

halloween-candy-project

Chloe J. Welch

10/29/2021

For this project, we will perform an exploratory analysis on Halloween candy.

First, we will begin by importing our candy data.

```
url <- "https://raw.githubusercontent.com/fivethirtyeight/data/master/candy-power-ranking/candy-data.csv"
candy_file <- read.csv(url)
candy = data.frame(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
```

Time to look at favorite candies! For a given candy this value is the percentage of people who prefer this candy over another randomly chosen candy from the dataset (what 538 term a matchup). Higher values indicate a more popular candy.

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

```
## [1] 81.64291
```

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

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

```
## [1] 67.03763
```

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
```






Now, let's see a quick overview of this dataset.








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

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?

Yes. The “percent” columns (10-12) seem to be on a 0-100 scale, while the rest of the columns appear to be on a 0-1 scale.

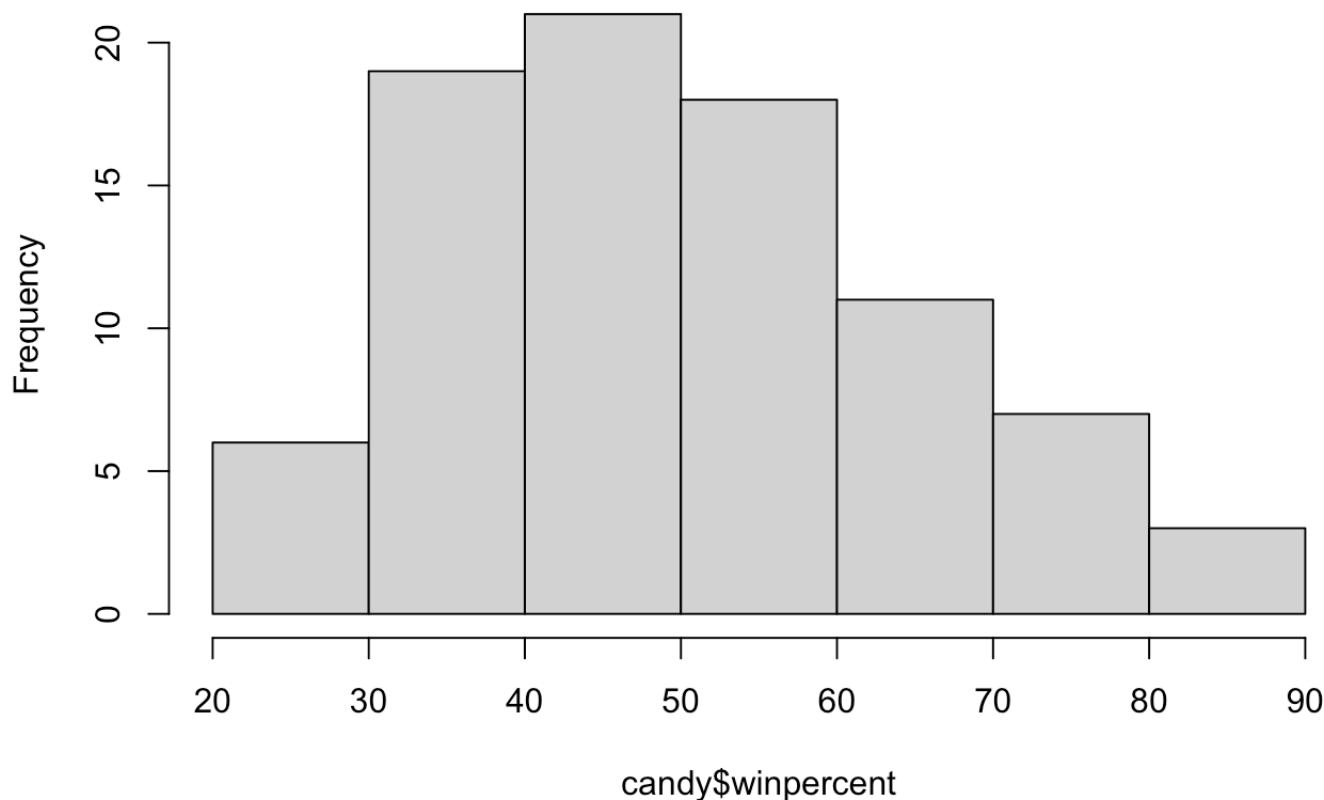
Q7. What do you think a zero and one represent for the candy\$chocolate column?

The “zero” means the candy is not “chocolate”, and the “one” means it is not.

Let’s make a histogram.

Q8. Plot a histogram of winpercent values.

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

Histogram of candy\$winpercent

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 is below 50%.

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

```
mean(candy$winpercent[as.logical(candy$chocolate)])
```

```
## [1] 60.92153
```

```
mean(candy$winpercent[as.logical(candy$fruit)])
```

```
## [1] 44.11974
```

```
as.logical(candy$chocolate)
```

```
## [1] TRUE TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE FALSE TRUE FALSE
## [13] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE TRUE
## [25] TRUE TRUE FALSE TRUE TRUE FALSE FALSE FALSE TRUE TRUE FALSE TRUE
## [37] TRUE TRUE TRUE TRUE TRUE FALSE TRUE TRUE FALSE FALSE FALSE FALSE TRUE
## [49] FALSE FALSE FALSE TRUE TRUE TRUE TRUE FALSE TRUE FALSE FALSE TRUE
## [61] FALSE FALSE TRUE FALSE TRUE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
## [73] FALSE FALSE TRUE TRUE TRUE TRUE FALSE TRUE FALSE FALSE FALSE FALSE
## [85] TRUE
```

Q12. Is this difference statistically significant?

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

```
##
## 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
```

Yes. The difference is statistically significant as the results yield a p-value of less than 0.05.

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

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

```
##               chocolate fruity caramel peanutyalmondy 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
##               crispedricewafer hard bar pluribus sugarpercent pricepercent
## 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
```

```
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(winpercent) %>% head(5)
```

```
##                chocolate fruity caramel peanutyalmondy 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
##                crispedricewafer hard bar pluribus sugarpercent pricepercent
## 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),], n=5)
```

```
##                chocolate fruity caramel peanutyalmondy 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
##                crispedricewafer hard bar pluribus sugarpercent pricepercent
## 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
```



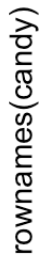
```
library("dplyr")
candy %>% arrange(desc(winpercent)) %>% head(5)
```

