

# David Griesel

## Data Analyst Portfolio

Hamburg, Germany

+49 157 534 08870

davidgrieselde@icloud.com

### About Me

---

I'm a data analyst with a background in accounting and auditing. I started out supporting small businesses with financial reporting and analysis before moving into public sector audit, where I worked with large, complex datasets across a range of government clients.

Over time, I began enhancing audit processes through the expanded use of data analytics and computer-assisted audit techniques (CAATs), which led to a dedicated role in data analytics. There, I introduced custom scripts to streamline workflows, reduce repetitive tasks, and improve efficiency.

Today, I apply the same principles to data analytics projects — uncovering insights, solving challenges, and continuing to explore new tools and approaches.

### Links

---

<https://www.linkedin.com/in/davidgriesel>

<https://github.com/davidgriesel>

<https://public.tableau.com/app/profile/david.griesel>

<https://davidgriesel.com>

### Tools

---



# Projects

## Medical Staffing Plan

---

Analysed historical influenza and demographic data to develop a staffing plan for a medical staffing agency by identifying high-risk states and forecasting seasonal flu trends.

## Online Grocery Store

---

Conducted an exploratory analysis of an online grocery store's data to gain insights into their customer base and purchasing behaviour, with the goal of optimising their marketing strategy.

## Coffee Quality

---

Performed statistical analysis on scoring and production data to investigate the factors influencing coffee quality scores



# Medical Staffing Plan

## Goal

---

Inform the timing and spatial distribution of a medical agency's staffing plan by analysing historical influenza and demographic data

## Process Flow

---

**Defining the Business Requirements** – Project Questions | Privacy & Ethics Considerations

**Designing a Data Research Project** – Research Hypothesis | Data Requirements

**Describing the Data** – Source | Collection Method | Contents | Limitations | Relevance

**Data Profiling and Integrity Assessment** – Data Types | Accuracy & Consistency | Cleaning | Summary Stats

**Data Quality Assessment** – Completeness | Uniqueness | Timeliness | Cleaning

**Data Transformation & Integration**

**Statistical Analysis** – Variance | Standard Deviation | Relationships | Correlation

**Statistical Hypothesis Testing**

**Compiling an Interim Report**

**Data Visualisation & Spatial Analysis**

**Writing a Narrative & Creating an Interactive Storyboard**

## Tools

---

 – Excel

 – Tableau

## Data

---

**Number of influenza-related deaths** – Number of deaths by location, time, and age. (CDC)

**Influenza vaccination rates in children** – Number of vaccinations by location, time, and age. (CDC)

**Population** – Number of people by location, time, age, and gender. (US Census Bureau)

## Links

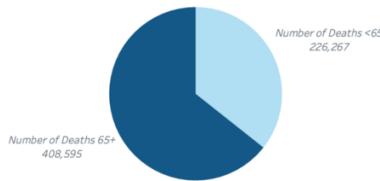
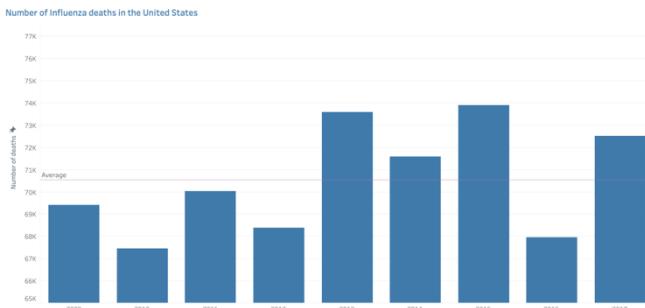
---

[GitHub Repository](#)

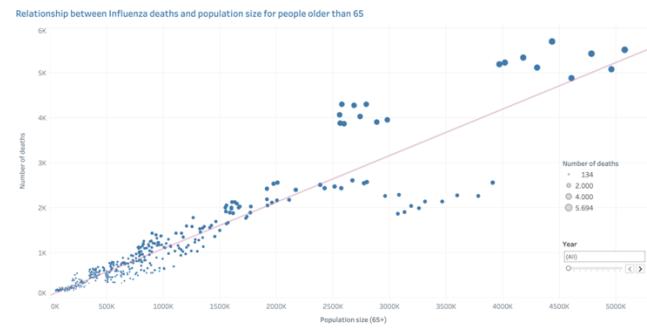
[Tableau Storyboard](#)

