# Liquor Sales in Iowa, USA 2016-2019

**Initial dataset and questions to answer**

The initial dataset contained the Liquor sales in the state of Iowa (USA) between the years 2012-2020 however we have been asked to find:

- the most popular item per zipcode and
- the percentage of sales per store

in the period between 2016-2019.

**Description of steps taken**

A SQL query was used to extract all the columns between 2016 – 2019 in a CSV format from Workbench. Then, utilising Python and the Pandas library, the CSV file was imported and a dataframe was created to be further analysed.

I printed the information about the dataframe which contains 24 columns and 74 entries, 0 to 73.  The data types are object, integers and float and there also some missing values on the variables 'store_location' (9) and 'category_name'(6). However, to answer the questions I kept only the following 7 columns: 'invoice_and_item_number', 'store_number', 'store_name', 'zip_code', 'item_description', 'bottles_sold', 'sale_dollars'.

For the first question, I grouped the data based on zip code and item description, I found the sum of the 'bottles_sold' for each zipcode, sorted the values in descending order and printed the first row.
For the second question, at first, I found the total of the 'sale_dollars', I created a new column called '%_of_sales' and calculated the percentage with the following script:

```
df["%_of_sales"] =(df['sale_dollars'] / sum_of_sales) * 100
```

Then, I grouped the data by 'store_number', found the sum of '%_of_sales' for each row and sorted the values in descending order.

For the visualization of the data, I opened the CSV file in Tableau, created the graphs for the questions and then I created a dashboard with the two graphs. A link has been attached:
https://public.tableau.com/views/Finalassignment_16744901562800/Dashboard1?:language=en-US&publish=yes&:display_count=n&:origin=viz_share_link
The first graph is a horizontal bar chart showing the sum of bottles sold in each zip code and the item description.
The second graph is also a horizontal bar chart showing the percentage of sales for each store. To calculate the percentage, I converted the 'sale_dollars' variable to measure, I chose the measure sum and from the quick table calculation I chose the percent of total.

**Challenges I faced**

During the initial steps (importing the file and analysing the data) I did not face any significant difficulty as I was using google and the provided material to find everything I needed.
The main challenge I faced was plotting the data using the matplotlib. Even though I was searching on google I kept getting errors or wrong values. Specifically, I had difficulty to plot the aggregated variable 'bottles_sold' and for the second graph I was trying to create an horizontal bar chart however, the bars were displayed as lines and only 3 of them were shown on the graph.

# Liquor Sales in Iowa, USA 2016-2019

## Most popular item sold based on zipcode

| item_description | zip_code | bottles_sold |
|---|---|---|
| Juarez Gold Dss | 52314 | 1,560 |
| Tortilla Gold Dss | 50320 | 72 |
| | 50702 | 768 |
| | 51360 | 48 |
| Juarez Triple Sec | 50314 | 240 |
| | 51106 | 144 |
| | 51501 | 48 |
| | 52240 | 60 |
| | 52556 | 6 |
| Member's Mark Spiced Ru.. | 50010 | 288 |
| Kahlua Coffee Liqueur | 51106 | 240 |
| Montezuma Triple Sec | 52402 | 216 |
| Pinnacle Peach w/ Punch .. | 50703 | 180 |
| Di Amore Quattro Orange | 50320 | 120 |
| Hennessy VS | 50158 | 24 |
| | 50316 | 48 |
| | 50703 | 24 |
| | 52601 | 24 |
| Kapali Coffee Liqueur | 51106 | 84 |
| | 52627 | 36 |

*bottles_sold* (x-axis: 0 to 1600)

## Percentage of sales per store

| store_number | % of Total Sum of sale_dollars |
|---|---|
| 5102 | |
| 3447 | 10.84% |
| 3494 | 8.49% |
| 2633 | 6.73% |
| 3524 | 6.39% |
| 9001 | 6.06% |
| 2619 | 5.50% |
| 4829 | 4.21% |
| 4971 | 3.75% |
| 2571 | 3.25% |
| 2562 | 2.55% |
| 2665 | 2.48% |
| 4153 | 2.20% |
| 5146 | 2.20% |
| 4136 | 2.12% |
| 2515 | 1.47% |
| 2538 | 1.32% |

*% of Total Sum of sale_dollars*