

# Lab 9- Halloween Candy Mini Project

Trevor Hoang (A16371830)

## Table of contents

Importing Candy Data . . . . .	1
Favorite Candy . . . . .	2
Overall Candy Rankings . . . . .	7
Time to add some color . . . . .	11
Taking a look at price percent . . . . .	13
Exploring the correlation structure . . . . .	17
Principal Component Analysis . . . . .	19

Today we will examine data from 538 on common Halloween candy. In particular we will use ggplot, dplyr and PCA to make sense of this multivariate dataset.

## Importing Candy Data

```
candy = read.csv("candy-data.csv", 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
```

How many chocolate candy are there in this dataset?

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

```
[1] 37
```

## Favorite Candy

Finding the favorite candy by percentage using the winpercent statistic

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

```
candy["Nerds", "winpercent"]
```

```
[1] 55.35405
```

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
```

Using 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

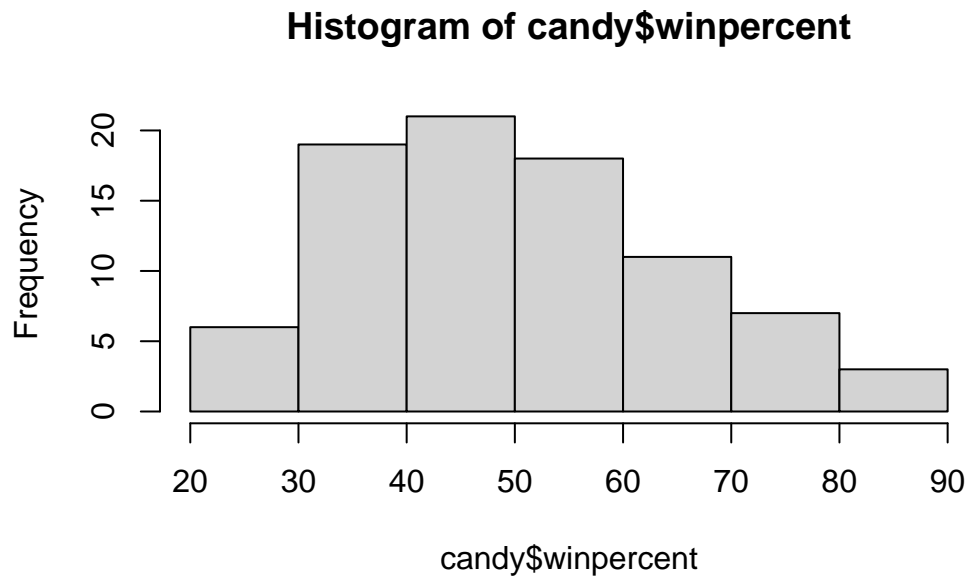
skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
chocolate	0	1	0.44	0.50	0.00	0.00	0.00	1.00	1.00	
fruity	0	1	0.45	0.50	0.00	0.00	0.00	1.00	1.00	
caramel	0	1	0.16	0.37	0.00	0.00	0.00	0.00	1.00	
peanutyalmondy	0	1	0.16	0.37	0.00	0.00	0.00	0.00	1.00	
nougat	0	1	0.08	0.28	0.00	0.00	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? It looks like the **winpercent** column is on a different scale than the others (0-100% rather than 0-1). I will need to scale this dataset before analysis like PCA

Q7. What do you think a zero and one represent for the candy\$chocolate column?  
I think that a 1 represents the candy being considered a chocolate while a 0 means that it is not considered a chocolate.

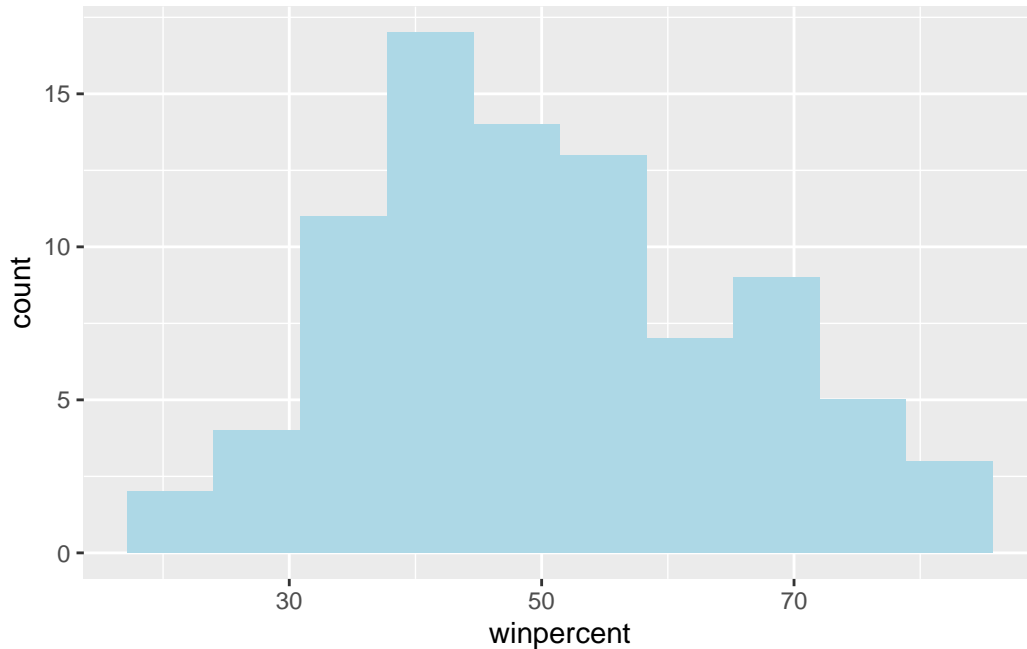
Q8. Plot a histogram of winpercent values