# Analysis & Insights

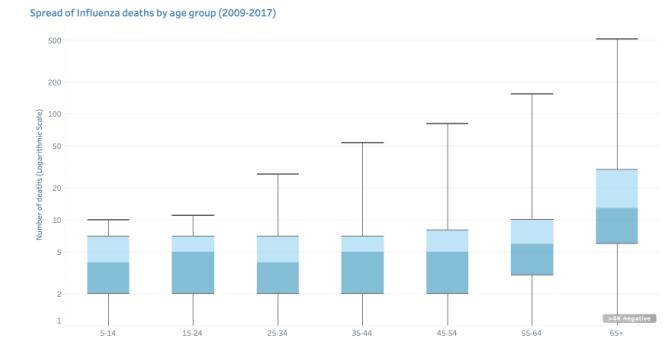
## Spatial Distribution & Risk Classification



Dataset		Integrated Data Set	
Sample or Population?	Sample	> 30 - Central Limit Theorem applies	
Variable		Population 65+	Influenza Deaths 65+
Number of Records	459	459	
Variance	786799847780,66	651954,62	
Standard Deviation	887017,39	807,44	
Mean	806989,68	1263,72	
-2SD	-967045,10	-351,15	
+2SD	2581024,46	2878,60	
Outliers below -2SD	0	0	
Outliers above +2SD	29	18	
Outlier Percentage (Empirical Rule)	6%	4%	
Pearson's Correlation Coefficient		0,84	
Relationship		Strong Relationship	
Interpretation		The more people there are aged 65 and older, the higher the number of influenza-related deaths will be.	

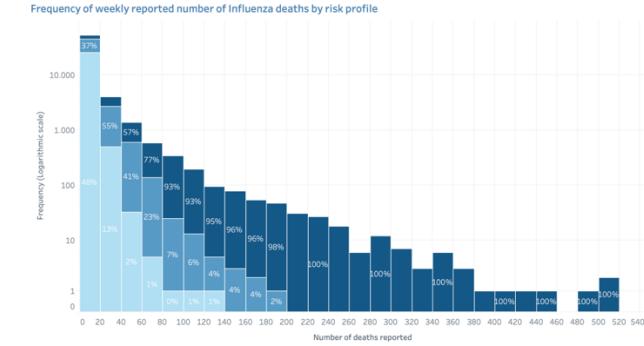
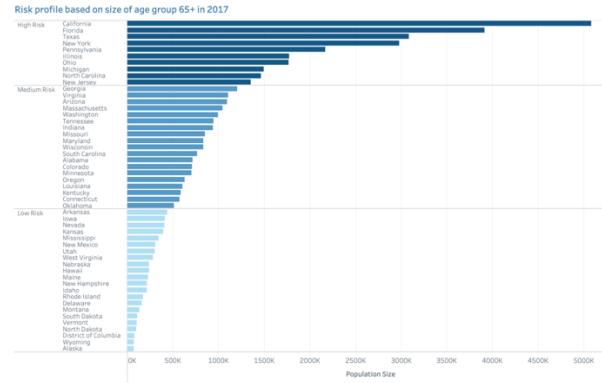
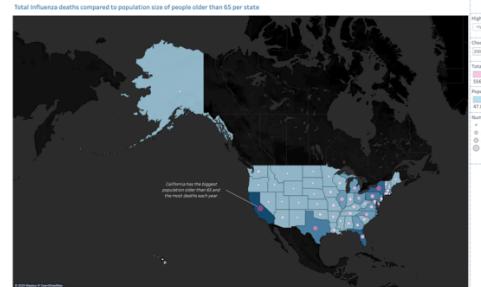


