

class 10: Halloween Candy Mini Project

Fu-Hsuan Ko

Table of contents

Background	2
Q1. How many different candy types are in this dataset?	3
Q2. How many fruity candy types are in the dataset?	3
Q3. What is your favorite candy in the dataset and what is it's winpercent value?	3
Q4. What is the winpercent value for "Kit Kat"?	4
Q5. What is the winpercent value for "Tootsie Roll Snack Bars"?	4
Q6. Is there any variable/column that looks to be on a different scale to the majority of the other columns in the dataset?	5
Q7. What do you think a zero and one represent for the candy\$chocolate column?	5
Q8. Plot a histogram of winpercent values	5
Q9. Is the distribution of winpercent values symmetrical?	6
Q10. Is the center of the distribution above or below 50%?	6
Q11. On average is chocolate candy higher or lower ranked than fruit candy?	6
Q12. Is this difference statistically significant?	7
3. Overall Candy Rankings	7
Q13. What are the five least liked candy types in this set?	8
Q14. What are the top 5 all time favorite candy types out of this set?	9

Q15. Make a first barplot of candy ranking based on winpercent values.	10
Barplot	10
Q16. This is quite ugly, use the reorder() function to get the bars sorted by winpercent?	12
Q17. What is the worst ranked chocolate candy?	16
Q18. What is the best ranked fruity candy?	16
4. Taking a look at pricepercnet	16
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?	18
Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?	18
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().	19
Q22. Examining this plot what two variables are anti-correlated (i.e. have minus values)?	21
Q23. Similarly, what two variables are most positively correlated?	21
PCA: Principal Component Analysis	22
Q24. What original variables are picked up strongly by PC1 in the positive direction? Do these make sense to you?	24

Background

In this mini-project we will examine 538 Halloween Candy data. What is your favourite candy? What is nougat anyway? And how do you say it in America?

```
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
```

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

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

```
[1] 81.64291
```

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
```

```
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_ratio	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	

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
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?

winepercent

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

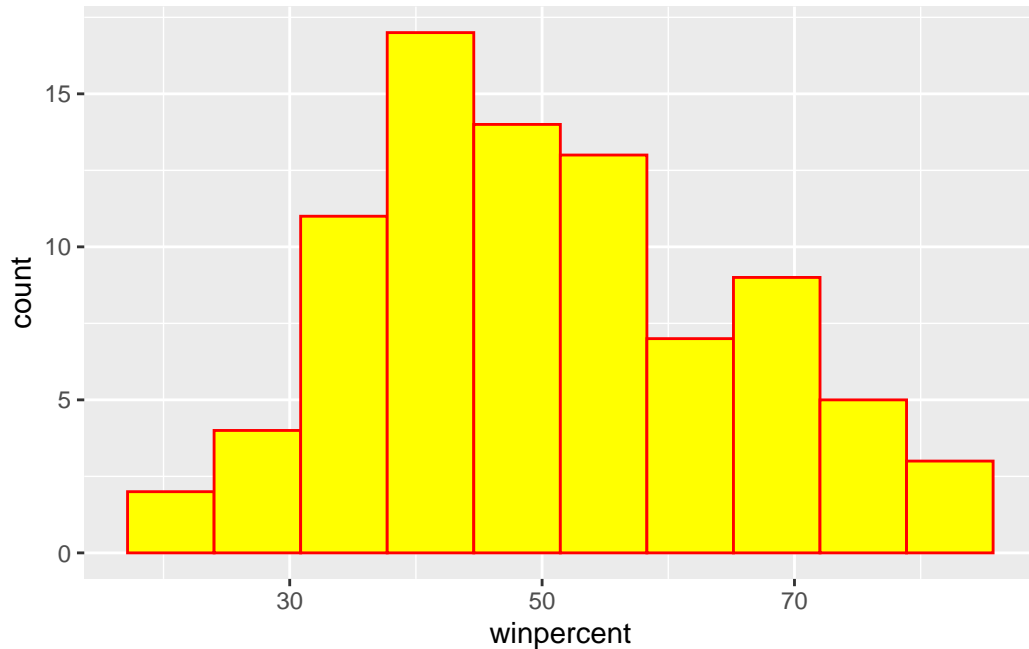
```
candy$chocolate
```

```
[1] 1 1 0 0 0 1 1 0 0 0 1 0 0 0 0 0 0 0 0 0 1 1 1 1 0 1 1 0 0 0 1 1 0 1 1 1
[39] 1 1 1 0 1 1 0 0 0 1 0 0 0 1 1 1 1 0 1 0 0 1 0 0 1 0 1 1 0 0 0 0 0 0 0 1 1
[77] 1 1 0 1 0 0 0 0 1
```

Zero means this candy is not in the type of chocolate, while one represents the candy is in the chocolate type

Q8. Plot a histogram of winpercent values

```
library(ggplot2)
ggplot(candy, aes(winpercent))+
  geom_histogram(bins=10, col="red", fill="yellow")
```



Q9. Is the distribution of winpercent values symmetrical?

