

# Class 9: Halloween Candy Mini Project

Alexis Galano (PID: A17628362)

```
candy_file <- read.csv("candy-data.txt")
candy = read.csv("candy-data.txt", row.names=1)
head(candy)
```

	chocolate	fruity	caramel	peanutyalmondy	nougat	crispedrice	wafer
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

- A1. There are 85 different candy types.

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

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

[1] 38

- A2. There are 38 fruity candy types.

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

[1] 37

```
View(candy)
```

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

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

[1] 81.64291

- A3. My favorite candy is Twix with a winpercent of 81.64%.

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

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

[1] 76.7686

- A4. The winpercent for Kit Kat is 76.77%

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

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

[1] 49.6535

- A5. The winpercent for Tootsie Roll Snack Bars is 49.65%

Q. What is the least liked candy in the dataset? - lowest winpercent

```
min(candy$winpercent)
```

[1] 22.44534

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

- A. The least liked candy is Nik L Nip with a 22.44% winpercent.

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

```
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?

- A6. There is a histogram column and an n\_missing column that look to be different than the other columns.

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

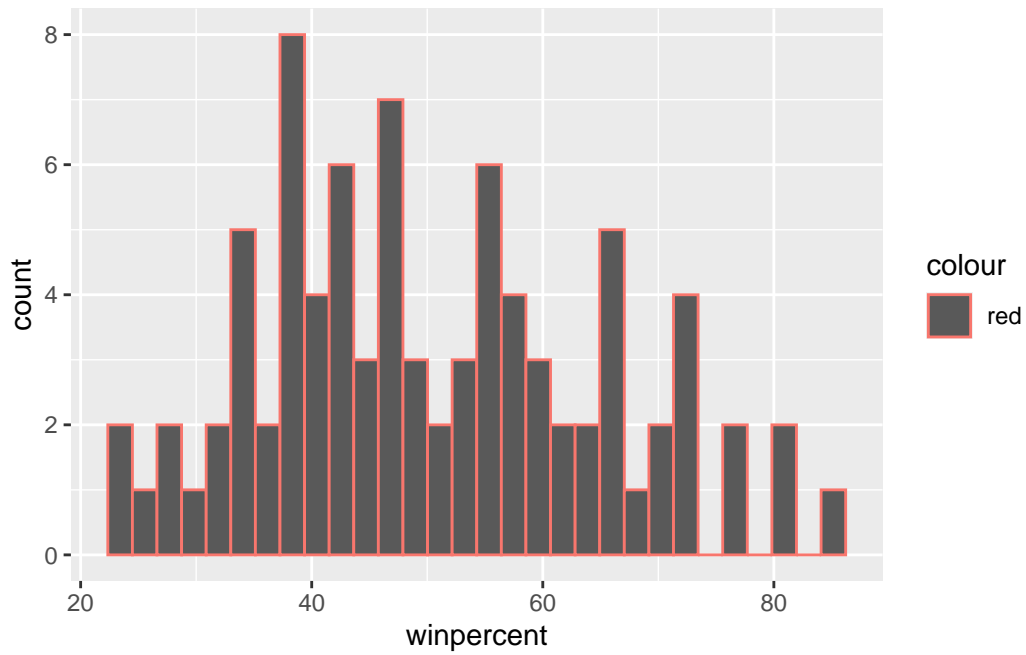
- A7. 0 in the n\_missing column for candy\$chocolate means there are no missing numeric values for the chocolate data. The 1 in the complete\_rate column means that all the values in the candy data frame are relevant and included in the statistics.

Q8. Plot a histogram of winpercent values

```
library(ggplot2)