Statistical Hypothesis Testing		
Research hypothesis	Individuals aged 65 and older have a higher influenza-related death rate than Population (65+ vs. <65)	
Independent Variable	Death Rate	
Dependent Variable	Death Rate	
H <sub>0</sub> (Null)	Death rate for people 65+ is ≤ death rate for people under 65	
H <sub>1</sub> (Alternative)	Death rate for people 65+ is > death rate for people under 65	
Type of Test	One-Tailed Test (Null hypothesis -> equal to, or lower than / one direction of)	
Significance Level ( $\alpha$ )	0,0500	
p-value	0,0000	
Assessment	The p-value of 0.00 is below the significance level of 0.05, indicating a statistically significant difference.	
t-Test: Two-Sample Assuming Unequal Variances		
	Death Rate 0-64	
	Death Rate 65+	
Mean	0,00059714382	0,00395811359
Variance	0,00000040031	0,00002543436
Observations	459	459
Hypothesized Mean Difference	0	
df	472	
t Stat	-14,16672411596	
P(T<=t) one-tail	0,00000000000	
t Critical one-tail	1,64808833557	
P(T<=t) two-tail	0,00000000000	
t Critical two-tail	1,96500267646	

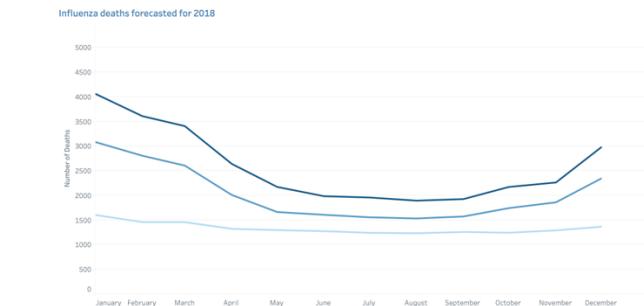
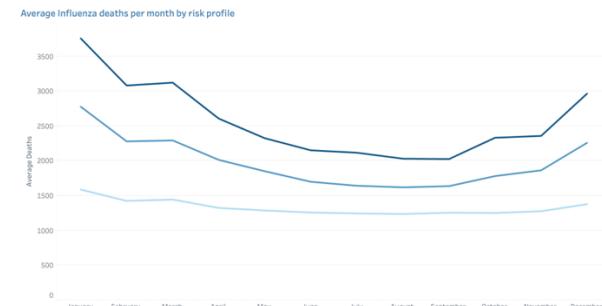
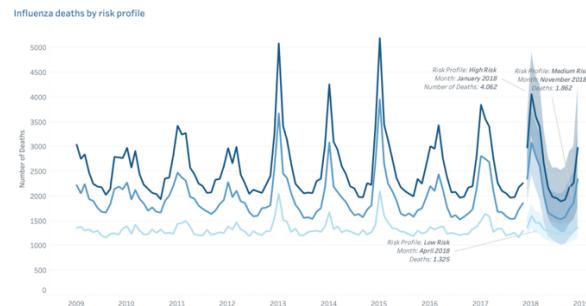


# Analysis & Insights

## Spatial Distribution & Risk Classification



## Seasonality & Forecasting





## Reflection

### Successes

The project provided actionable insights into the spatial distribution of vulnerable populations, enabling the development of a risk profile to guide regional staff allocations across states. Additionally, it offered valuable understanding of influenza seasonality, which supported the strategic timing of staff deployments to meet anticipated healthcare demands.

### Challenges

The available data lacked sufficient detail to account for all known factors in the analysis. The suppression of all records for children under 5 limited the ability to assess risk in this age group, highlighting how data privacy laws can restrict public health analyses when key demographic groups are excluded.

## Recommendations

It is recommended that staffing resources be prioritised for high-risk states, with allocations scaled according to the size of their elderly populations. Additionally, deployment should begin ahead of the seasonal peak — ideally starting in November — to ensure adequate coverage throughout the critical period from December through April. This approach supports a targeted and data-informed strategy to minimise influenza-related mortality and optimise medical staffing efficiency.

