

Class 9: Candy Mini-Project

Jacob Hizon A17776679

Table of contents

Importing candy data	1
Exploratory analysis	4
Overall candy rankings	7
Taking a look at pricepercent	12
Exploring the correlation structure	14
Principal Component Analysis	15

Importing candy data

```
candy_file <- "candy-data.csv"
```

```
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?

There are 85 different types of candy in this dataset

```
nrow(candy)
```

```
[1] 85
```

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

There are 38 types of fruity candy.

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

```
[1] 38
```

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

```
[1] 81.64291
```

```
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 |>
  filter(row.names(candy)=="Twix") |>
  select(winpercent)
```

```
      winpercent
Twix    81.64291
```

Q3. What is your favorite candy (other than Twix) in the dataset and what is it's winpercent value?

```
candy["100 Grand", ]$winpercent
```

```
[1] 66.97173
```

66.98%

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_vari- able	n_miss- ing	com- plete_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	
peanutyal- mondy	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	
crispedrice- wafer	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 winpercent.

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

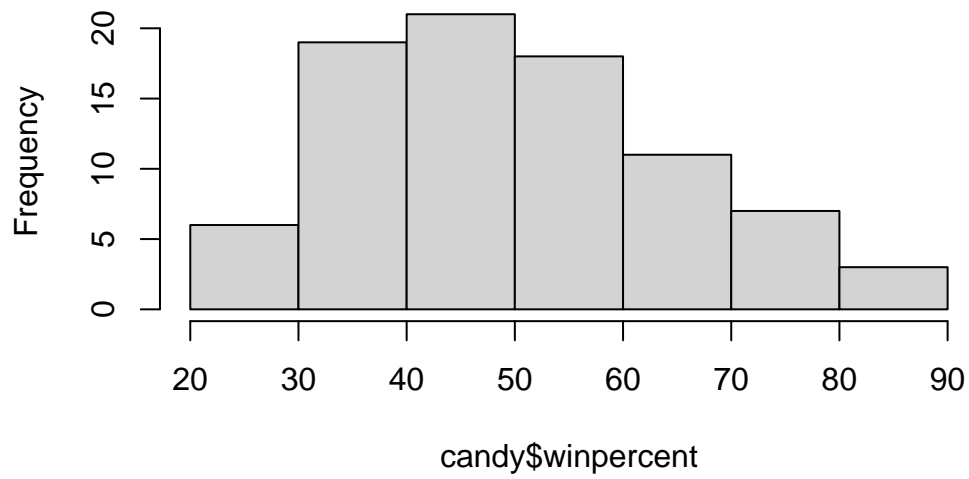
It represents zero not a chocolate, and a 1 represents true or is a chocolate candy.

Exploratory analysis

Q8. Plot a histogram of winpercent values using both base R and ggplot2

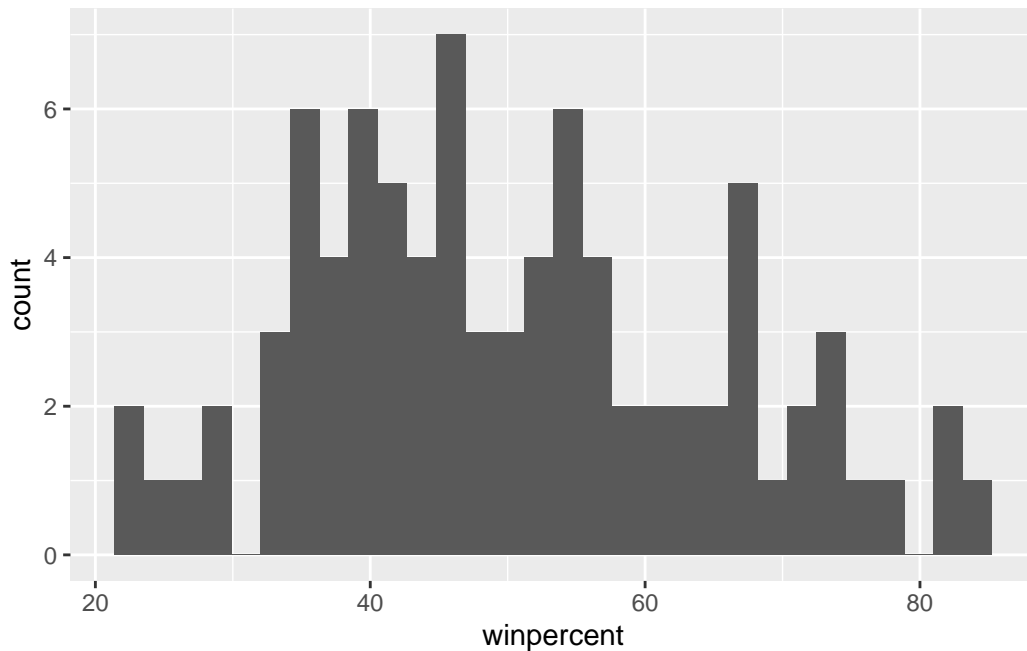
```
hist(candy$winpercent)
```

Histogram of candy\$winpercent



