

Class 9: Candy Analysis Mini Project

Dani Baur (A16648266)

In today's class we will examine some data about candy from the 538 website.

1. Importing Candy Data

```
candy_file <- "candy-data.txt"
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?

```
nrow(candy)
```

```
[1] 85
```

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

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

```
[1] 38
```

2. What is your favorite candy?

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
```

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

Table 1: Data summary

Name	candy
Number of rows	85
Number of columns	12
<hr/>	
Column type frequency:	
numeric	12
<hr/>	

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?

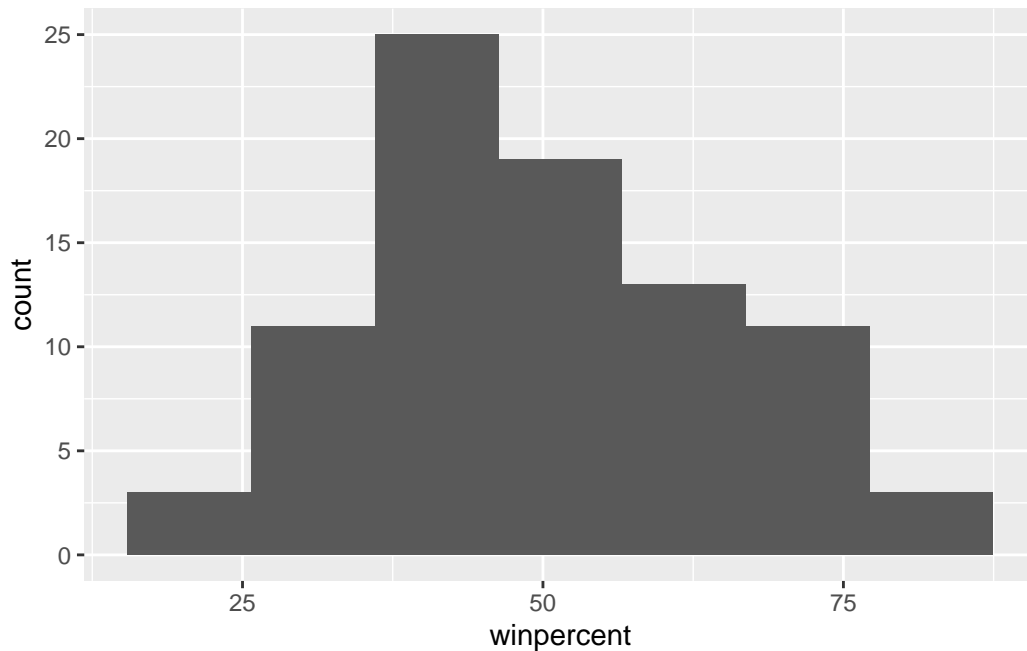
The winpercent is significantly larger than the other values such as sugarpercent and pricepercent.

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

A zero is similar to a “FALSE” logical which means that there is not any chocolate in the candy. A one is similar to a “TRUE” logical which means that there is chocolate in the candy.

Q8. Plot a histogram of winpercent values.

```
library(ggplot2)
ggplot(candy) +
  aes(winpercent) +
  geom_histogram(bins=7)
```



Q9. Is the distribution of winpercent values symmetrical?

No, the distribution is not symmetrical.

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

The center of the distribution is below 50%.

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

```
[1] 50.31676
```

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

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

```
[1] TRUE
```

On average, chocolate candy was ranked higher than fruit candy.

Q12. Is this different statistically significant?

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

Welch Two Sample t-test

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

The order function allows you to sort a dataset based on any value you want: `x[order(dataset)]`.

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

```
inds <- order(candy$winpercent)
head(candy[inds,], 5)
```

	chocolate	fruity	caramel	peanut	almondy	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 five all time favorite candy types out of this set?

```
tail(candy[inds,], 5)
```

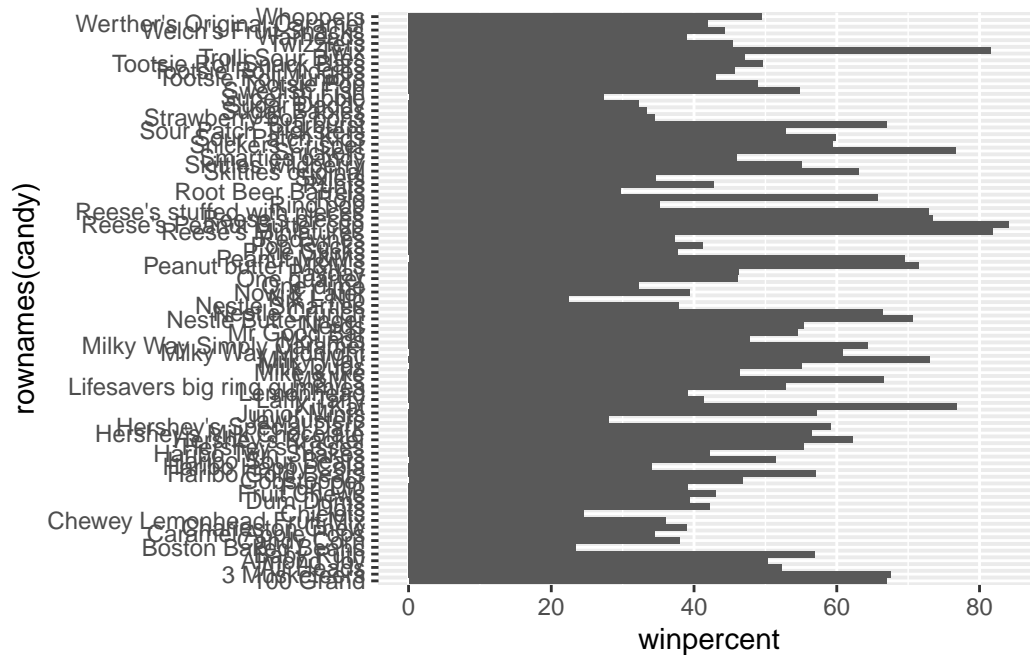
	chocolate	fruity	caramel	peanut	almond	nougat
Snickers	1	0	1		1	1
Kit Kat	1	0	0		0	0
Twix	1	0	1		0	0
Reese's Miniatures	1	0	0		1	0
Reese's Peanut Butter cup	1	0	0		1	0