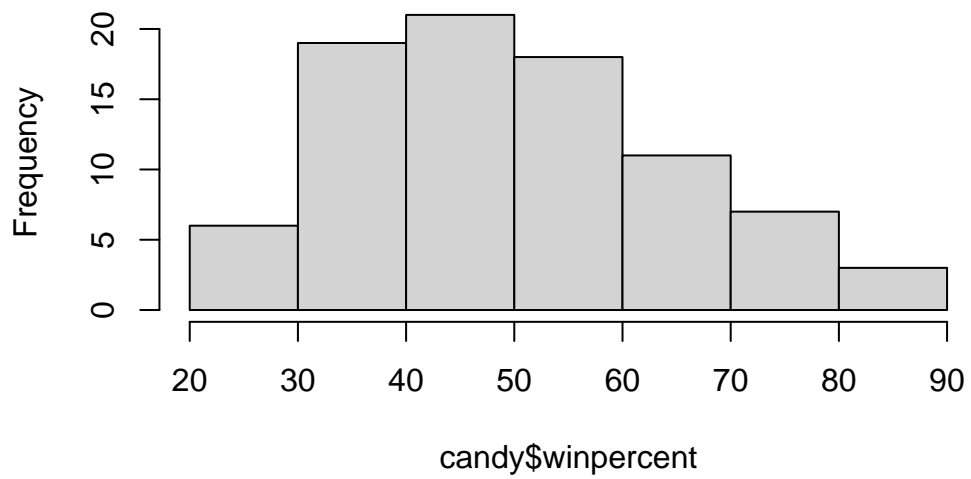
ggplot(candy) +
  aes(winpercent, col="red") +
  geom_histogram()
```

`stat\_bin()` using `bins = 30`. Pick better value with `binwidth`.



```
hist(candy$winpercent, breaks=8)
```

**Histogram of candy\$winpercent**



Q9. Is the distribution of winpercent values symmetrical?

- A9. No the distribution of winpercent values is not symmetrical.

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

- A10. The center of distribution is below 50%.

Q11. On average is chocolate candy higher or lower ranked than fruit candy? First find all chocolate candy and their \$winpercent values. Next summarize these values into one number. Then do the same for fruit candy and compare the numbers.

```
choc.inds <- as.logical(candy$chocolate)
choc.win <- candy[choc.inds,]$winpercent
mean(choc.win)
```

```
[1] 60.92153
```

```
fruit.inds <- as.logical(candy$fruity)
fruit.win <- candy[fruit.inds,]$winpercent
mean(fruit.win)
```

```
[1] 44.11974
```

- A11. Chocolate candy is higher ranked on average compared to fruit candy with a difference of about 16%.

Q12. Is this difference statistically significant?

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

Welch Two Sample t-test

```
data:  choc.win and fruit.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
```

- A12. The difference in the means is not equal to 0 and so they are statistically significant. The p-value is 2.871e-08

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

- A13. The bottom 5 are Nik L Nip, Boston Baked Beans, Chiclets, Super Bubble, Jawbusters, Rootbeer barrels.

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