```
hist(candy$winpercent)
```



```
library(ggplot2)

ggplot(candy)+
  aes(winpercent)+
  geom_histogram(bins=10, fill="lightblue")
```



Q9. Is the distribution of winpercent values symmetrical? No it looks like the histogram is skewed left

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

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

- Step 1: Find all “chocolate” candy
- Step 2: Find their “winpercent” values
- Step 3: Summarize these values
- Step 4: Find all “Fruity” candy
- Step 5: Find their “winpercent” values
- Step 6: Summarize these values
- Step 7: Compare the 2 summary values

Step 1: Find all chocolate candy

```
choc.inds = candy$chocolate == 1
```

Step 2: Find their winpercent values

```
choc.win = candy[choc.inds,]$winpercent
```

Step 3: Summarize these values

```
choc.mean = mean(choc.win)
```

Step 4: Find all “Fruity” candy

```
fruity.inds = candy$fruity == 1
```

Step 5: Find their “winpercent” values

```
fruity.wins = candy[fruity.inds,]$winpercent
```

Step 6: Summarize these values

```
fruit.mean = mean(fruity.wins)
```

Step 7: Compare the 2 summary values Chocolate has the higher mean

```
choc.mean
```

```
[1] 60.92153
```

```
fruit.mean
```

```
[1] 44.11974
```

Q12. Is this difference statistically significant?

```
t.test(choc.win, fruity.wins)
```

Welch Two Sample t-test

data: choc.win and fruity.wins

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

## Overall Candy Rankings

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

```
#Sort is not that useful - it just sorts values  
sort(candy$winpercent)
```

```
[1] 22.44534 23.41782 24.52499 27.30386 28.12744 29.70369 32.23100 32.26109  
[9] 33.43755 34.15896 34.51768 34.57899 34.72200 35.29076 36.01763 37.34852  
[17] 37.72234 37.88719 38.01096 38.97504 39.01190 39.14106 39.18550 39.44680  
[25] 39.46056 41.26551 41.38956 41.90431 42.17877 42.27208 42.84914 43.06890  
[33] 43.08892 44.37552 45.46628 45.73675 45.99583 46.11650 46.29660 46.41172  
[41] 46.78335 47.17323 47.82975 48.98265 49.52411 49.65350 50.34755 51.41243  
[49] 52.34146 52.82595 52.91139 54.52645 54.86111 55.06407 55.10370 55.35405  
[57] 55.37545 56.49050 56.91455 57.11974 57.21925 59.23612 59.52925 59.86400  
[65] 60.80070 62.28448 63.08514 64.35334 65.71629 66.47068 66.57458 66.97173  
[73] 67.03763 67.60294 69.48379 70.73564 71.46505 72.88790 73.09956 73.43499  
[81] 76.67378 76.76860 81.64291 81.86626 84.18029
```

```
x = c(10, 1, 100)  
order(x)
```

```
[1] 2 1 3
```

```
x[order(x)]
```

```
[1] 1 10 100
```

The `order()` function tells us how to arrange the elements of the input to make them sortable - i.e. how to order them

We can determine the order to winpercent to make them sorted and use that order to arrange the whole dataset.