Yes

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

```
median(candy$winpercent)
```

```
[1] 47.82975
```

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

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

```
[1] 60.92153
```

And for fruit candy...

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

```
[1] 44.11974
```

Q12. Is this difference statistically significant?

Yes

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

Welch Two Sample t-test

```
data: chocolate.win and fruity.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
```

3. Overall Candy Rankings

The base R “sort()” and “order()” functions are very useful!

```
x <- c(5,1,2,6)
sort(x, decreasing = T)
```

```
[1] 6 5 2 1
```

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

```
[1] 1 2 5 6
```

```
y <- c("berry", "alice", "chandra")
y
```

```
[1] "berry" "alice" "chandra"
```

```
sort(y)
```

```
[1] "alice" "berry" "chandra"
```

```
order(y)
```

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

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

First I want to order/arrange the whole dataset by winpercent values

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

```
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

```
head(arrange(candy, desc(winpercent)), n=5)
```

	chocolate	fruity	caramel	peanut	almondy	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					

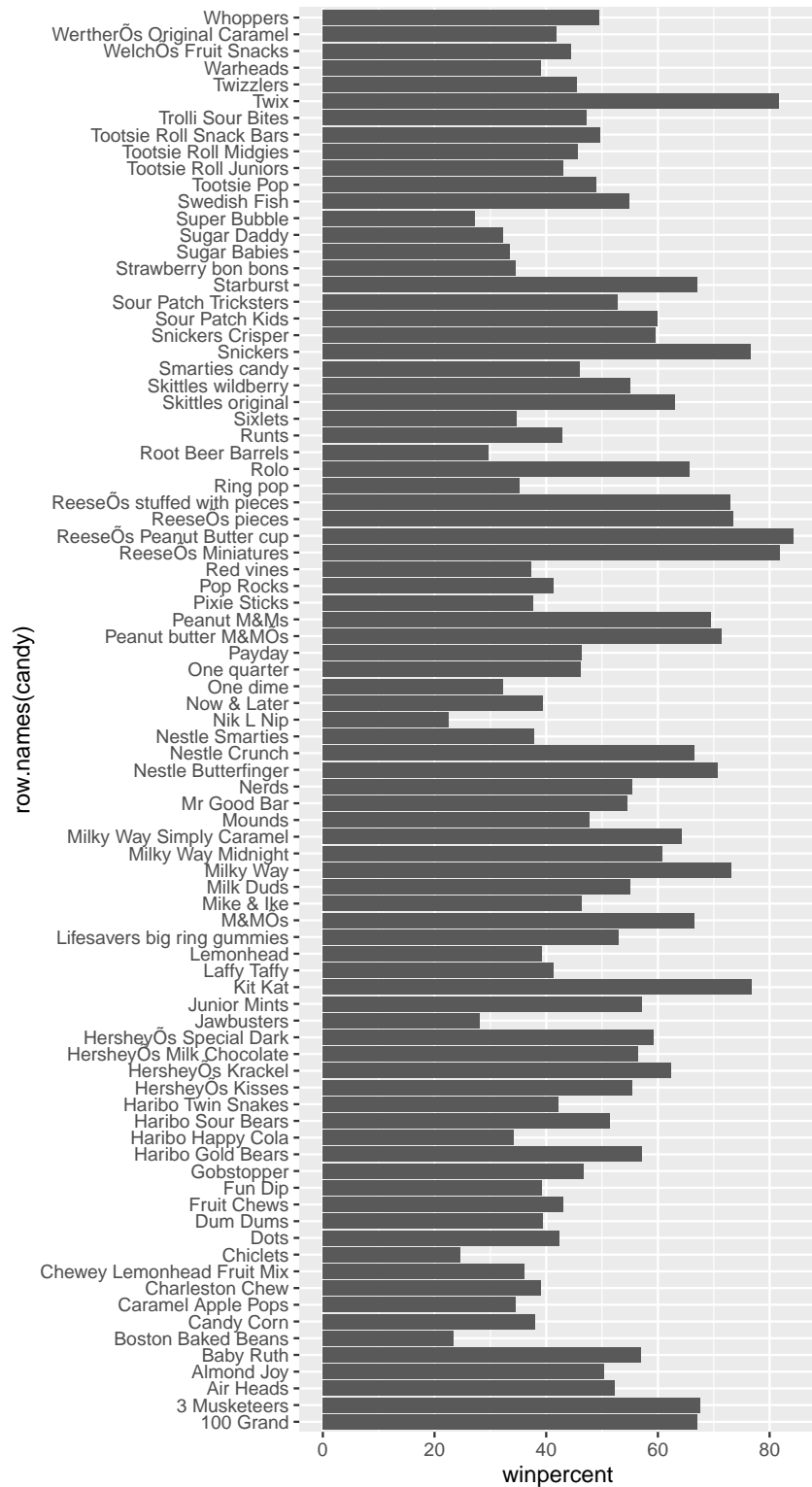
Reese's Peanut Butter cup	0.651	84.18029
Reese's Miniatures	0.279	81.86626
Twix	0.906	81.64291
Kit Kat	0.511	76.76860
Snickers	0.651	76.67378

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

Barplot

The default barplot, made with “geom_col” has the bars in the order they are in the dataset...