```
##               chocolate fruity caramel peanutyalmondy 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
##               crispedricewafer hard bar pluribus sugarpercent
## 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
##               pricepercent 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
```

Next, let's make a barplot.

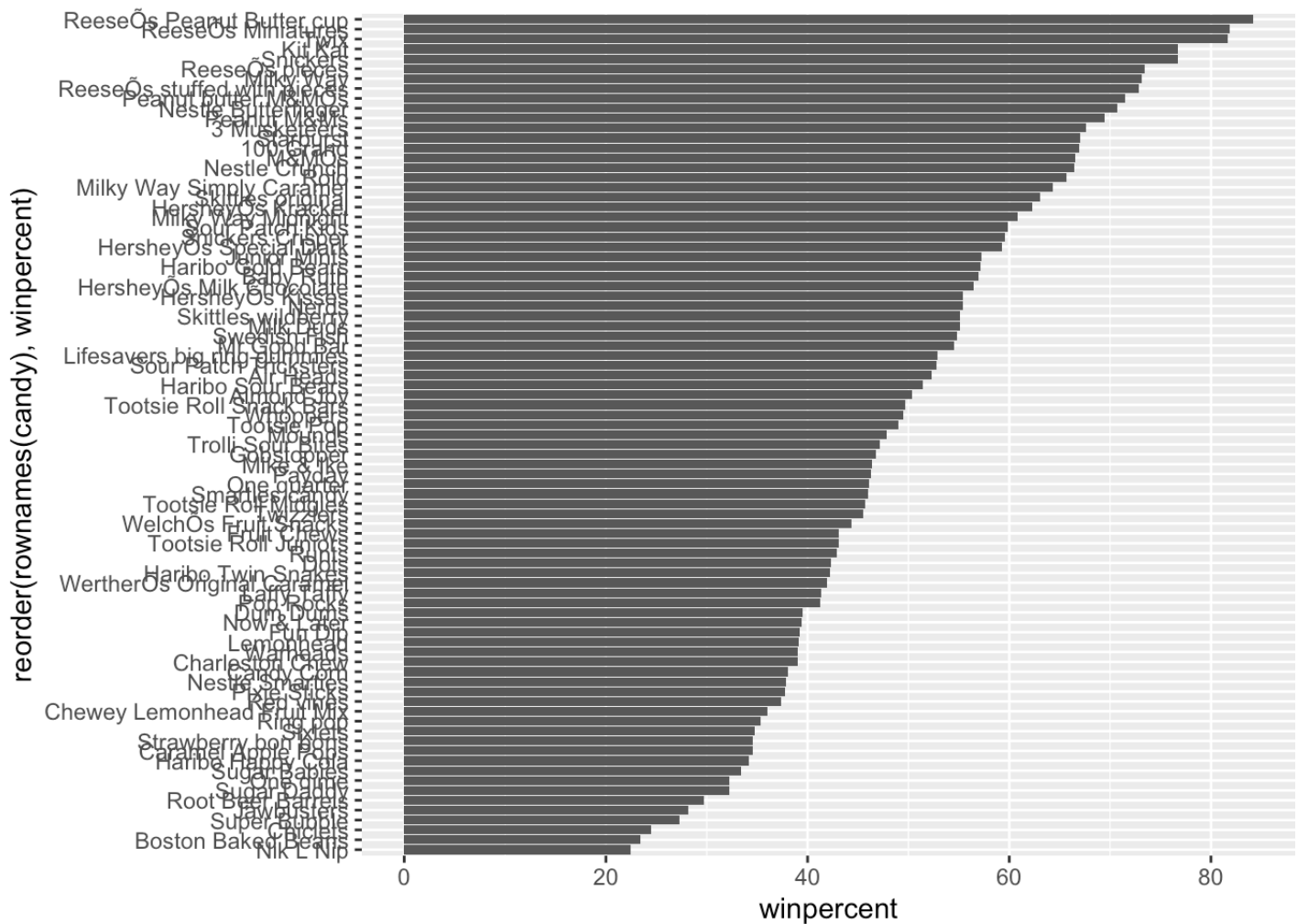
Q15. Make a first barplot of candy ranking based on winpercent values.

```
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`?

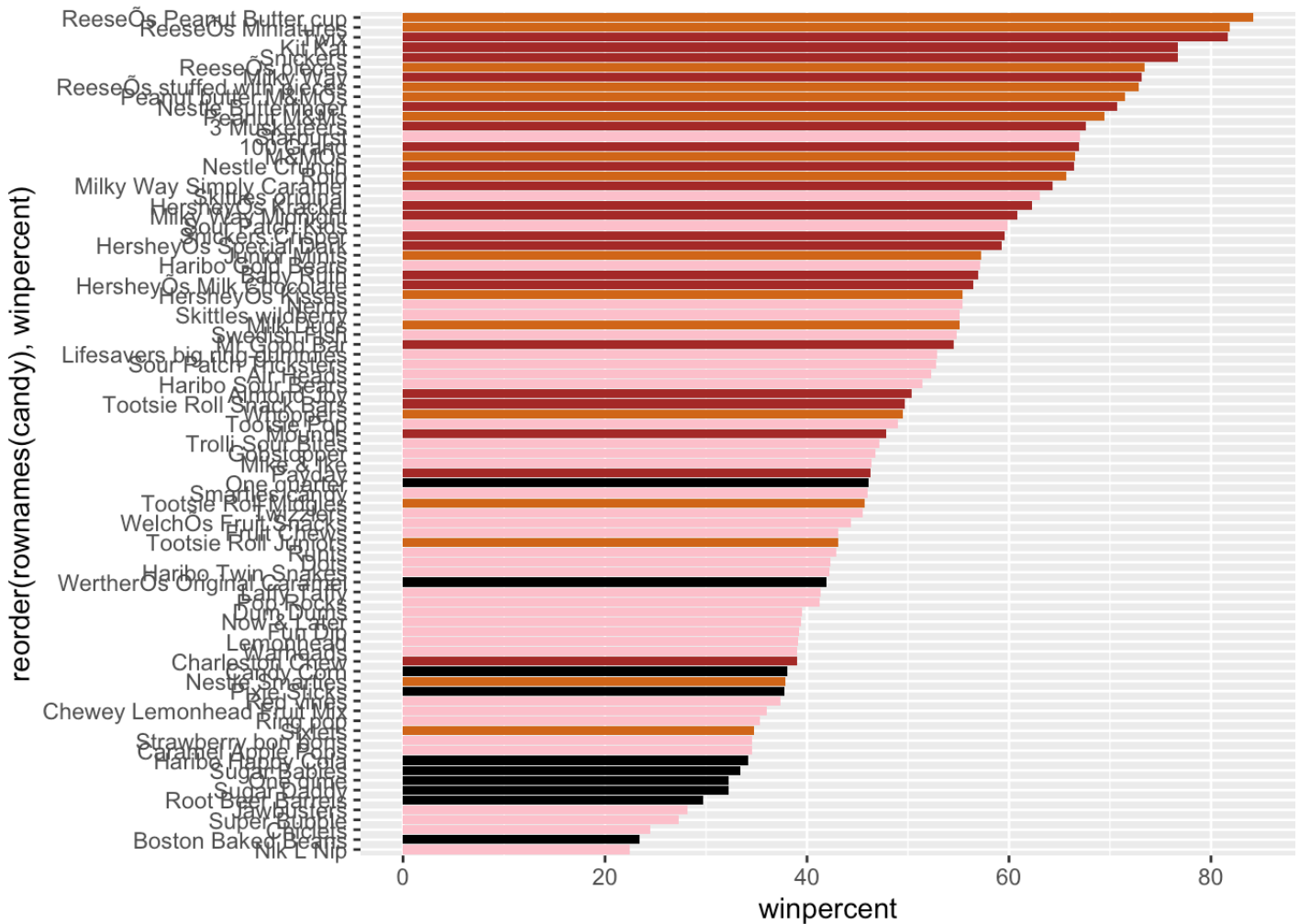
```
library(ggplot2)
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy), winpercent)) +
  geom_col()
```