```
ord.inds = order(candy$winpercent)  
head(candy[ord.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

	win	percent
Nik L Nip	22.44	534
Boston Baked Beans	23.41	782
Chiclets	24.52	499
Super Bubble	27.30	386
Jawbusters	28.12	744
Root Beer Barrels	29.70	369

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

```
tail(candy[ord.inds,])
```

	chocolate	fruity	caramel	peanut	almond	nougat
Reese's pieces	1	0	0		1	0
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
Reese's pieces			0	0	0	1		0.406
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	percent
Reese's pieces	0.651		73.43	499



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

```
ord.inds = order(candy$winpercent, decreasing =T)
head(candy[ord.inds,])
```

	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
Reese's pieces	1	0	0		1	0

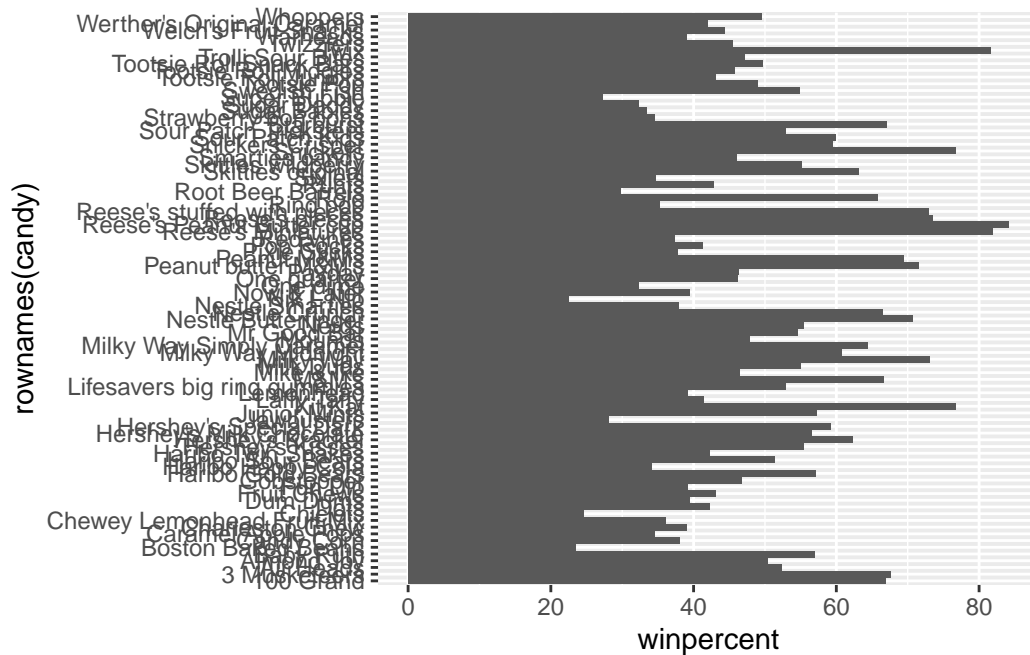
	crisped	rice	wafer	hard	bar	pluribus	sugar	percent
Reese's Peanut Butter cup		0	0	0		0		0.720
Reese's Miniatures		0	0	0		0		0.034
Twix		1	0	1		0		0.546
Kit Kat		1	0	1		0		0.313
Snickers		0	0	1		0		0.546
Reese's pieces		0	0	0		1		0.406

	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
Reese's pieces	0.651		73.43499

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

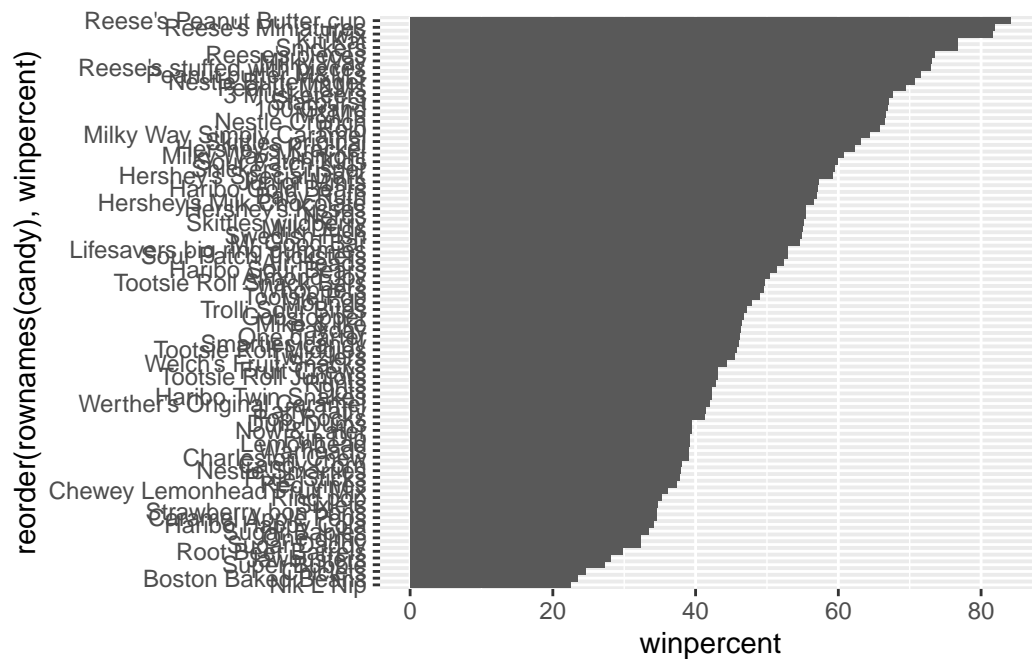
```
ggplot(candy) +
  aes(winpercent, rownames(candy)) +
  geom_col()
```



