# Class 9 Halloween Mini Project

Katelyn (PID: A16682595)

## Importing Candy Data

```
candy_data <- "candy-data.csv"
candy <- read.csv(candy_data, 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 fruit candy types are in the dataset?*

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

```
[1] 38
```

## What is Your Favorite Candy?

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

```r
candy["Whoppers",]$winpercent
```

```
[1] 49.52411
```

*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
# Loading and using the skimr package
library("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_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?*
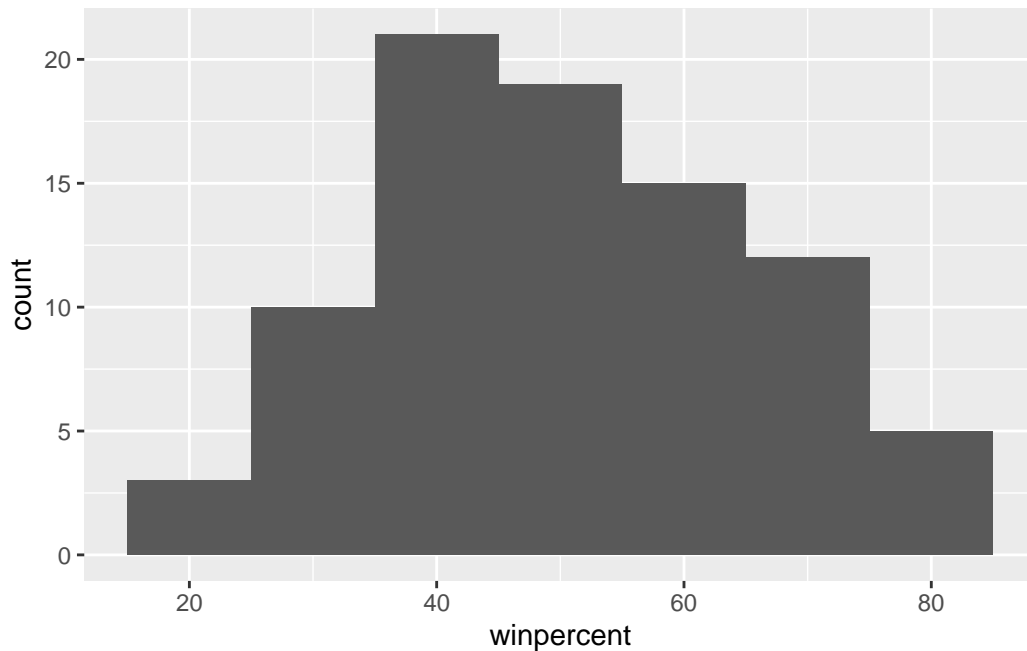
Win percent goes all the way up to 82.18, but all other columns are binary(either 0 or 1) or between 0 and 1.

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

0 = "no", 1 = "yes"

*Q8. Plot a histogram of winpercent values*

```
library(ggplot2)
ggplot(candy, aes(winpercent)) +
  geom_histogram(binwidth = 10)
```

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

Not quite. It looks a little skewed.

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

The center of the distribution looks below 50%.

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

```r
# Turning chocolate into T/F values
choc.Inds <- as.logical(candy$chocolate)
# Identifying the win percents of chocolate candy
choc.win <- candy[choc.Inds,]$winpercent
# Taking the average win percent of chocolate
mean(choc.win)
```

```
[1] 60.92153
```

```r
# Doing the same for fruit candy
fruit.Inds <- as.logical(candy$fruity)
fruit.win <- candy[fruit.Inds,]$winpercent
mean(fruit.win)
```

```
[1] 44.11974
```

```
#Is chocolate candy higher ranked than fruit candy on average?
mean(choc.win) > mean(fruit.win)
```

```
[1] TRUE
```

*Q12. Is this difference statistically significant?*

```
t.test(choc.win, fruit.win)
```

```
	Welch Two Sample t-test

data:  choc.win and fruit.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
```

Yes. The p-value is really small, indicating the difference is significant.

## Overall Candy Rankings

*Q13. What are the five least liked candy types in 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
```

```
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.*

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



*Q16. Use the `reorder()` function to get the bars sorted by `winpercent`*

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