Next, let's setup a color vector (that signifies candy type) that we can then use for some future plots. We start by making a vector of all black values (one for each candy). Then we overwrite chocolate (for chocolate candy), brown (for candy bars) and red (for fruity candy) values.

```
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 candy appears to be Sixlets.

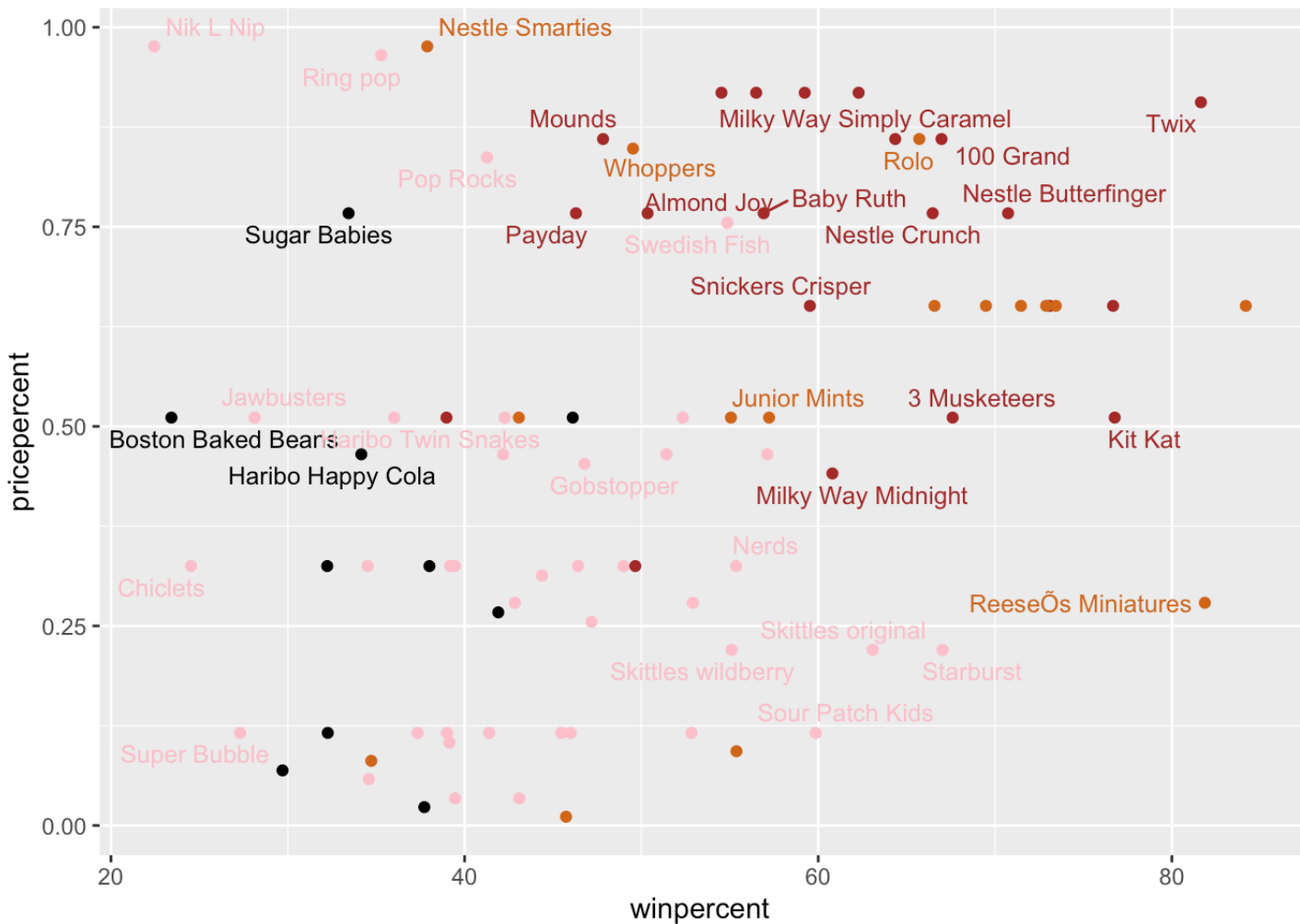
Q18. What is the best ranked fruity candy?

The best ranked fruit candy appears to be Starburst (I agree!!!)

Now, we are going to look at the best candy for least amount of money.

```
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 = 5)
```

```
## 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 Reese's Miniatures are the highest ranked in terms of winpercent for the least money.