```
inds <- order(candy$winpercent)
tail(candy[inds,], n=5)
```

	chocolate	fruity	caramel	peanut	almond	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

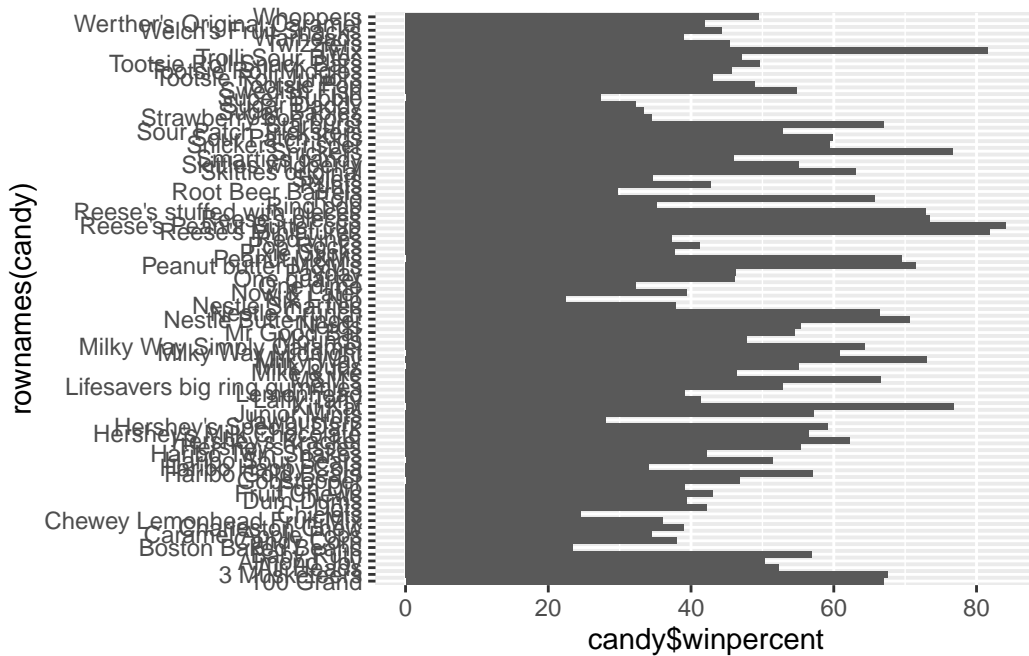
Kit Kat	1	0	1	0	0.313
Twix	1	0	1	0	0.546
Reese's Miniatures	0	0	0	0	0.034
Reese's Peanut Butter cup	0	0	0	0	0.720

	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

- The top 5 are Snickers, Kit Kat, Twix, Reese's Miniatures, Reese's Peanut Butter Cup.

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

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



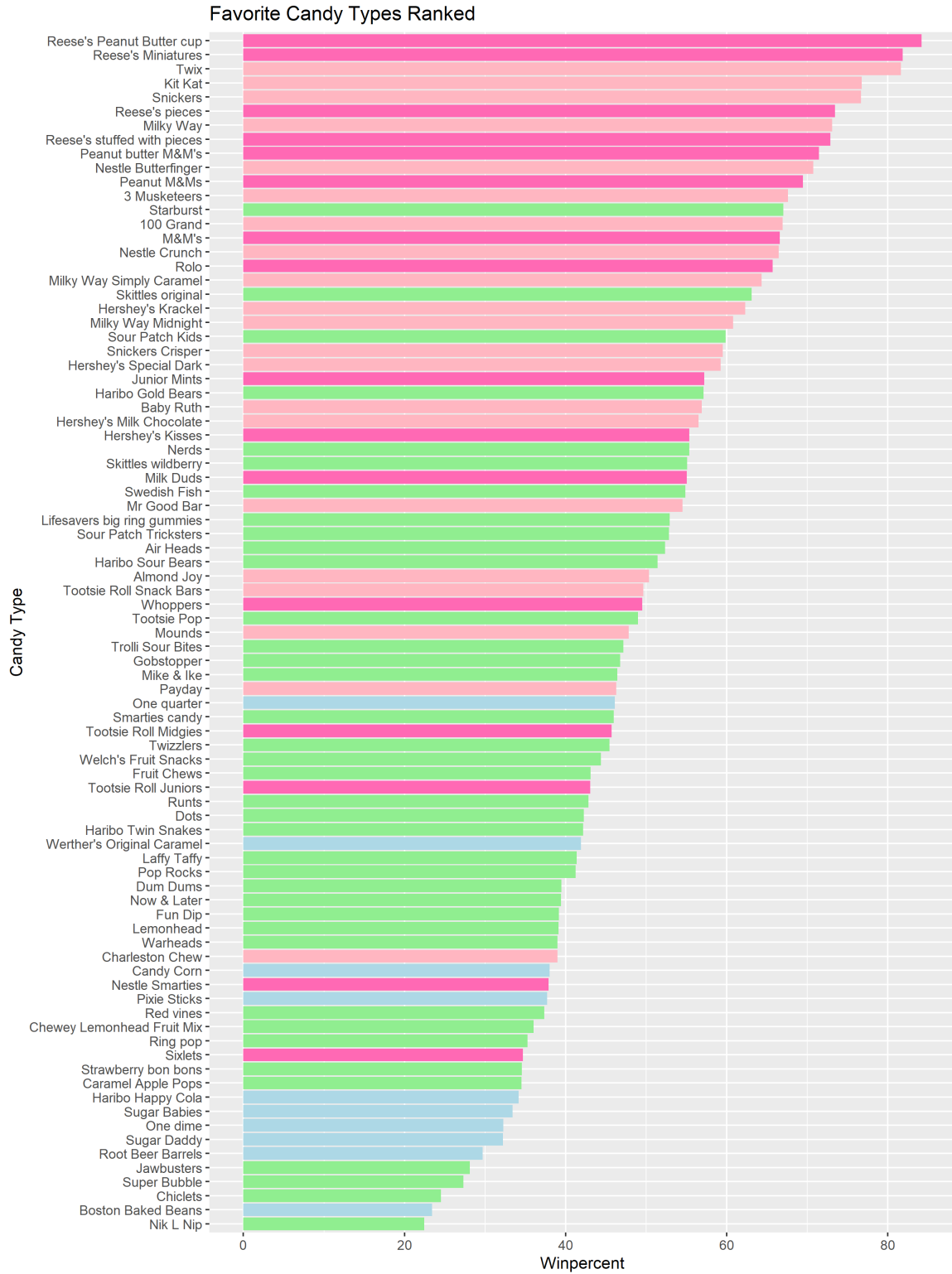
Q16. This is quite ugly, use the reorder() function to get the bars sorted by winpercent?

Add some color to our ggplot. We need to make a custom color vector.