### Moving Forward

Monitor the impact of the 2018 deployment strategy to evaluate its effectiveness and inform the next planning cycle, ensuring continuous improvement in staff allocation strategies.



# Online Grocery Store

## Goal

---

Perform an exploratory analysis of customer demographics and transactional data to gain insights into ordering behaviours, thereby enhancing customer profiling and optimising marketing strategies through targeted product placement.

## Process Flow

---

### Conducted Descriptive Exploratory Data Analysis

**Data Wrangling** – Data Types | Accuracy | Drop & Rename Columns

**Data Consistency Checks** – Missing Values | Duplicates | Mixed-Type variables

### Merging Dataframes

### Subsetting & Deriving New Variables

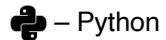
### Grouping & Aggregating Data

### Data Visualisations & Analysis

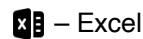
### Compiling and Interim Report

## Tools

---



– Python



– Excel

## Data

---

**The Instacart Online Grocery Shopping Dataset** – Product, department, and order details. (Instacart)

**Supplementary Customer dataset** – Customer details and demographic information. (CareerFoundry)

## Links

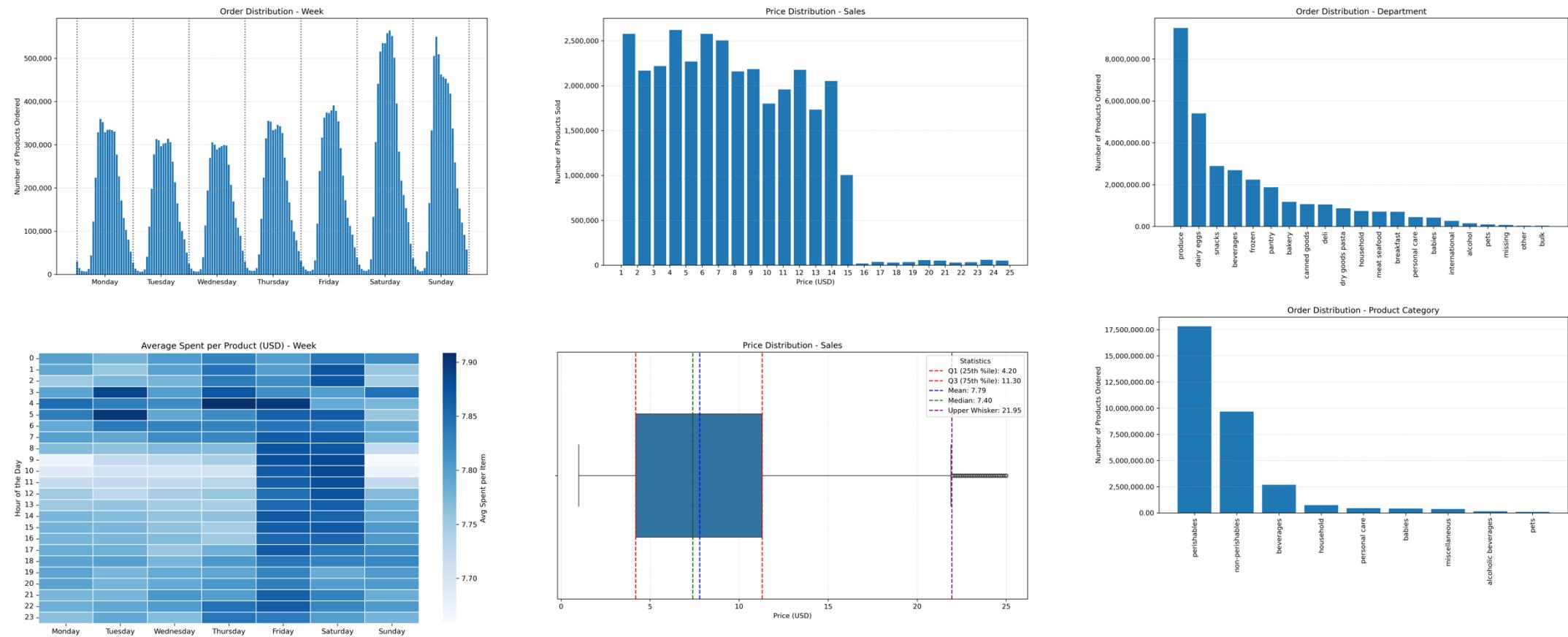
---

[Kaggle](#)