```
library(ggplot2)
ggplot(candy, aes(winpercent)) +
  geom_histogram()
```

`stat_bin()` using `bins = 30`. Pick better value `binwidth`.



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%?

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

```
[1] 47.82975
```

Below 50%.

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

```
chocwin <- candy$winpercent[as.logical(candy$chocolate)]
fruitwin <- candy$winpercent[as.logical(candy$fruity)]

t.test(chocwin, fruitwin)
```

Welch Two Sample t-test

```
data:  chocwin and fruitwin
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
```

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

```
mean(chocwin)
```

```
[1] 60.92153
```

```
mean(fruitwin)
```

```
[1] 44.11974
```

Chocolate candies have a higher average winpercent than fruity candies.

Q12. Is this difference statistically significant?

Yes this difference is statistically significant because the p-value is extremely low. Therefore, we can say that the difference in winpercent between chocolate and fruity candies is not due to random chance.

Overall candy rankings

```
order(candy$winpercent)
```

```
[1] 45  8 13 73 27 58 72  3 71 20 10 70 60 56 12 51 49 63  9 11 82 31 17 46 15
[26] 50 30 84 22 14 59 76 16 83 81 77 64  4 47 35 18 79 40 75 85 78  6 21  5 68
[51] 32 41 74 36 62 42 23 25  7 19 28 26 66 67 38 24 61 39 57 44 34  1 69  2 48
[76] 43 33 55 37 54 65 29 80 52 53
```

```
candy %>%
  arrange(winpercent) %>%
  head(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								

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

Nik L Nip, Boston Baked Beans, Chiclets, Super Bubble, Jawbusters

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

Reese's Peanut Butter cup, Reese's Minatures, Twix, Kit kat, snickers.

```
candy %>%
  arrange(desc(winpercent)) %>%
  head(5)
```

	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	
	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

	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

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

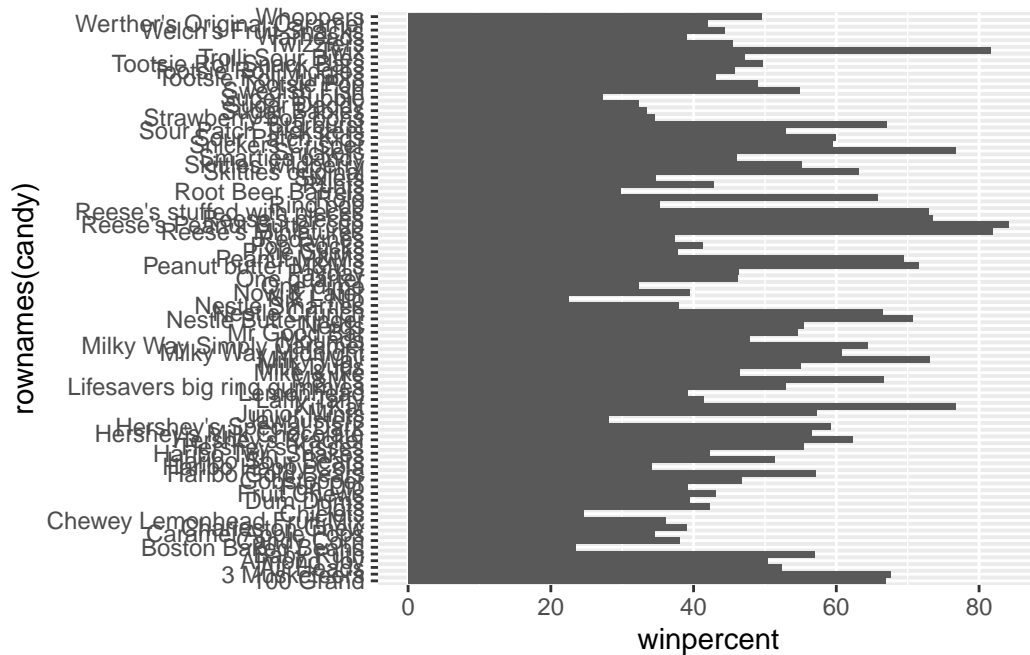
	chocolate	fruity	caramel	peanutyalmondy	nougat
Nik L Nip	0	1	0	0	0
Boston Baked Beans	0	0	0	1	0
Chiclets	0	1	0	0	0
Super Bubble	0	1	0	0	0
Jawbusters	0	1	0	0	0