```
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy), winpercent)) +
  geom_col(fill=my_cols) +
  labs( x='Winpercent',
        y='Candy Type',
        title='Favorite Candy Types Ranked')
```





> Q17. What is the worst ranked chocolate candy?

- A17. Sixlets

Q18. What is the best ranked fruity candy?

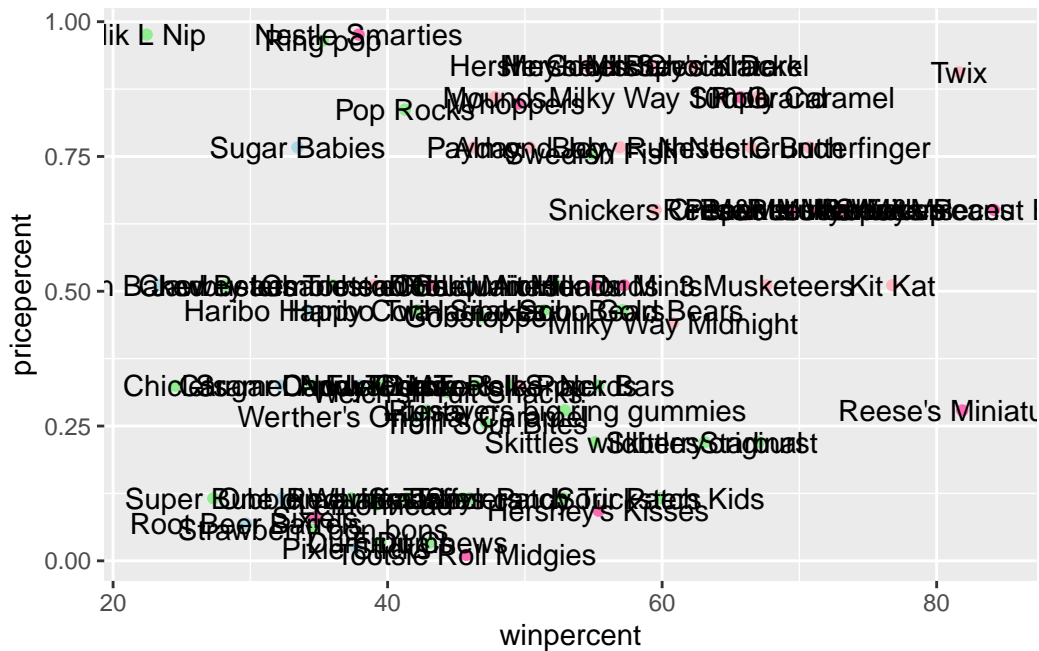
- A18. Starburst