```
# Adding a color vector, mycols
mycols <- rep("black", nrow(candy))
mycols[as.logical(candy$chocolate)] <- "chocolate"
mycols[as.logical(candy$bar)] <- "brown"
mycols[as.logical(candy$fruity)] <- "pink"
```

```
# Generating a colored graph
ggplot(candy) +
  aes(winpercent, reorder(rownames(candy), winpercent)) +
  geom_col(fill = mycols)
```

*Q17. What is the worst ranked chocolate candy?*

Sixlets

*Q18. What is the best ranked fruity candy?*

Starburst

## Taking a Look at Pricepercent

```
library(ggrepel)

# How about a plot of price vs win
ggplot(candy) +
  aes(winpercent, pricepercent, label=rownames(candy)) +
  geom_point(col=mycols) +
  geom_text_repel(col=mycols, 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 miniatures are the most bang for your buck.

**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
```

Nik L Nips are the least popular and one of the most expensive.

```
# Pricepercent ggplot:
ggplot(candy) +
```

```
aes(pricepercent, reorder(rownames(candy), pricepercent)) +
geom_point() +
geom_segment(aes(yend = reorder(rownames(candy), pricepercent), xend = 0), col = "grey40
```



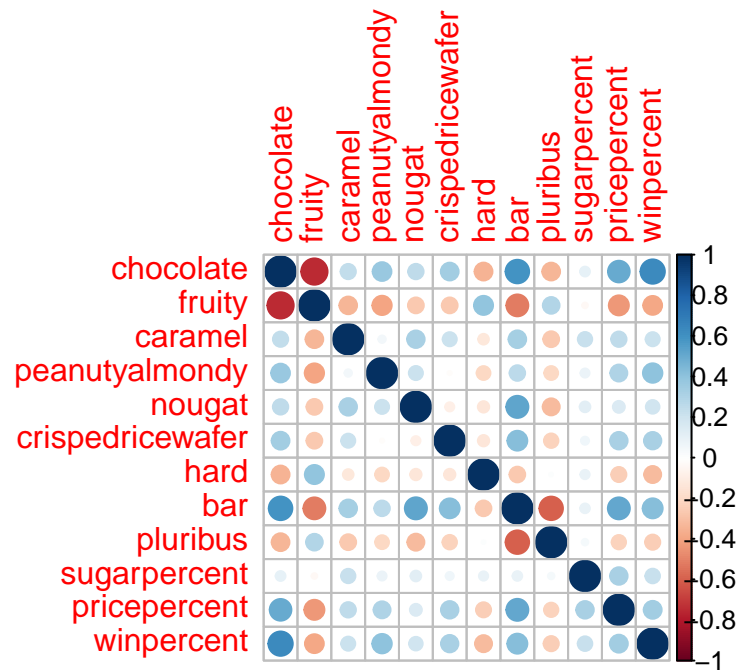## Exploring the Correlation Structure

```
library(corrplot)
```

corrplot 0.92 loaded

```
#calculate correlation
cij <- cor(candy)
#plot correlation
corrplot(cij)
```

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

Chocolate and fruity are strongly anti-correlated.

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

Chocolate being a bar and having a high win percent are most positively correlated.

## Principal Component Analysis