[GitHub Repository](#)

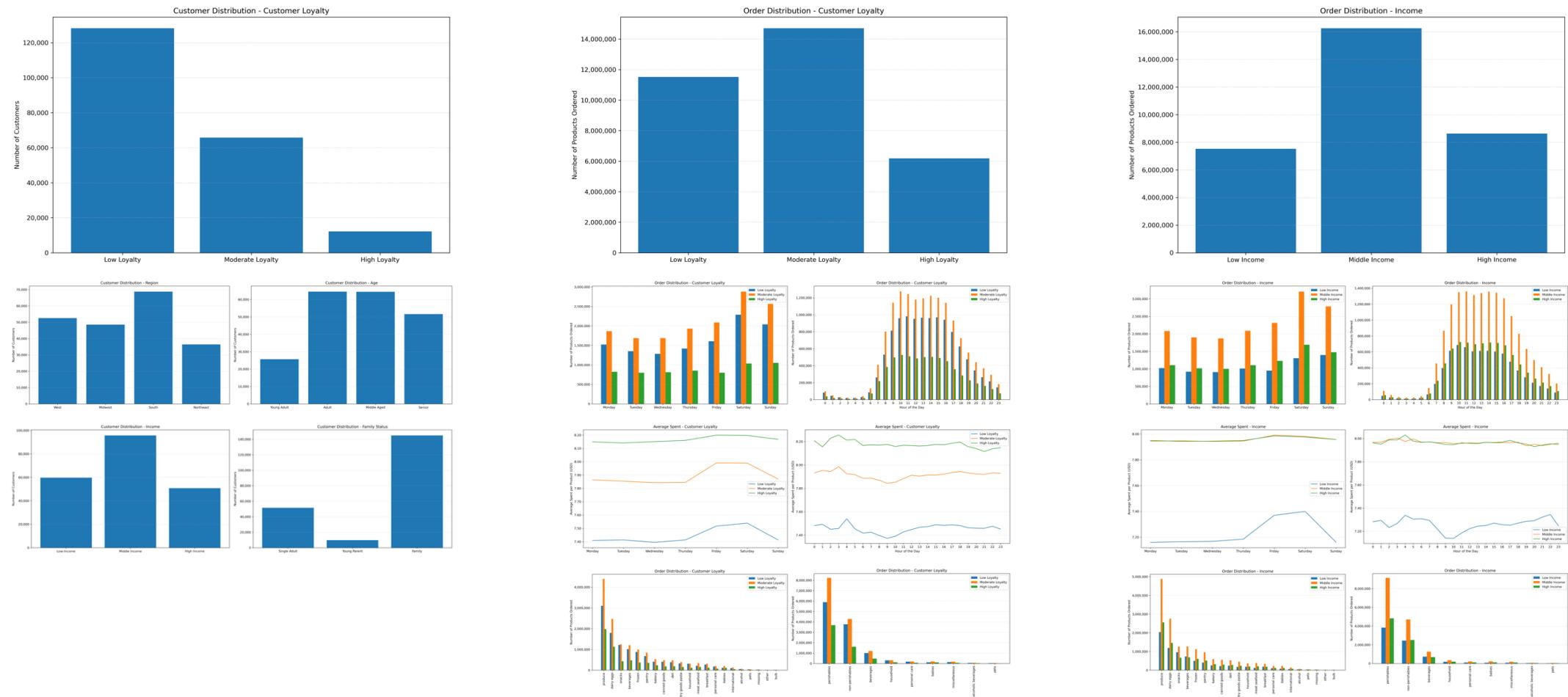
# Analysis & Insights

## Ordering Trends



# Analysis & Insights

## Ordering Trends





## Reflection

### Successes

The project provided actionable insights into ordering behaviour across customer profiles, enabling a more targeted marketing strategy through better product placement.