	crisp	rice	wafers	hard	bar	pluribus	sugar
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
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	

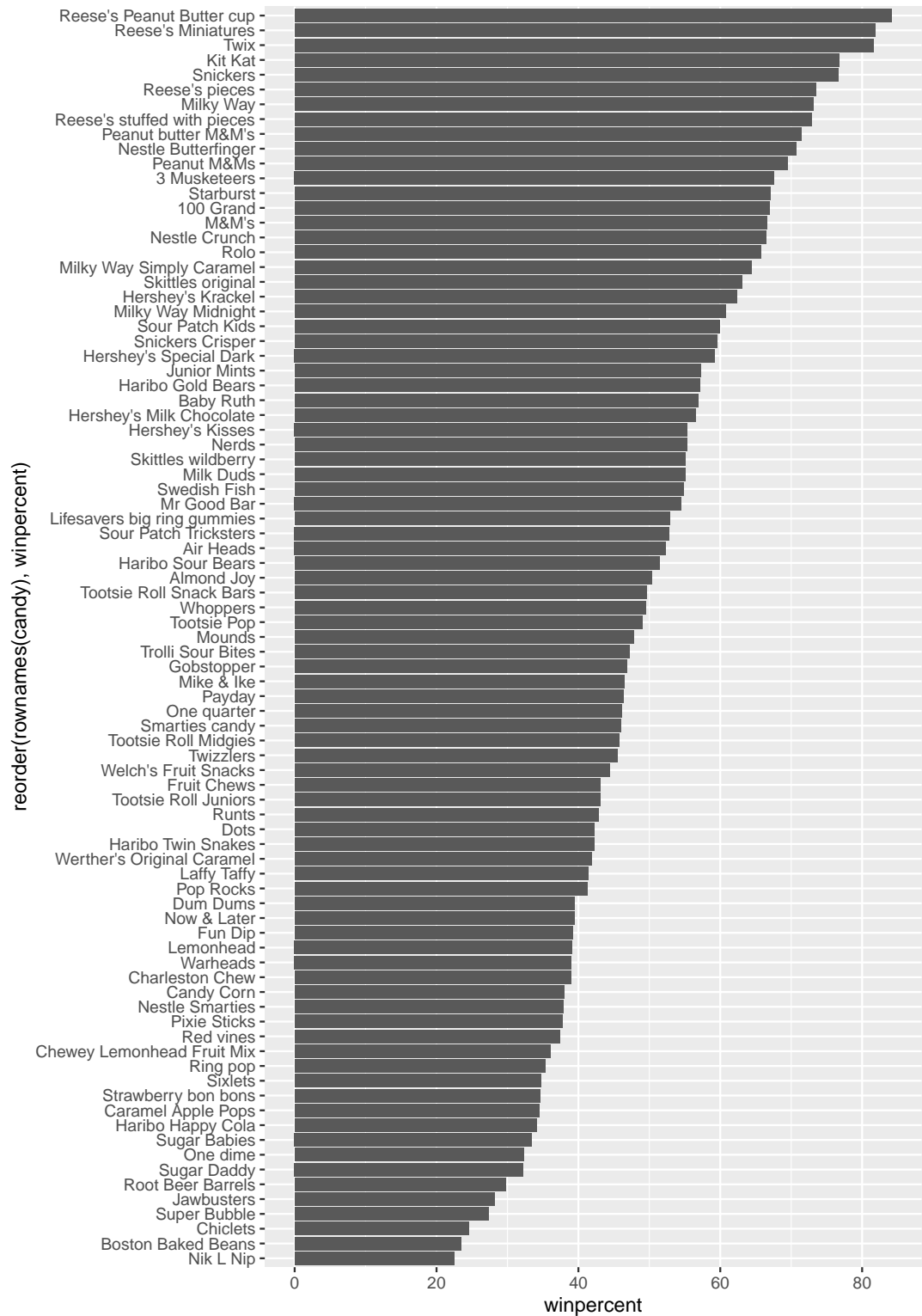
Q15. Make a first barplot of candy ranking based on winpercent values.

```
library(ggplot2)
ggplot(candy) +
  aes(winpercent, rownames(candy)) +
  geom_col()
```



