# Case study: *gummy-like* or *cookie-like* candies?

## Problem statement

- We want to create a brand-new candy.
- We have a **market survey** on **customer sentiment** for competitor brands of candies.
- We want our analysis to be driven by the **maximization of customer preference**.
- We need to choose what the new candy will be: a **gummy** or a **cookie-based sweet**?

## Recommendation

Results from data analysis suggest that we should opt for a **cookie-based sweet**.

# Assumptions and procedure

**The market survey does not specify which candies are gummy-like or cookie-like.**

So, first of all we need to create two clusters (based on the properties listed in the marked survey).
**We assume that gummies and cookies can be classified as following:**

| gummy-like |
| --- |
| • **Not hard:** they are very different from hard candies.<br>• **Fruity**: often, gummies have fruity flavours.<br>• **Pluribus**: often, gummies come in packages which contains many of the same type. |

| cookie-like |
| --- |
| • **Not hard:** they are very different from hard candies.<br>• **Crispedricewafer**: of course they must have a crispy texture, similar to cookies. |

# Results of the analysis

So far, the likeability of a candy is indicated by a continuous value (*winpercent*). We want to simplify things:

> **We say that a candy is a _winner_ if its winpercent is > 50**

This means that, in a random match with another candy, a winner candy has more than 50% the chance of being preferred by a customer.

What percentage of *cookie-like* candies in our dataset are winners?
What about *gummy-like* candies?

```
Percentage of cookie-like candies that are winners:  85.71428571428571
Percentage of gummmy-like candies that are winners:  42.10526315789473
```

# Results of the analysis

Let's now look at flavours and physical features of candies.



**Some features are heavily correlated** (see for example chocolate and fruity).
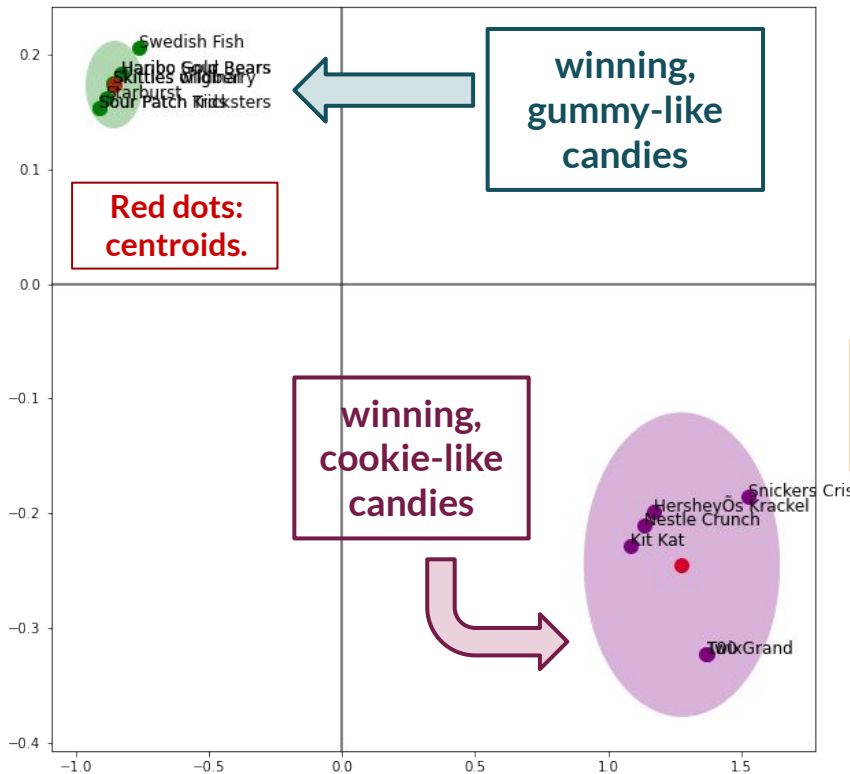
There most likely are **redundant variables**.

This suggest an approach of **dimensionality reduction** before trying to extract value from the data.

# Results of the analysis

We reduce the dimensionality of the dataset by applying **Principal Components Analysis**. Now every candy can be described by **only two variables**, so we can plot them on a cartesian grid.



winning, gummy-like candies

Red dots: centroids.

winning, cookie-like candies

We can expect with 95% confidence that a new cookie-like/gummy-like candy will fall inside the purple/green ellipse.

## Linear regression

**Q:** What is the expected likeability for the **mean** gummy-like or cookie-like candy?

| *gummy-like* | *cookie-like* |
|:---:|:---:|
| **42.52%** | **62.02%** |