```
ggplot(candy)+  
  aes(winpercent, row.names(candy))+  
  geom_col()
```



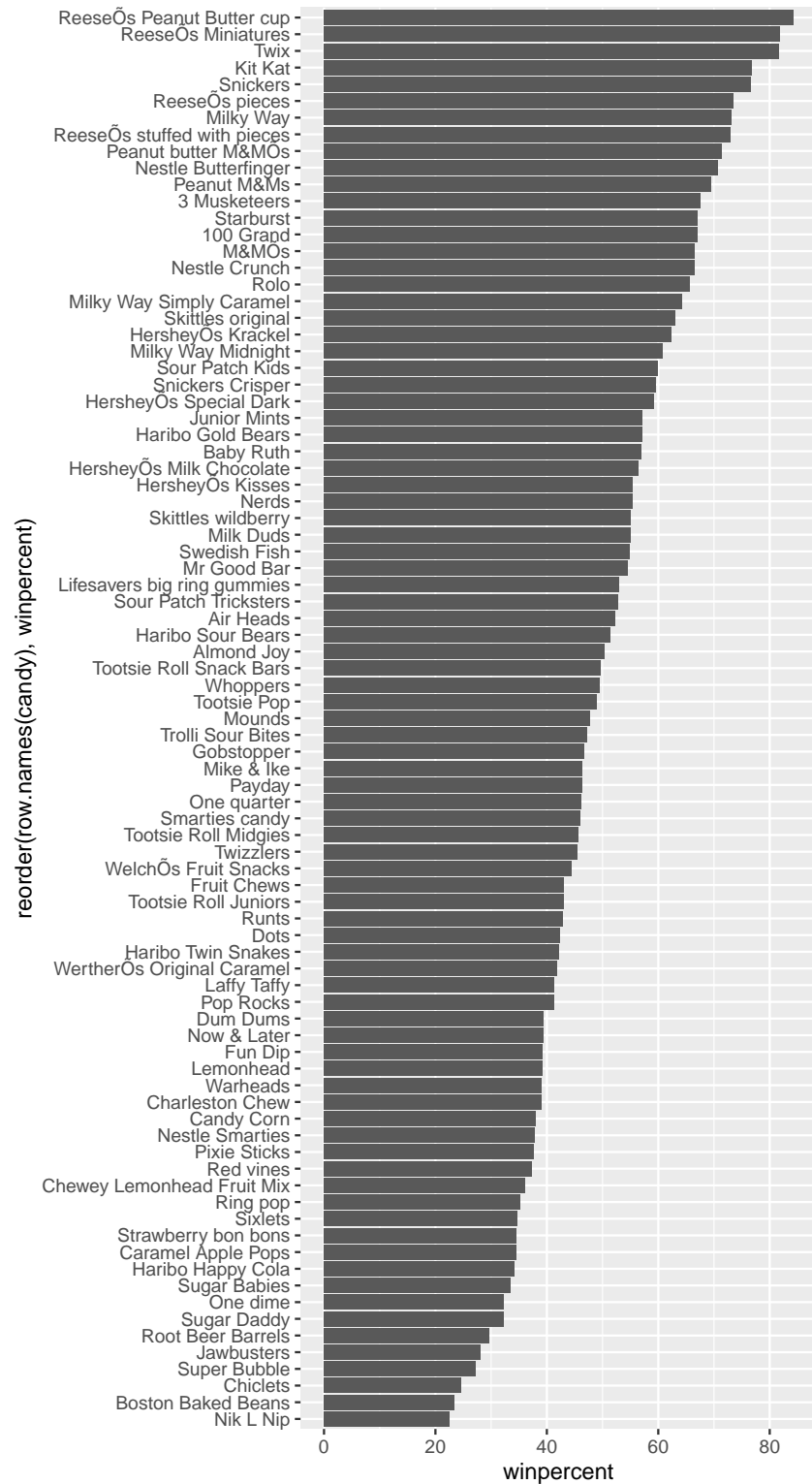
Q16. This is quite ugly, use the `reorder()` function to get the bars sorted by winpercent?

