```{r}
---
title: "Class 10"
format: html
---

1. Importing Data
```{r}
candy_file <- "https://raw.githubusercontent.com/fivethirtyeight/data/master/candy-power-ranking/candy-data.csv"
candy = read.csv(candy_file)
head(candy)
```
```

Description: df [6 × 13]

| | competitorname <chr> | chocolate <int> | fruity <int> | caramel <int> | peanutyalmondy <int> | nougat <int> | crispedricewafer <int> | hard <int> | bar <int> |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 100 Grand | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 2 | 3 Musketeers | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| 3 | One dime | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | One quarter | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | Air Heads | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | Almond Joy | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |

6 rows | 1–10 of 13 columns

Q1. How many different candy types are in this dataset?
A1: 9 types

Q2. How many fruity candy types are in the dataset?
A2: 37 types

2. What is your favorate candy?
```{r}
candy["Kit Kat", ]$winpercent
```

[1] NA

Q3. What is your favorite candy in the dataset and what is it's winpercent value?
A3: Almond Joy Winpercent: 50.347545

Q4. What is the winpercent value for "Kit Kat"?
A4: 76.7686

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

```
31  Q5. What is the winpercent value for "Tootsie Roll Snack Bars"?
32  A5: 49.653503
33
34  SKIMR
35 ▾ ```{r}
36  install.packages("skimr")
37  library("skimr")
38  skim(candy)
39 ▴ ```
```

R Console    one_skim_df 1 × 8    one_skim_df 12 × 11

A tibble: **12 × 11**

| | skim_variable <chr> | n_missing <int> | complete_rate <dbl> | mean <dbl> | sd <dbl> | p0 <dbl> | p25 <dbl> | p50 <dbl> | p75 <dbl> |
|---|---|---|---|---|---|---|---|---|---|
| 1 | chocolate | 0 | 1 | 0.43529412 | 0.4987379 | 0.00000 | 0.00000 | 0.00000 | 1.000 |
| 2 | fruity | 0 | 1 | 0.44705882 | 0.5001400 | 0.00000 | 0.00000 | 0.00000 | 1.000 |
| 3 | caramel | 0 | 1 | 0.16470588 | 0.3731162 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 4 | peanutyalmondy | 0 | 1 | 0.16470588 | 0.3731162 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 5 | nougat | 0 | 1 | 0.08235294 | 0.2765332 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 6 | crispedricewafer | 0 | 1 | 0.08235294 | 0.2765332 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 7 | hard | 0 | 1 | 0.17647059 | 0.3834825 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 8 | bar | 0 | 1 | 0.24705882 | 0.4338609 | 0.00000 | 0.00000 | 0.00000 | 0.000 |
| 9 | pluribus | 0 | 1 | 0.51764706 | 0.5026540 | 0.00000 | 0.00000 | 1.00000 | 1.000 |
| 10 | sugarpercent | 0 | 1 | 0.47864705 | 0.2827779 | 0.01100 | 0.22000 | 0.46500 | 0.732 |

1–10 of 12 rows | 1–10 of 11 columns      Previous 1 2 Next

```
40  Q6. Is there any variable/column that looks to be on a different scale to the majority of the other columns in the dataset?
41  A6: Pluribus average values significantly greater
42
43  Q7. What do you think a zero and one represent for the candy$chocolate column?
44  A7: A "1" means "yes" the candy is chocolate, "0" means "No"
45
46  Q8. Plot a histogram of winpercent values
47 ▾ ```{r}
48  m <- ggplot(candy_file, aes(x = winpercent)) +
49      geom_histogram()
50 ▴ ```
51  Q9. Is the distribution of winpercent values symmetrical?
52  A9: No
53  Q10. Is the center of the distribution above or below 50%?
54  A10: Above
55  Q11. On average is chocolate candy higher or lower ranked than fruit candy?
```

```
55  Q11. On average is chocolate candy higher or lower ranked than fruit candy?
56  A11: Higher
57  Q12. Is this difference statistically significant?
58
59  3. Overall Candy Rankings
60  ```{r}
61  tail(candy[order(candy$winpercent),], n=5)
62  ```
```

Description: df [5 × 13]

| | competitorname<br><chr> | chocolate<br><int> | fruity<br><int> | caramel<br><int> | peanutyalmondy<br><int> | nougat<br><int> | crispedricewafer<br><int> | hard<br><int> | bar<br><int> |
|---|---|---|---|---|---|---|---|---|---|
| 65 | Snickers | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 |
| 29 | Kit Kat | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| 80 | Twix | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 52 | ReeseÕs Miniatures | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 53 | ReeseÕs Peanut Butter cup | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

