# Class 10: Halloween Mini-Project

Xiaoyan Wang(A16454055)

**Exploratory Analysis of Halloween Candy**

```
candy_file<- "https://raw.githubusercontent.com/fivethirtyeight/data/master/candy-power-rank:
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

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

```r
candy["Skittles original", ]$winpercent
```

```
[1] 63.08514
```

- **Q4**. What is the `winpercent` value for "Kit Kat"?

```r
candy["Kit Kat", ]$winpercent
```

```
[1] 76.7686
```

- **Q5**. What is the `winpercent` value for "Tootsie Roll Snack Bars"?

```r
candy["Tootsie Roll Snack Bars", ]$winpercent
```

```
[1] 49.6535
```

```r
library(dplyr)
```

```
Warning: package 'dplyr' was built under R version 4.3.3
```

```
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
```

```r
candy |>
  filter(rownames(candy)=="Haribo Happy Cola")|>
  select(winpercent)
```

```
                  winpercent
Haribo Happy Cola   34.15896
```

Q: Find Fruity candy that have a winpercent >= 50

```
candy |>
  filter(winpercent> 50)|>
  filter(fruity==1)
```

|                           | chocolate | fruity | caramel | peanutyalmondy | nougat |
|---------------------------|-----------|--------|---------|----------------|--------|
| Air Heads                 | 0         | 1      | 0       | 0              | 0      |
| Haribo Gold Bears         | 0         | 1      | 0       | 0              | 0      |
| Haribo Sour Bears         | 0         | 1      | 0       | 0              | 0      |
| Lifesavers big ring gummies | 0       | 1      | 0       | 0              | 0      |
| Nerds                     | 0         | 1      | 0       | 0              | 0      |
| Skittles original         | 0         | 1      | 0       | 0              | 0      |
| Skittles wildberry        | 0         | 1      | 0       | 0              | 0      |
| Sour Patch Kids           | 0         | 1      | 0       | 0              | 0      |
| Sour Patch Tricksters     | 0         | 1      | 0       | 0              | 0      |
| Starburst                 | 0         | 1      | 0       | 0              | 0      |
| Swedish Fish              | 0         | 1      | 0       | 0              | 0      |

|                           | crispedricewafer | hard | bar | pluribus | sugarpercent |
|---------------------------|------------------|------|-----|----------|--------------|
| Air Heads                 | 0                | 0    | 0   | 0        | 0.906        |
| Haribo Gold Bears         | 0                | 0    | 0   | 1        | 0.465        |
| Haribo Sour Bears         | 0                | 0    | 0   | 1        | 0.465        |
| Lifesavers big ring gummies | 0              | 0    | 0   | 0        | 0.267        |
| Nerds                     | 0                | 1    | 0   | 1        | 0.848        |
| Skittles original         | 0                | 0    | 0   | 1        | 0.941        |
| Skittles wildberry        | 0                | 0    | 0   | 1        | 0.941        |
| Sour Patch Kids           | 0                | 0    | 0   | 1        | 0.069        |
| Sour Patch Tricksters     | 0                | 0    | 0   | 1        | 0.069        |
| Starburst                 | 0                | 0    | 0   | 1        | 0.151        |
| Swedish Fish              | 0                | 0    | 0   | 1        | 0.604        |

|                           | pricepercent | winpercent |
|---------------------------|--------------|------------|
| Air Heads                 | 0.511        | 52.34146   |
| Haribo Gold Bears         | 0.465        | 57.11974   |
| Haribo Sour Bears         | 0.465        | 51.41243   |
| Lifesavers big ring gummies | 0.279      | 52.91139   |
| Nerds                     | 0.325        | 55.35405   |
| Skittles original         | 0.220        | 63.08514   |
| Skittles wildberry        | 0.220        | 55.10370   |
| Sour Patch Kids           | 0.116        | 59.86400   |
| Sour Patch Tricksters     | 0.116        | 52.82595   |
| Starburst                 | 0.220        | 67.03763   |
| Swedish Fish              | 0.755        | 54.86111   |