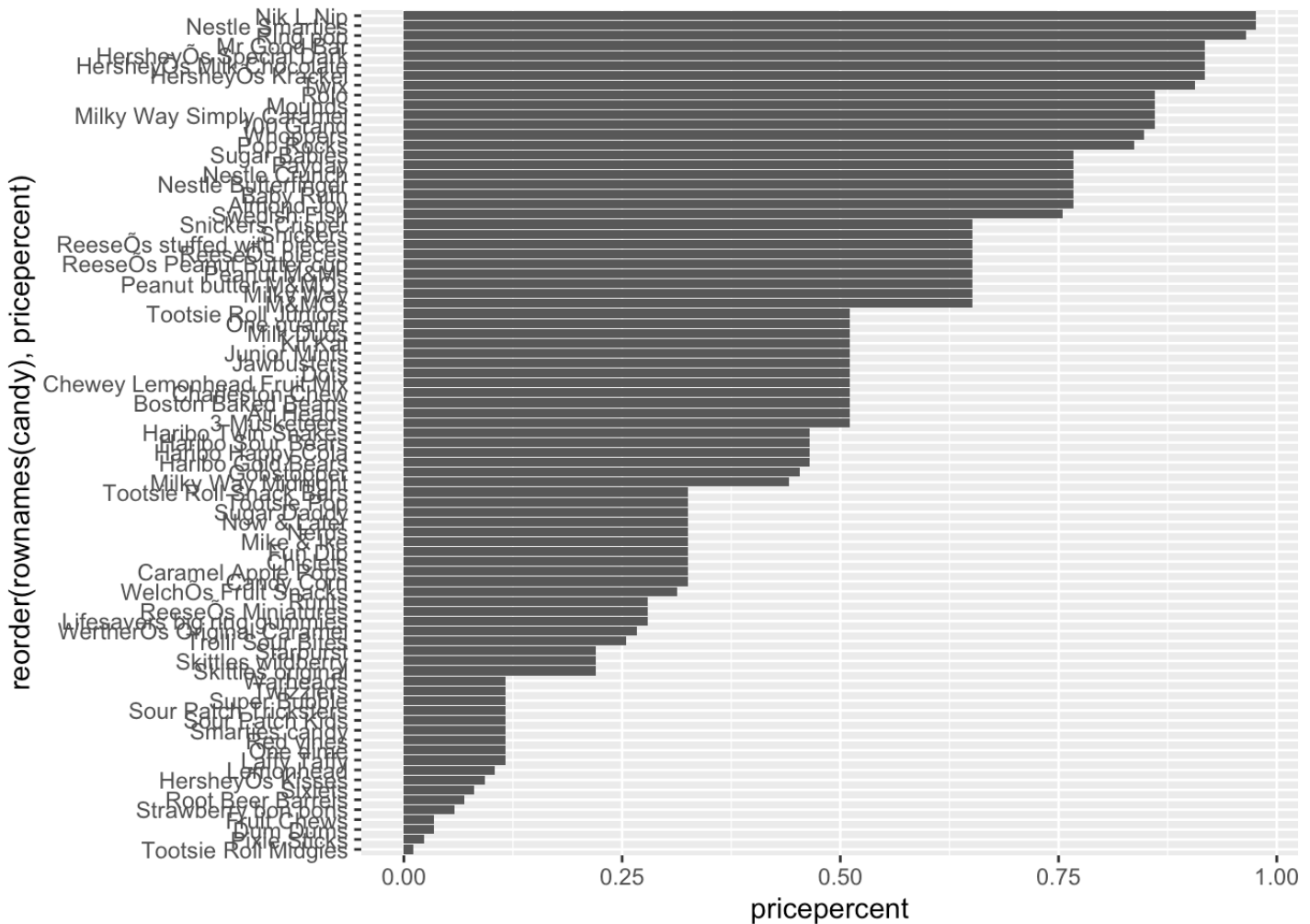
Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?

```
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()`.

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



Next, we will make a “lollipop” chart of pricepercent.

```
```r
ggplot(candy) +
 aes(pricepercent, reorder(rownames(candy), pricepercent)) +
 geom_segment(aes(yend = reorder(rownames(candy), pricepercent), xend = 0), col="gray40") +
 geom_point()
```
```

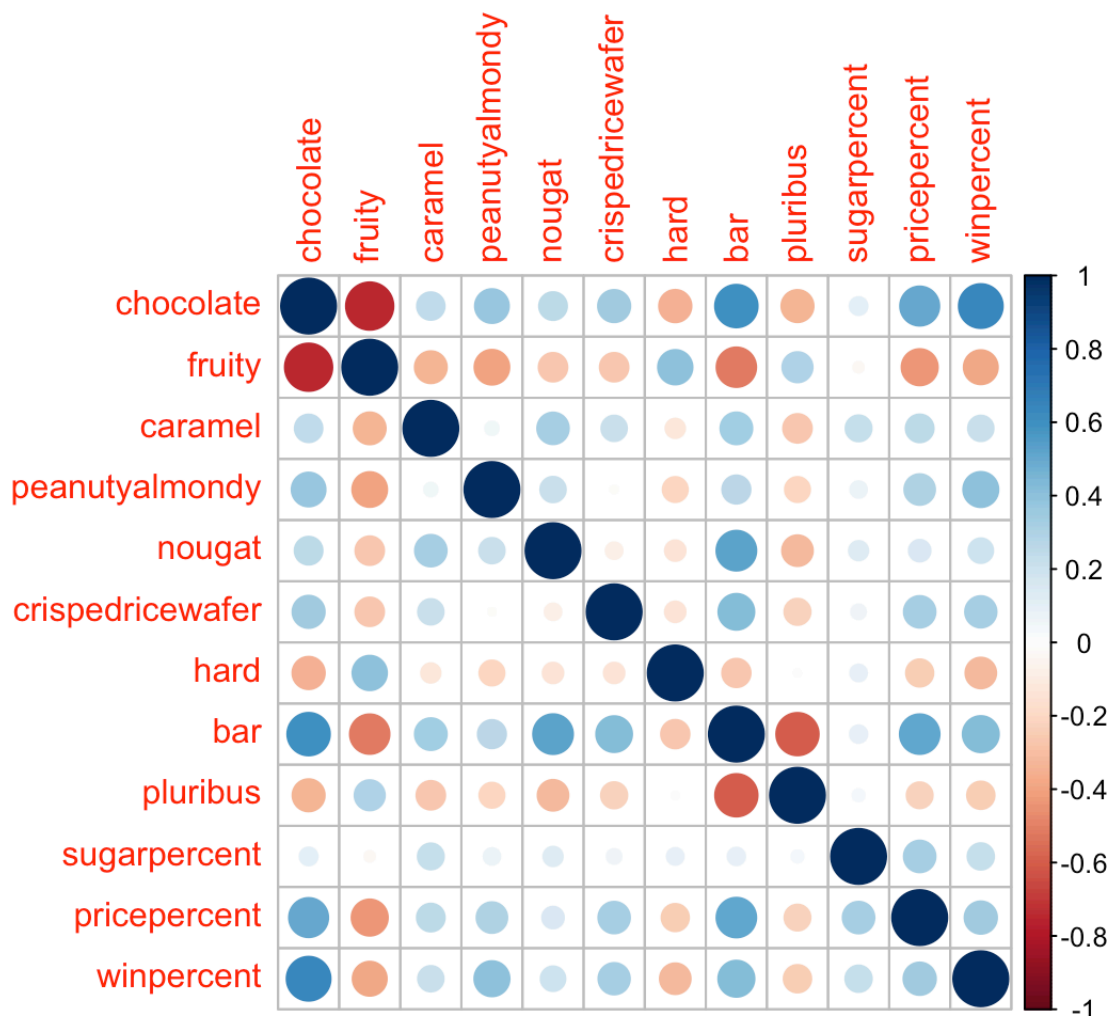

One of the most interesting aspects of this chart is that a lot of the candies share the same ranking, so it looks like quite a few of them are the same price.

Now that we've explored the dataset a little, we'll see how the variables interact with one another. We'll use correlation and view the results with the `corrplot` package to plot a correlation matrix.