5 rows | 1–10 of 13 columns

```
63  Q13. What are the five least liked candy types in this set?
64  A13: Least to most: Nik L Nip, Chiclets, Super Bubble, and Jawbusters
65  Q14. What are the top 5 all time favorite candy types out of this set?
66  Q14: Low to high: Snickers, Kit Kat, Twix, Reese's Miniatures, Reese's Peanut Butter Cup
67
68  Q15. Make a first barplot of candy ranking based on winpercent values.
69  ```{r}
70  library()
71
72  ggplot(candy) +
73    aes(winpercent, reorder(rownames(candy),winpercent)) +
74    geom_col(fill=my_cols)
75  ```
76  Q17. What is the worst ranked chocolate candy?
77  A17: Sixlets
78  Q18. What is the best ranked fruity candy?
79  A18: Starbursts
80
81  4. Taking a look at pricepercent
82  ```{r}
83  library(ggrepel)
84
85  # How about a plot of price vs win
86  ggplot(candy) +
87    aes(winpercent, pricepercent, label=rownames(candy)) +
88    geom_point(col=my_cols) +
89    geom_text_repel(col=my_cols, size=3.3, max.overlaps = 5)
```

91 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?
92 A19: Fruity Candy
93 Q20. What are the top 5 most expensive candy types in the dataset and of these which is the least popular?
94 A20: Most expensive: Nik n Lip, Nestle Smarties, Ring pop, Hershey's Krackel, Hershey's Milk Chocolate. Nik n Lip is least
   popular

```r
95 ```{r}
96 ord <- order(candy$pricepercent, decreasing = TRUE)
97 head( candy[ord,c(11,12)], n=5 )
98 ```
```

Description: df [5 × 2]

| | sugarpercent<br><dbl> | pricepercent<br><dbl> |
|---|---|---|
| 45 | 0.197 | 0.976 |
| 63 | 0.267 | 0.976 |
| 56 | 0.732 | 0.965 |
| 24 | 0.430 | 0.918 |
| 25 | 0.430 | 0.918 |

5 rows

99
100 5 Exploring the correlation structure

```r
101 ```{r}
102 install.packages(corrplot)
103 library(corrplot)
104 ```
```

105
106 Q22. Examining this plot what two variables are anti-correlated (i.e. have minus values)?
107 A22: Anti is chocolatey and fruity candy
108 Q23. Similarly, what two variables are most positively correlated?
109 A23: Most positively correlated is Nougat & Bar
110
111 6. Principal Component Analysis

```r
112 ```{r}
113 plot(pca$x[,1:2], col=my_cols, pch=16)
114 ```
```

115

```r
116 ```{r}
117 p <- ggplot(my_data) +
118       aes(x=PC1, y=PC2,
119           size=winpercent/100,
120           text=rownames(my_data),
121           label=rownames(my_data)) +
122       geom_point(col=my_cols)
123
```

| | <dbl> | <dbl> |
|---|---|---|
| 45 | 0.197 | 0.976 |
| 63 | 0.267 | 0.976 |
| 56 | 0.732 | 0.965 |
| 24 | 0.430 | 0.918 |
| 25 | 0.430 | 0.918 |

5 rows

```
99
100  5 Exploring the correlation structure
101 ▾ ```{r}
102  install.packages(corrplot)
103  library(corrplot)
104 ▴ ```
105
106  Q22. Examining this plot what two variables are anti-correlated (i.e. have minus values)?
107  A22: Anti is chocolatey and fruity candy
108  Q23. Similarly, what two variables are most positively correlated?
109  A23: Most positively correlated is Nougat & Bar
110
111  6. Principal Component Analysis
112 ▾ ```{r}
113  plot(pca$x[,1:2], col=my_cols, pch=16)
114 ▴ ```
115
116 ▾ ```{r}
117  p <- ggplot(my_data) +
118        aes(x=PC1, y=PC2,
119           size=winpercent/100,
120           text=rownames(my_data),
121           label=rownames(my_data)) +
122        geom_point(col=my_cols)
123
124 ▴ ```
125
126 ▾ ```{r}
127  library(ggrepel)
128
129  p + geom_text_repel(size=3.3, col=my_cols, max.overlaps = 7)  +
130    theme(legend.position = "none") +
131    labs(title="Halloween Candy PCA Space",
132        subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brown), fruity (red), other (black)",
133        caption="Data from 538")
134 ▴ ```
135
136  Q24. What original variables are picked up strongly by PC1 in the positive direction? Do these make sense to you?
137  A24: Fruity and Hard are picked up strongly in the positive direction. Generally fruity candies are harder to chew.
138
```