```r
top.candy<- candy[candy$winpercent >50,]
top.candy[top.candy$fruity == 1,]
```

|                              | chocolate | fruity | caramel | peanutyalmondy | nougat |
|------------------------------|-----------|--------|---------|----------------|--------|
| Air Heads                    | 0         | 1      | 0       | 0              | 0      |
| Haribo Gold Bears            | 0         | 1      | 0       | 0              | 0      |
| Haribo Sour Bears            | 0         | 1      | 0       | 0              | 0      |
| Lifesavers big ring gummies  | 0         | 1      | 0       | 0              | 0      |
| Nerds                        | 0         | 1      | 0       | 0              | 0      |
| Skittles original            | 0         | 1      | 0       | 0              | 0      |
| Skittles wildberry           | 0         | 1      | 0       | 0              | 0      |
| Sour Patch Kids              | 0         | 1      | 0       | 0              | 0      |
| Sour Patch Tricksters        | 0         | 1      | 0       | 0              | 0      |
| Starburst                    | 0         | 1      | 0       | 0              | 0      |
| Swedish Fish                 | 0         | 1      | 0       | 0              | 0      |

|                              | crispedricewafer | hard | bar | pluribus | sugarpercent |
|------------------------------|------------------|------|-----|----------|--------------|
| Air Heads                    | 0                | 0    | 0   | 0        | 0.906        |
| Haribo Gold Bears            | 0                | 0    | 0   | 1        | 0.465        |
| Haribo Sour Bears            | 0                | 0    | 0   | 1        | 0.465        |
| Lifesavers big ring gummies  | 0                | 0    | 0   | 0        | 0.267        |
| Nerds                        | 0                | 1    | 0   | 1        | 0.848        |
| Skittles original            | 0                | 0    | 0   | 1        | 0.941        |
| Skittles wildberry           | 0                | 0    | 0   | 1        | 0.941        |
| Sour Patch Kids              | 0                | 0    | 0   | 1        | 0.069        |
| Sour Patch Tricksters        | 0                | 0    | 0   | 1        | 0.069        |
| Starburst                    | 0                | 0    | 0   | 1        | 0.151        |
| Swedish Fish                 | 0                | 0    | 0   | 1        | 0.604        |

|                              | pricepercent | winpercent |
|------------------------------|--------------|------------|
| Air Heads                    | 0.511        | 52.34146   |
| Haribo Gold Bears            | 0.465        | 57.11974   |
| Haribo Sour Bears            | 0.465        | 51.41243   |
| Lifesavers big ring gummies  | 0.279        | 52.91139   |
| Nerds                        | 0.325        | 55.35405   |
| Skittles original            | 0.220        | 63.08514   |
| Skittles wildberry           | 0.220        | 55.10370   |
| Sour Patch Kids              | 0.116        | 59.86400   |
| Sour Patch Tricksters        | 0.116        | 52.82595   |
| Starburst                    | 0.220        | 67.03763   |
| Swedish Fish                 | 0.755        | 54.86111   |

Quick overview of a given dataset:

```
#install.packages("skimr")
library("skimr")
```

Warning: package 'skimr' was built under R version 4.3.3

```
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 | |

Looks like the "winpercent" variable or column is masured on a different scale than everything else. I will need to scale my data before doing any analysis like PCA etc.

- **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 in 0-100 range, representing a percentage, whereas other columns from the dataset have a range of 0-1.

- **Q7**. What do you think a zero and one represent for the `candy$chocolate` column?
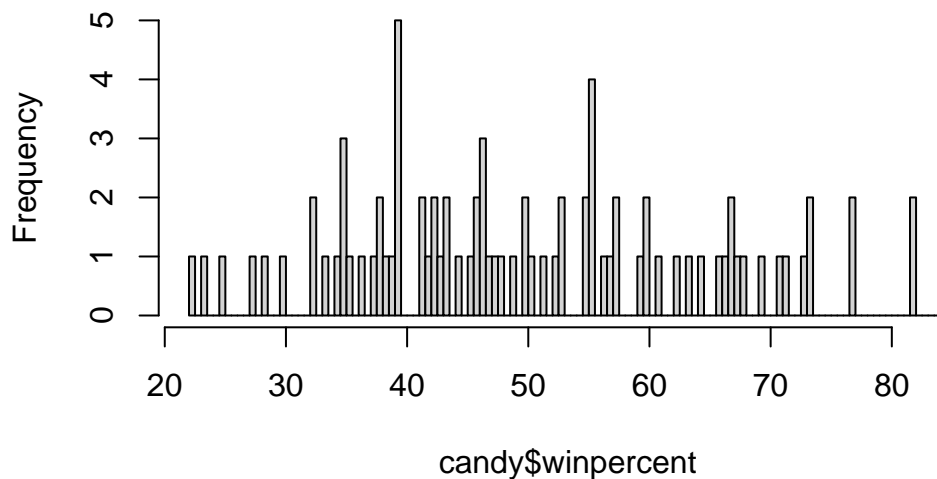
  The 0 means this candy is not chocolate, and 1 means this candy is/contains chocolate