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

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




```
ggsave("mybarplot.png", height=10)
```

Saving 5.5 x 10 in image

Add my custom colors to my barplot

```
my_cols=rep("pink3", nrow(candy))
my_cols[as.logical(candy$chocolate)] = "peru"
my_cols[as.logical(candy$bar)] = "wheat2"
my_cols[as.logical(candy$fruity)] = "palegreen3"
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy),winpercent)) +
  geom_col(fill=my_cols)
```

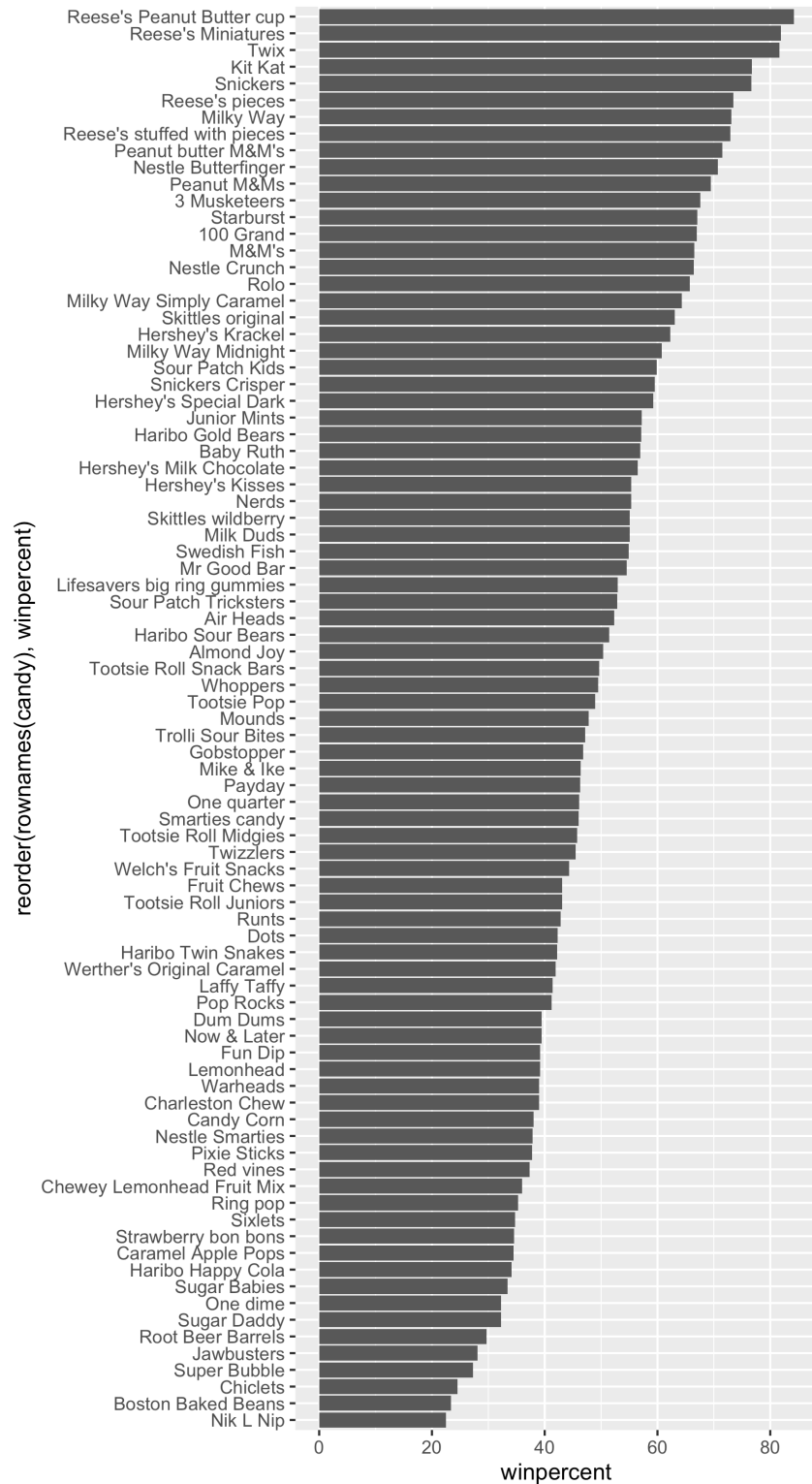
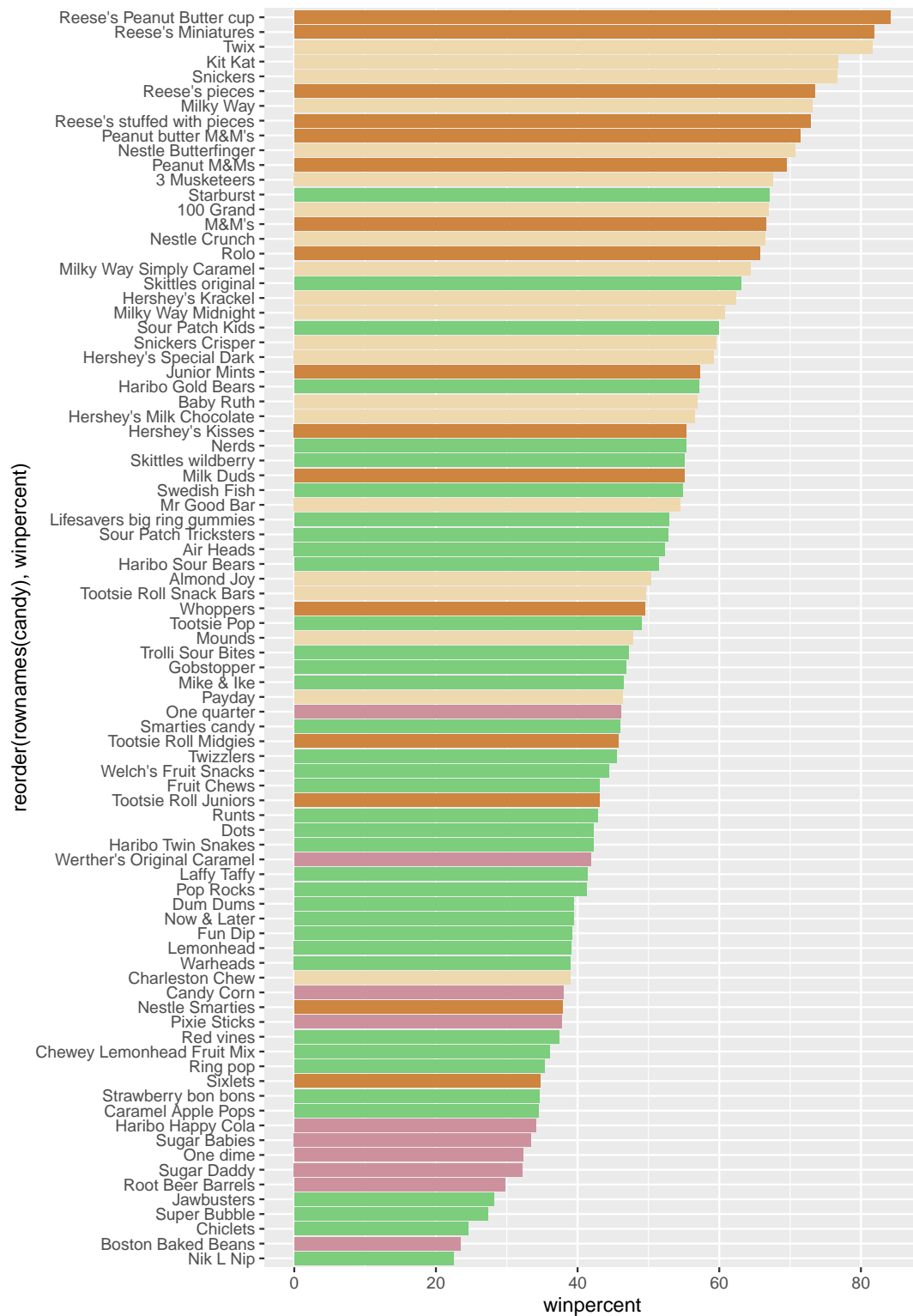