### Challenges

Generating visualisations in Python is cumbersome and creating clear and intuitive charts required significant customisation.

## Recommendations

To optimise engagement and revenue, ads should run during low-order periods — specifically between 4 p.m. and 7 a.m. Product promotions should align with spending behaviour: mid-priced items in the evening, high-priced products in the early morning and late on Thursday and Friday evenings, and premium items between 3–6 a.m., when spending peaks. Focus should be placed on food and beverage categories, especially fresh produce, dairy & bakery items.

Targeting should also consider customer profiles:

- Moderate loyalty customers respond well to a mix of staple and premium products.
- High loyalty customers prefer healthier, premium options.
- Low loyalty customers are more price sensitive and respond to budget-friendly goods.

Income segmentation can further refine targeting, with premium products appealing to high-income shoppers and essentials resonating more with low-income groups.

### Moving Forward

Create an interactive Tableau storyboard to visually present key findings, enabling stakeholders to explore ordering trends, customer profiles, and spending behaviours with filters, dashboards, and drill-down insights.



# Coffee Quality

## Goal

---

Analyse coffee quality across countries, explore predictive relationships between quality measures, and examine the link between quality scores and demand.

## Process Flow

---

**Define Project Requirements** – Define Project Deliverables & Data Requirements

**Data Sourcing & Preparation** – Source Prepare Open Source Data

**Conduct Exploratory Data Analysis** – Visualising Relationships & Developing a Hypothesis

**Geospatial Analysis** – Integrate & Analyse Spatial Data with Choropleth Maps using Shapefiles

**Supervised Machine Learning** – Apply Regression Techniques for Predictive Analysis

**Unsupervised Machine Learning** – Apply k-means Clustering using Elbow Technique for pattern discovery

**Time Series Analysis** – Analysing Temporal Data through Decomposition, Stationary Tests, and Transformations

**Data Visualisation & Storytelling** – Communicate Insights through Interactive Dashboard.

## Tools

---

 – Python

 – Tableau

## Data

---

**Coffee Quality with Locations of Origin** – Quality scores of coffee beans by country. (Coffee Quality Institute)

**Average Production Per Year** – Consumption, import, export, and production metrics over time. (Coffee Quality Institute)