	crispedricewafer	hard bar	pluribus	sugarpercent	pricepercent
Nik L Nip	0	0	0	1	0.197
Boston Baked Beans	0	0	0	1	0.313
Chiclets	0	0	0	1	0.046
Super Bubble	0	0	0	0	0.162
Jawbusters	0	1	0	1	0.093

	winpercent
Nik L Nip	22.44534
Boston Baked Beans	23.41782
Chiclets	24.52499
Super Bubble	27.30386
Jawbusters	28.12744

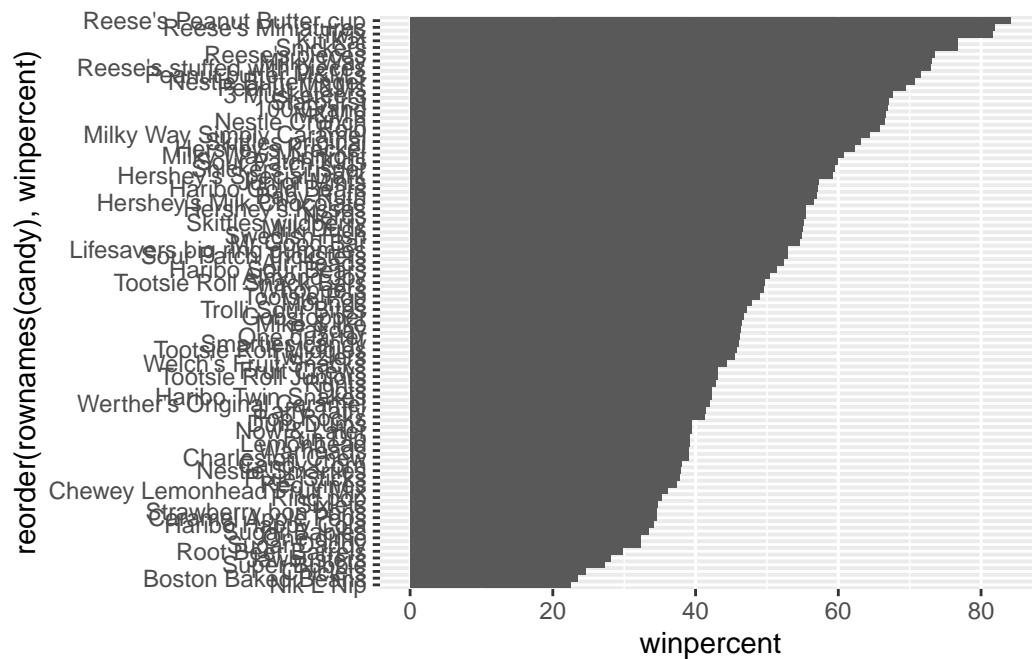
Q15. Make a first barplot of candy ranking based on winpercent values.