```
candy$pricepercent
```

```
[1] 0.860 0.511 0.116 0.511 0.511 0.767 0.767 0.511 0.325 0.325 0.511 0.511
[13] 0.325 0.511 0.034 0.034 0.325 0.453 0.465 0.465 0.465 0.465 0.093 0.918
[25] 0.918 0.918 0.511 0.511 0.511 0.116 0.104 0.279 0.651 0.651 0.325 0.511
[37] 0.651 0.441 0.860 0.860 0.918 0.325 0.767 0.767 0.976 0.325 0.767 0.651
[49] 0.023 0.837 0.116 0.279 0.651 0.651 0.651 0.965 0.860 0.069 0.279 0.081
[61] 0.220 0.220 0.976 0.116 0.651 0.651 0.116 0.116 0.220 0.058 0.767 0.325
[73] 0.116 0.755 0.325 0.511 0.011 0.325 0.255 0.906 0.116 0.116 0.313 0.267
[85] 0.848
```

If we want to see what is good candy to buy in terms of winpercent and pricepercent we can plot these two variables and then see the best candy for the least amount of money.

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



To avoid the overplotting of these labels we can use an add on package called ggrepel

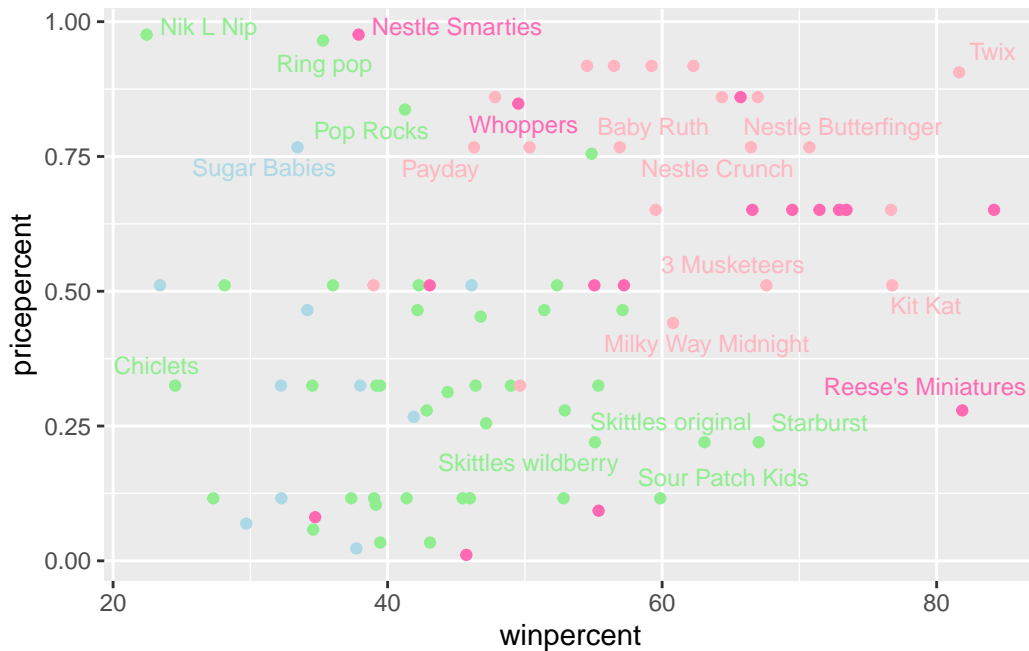
```
#install.packages("ggrepel")
```

```
library(ggrepel)
library(ggplot2)
```

Play with the 'max.overlaps' parameter to 'geom\_text\_repel()'

