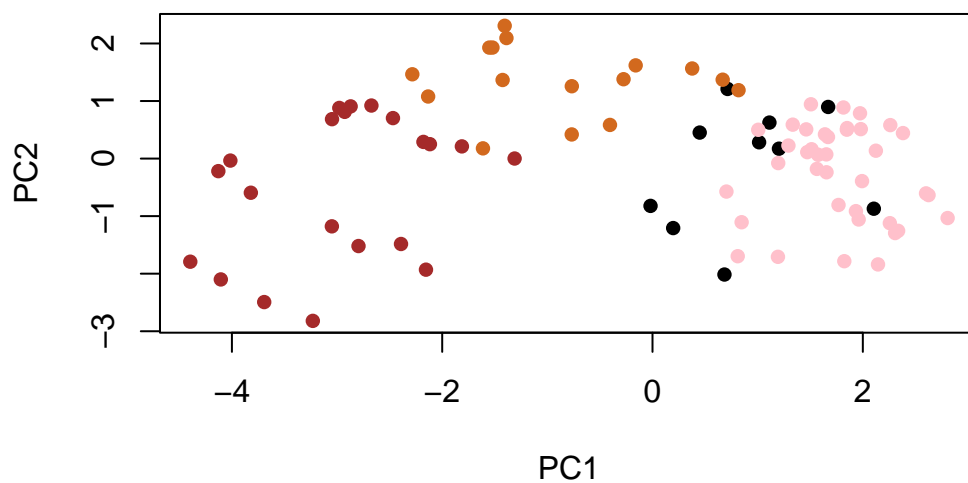
	PC1	PC2	PC3	PC4	PC5	PC6	PC7
Standard deviation	2.0788	1.1378	1.1092	1.07533	0.9518	0.81923	0.81530
Proportion of Variance	0.3601	0.1079	0.1025	0.09636	0.0755	0.05593	0.05539
Cumulative Proportion	0.3601	0.4680	0.5705	0.66688	0.7424	0.79830	0.85369