```
library(corrplot)
```

```
## corrplot 0.90 loaded
```

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



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

Variables include chocolate/fruity, bar/pluribus, and fruity/bar.

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

Variables include chocolate/bar, chocolate/pricepercent, and chocolate/winpercent.

Finally, it's time to apply PCA!

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

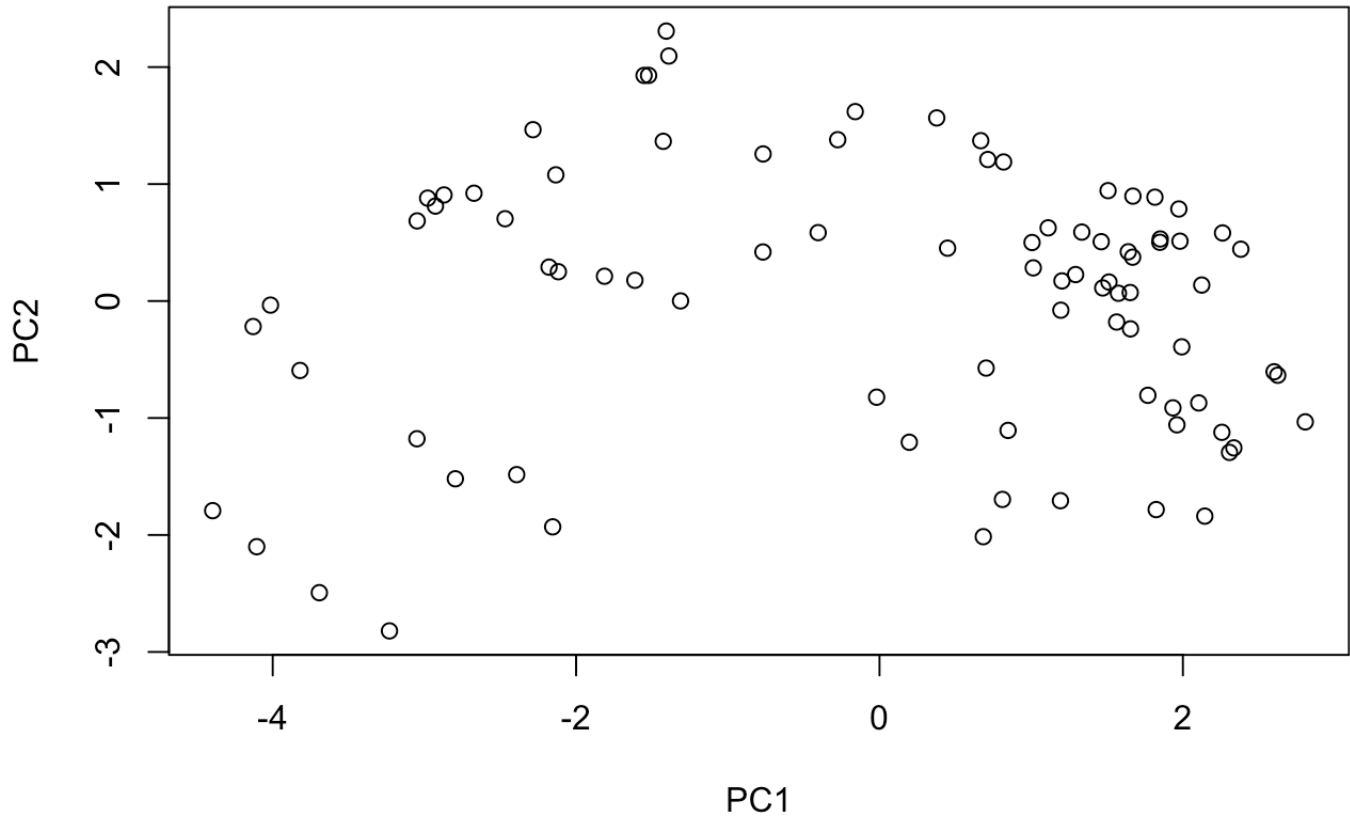
```
## Importance of components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  2.0788 1.1378 1.1092 1.07533 0.9518 0.81923 0.81530
## Proportion of Variance 0.3601 0.1079 0.1025 0.09636 0.0755 0.05593 0.05539
## Cumulative Proportion 0.3601 0.4680 0.5705 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
```