```
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(max.overlaps=5, size=3.3, col=my_cols)
```

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?

- A19. I would say Reese's Miniatures is the highest ranked for the lowest price.

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

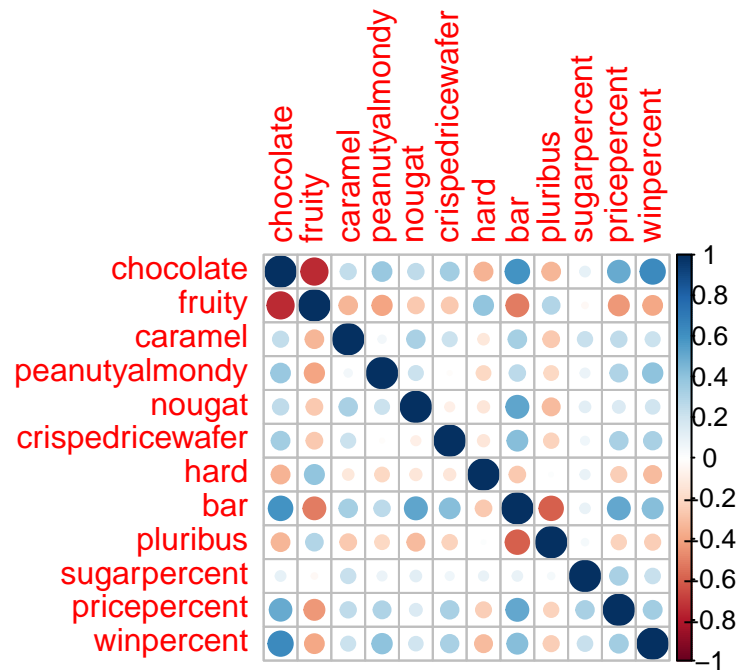
- A20. Nik L Nip, Ring pop, Nestle Smarties, Mr. Goodbar, and Hersheys are the most expensive. Of those Nik L Nip is the least popular.

## Exploring the correlation structure

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

- A22. Chocolate and Fruity are two variables that are anti-correlated because they are highly UNLIKELY to be combined together in a candy. Also it's unlikely that a fruity candy is in a bar form.

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

- A23. Chocolate and bar are two variables that are positively correlated meaning chocolate can typically be found in bar form.

## Principal Component Analysis

This function for this is called 'prcomp()' and here we know we need to scale our data with the 'scale=TRUE' argument

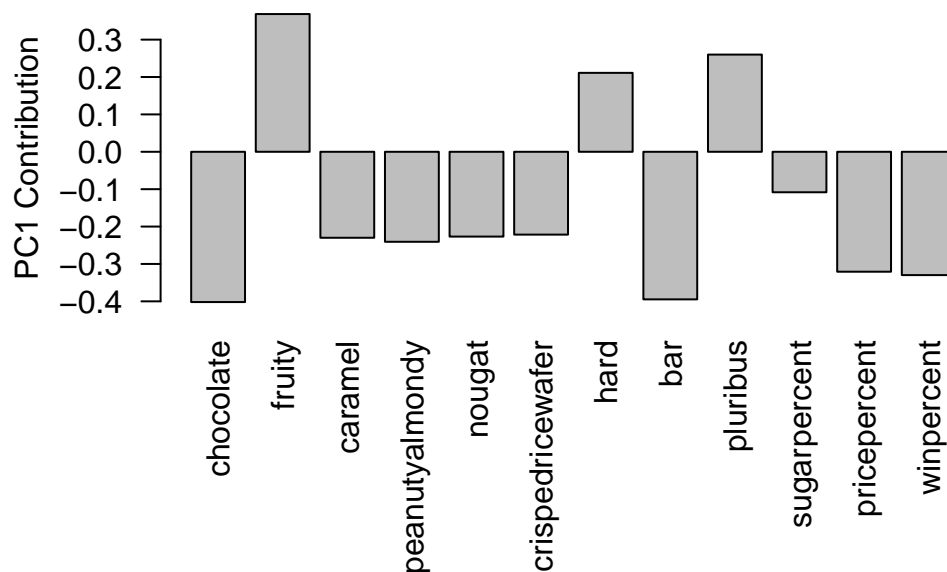
```
pca <- prcomp(candy, scale=TRUE)
summary(pca)
```

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		

```
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?