Figure 1: exported image that is a little bit bigger so I can read it!



Q17. What is the worst ranked chocolate candy?

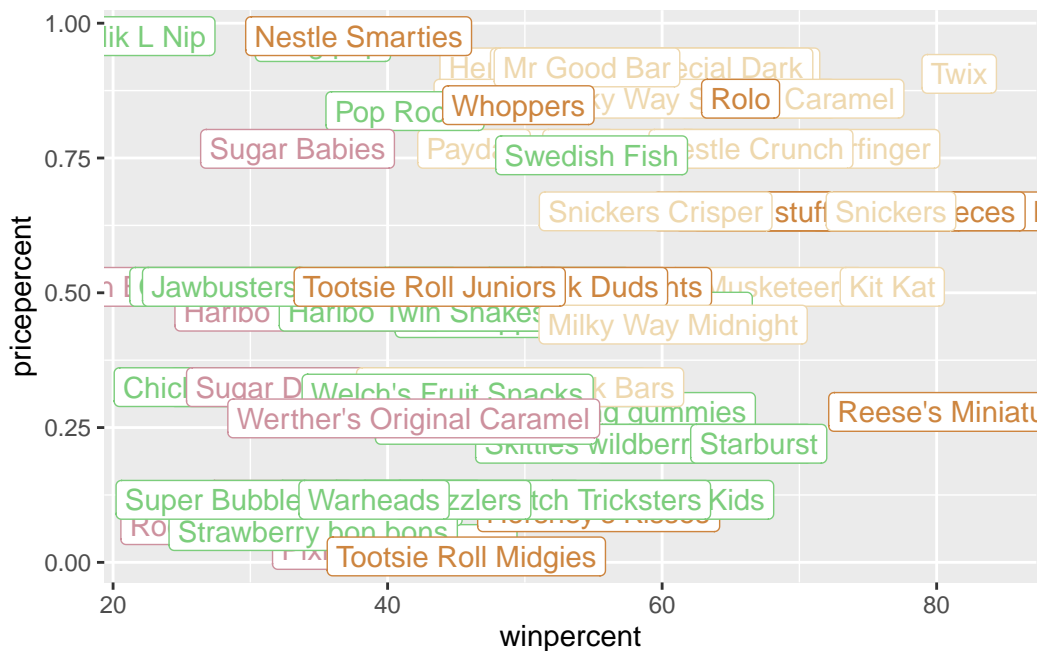
The worst ranked chocolate candy is Sixlets.