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

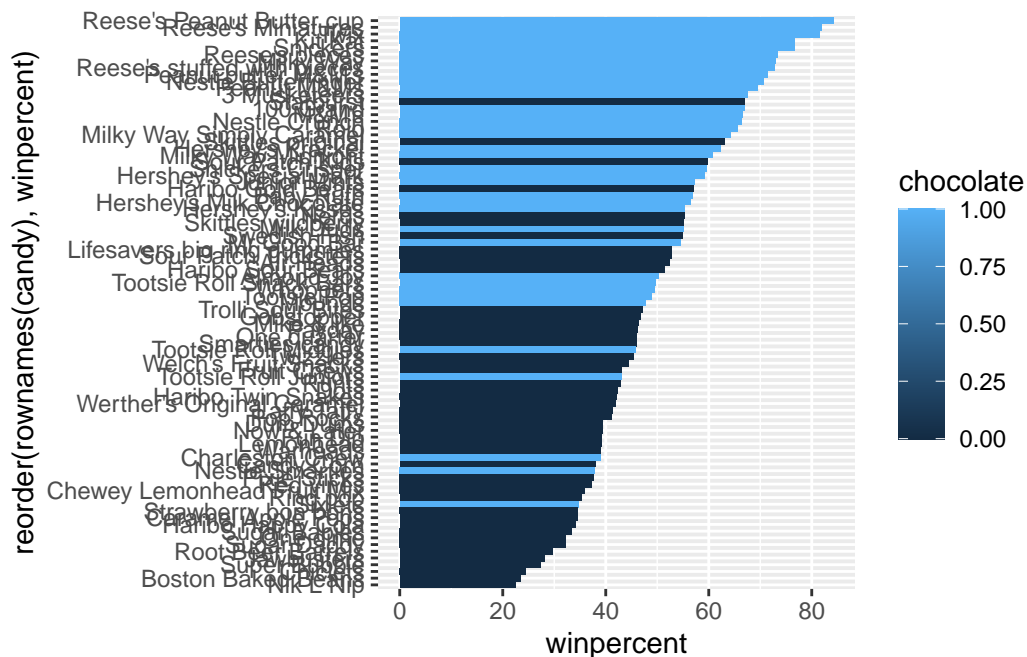
Let's rearrange these

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



**Time to add some color**

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



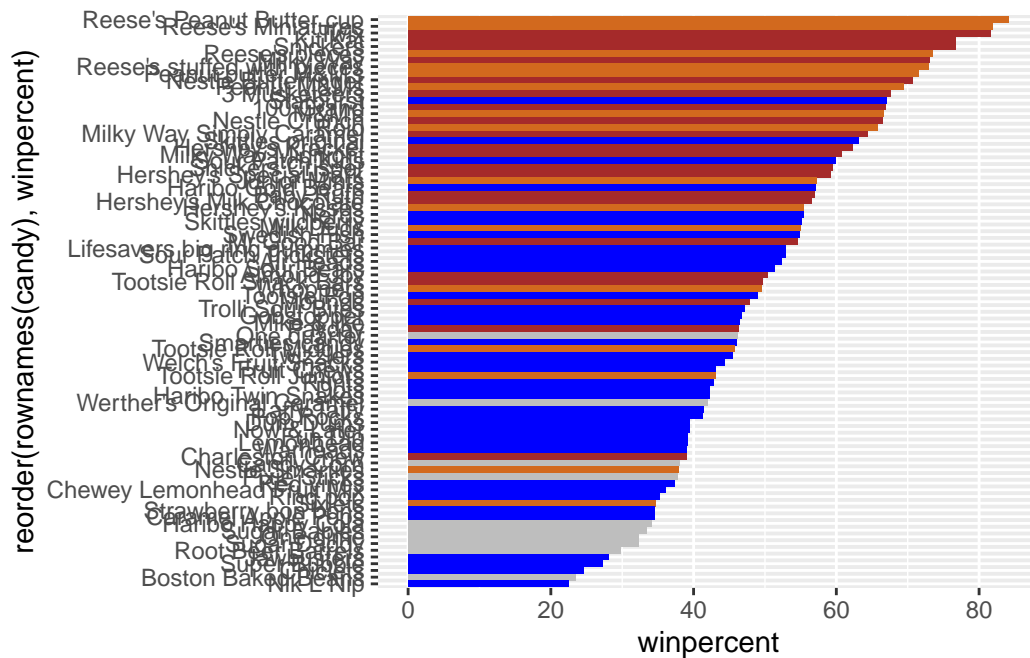
Need to make our own separate color vector where we can spell out exactly what candy is colored a particular color.

```
mycols = rep("gray", nrow(candy))
mycols[candy$chocolate == 1] = "chocolate"
mycols[candy$bar == 1] = "brown"
mycols[candy$fruity == 1] = "blue"
mycols
```

```
[1] "brown"    "brown"    "gray"     "gray"     "blue"     "brown"
[7] "brown"    "gray"     "gray"     "blue"     "brown"    "blue"
[13] "blue"     "blue"     "blue"     "blue"     "blue"     "blue"
[19] "blue"     "gray"     "blue"     "blue"     "chocolate" "brown"
[25] "brown"    "brown"    "blue"     "chocolate" "brown"    "blue"