```
p <- ggplot(candy)+  
  aes(winpercent, reorder(row.names(candy), winpercent))+  
  geom_col()
```

```
ggsave("mybarplot.png", p)
```

Saving 5.5 x 3.5 in image

```
p
```



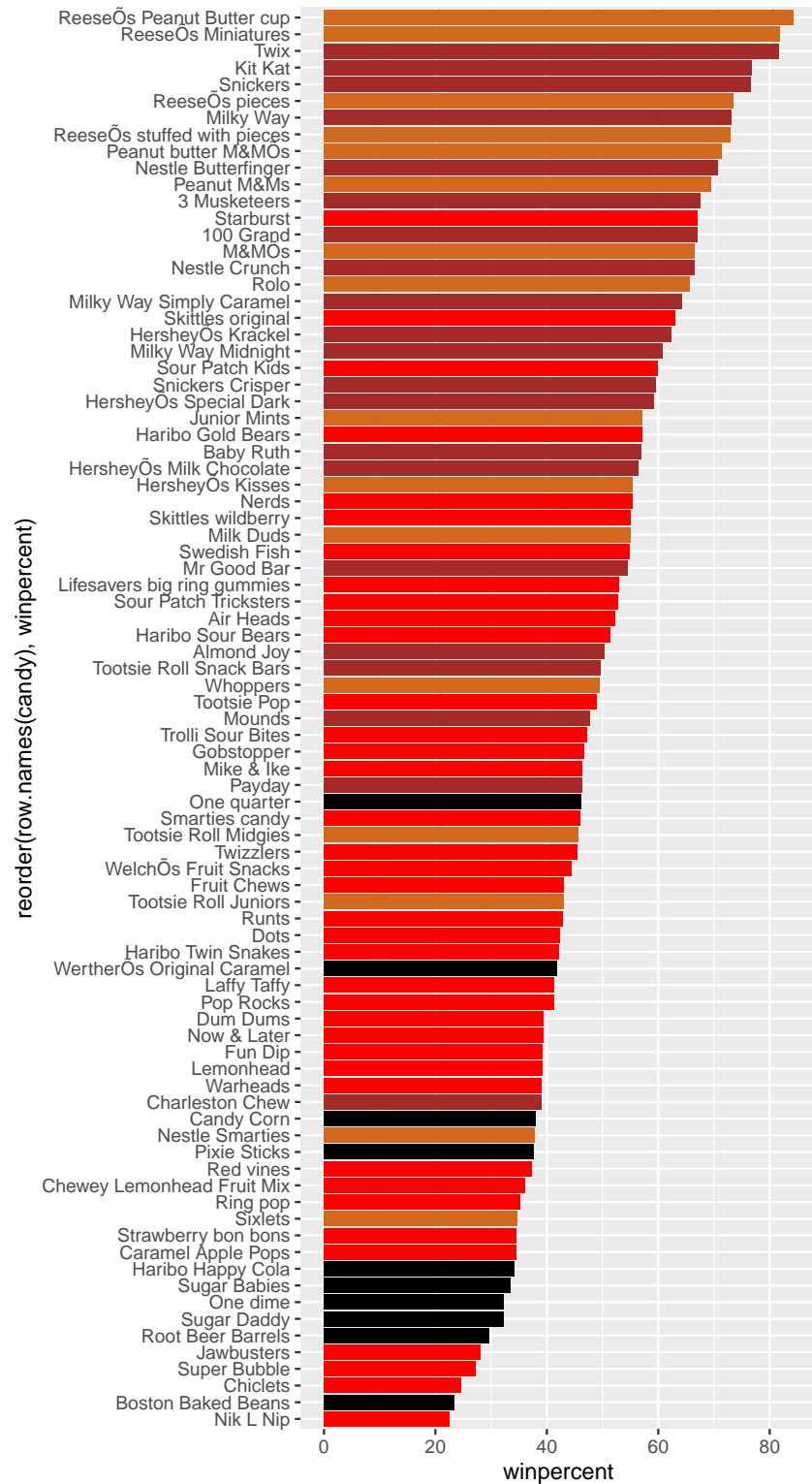
Let's setup a color vector (that signifies candy type) that we can then use for some future bar plots (one for each candy). Then we overwrite chocolate (for chocolate candy), brown (for candy bars) and red (for fruity candy)

```
my_cols <- rep("black", nrow(candy))
#my_cols
my_cols[as.logical(candy$chocolate)] <- "chocolate"
my_cols[as.logical(candy$bar)] <- "brown"
my_cols[as.logical(candy$fruity)] <- "red"
my_cols
```

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

Now I can use this vector to color up my barplot

```
ggplot(candy)+
  aes(winpercent, reorder(row.names(candy), winpercent))+
  geom_col(fill=my_cols)
```



Now, for the first time, using this plot we can answer questions like:

Q17. What is the worst ranked chocolate candy?