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



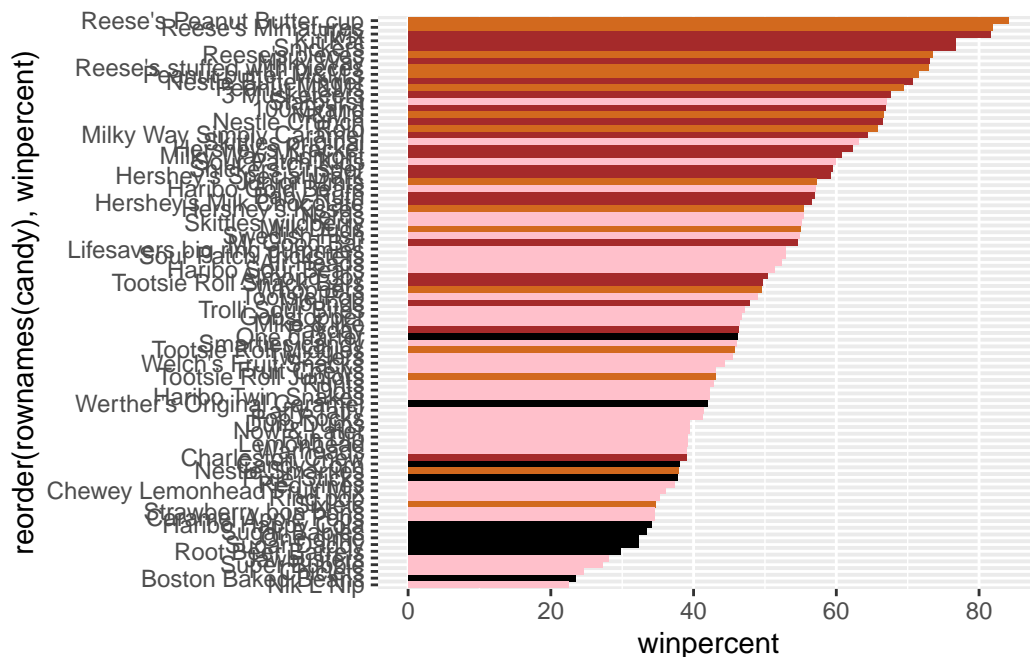
Q16. This is quite ugly, use the `reorder()` function to get the bars sorted by winpercent?

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



```
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?

Nik L Lip

Q18. What is the best ranked fruity candy?

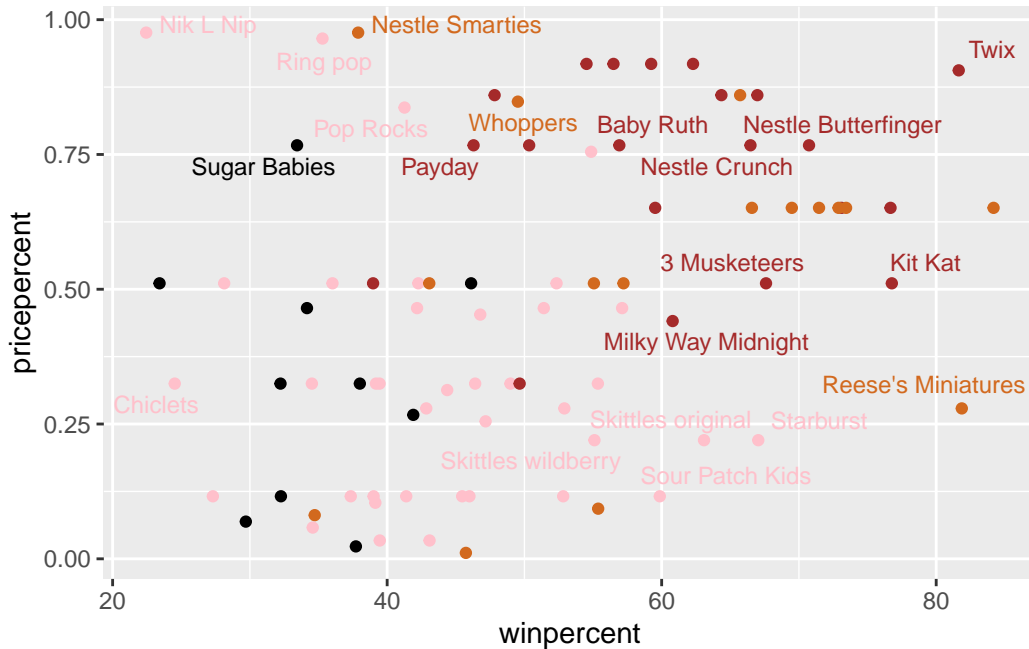
Starburst

Taking a look at pricepercent

```
library(ggrepel)

# How about a plot of win vs price
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: 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?