[31] "blue"     "blue"     "chocolate" "chocolate" "blue"     "chocolate"
[37] "brown"    "brown"    "brown"    "brown"    "brown"    "blue"
[43] "brown"    "brown"    "blue"     "blue"     "brown"    "chocolate"
[49] "gray"     "blue"     "blue"     "chocolate" "chocolate" "chocolate"
[55] "chocolate" "blue"     "chocolate" "gray"     "blue"     "chocolate"
[61] "blue"     "blue"     "chocolate" "blue"     "brown"    "brown"
[67] "blue"     "blue"     "blue"     "blue"     "gray"     "gray"
[73] "blue"     "blue"     "blue"     "chocolate" "chocolate" "brown"
[79] "blue"     "brown"    "blue"     "blue"     "blue"     "gray"
```

```
[85] "chocolate"
```

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



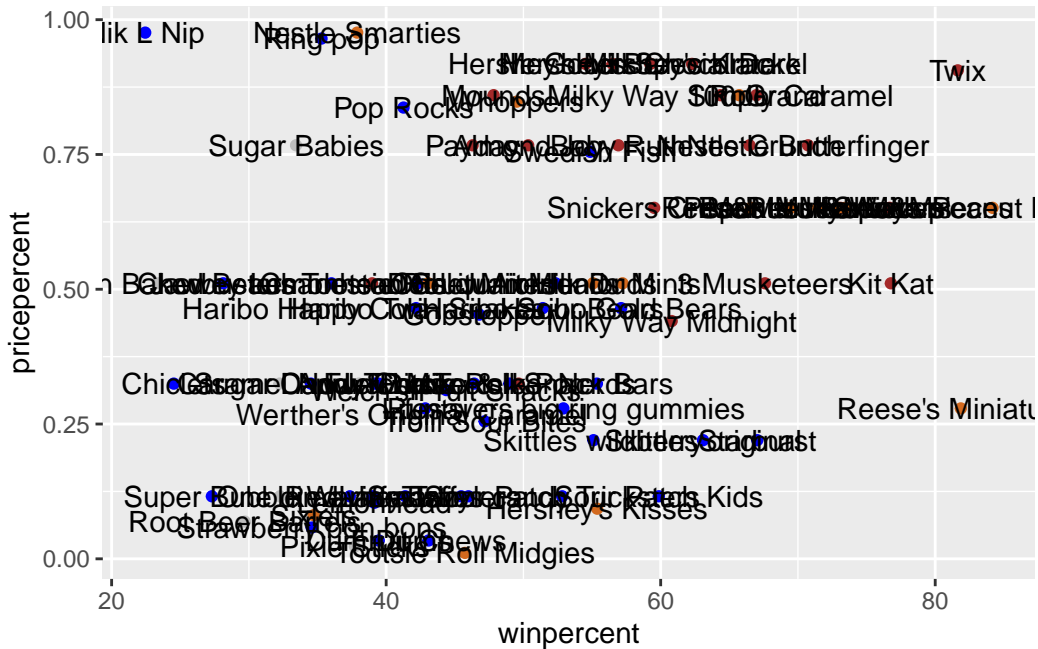
Q17. What is the worst ranked chocolate candy? The worst ranked chocolate candy is sixlet