```
pca.false <- prcomp(candy, scale=FALSE)
summary(pca.false)
```

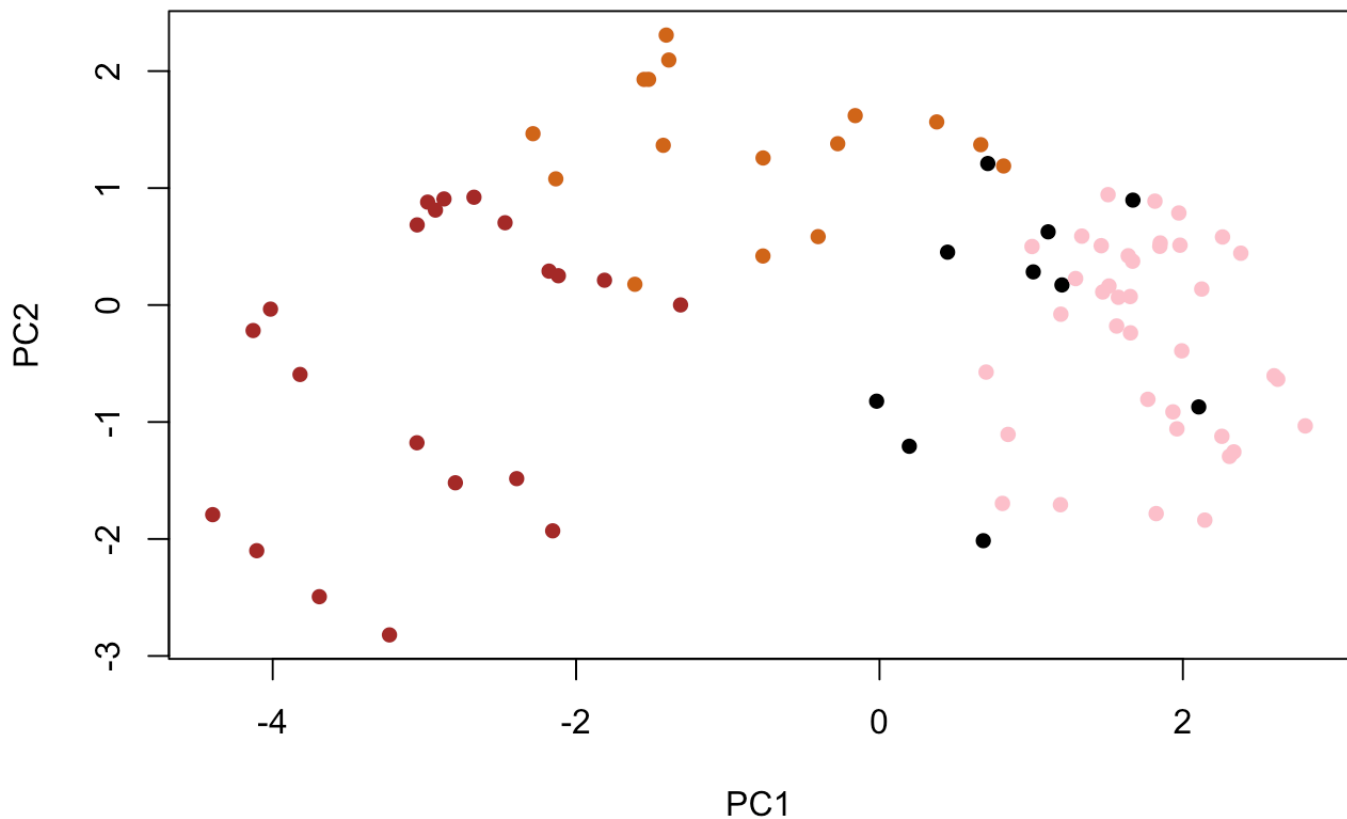
```
## Importance of components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation 14.7231 0.70241 0.47762 0.37292 0.34641 0.33614 0.30748
## Proportion of Variance 0.9935 0.00226 0.00105 0.00064 0.00055 0.00052 0.00043
## Cumulative Proportion 0.9935 0.99574 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
```

Now, we can plot our main PCA score plot of PC1 vs PC2.

```
plot(pca$x[, 1:2])
```



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

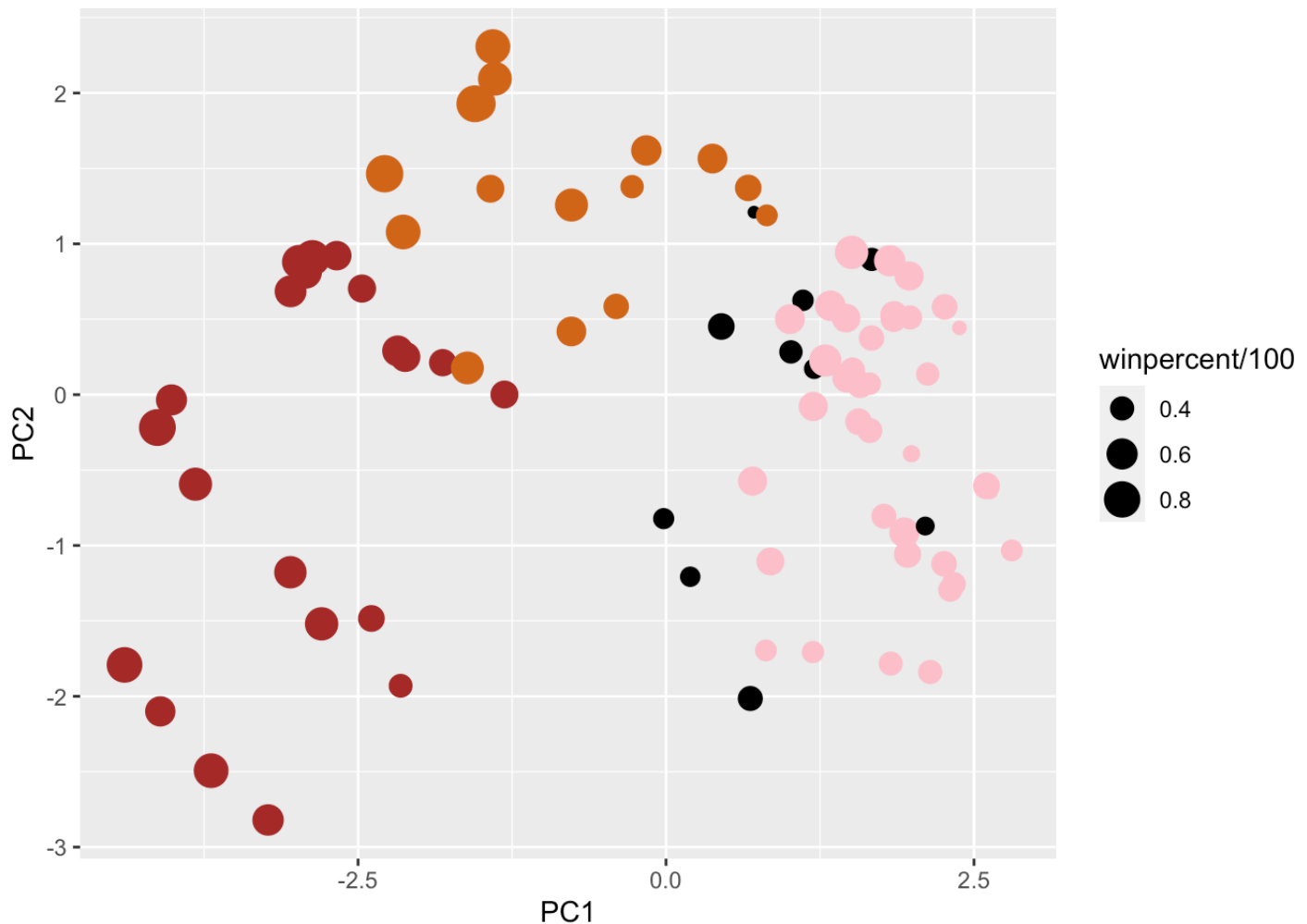


Make a new data-frame with our PCA results and candy data.