Reese's minatures

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

The tope 5 most expensive candy and the least popular are Nik L Nip, Nestle smarties, Ring pop, Hershey's Krackel, and hershey's milk chocolate respectively.

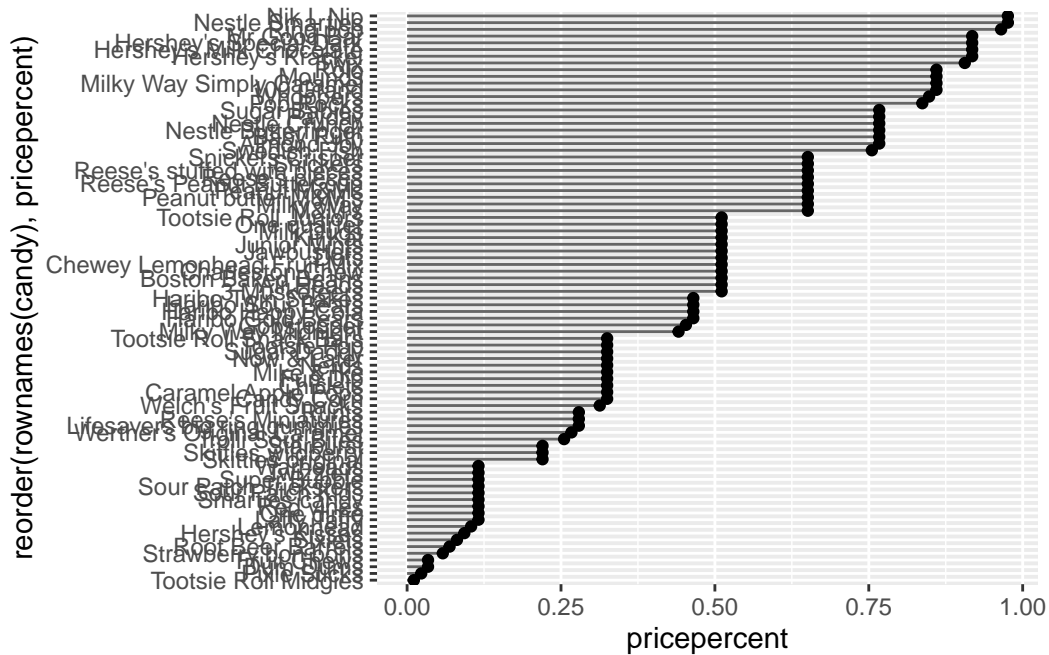
```
ord <- order(candy$pricepercent, decreasing = TRUE)
head( candy[ord,c(11,12)], n=5 )
```

	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()
```

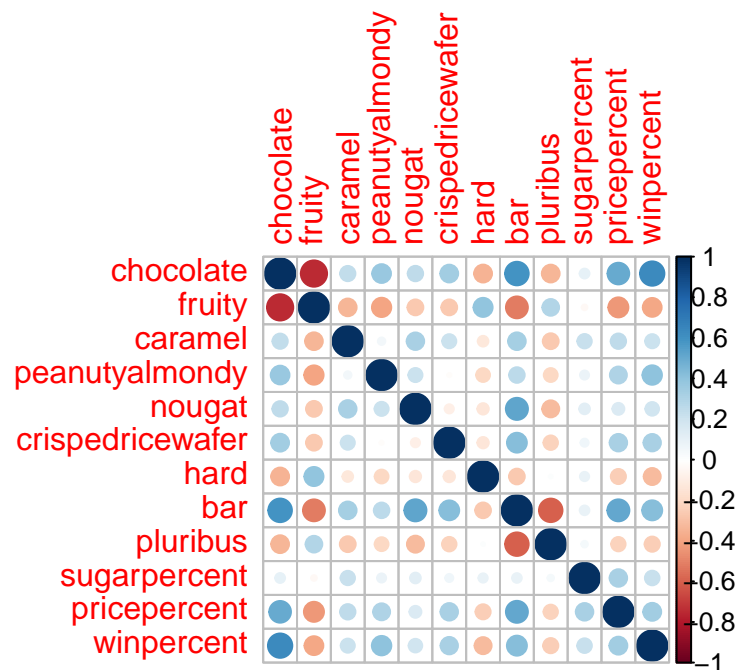


Exploring the correlation structure

```
library(corrplot)
```

```
corrplot 0.95 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 win percent.