Q18. What is the best ranked fruity candy? The best ranked fruity candy is starburts

## Taking a look at price percent

Make a plot of winpercent (x-axis) vs pricepercent (y-axis)

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

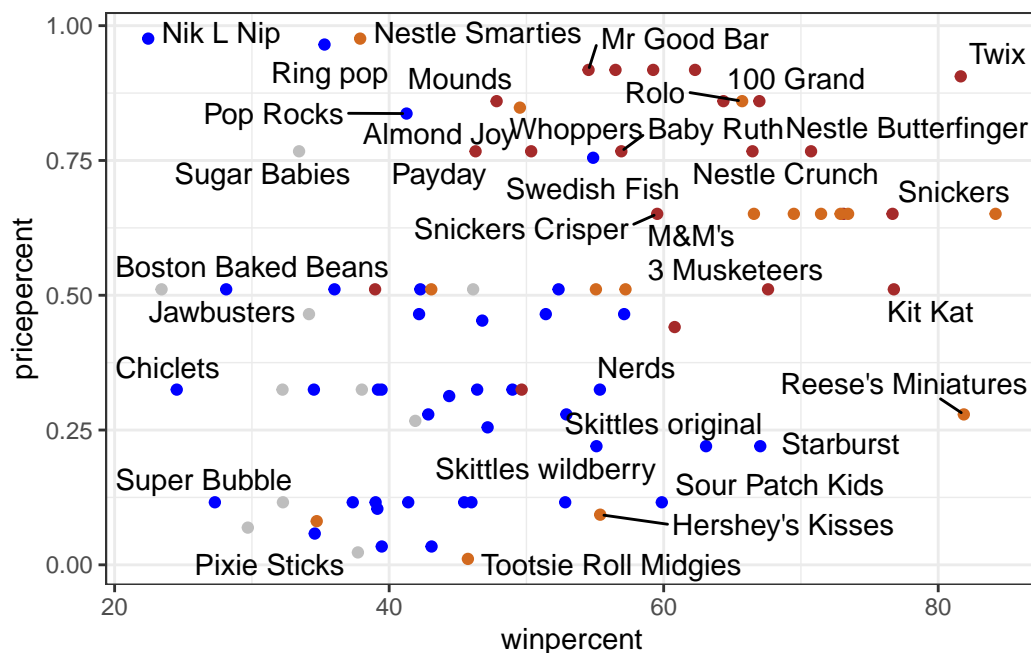


To avoid the overplotting of the text labels we can use the add on package **ggrepel**

```
library(ggrepel)

ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=mycols) +
  geom_text_repel(max.overlaps=10) +
  theme_bw()
```

Warning: ggrepel: 50 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? The highest ranked candy with for the least money is reese's miniatures

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

```
ord.inds = order(candy$pricepercent, decreasing =T)
head(candy[ord.inds,])
```

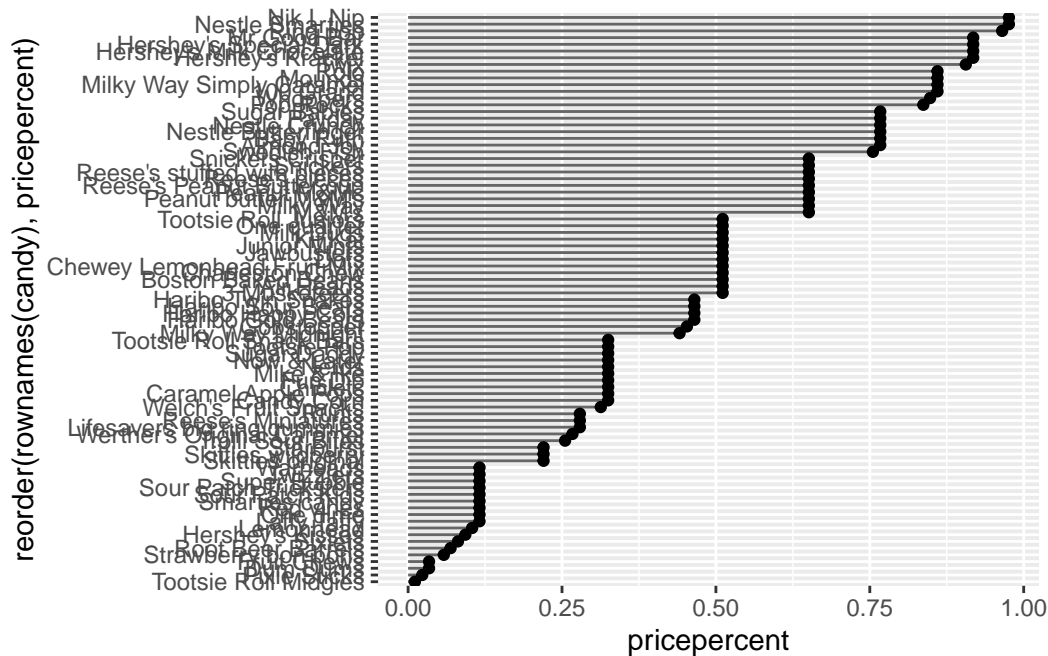
	chocolate	fruity	caramel	peanut	almond	nougat		
Nik L Nip	0	1	0		0	0		
Nestle Smarties	1	0	0		0	0		
Ring pop	0	1	0		0	0		
Hershey's Krackel	1	0	0		0	0		
Hershey's Milk Chocolate	1	0	0		0	0		
Hershey's Special Dark	1	0	0		0	0		
	crisped	rice	wafer	hard	bar	pluribus	sugar	percent
Nik L Nip		0	0	0		1		0.197
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