```
library(dplyr)
candy%>%
  filter(chocolate==TRUE)%>%
  arrange(winpercent)%>%
  .[1,]
```

	chocolate	fruity	caramel	peanutyalmondy	nougat	crispedricewafer	hard
Sixlets	1	0	0	0	0		0
	bar	pluribus	sugarpercent	pricepercent	winpercent		
Sixlets	0	1	0.22	0.081	34.722		

Q18. What is the best ranked fruity candy?

```
candy%>%
  filter(fruity==TRUE)%>%
  arrange(desc(winpercent))%>%
  .[1,]
```

	chocolate	fruity	caramel	peanutyalmondy	nougat	crispedricewafer	hard
Starburst	0	1	0	0	0		0
	bar	pluribus	sugarpercent	pricepercent	winpercent		
Starburst	0	1	0.151	0.22	67.03763		

4. Taking a look at pricepercnet

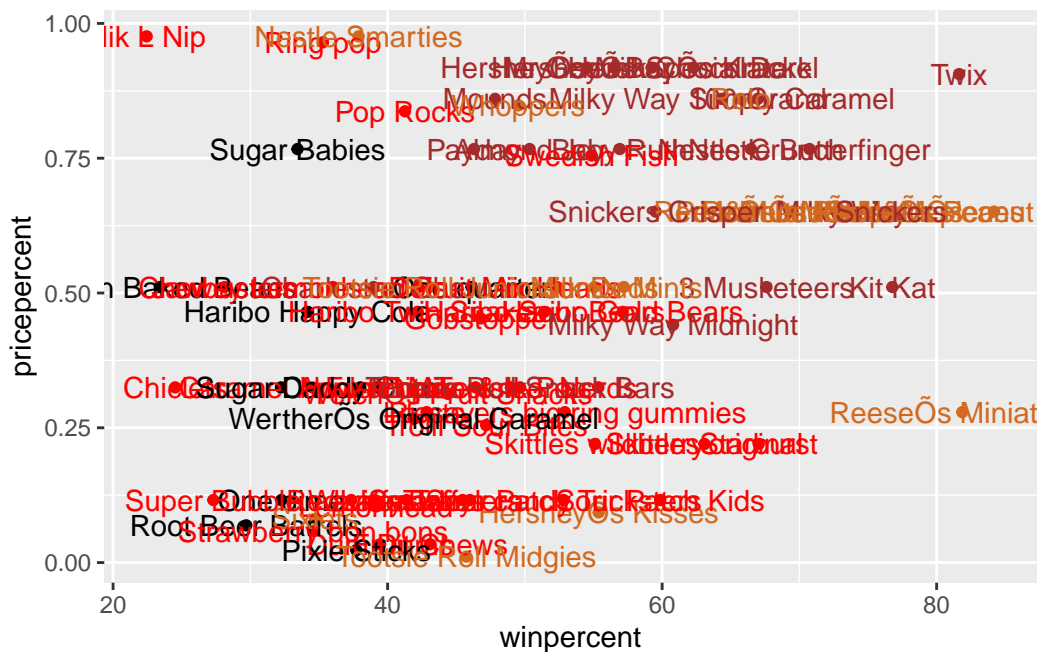
What about value for money? What is the best candy for the least money?

One way to get at this would be to make a plot of “winpercent” as the “pricepercent” values

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