```
pca <- prcomp(candy, scale = T)
summary(pca)
```
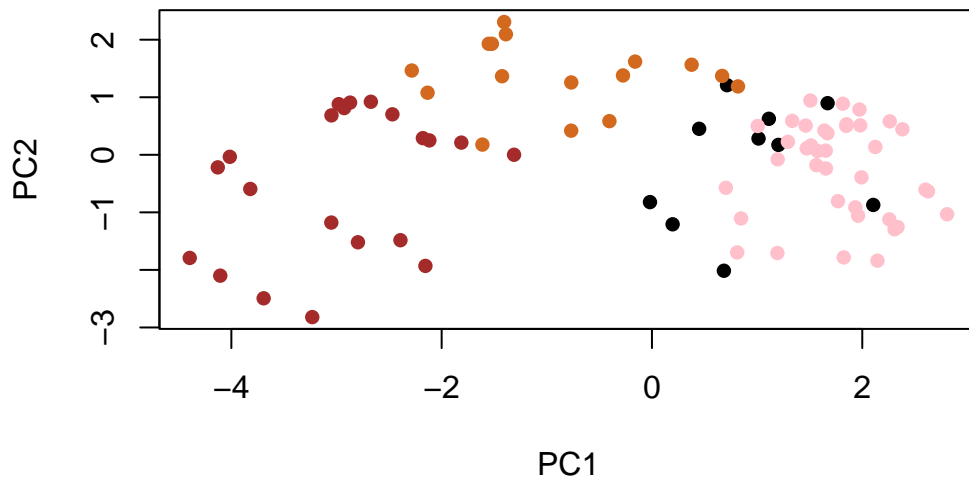
```
Importance of components:
                          PC1    PC2    PC3     PC4    PC5     PC6     PC7
Standard deviation     2.0788 1.1378 1.1092 1.07533 0.9518 0.81923 0.81530
Proportion of Variance 0.3601 0.1079 0.1025 0.09636 0.0755 0.05593 0.05539
Cumulative Proportion  0.3601 0.4680 0.5705 0.66688 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
```
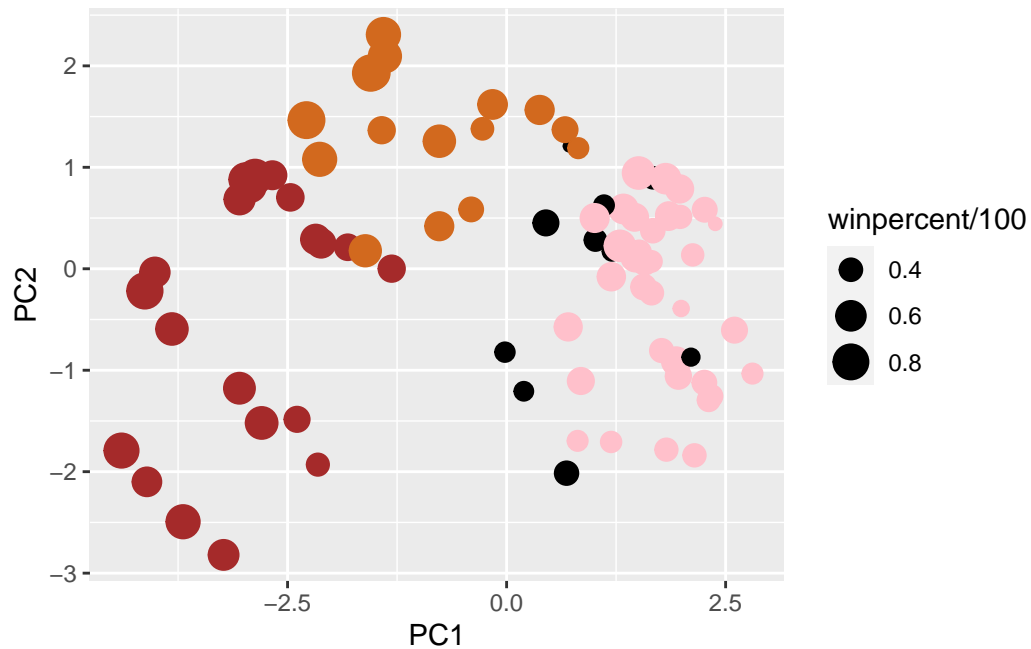
```
# Comparing PC1 with PC2 in basic R:
plot(pca$x[,1:2], col = mycols, pch = 16)
```



```
# Make a new dataframe for ggplot to read
my_data <- cbind(candy, pca$x[,1:3])