Hershey's Special Dark	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			
Hershey's Special Dark	0.918	59.23612			

Nik L Nip is the least popular

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()`.

```
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

Now that we have explored the dataset a little, we will see how the variables interact with one another.

First we will use correlation and view the results with the **corrplot** package to plot a correlation matrix

```
cij = cor(candy)
cij
```

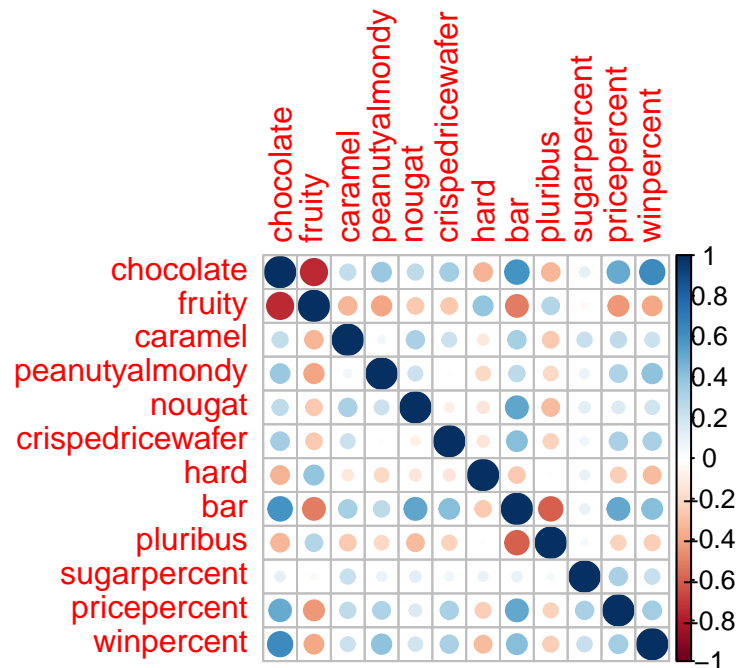
	chocolate	fruity	caramel	peanutyalmondy	nougat
chocolate	1.0000000	-0.74172106	0.24987535	0.37782357	0.25489183
fruity	-0.7417211	1.00000000	-0.33548538	-0.39928014	-0.26936712
caramel	0.2498753	-0.33548538	1.00000000	0.05935614	0.32849280
peanutyalmondy	0.3778236	-0.39928014	0.05935614	1.00000000	0.21311310
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

```
library(corrplot)
```

```
corrplot 0.95 loaded
```

```
corrplot(cij)
```



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

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

## Principal Component Analysis

Let's apply PCA using the `prcomp()` function to our candy dataset remembering to set the `scale=TRUE` argument.

```
pca = prcomp(candy, scale=T)
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

```
attributes(pca)
```