- **Q8**. Plot a histogram of `winpercent` values

  We can do this in few ways. e.g. the "base" R 'hist()' function or with 'ggplot'.

```
hist(candy$winpercent, breaks = 100)
```

## Histogram of candy$winpercent



```
library(ggplot2)
```

```
Warning: package 'ggplot2' was built under R version 4.3.3
```

```
ggplot(candy)+
  aes(winpercent)+
  geom_histogram(binwidth =8)
```

- **Q9**. Is the distribution of `winpercent` values symmetrical?

  No

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

```
summary(candy$winpercent)
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  22.45   39.14   47.83   50.32   59.86   84.18
```

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

```
candy |>
  filter(as.logical(fruity))
```

|                          | chocolate | fruity | caramel | peanutyalmondy | nougat |
|--------------------------|-----------|--------|---------|----------------|--------|
| Air Heads                | 0         | 1      | 0       | 0              | 0      |
| Caramel Apple Pops       | 0         | 1      | 1       | 0              | 0      |
| Chewey Lemonhead Fruit Mix | 0       | 1      | 0       | 0              | 0      |
| Chiclets                 | 0         | 1      | 0       | 0              | 0      |
| Dots                     | 0         | 1      | 0       | 0              | 0      |
| Dum Dums                 | 0         | 1      | 0       | 0              | 0      |
| Fruit Chews              | 0         | 1      | 0       | 0              | 0      |
| Fun Dip                  | 0         | 1      | 0       | 0              | 0      |
| Gobstopper               | 0         | 1      | 0       | 0              | 0      |

| | | | | | |
|---|---|---|---|---|---|
| Haribo Gold Bears | 0 | 1 | 0 | 0 | 0 |
| Haribo Sour Bears | 0 | 1 | 0 | 0 | 0 |
| Haribo Twin Snakes | 0 | 1 | 0 | 0 | 0 |
| Jawbusters | 0 | 1 | 0 | 0 | 0 |
| Laffy Taffy | 0 | 1 | 0 | 0 | 0 |
| Lemonhead | 0 | 1 | 0 | 0 | 0 |
| Lifesavers big ring gummies | 0 | 1 | 0 | 0 | 0 |
| Mike & Ike | 0 | 1 | 0 | 0 | 0 |
| Nerds | 0 | 1 | 0 | 0 | 0 |
| Nik L Nip | 0 | 1 | 0 | 0 | 0 |
| Now & Later | 0 | 1 | 0 | 0 | 0 |
| Pop Rocks | 0 | 1 | 0 | 0 | 0 |
| Red vines | 0 | 1 | 0 | 0 | 0 |
| Ring pop | 0 | 1 | 0 | 0 | 0 |
| Runts | 0 | 1 | 0 | 0 | 0 |
| Skittles original | 0 | 1 | 0 | 0 | 0 |
| Skittles wildberry | 0 | 1 | 0 | 0 | 0 |
| Smarties candy | 0 | 1 | 0 | 0 | 0 |
| Sour Patch Kids | 0 | 1 | 0 | 0 | 0 |
| Sour Patch Tricksters | 0 | 1 | 0 | 0 | 0 |
| Starburst | 0 | 1 | 0 | 0 | 0 |
| Strawberry bon bons | 0 | 1 | 0 | 0 | 0 |
| Super Bubble | 0 | 1 | 0 | 0 | 0 |
| Swedish Fish | 0 | 1 | 0 | 0 | 0 |
| Tootsie Pop | 1 | 1 | 0 | 0 | 0 |
| Trolli Sour Bites | 0 | 1 | 0 | 0 | 0 |
| Twizzlers | 0 | 1 | 0 | 0 | 0 |
| Warheads | 0 | 1 | 0 | 0 | 0 |
| Welch's Fruit Snacks | 0 | 1 | 0 | 0 | 0 |

