



# Exploring Avocados

March 2, 2022

## The Dataset

The data was pulled from <https://www.kaggle.com/neuromusic/avocado-prices> Some relevant columns in the dataset are shown below:

- Date - The date of the observation
- AveragePrice - the average price of a single avocado
- type - conventional or organic
- year - the year
- likes - Number of likes of the Talk
- Region - the city or region of the observation
- Total Volume - Total number of avocados sold

Many thanks to the Hass Avocado Board for sharing this data!!

1. What three regions have the highest mean Average Price for conventional avocados?
2. What three regions have the lowest mean total volume of conventional avocados sold?
3. Create a pairwise boxplot for the average price of organic avocados sold across the years. (Remember you will need to 'factor' years in your ggplot). Looking at the plot is there any significant difference in average price across years?
4. Visualize and describe the distribution for average price of organic avocados in New York. Be sure to include a thorough description.
5. Visualize the distribution for total volume of conventional avocados sold in Houston in the year 2015. Be sure to include a thorough description.
6. Visualize the average price of conventional avocados vs. organic avocados in California using a pairwise boxplot. Does there appear to be a significant in average price across the different types?
7. Visualize the total volume of avocados sold in Houston vs. California using a pairwise boxplot. Does there appear to be a significant difference in total volume across the two regions?

### Solution:

```
1  ## Load Packages
2  library(ggplot2)
3  library(tidyverse)
4
5  ##### Load the dataset
```

```

6 avo <- read.csv('/data/datasets/avocado.csv')
7 glimpse(avo)
8
9 ##### Question 1 Find top three regions by average price conventional
10 avo %>% filter(type=='conventional') %>%
11   group_by(region) %>% summarise(average = mean(AveragePrice)) %>% arrange(desc(
12     average))
13   # Highest prices are in Hartford, NewYork, SanFran
14 ##### Question 2 Find bottom three regions by total volume sold conventional
15 avo %>% filter(type=='conventional') %>%
16   group_by(region) %>%
17   summarise(average = mean(Total.Volume)) %>% arrange(average)
18   # Lowest volume in Syracuse, Boise, and Spokane
19
20 ##### Question 3 Pairwise boxplot for organic across years
21 avo %>% filter(type == 'organic') %>%
22   ggplot(aes(y=AveragePrice, fill=factor(year))) +
23   geom_boxplot() +
24   theme_classic()
25
26 ##### Question 4 Average price of organic avocados in New York
27 avo %>% filter(type == 'organic' & region == 'NewYork') %>%
28   ggplot(aes(x=AveragePrice)) +
29   geom_histogram(fill='green', col='black')+
30   theme_classic()+
31   labs(x='Average price of avocados', title='Distribution for the average price of
32     avocados')
33 ##### Question 5 Total volume for conventional avocados in Houston 2015
34 avo %>% filter(region == 'Houston' & type == 'conventional' & year==2015) %>%
35   ggplot(aes(x=Total.Volume))+
36   geom_histogram(fill='darkgreen', col='black')+
37   labs(x='Total Volume Sold')
38
39 ##### Question 6 Average price for conventional vs. organic avocados in California
40 avo %>% filter(region=='California') %>%
41   ggplot(aes(y=AveragePrice, fill=type)) +
42   geom_boxplot()
43
44 ##### Question 7 Total volume of organic sold in Houston vs. California
45 avo %>% filter(type == 'organic') %>%
46   filter(region == 'Houston' | region == 'California') %>%
47   ggplot(aes(y=Total.Volume, fill=region)) +
48   geom_boxplot() +
49   labs(y='Total volume sold')+
50   theme_classic()
51
52

```