Principal Component Analysis

```
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

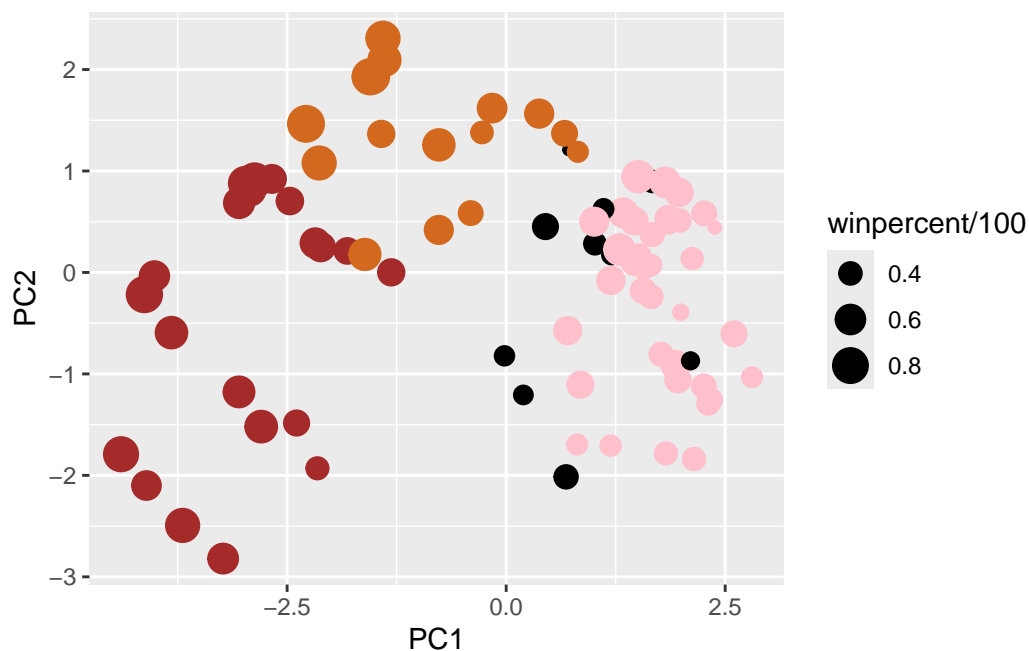
	PC8	PC9	PC10	PC11	PC12
Standard deviation	0.74530	0.67824	0.62349	0.43974	0.39760
Proportion of Variance	0.04629	0.03833	0.03239	0.01611	0.01317
Cumulative Proportion	0.89998	0.93832	0.97071	0.98683	1.00000

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

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)
```