| | crispedricewafer | hard | bar | pluribus | sugarpercent |
|---|---|---|---|---|---|
| Air Heads | 0 | 0 | 0 | 0 | 0.906 |
| Caramel Apple Pops | 0 | 0 | 0 | 0 | 0.604 |
| Chewey Lemonhead Fruit Mix | 0 | 0 | 0 | 1 | 0.732 |
| Chiclets | 0 | 0 | 0 | 1 | 0.046 |
| Dots | 0 | 0 | 0 | 1 | 0.732 |
| Dum Dums | 0 | 1 | 0 | 0 | 0.732 |
| Fruit Chews | 0 | 0 | 0 | 1 | 0.127 |
| Fun Dip | 0 | 1 | 0 | 0 | 0.732 |
| Gobstopper | 0 | 1 | 0 | 1 | 0.906 |
| Haribo Gold Bears | 0 | 0 | 0 | 1 | 0.465 |
| Haribo Sour Bears | 0 | 0 | 0 | 1 | 0.465 |
| Haribo Twin Snakes | 0 | 0 | 0 | 1 | 0.465 |
| Jawbusters | 0 | 1 | 0 | 1 | 0.093 |

| | | | | | |
|---|---|---|---|---|---|
| Laffy Taffy | 0 | 0 | 0 | 0 | 0.220 |
| Lemonhead | 0 | 1 | 0 | 0 | 0.046 |
| Lifesavers big ring gummies | 0 | 0 | 0 | 0 | 0.267 |
| Mike & Ike | 0 | 0 | 0 | 1 | 0.872 |
| Nerds | 0 | 1 | 0 | 1 | 0.848 |
| Nik L Nip | 0 | 0 | 0 | 1 | 0.197 |
| Now & Later | 0 | 0 | 0 | 1 | 0.220 |
| Pop Rocks | 0 | 1 | 0 | 1 | 0.604 |
| Red vines | 0 | 0 | 0 | 1 | 0.581 |
| Ring pop | 0 | 1 | 0 | 0 | 0.732 |
| Runts | 0 | 1 | 0 | 1 | 0.872 |
| Skittles original | 0 | 0 | 0 | 1 | 0.941 |
| Skittles wildberry | 0 | 0 | 0 | 1 | 0.941 |
| Smarties candy | 0 | 1 | 0 | 1 | 0.267 |
| Sour Patch Kids | 0 | 0 | 0 | 1 | 0.069 |
| Sour Patch Tricksters | 0 | 0 | 0 | 1 | 0.069 |
| Starburst | 0 | 0 | 0 | 1 | 0.151 |
| Strawberry bon bons | 0 | 1 | 0 | 1 | 0.569 |
| Super Bubble | 0 | 0 | 0 | 0 | 0.162 |
| Swedish Fish | 0 | 0 | 0 | 1 | 0.604 |
| Tootsie Pop | 0 | 1 | 0 | 0 | 0.604 |
| Trolli Sour Bites | 0 | 0 | 0 | 1 | 0.313 |
| Twizzlers | 0 | 0 | 0 | 0 | 0.220 |
| Warheads | 0 | 1 | 0 | 0 | 0.093 |
| Welch's Fruit Snacks | 0 | 0 | 0 | 1 | 0.313 |