Q18. What is the best ranked fruity candy?

The best ranked fruity candy is Starburst.

4. Taking a Look at Pricepercent

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

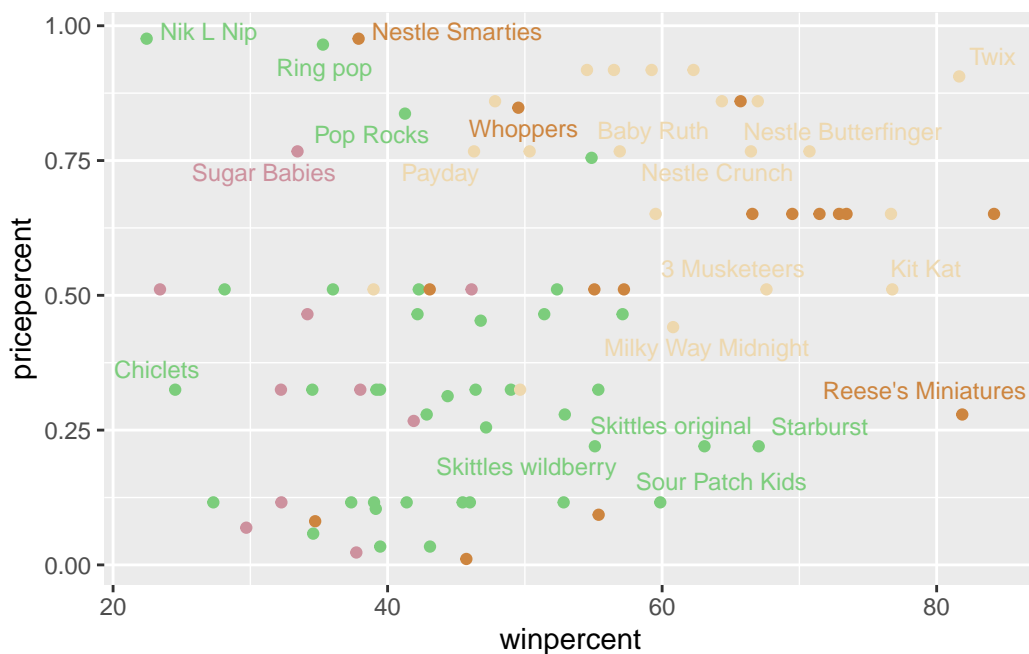


There are just too many labels in this above plot to be readable. Use `ggrepel` to make it more visually appealing.

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

```
geom_text_repel(col=my_cols, max.overlaps=5, size=3.3)
```

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?

Fruity candy has a good winpercent with a lower pricepercent.

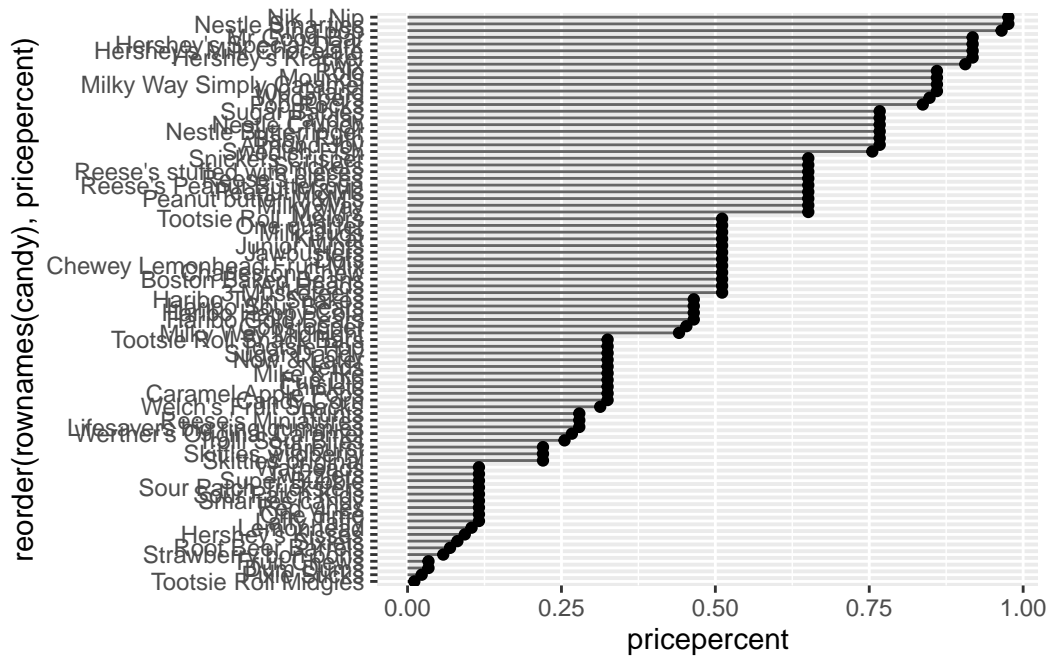
Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?

```
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. 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()`.