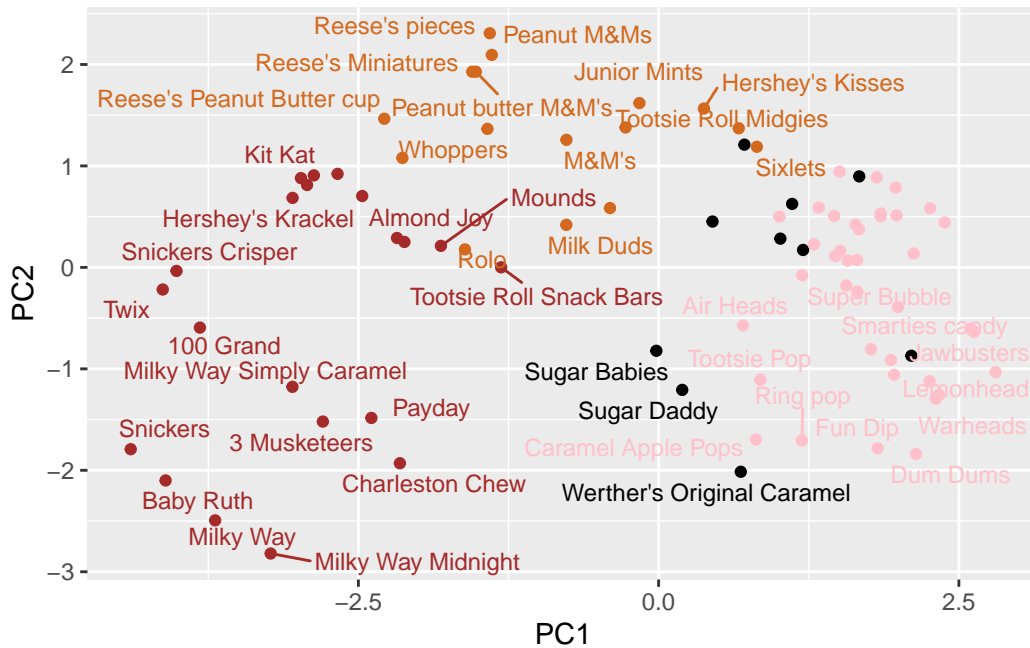
p



```
ggplot(pca$x) +
  aes(PC1, PC2, label = row.names(pca$x)) +
  geom_point(col = my_cols) +
  geom_text_repel(max.overlaps = 5, size = 3.3, col = my_cols)
```

Warning in geom_text_repel(max.overlaps = 5, size = 3.3, col = my_cols):
Ignoring unknown parameters: `max.overlaps`

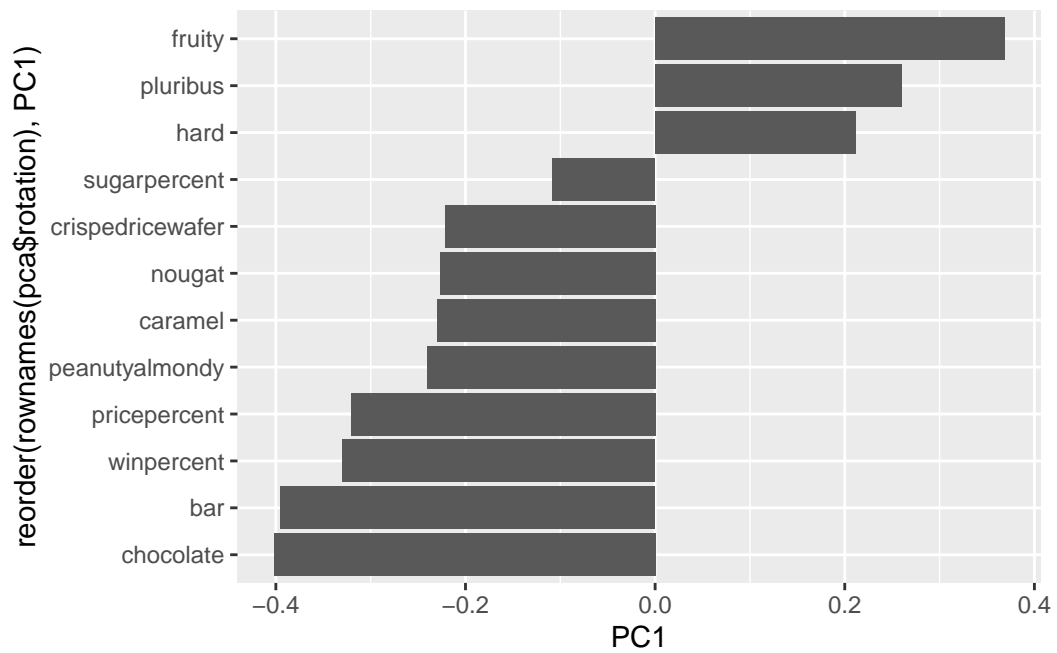
Warning: ggrepel: 42 unlabeled data points (too many overlaps). Consider increasing max.overlaps