| | pricepercent | winpercent |
|---|---|---|
| Air Heads | 0.511 | 52.34146 |
| Caramel Apple Pops | 0.325 | 34.51768 |
| Chewey Lemonhead Fruit Mix | 0.511 | 36.01763 |
| Chiclets | 0.325 | 24.52499 |
| Dots | 0.511 | 42.27208 |
| Dum Dums | 0.034 | 39.46056 |
| Fruit Chews | 0.034 | 43.08892 |
| Fun Dip | 0.325 | 39.18550 |
| Gobstopper | 0.453 | 46.78335 |
| Haribo Gold Bears | 0.465 | 57.11974 |
| Haribo Sour Bears | 0.465 | 51.41243 |
| Haribo Twin Snakes | 0.465 | 42.17877 |
| Jawbusters | 0.511 | 28.12744 |
| Laffy Taffy | 0.116 | 41.38956 |
| Lemonhead | 0.104 | 39.14106 |
| Lifesavers big ring gummies | 0.279 | 52.91139 |
| Mike & Ike | 0.325 | 46.41172 |

```
Nerds                           0.325   55.35405
Nik L Nip                       0.976   22.44534
Now & Later                     0.325   39.44680
Pop Rocks                       0.837   41.26551
Red vines                       0.116   37.34852
Ring pop                        0.965   35.29076
Runts                           0.279   42.84914
Skittles original               0.220   63.08514
Skittles wildberry              0.220   55.10370
Smarties candy                  0.116   45.99583
Sour Patch Kids                 0.116   59.86400
Sour Patch Tricksters           0.116   52.82595
Starburst                       0.220   67.03763
Strawberry bon bons             0.058   34.57899
Super Bubble                    0.116   27.30386
Swedish Fish                    0.755   54.86111
Tootsie Pop                     0.325   48.98265
Trolli Sour Bites               0.255   47.17323
Twizzlers                       0.116   45.46628
Warheads                        0.116   39.01190
Welch's Fruit Snacks            0.313   44.37552
```

```r
choc.candy <- candy %>% filter(as.logical(chocolate))
fruit.candy <- candy %>% filter(as.logical(fruity))
mean(choc.candy$winpercent, na.rm = TRUE)
```

```
[1] 60.92153
```

```r
mean(fruit.candy$winpercent, na.rm = TRUE)
```

```
[1] 44.11974
```

- **Q12**. Is this difference statistically significant?

```r
t.test(choc.candy$winpercent, fruit.candy$winpercent)
```

```
    Welch Two Sample t-test

data:  choc.candy$winpercent and fruit.candy$winpercent
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 p-value is small, indicating that we can reject the null hypothesis and state there is significant differences between the mean between winpercent of fruity candy and chocolate candy.

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

```
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?

```
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
```
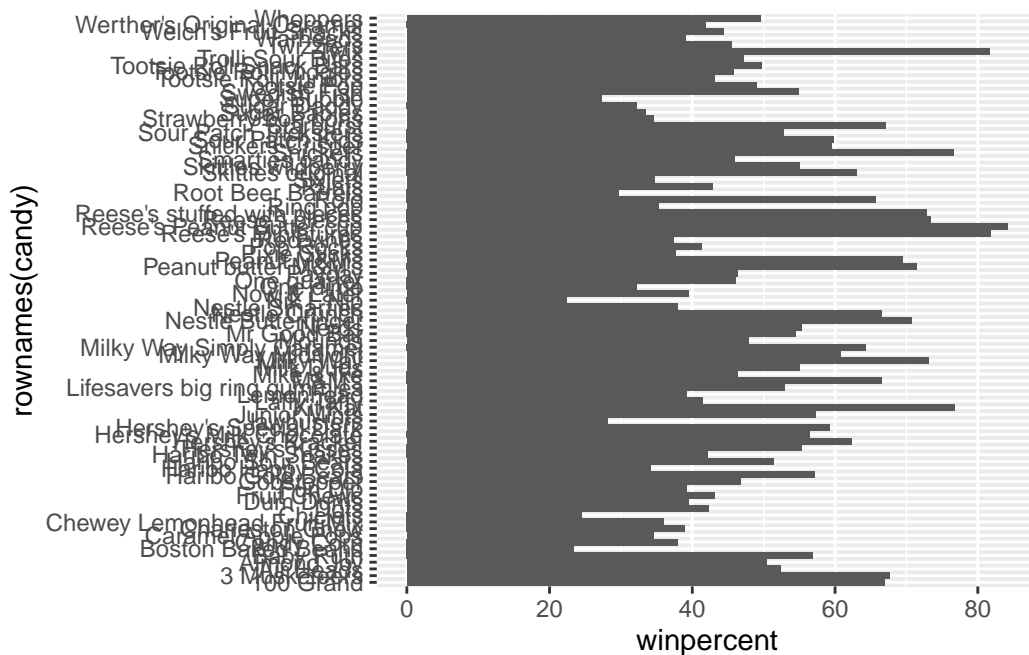
```
Kit Kat                                    1    0   1        0         0.313
Snickers                                   0    0   1        0         0.546
                                 pricepercent winpercent
Reese's Peanut Butter cup              0.651    84.18029
Reese's Miniatures                     0.279    81.86626
Twix                                   0.906    81.64291
Kit Kat                                0.511    76.76860
Snickers                               0.651    76.67378
```