```
# 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()
```

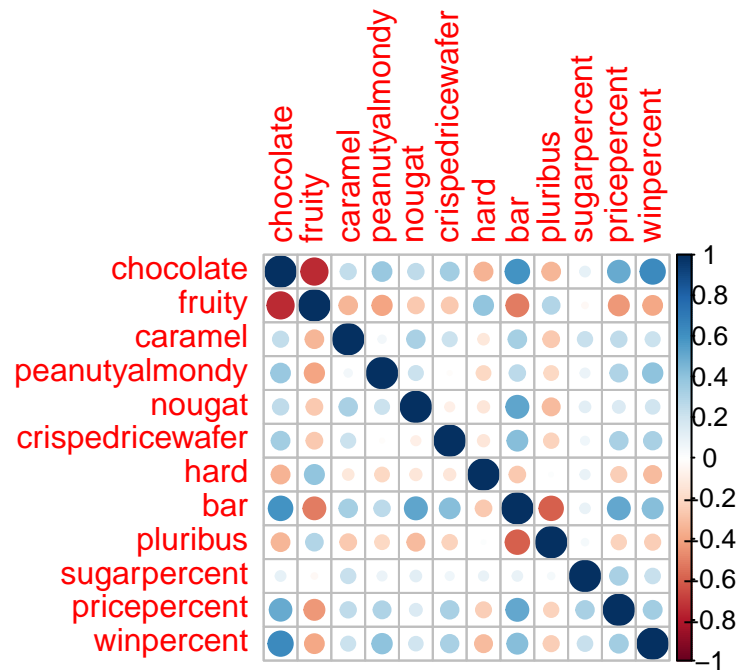


5. Exploring the Correlation Structure

```
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)?

Chocolate and fruit are not correlated to each other.

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

Chocolate is most likely a bar and expensive.

6. Principal Component Analysis

We will perform a PCA of the candy. Key question: do we need to scale the data before PCA?

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

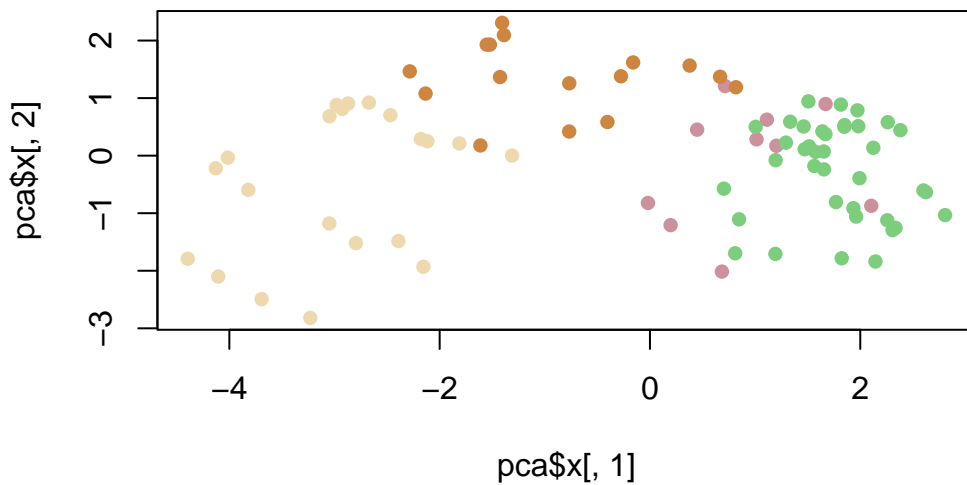
Importance of components:

	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

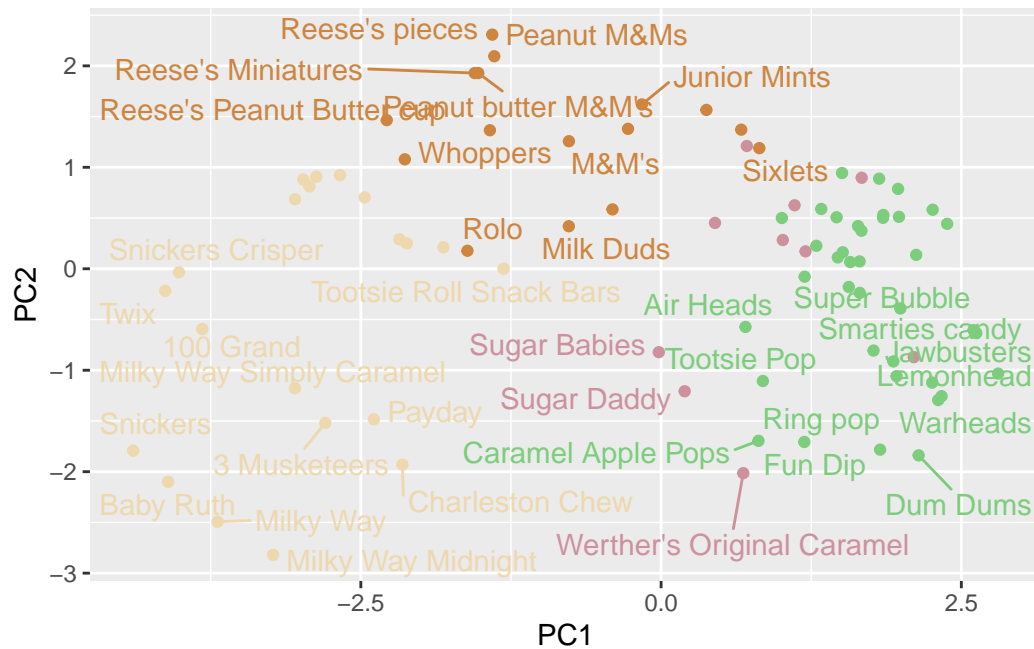
```
plot(pca$x[,1], pca$x[,2], col=my_cols, pch=16)
```



We can make a much nicer ggplot!

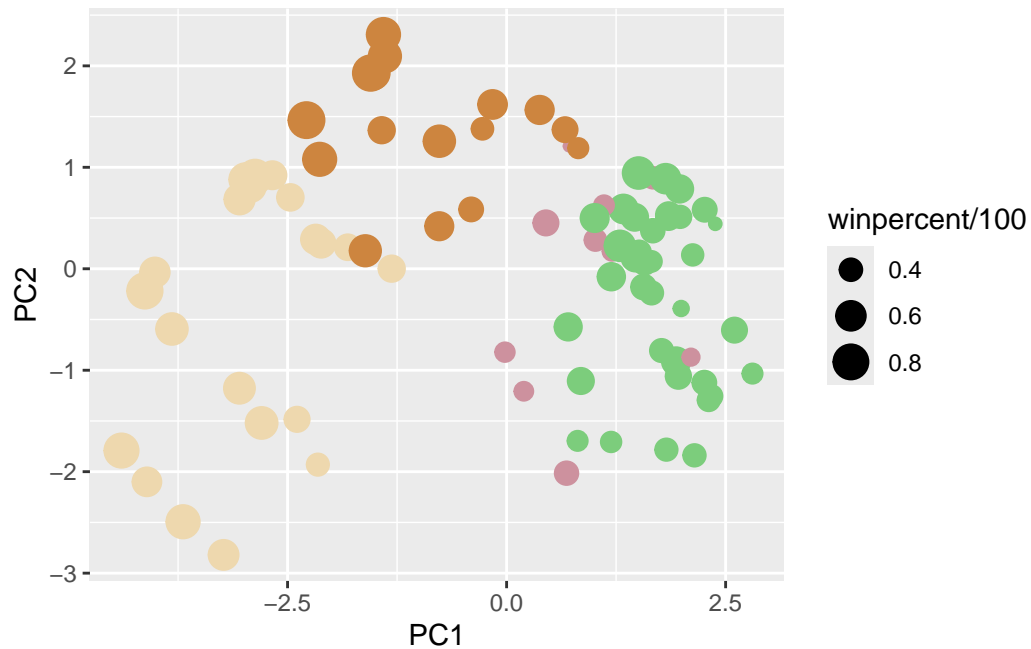
```
#Make a new data-frame with our PCA results and candy data
my_data <- cbind(candy, pca$x[,1:3])
ggplot(my_data) +
  aes(x=PC1, y=PC2, label=rownames(candy)) +
  geom_point(col=my_cols) +
  geom_text_repel(col=my_cols)
```

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



Make this a bit nicer:

```
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
```