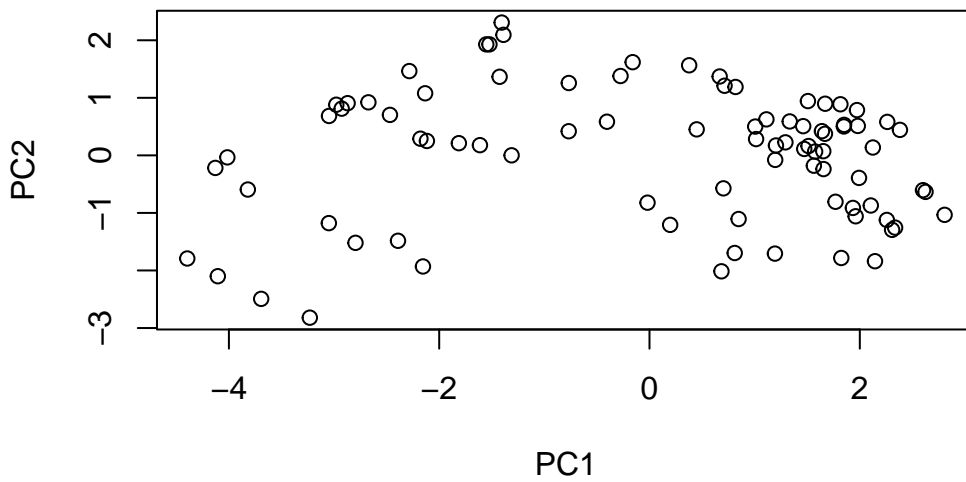
Q24. Complete the code to generate the loadings plot above. What original variables are picked up strongly by PC1 in the positive direction? Do these make sense to you? Where did you see this relationship highlighted previously?

Fruity, plurbis, and hard are picked up strongly by PC1 in the positive direction. We see this relationship in the correlation plot that are negatively correlated with chocolate.

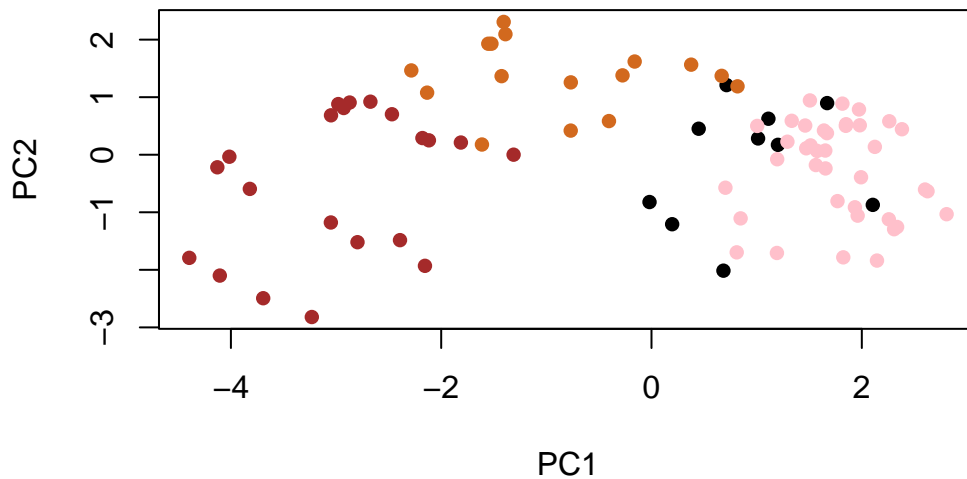
```
ggplot(pca$rotation) +
  aes(PC1,
    reorder(rownames(pca$rotation), PC1)) +
  geom_col()
```



```
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])
```

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

Q25. Based on your exploratory analysis, correlation findings, and PCA results, what combination of characteristics appears to make a “winning” candy? How do these different analyses (visualization, correlation, PCA) support or complement each other in reaching this conclusion?

Winning are usually chocolate and often containing nuts/caramel and are usually not fruity or hard. This was supported by three matrices exploratory plots which showed the overall distribution of popularity by win percentage, correlation matrix which elucidated how candy characteristics are related (e.g. price point and chocolate are positively correlated), and PCA which compressed and clustered the data into less dimensions and showed that nutty chocolates have higher winpercents.