- A24. PC1 represents the correlation structure in the candy data. In the positive direction we have the variables fruity, hard, and pluribus because they are variables that can be categorized together due to their correlation. Example: Skittles is fruity, has a slightly hard outer shell, and there are multiple pieces in a single pack.

```
pca <- prcomp(candy, SCALE=FALSE)
```

Warning: In prcomp.default(candy, SCALE = FALSE) :  
extra argument 'SCALE' will be disregarded

```
summary(pca)
```

Importance of components:

	PC1	PC2	PC3	PC4	PC5	PC6	PC7
Standard deviation	14.7231	0.70241	0.47762	0.37292	0.34641	0.33614	0.30748
Proportion of Variance	0.9935	0.00226	0.00105	0.00064	0.00055	0.00052	0.00043
Cumulative Proportion	0.9935	0.99574	0.99678	0.99742	0.99797	0.99849	0.99892

	PC8	PC9	PC10	PC11	PC12
Standard deviation	0.27417	0.23826	0.21435	0.18434	0.15331
Proportion of Variance	0.00034	0.00026	0.00021	0.00016	0.00011
Cumulative Proportion	0.99927	0.99953	0.99974	0.99989	1.00000

```
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)
```

```
library(ggrepel)
```

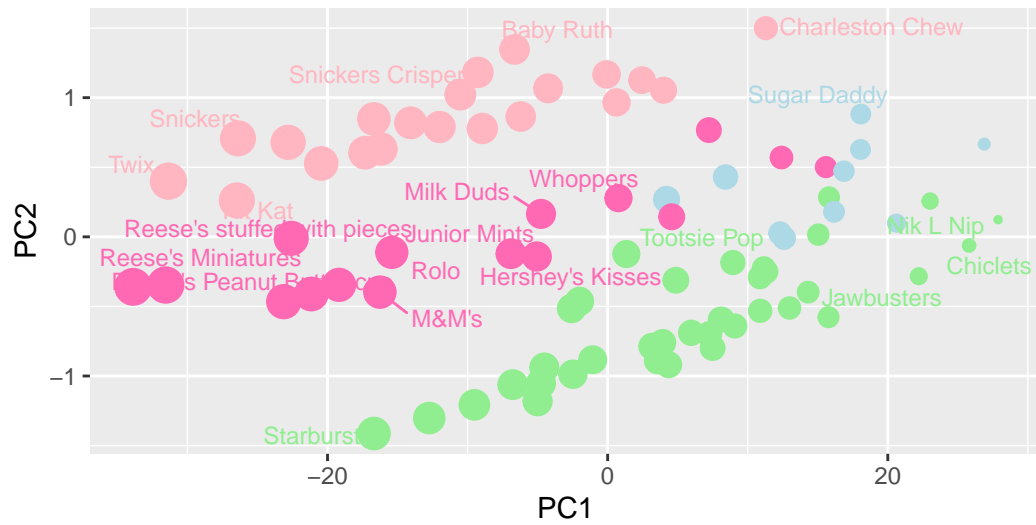
```
p + geom_text_repel(size=3.3, col=my_cols, max.overlaps = 7) +  
  theme(legend.position = "none") +  
  labs(title="Halloween Candy PCA Space",  
        subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brown)",  
        caption="Data from 538")
```

Warning: ggrepel: 64 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),



Data from 538

## LOADINGS PLOT

```
loadings <- as.data.frame(pca$rotation)

ggplot(loadings) +
  aes(PC1, reorder(rownames(loadings), PC1)) +
  geom_col()
```