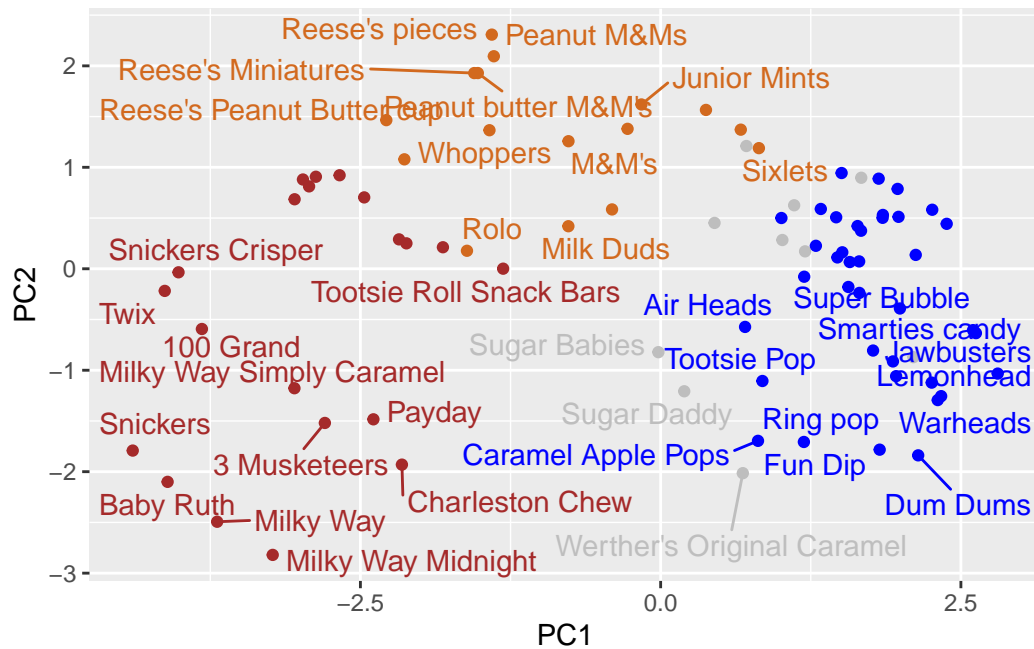
```
$names
[1] "sdev"      "rotation" "center"    "scale"     "x"

$class
[1] "prcomp"
```

Let's plot our main results as our PCA "score plot"

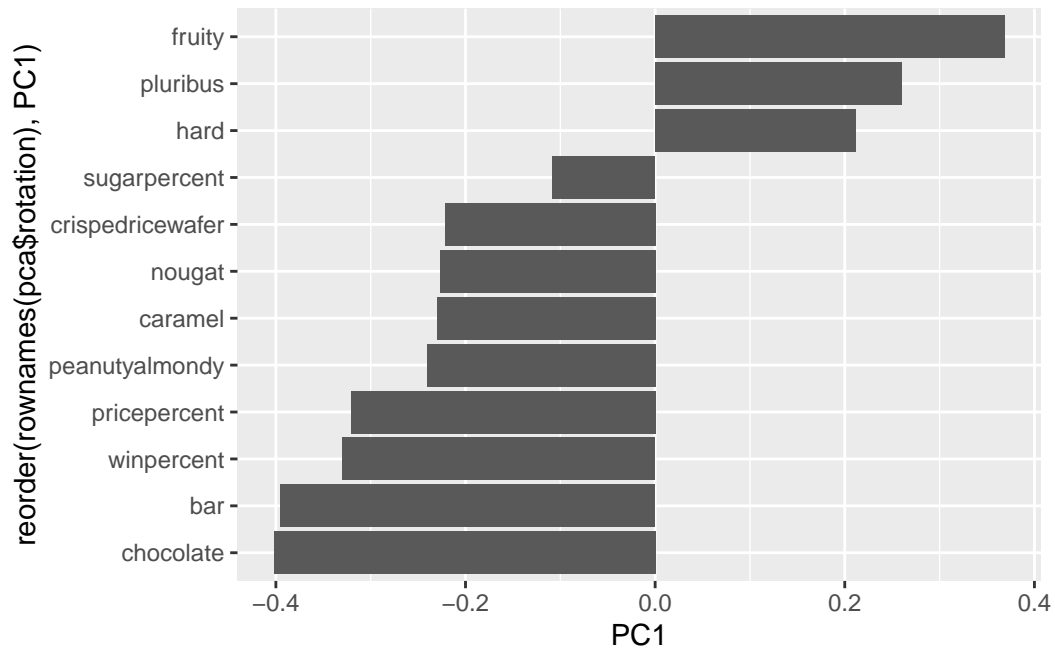
```
ggplot(pca$x) +
  aes(PC1, PC2, label=rownames(pca$x)) +
  geom_point(col=mycols) +
  geom_text_repel(col=mycols)
```

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



Finally let's look at how the original variables contribute to the PC's, start with PC1

```
ggplot(pca$rotation) +
  aes(PC1, reorder(rownames(pca$rotation), PC1)) +
  geom_col()
```



Q24. What original variables are picked up strongly by PC1 in the positive direction? Do these make sense to you? Fruitym pluribus and hard were strongly in the positive direction. This makes sense because they are generally correlated with one another.