```
geom_text(col=my_cols)
```

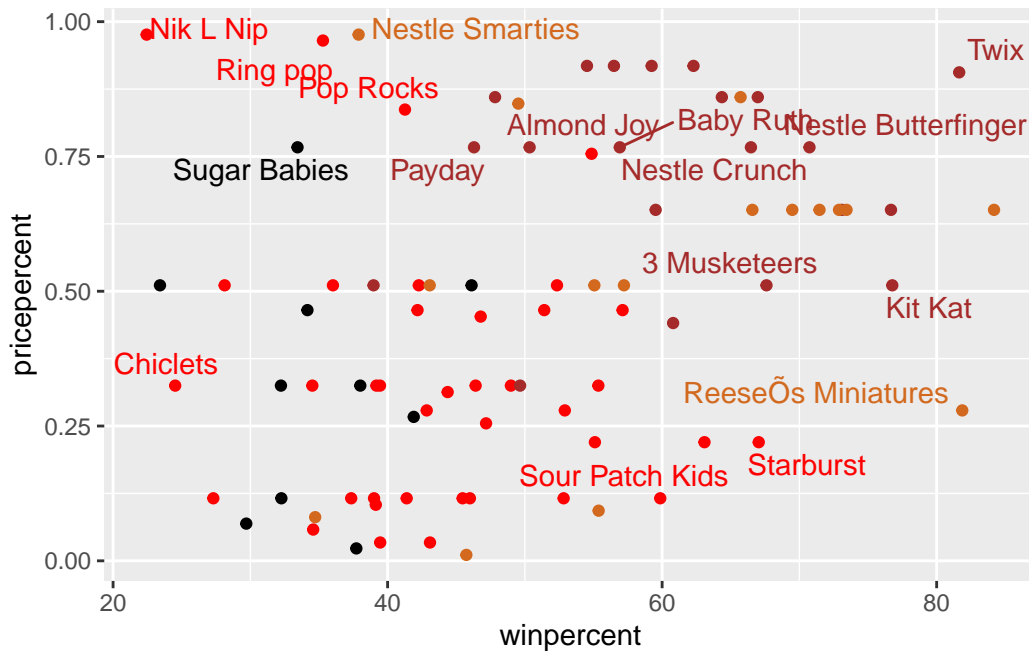


This plot sucks! I can not read the labels... We can use ggrepel package to help with this

```
library(ggrepel)

ggplot(candy)+
  aes(winpercent, pricepercent, label=rownames(candy))+
  geom_point(col=my_cols)+
  geom_text_repel(col=my_cols, max.overlaps = 7)
```

Warning: ggrepel: 68 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?

Reese's Miniatures

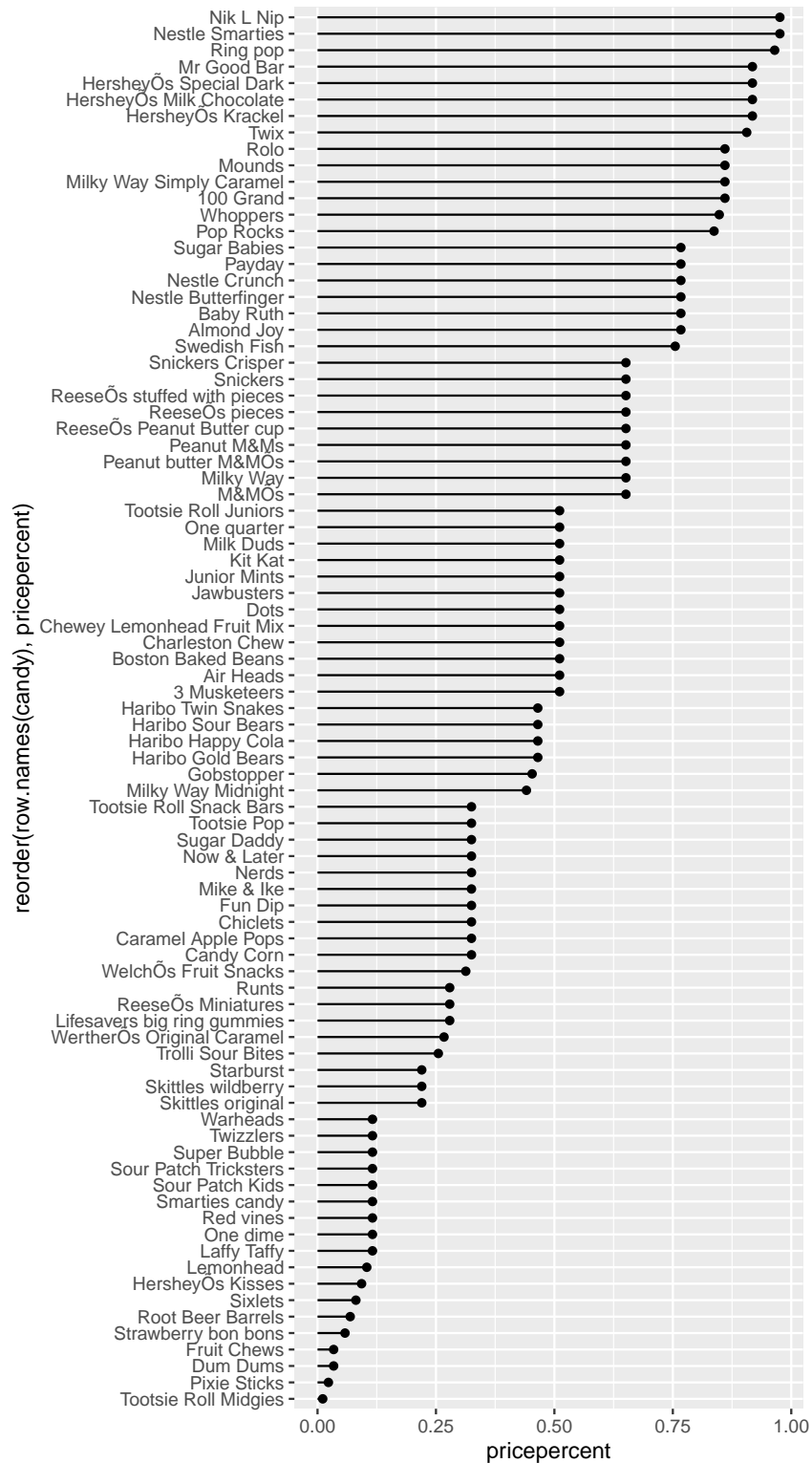
Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?