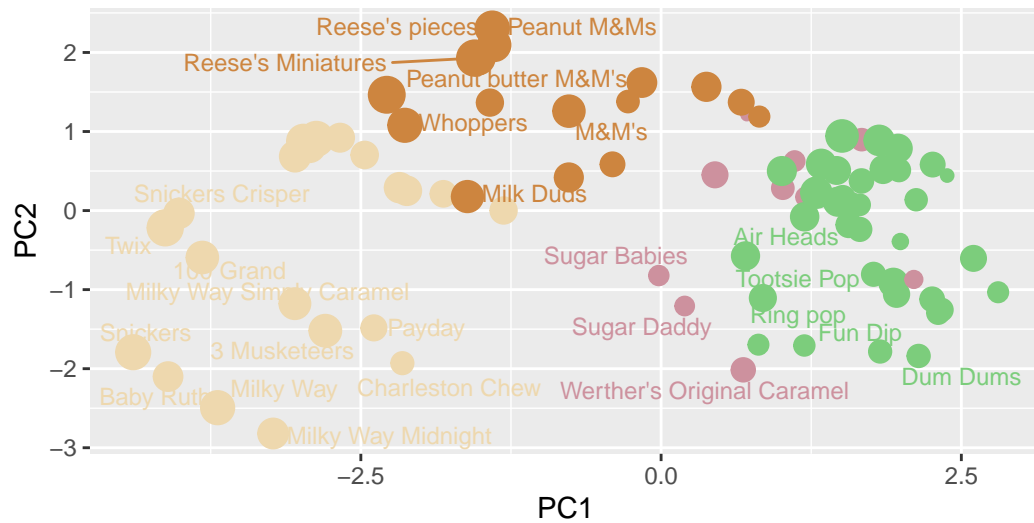
```
library(ggrepel)

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: 59 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),



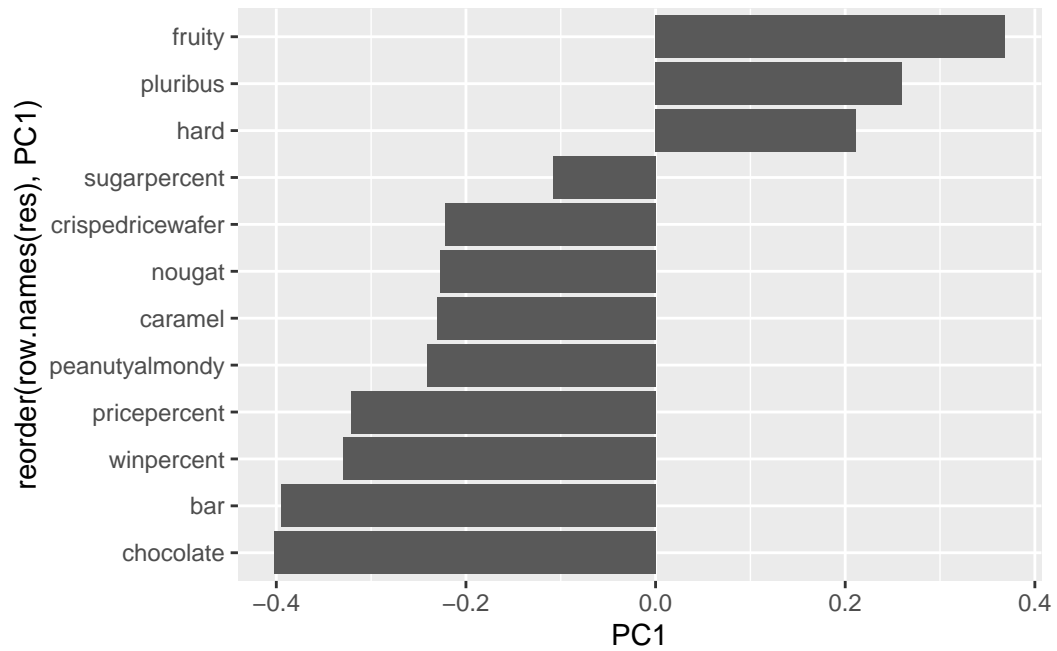
Data from 538

```
#library(plotly)
#ggplotly(p)
```

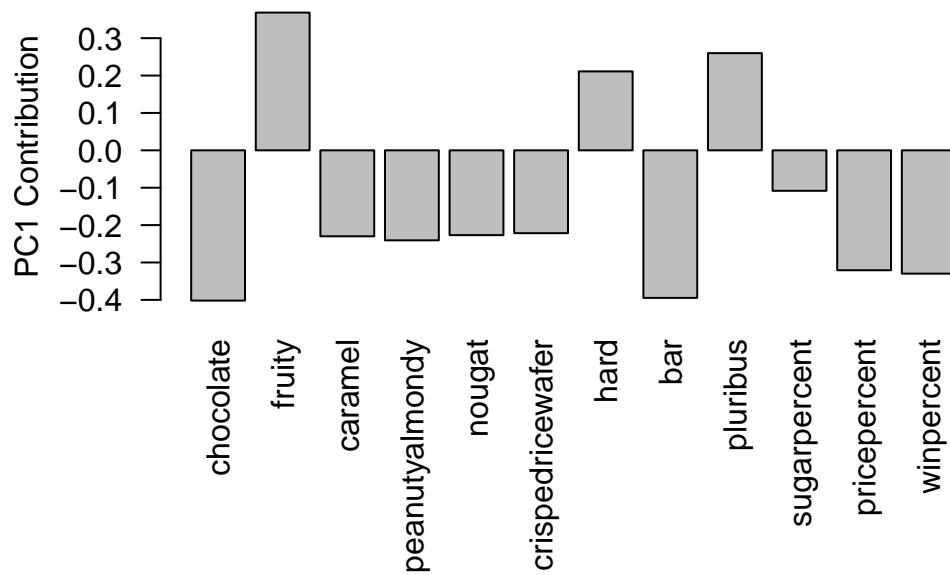
How do the original variables contribute to our PCs? For this we look at the loading component of our results object- i.e. the `pca$rotation` object.

Make a barplot with ggplot and order the bars by their value. Recall that you need a data.frame as input for ggplot.

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



```
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?

Fruit, pluribus and hard are all picked up in the positive direction and these do make sense based on the correlation structure in the dataset. If you are a fruity candy then you will tend to be hard and come in a packet with multiple candies in it (pluribus).