	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

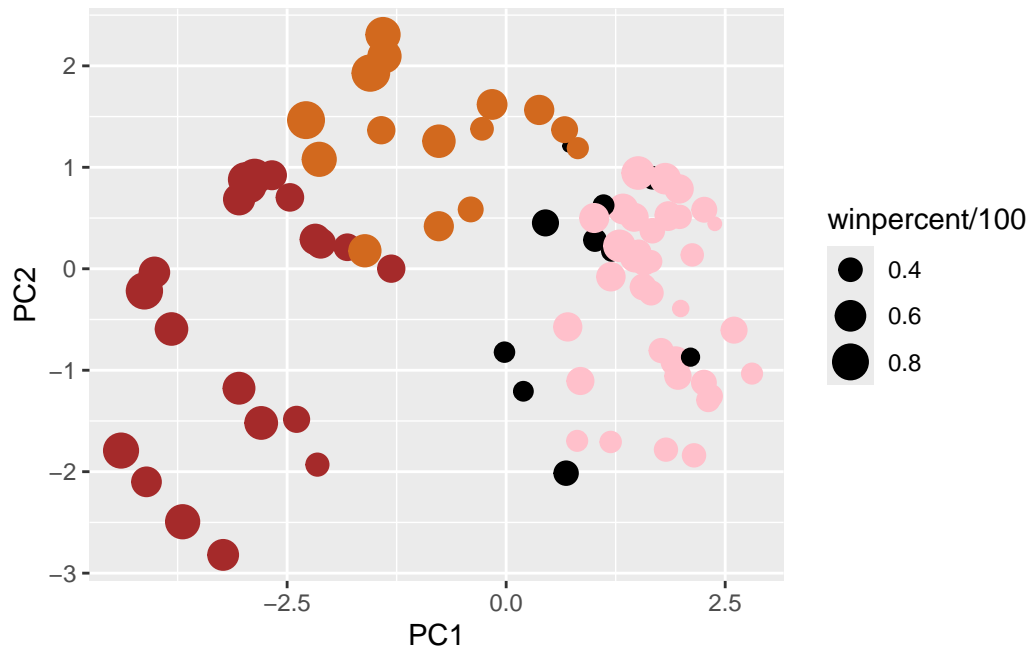
PC1vPC2

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



PCAggplot

```
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
```



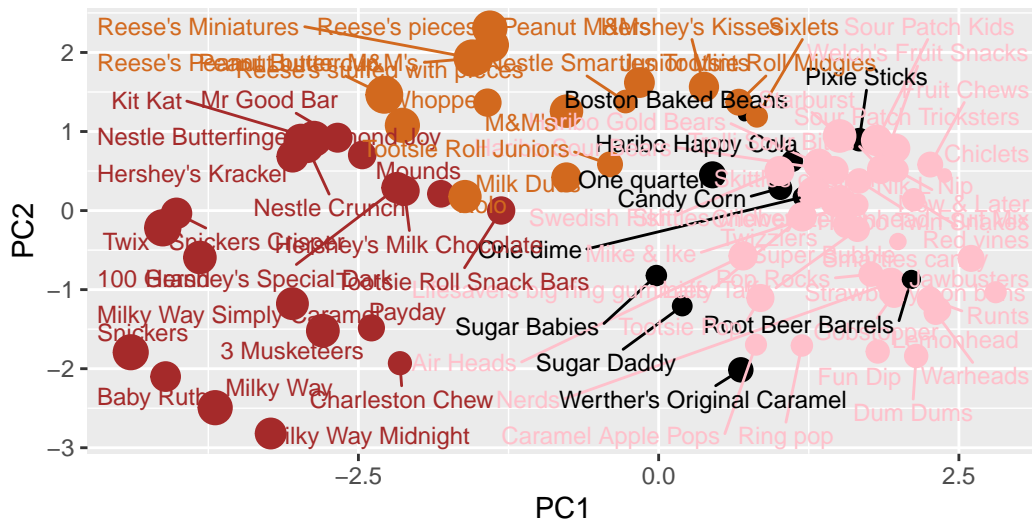
PCAggrepel

```
library(ggrepel)

p + geom_text_repel(size=3.3, col=my_cols, max.overlaps = 53) +
  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

PCAplotly

```
library(plotly)
```

Attaching package: 'plotly'

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

last_plot

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

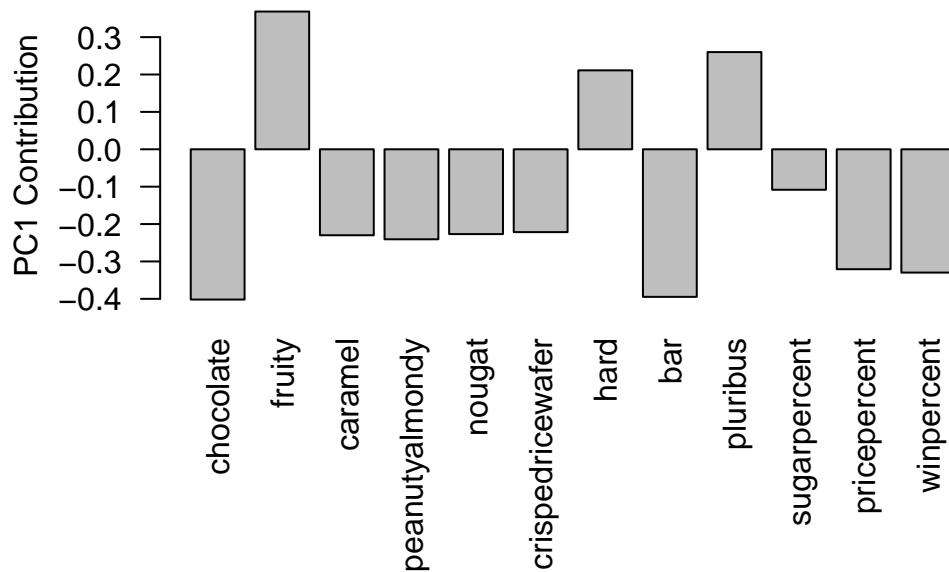
filter

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

layout

PCAlodings

```
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 strongly in the positive direction. Likely since they have lower win percentages it is easier to sell them in bulk compared to singularly.