```
ord <- 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()`.

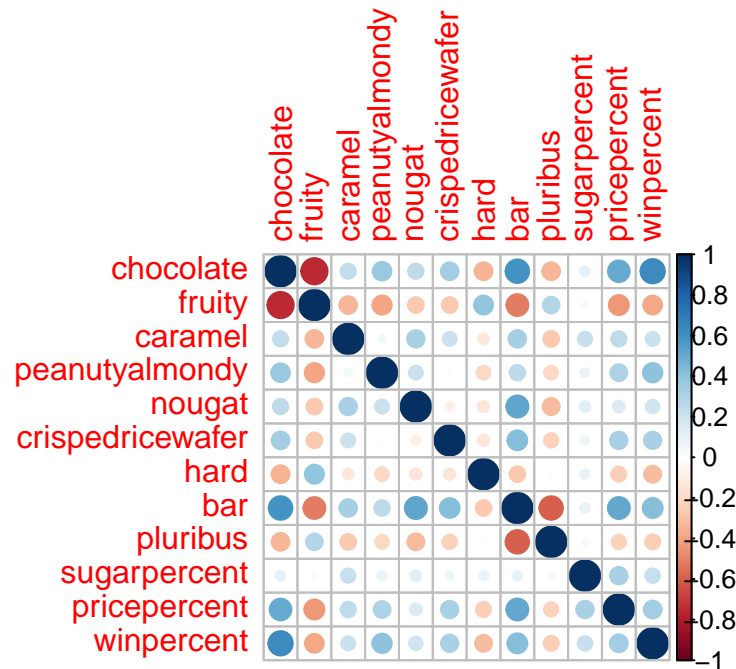
```
ggplot(candy)+  
  aes(pricepercent, reorder(row.names(candy), pricepercent))+  
  geom_point()+  
  geom_segment(aes(yend = reorder(row.names(candy), pricepercent), xend = 0))
```



```
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 fruity

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

Chocolate and winpercent

PCA: Principal Component Analysis

The main function that always there for us is “prcomp()” It has an important argument that is set to “scale=False”

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

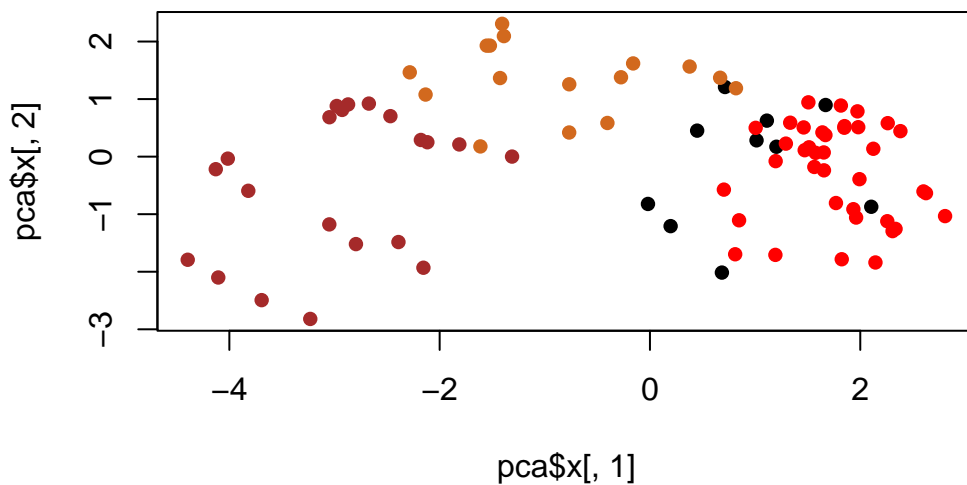
Importance of components:

	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

My PCA plot (a.k.a.) PC1 vs PC2 score plot.

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

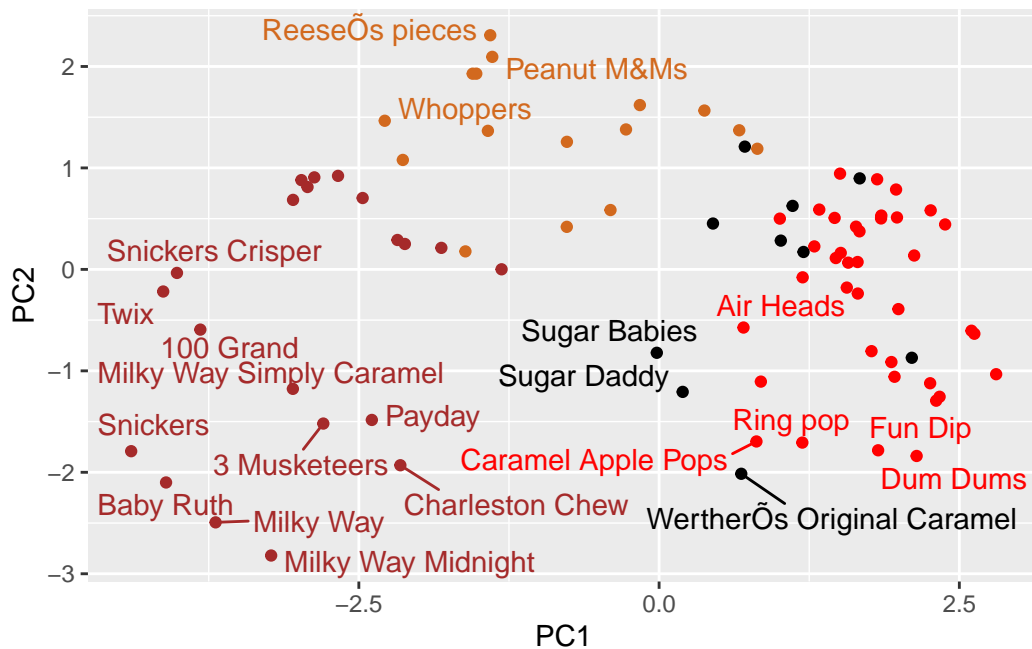


I will make a “nicer” plot with ggplot. ggplot only works with data.frames as input so I need to make one for it first...

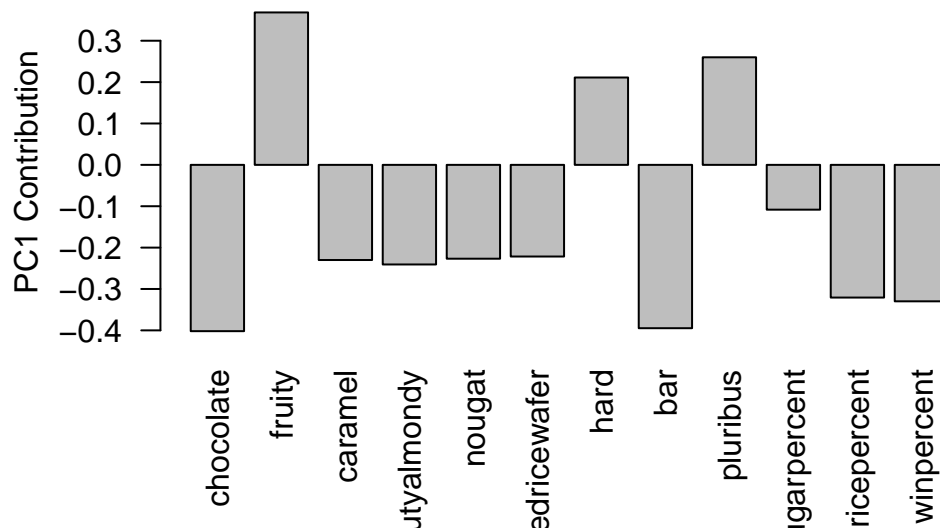
```
#Make a new data-frame with our PCA results and candy data
my_data <- cbind(candy, pca$x[,1:3])
```

```
ggplot(my_data)+
  aes(PC1, PC2, label=row.names(my_data))+
  geom_point(col=my_cols)+
  geom_text_repel(col=my_cols, max.overlaps = 7)
```

Warning: ggrepel: 63 unlabeled data points (too many overlaps). Consider increasing max.overlaps



```
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. Yes