# PC2 vs. PC1 in ggplot:
p <- ggplot(my_data, aes(PC1, PC2, size = winpercent/100, text=rownames(my_data),
            label=rownames(my_data))) +
  geom_point(col = mycols)
p
```
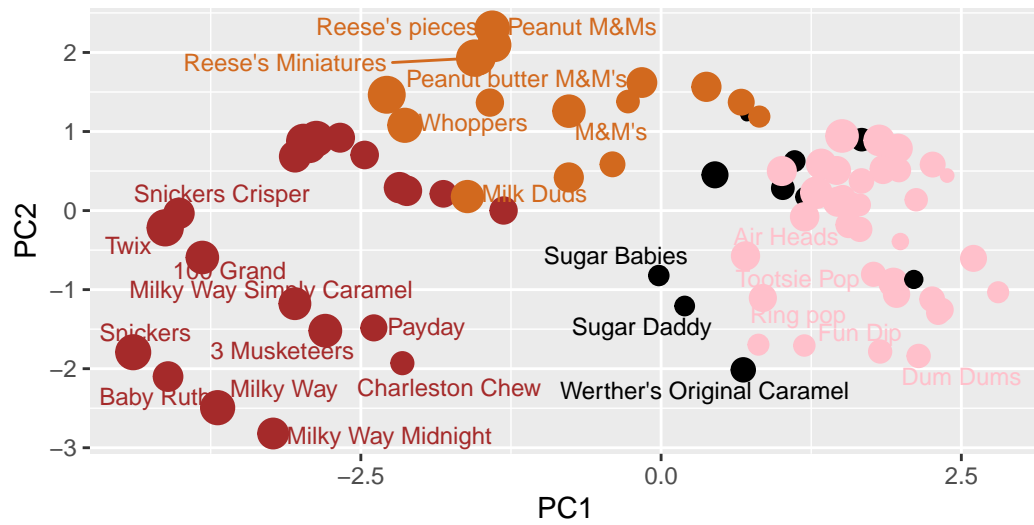
```
# Adding text labels with ggrepel:
p + geom_text_repel(size=3.3, col=mycols, max.overlaps = 7)  +
  theme(legend.position = "none") +
  labs(title="Halloween Candy PCA Space",
       subtitle="Colored by type: chocolate bar (dark brown), chocolate other (light brown
       caption="Data from 538")
```

Warning: ggrepel: 59 unlabeled data points (too many overlaps). Consider
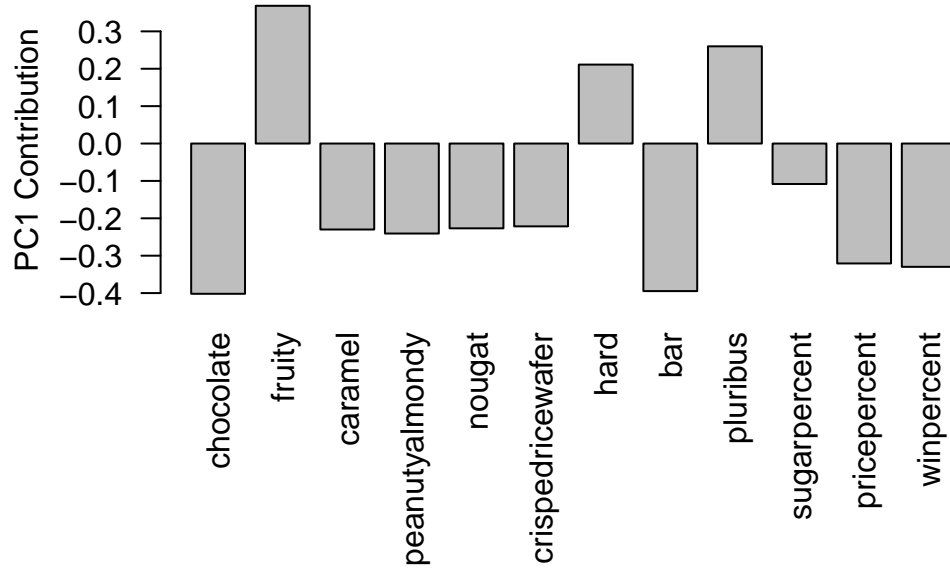increasing max.overlaps

## Halloween Candy PCA Space

Colored by type: chocolate bar (dark brown), chocolate other (light brown),

```
# Looking at our PCA loadings:
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 and pluribus. This makes sense because fruity candy tend to be hard and come in multiples.