```
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



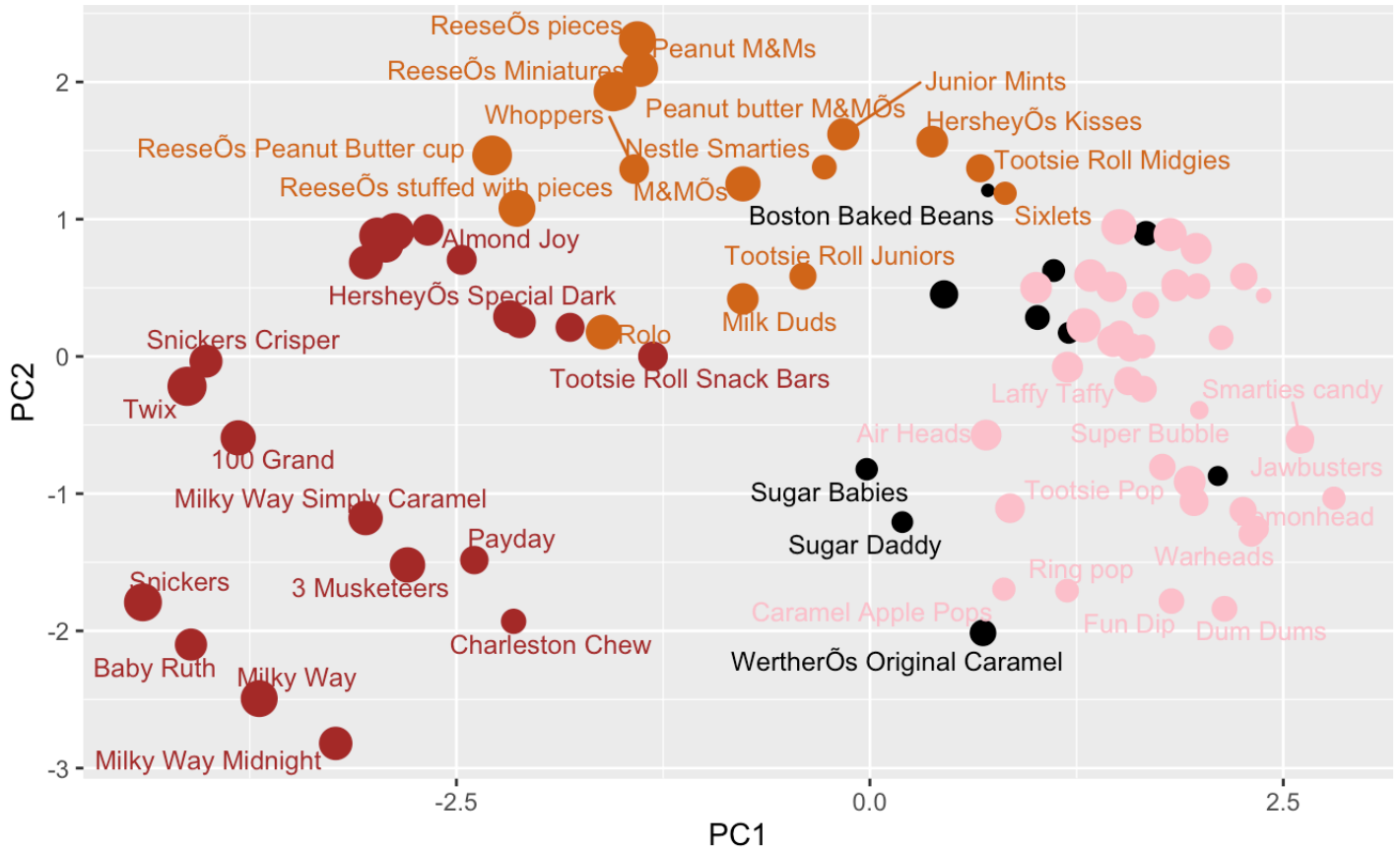
Again, we can use the `ggrepel` package and the function `ggrepel::geom_text_repel()` to label up the plot with non overlapping candy names like. We will also add a title and subtitle like so:

```
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), fruity (red), other (black)",
        caption="Data from 538")
```

```
## Warning: ggrepel: 39 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), fruity (red), other (black)



Data from 538

```
library(plotly)
```

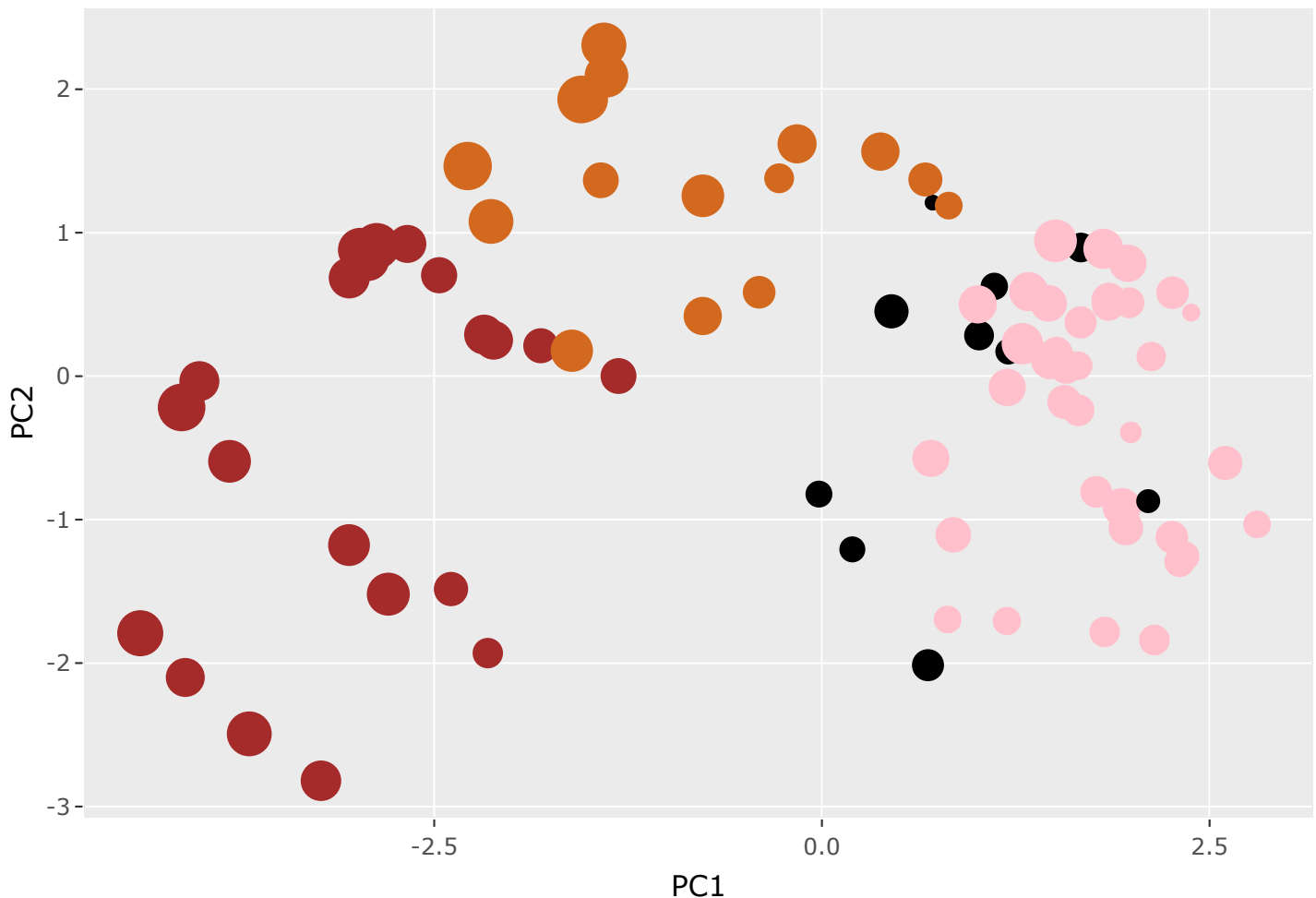
```
##
## Attaching package: 'plotly'
```