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

  lets do a barplot

```
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(x=winpercent,
      y= reorder(rownames(candy),winpercent),
      fill=chocolate)+
  geom_col()
```

I want a more custom color scheme where I can see both chocolate and bar and fruity etc. all from the one plot. To do this, we can roll our own color vector...

```
mycol<- rep("black",nrow(candy))
mycol[as.logical(candy$chocolate)] <- "chocolate"
mycol[as.logical(candy$bar)] <- "red"
mycol[as.logical(candy$fruity)] <- "pink"
mycol[row.names(candy)=="Skittles original"] <- "blue"
```
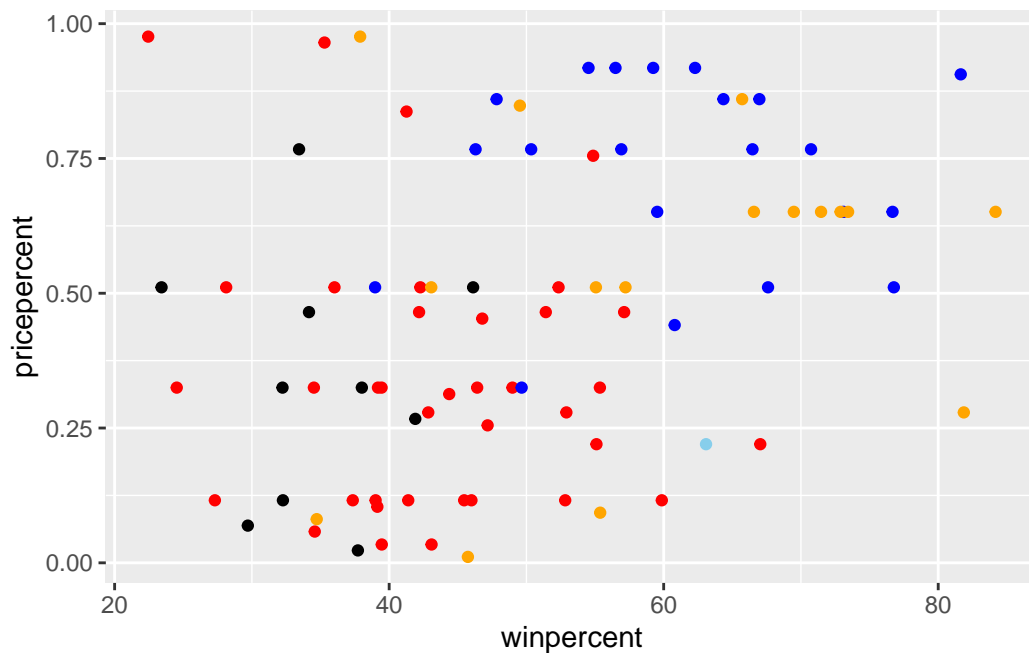
```
ggplot(candy)+
  aes(x=winpercent,
      y= reorder(rownames(candy),winpercent),
      fill=chocolate)+
  geom_col(fill=mycol)
```

plot of winpercent vs pricepercent to see what would be the candy to by

```r
mycol<- rep("black",nrow(candy))
mycol[as.logical(candy$chocolate)] <- "orange"
mycol[as.logical(candy$bar)] <- "blue"
mycol[as.logical(candy$fruity)] <- "red"
mycol[row.names(candy)=="Skittles original"] <- "skyblue"
```

```r
ggplot(candy)+
  aes(x= winpercent,
      y= pricepercent)+
  geom_point(col=mycol)
```

**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
```