```
## The following object is masked from 'package:ggplot2':
##
## last_plot
```

```
## The following object is masked from 'package:stats':
##
## filter
```

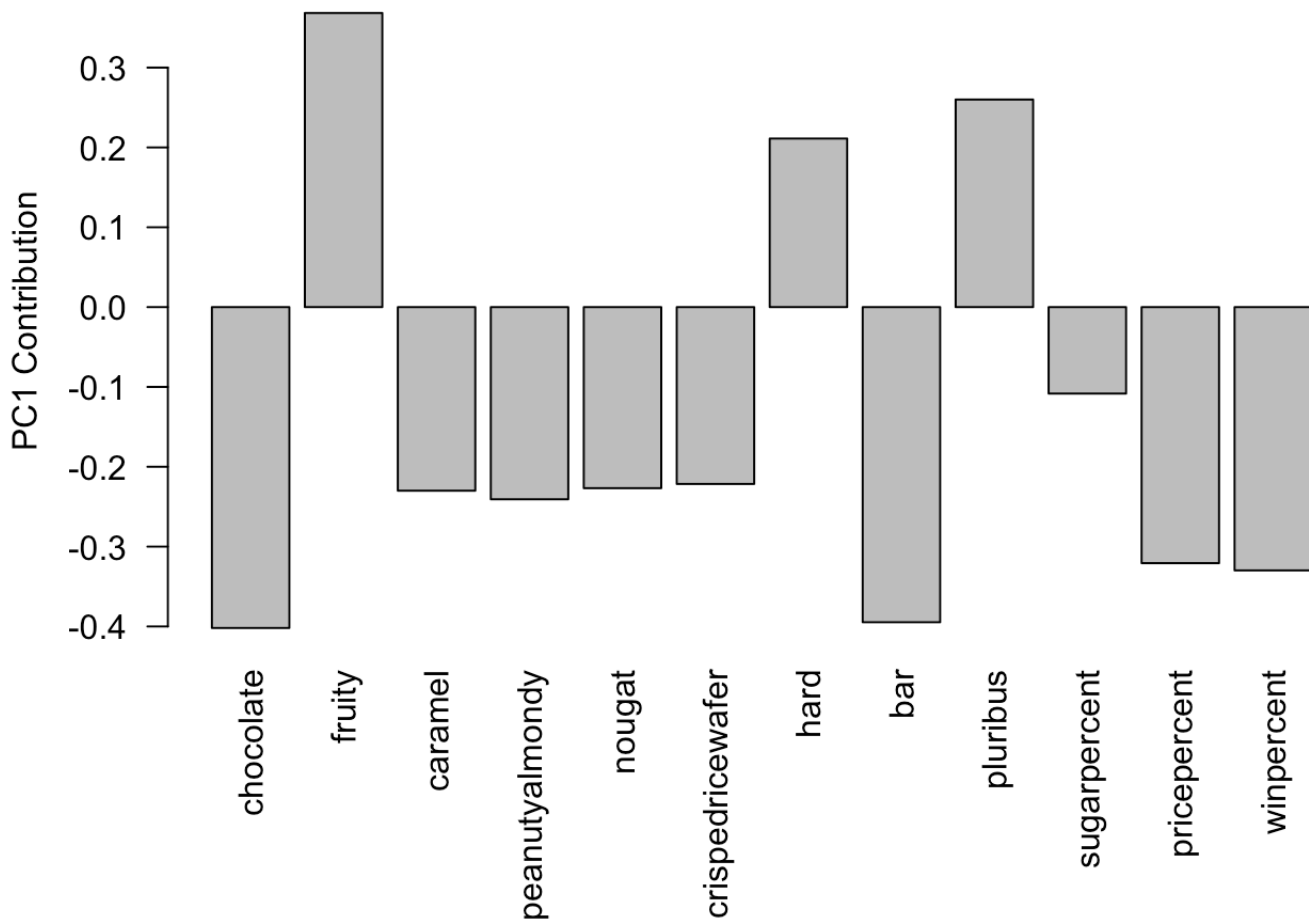
```
## The following object is masked from 'package:graphics':  
##  
## layout
```

```
ggplotly(p)
```



**Let's finish by taking a quick look at PCA
our loadings. Do these make sense to you?
Notice the opposite effects of chocolate
and fruity and the similar effects of
chocolate and bar (i.e. we already know
they are correlated).**

```
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/pluribus variables are picked up strongly by PC1 in the positive direction. These make sense because a lot of fruit candies are hard and come in packages with multiple candies in one. These variables also correlate with the data from the correlation plot we looked at earlier in this lab activity.