## Links

---

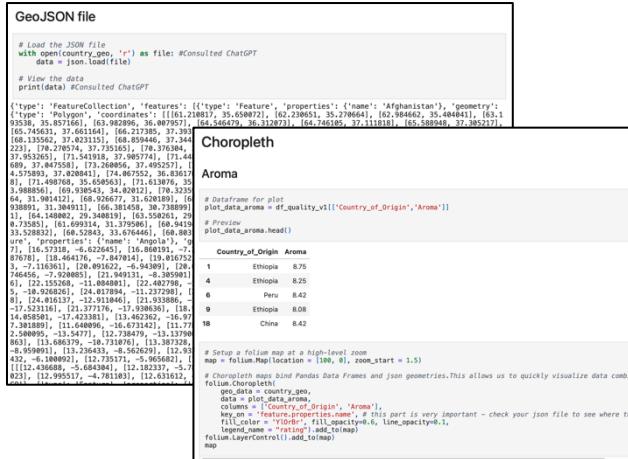
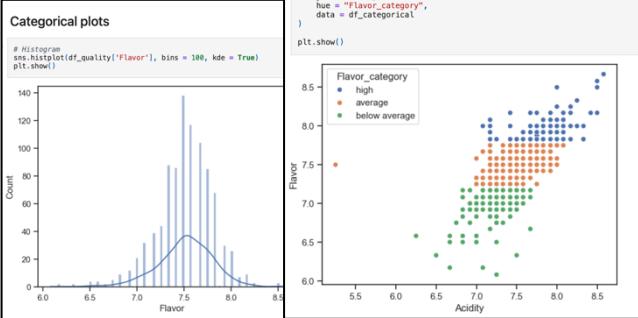
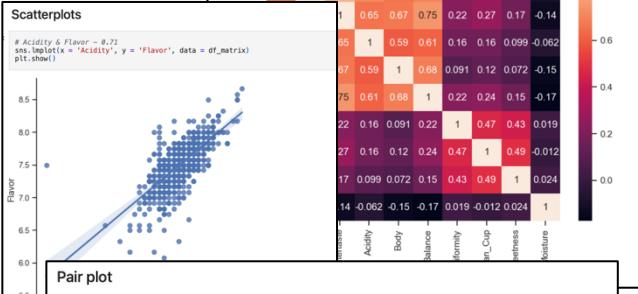
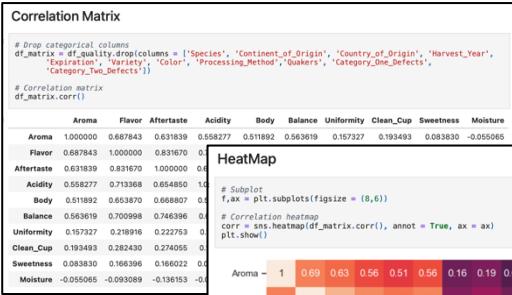
[Kaggle](#)

[GitHub Repository](#)

[Tableau Storyboard](#)

# Analysis & Insights

## Ordering Trends



# Analysis & Insights

## Ordering Trends

### Elbow technique

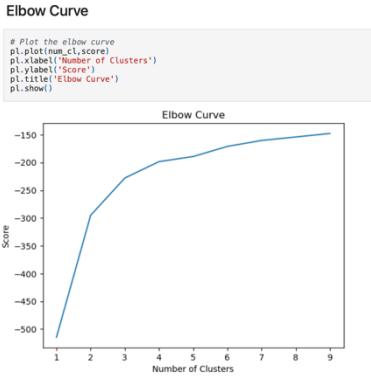
```
# Defines the range
num_cl = range(1, 10)

# Defines k-means clusters
kmeans = KMeans(n_clusters = i) for i in num_cl

# Create scores
score = [kmeans[i].fit(df_matrix).score(df_matrix) for i in range(len(kmeans))]

# View scores
score

[-514.793415919853,
-295.1728474245013,
-228.0014291244987,
-198.50689513746968,
-169.4213326131645,
-171.89118326131645,
-168.4213326131645,
-154.1394772245893,
-147.7822725919875]
```



### k-means clustering

```
# Create the k-means object
kmeans = KMeans(n_clusters = 3)

# Fit the object
kmeans.fit(df_matrix)

KMeans(n_clusters=3)
```

### New column

```
# Create 'clusters' variable
df_matrix['clusters'] = kmeans.fit_predict(df_matrix)

# Preview
df_matrix.head()

Aroma Flavor Aftertaste Acidity Body Balance clusters
1 8.75 8.50 8.58 8.42 8.42 1
4 8.25 8.50 8.25 8.50 8.42 8.33 1
6 8.42 8.50 8.33 8.50 8.25 8.25 1
9 8.08 8.58 8.50 8.50 7.67 8.42 1
18 8.42 8.25 8.08 8.17 7.92 8.00 1

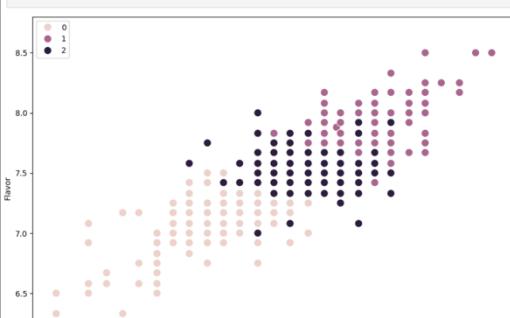
# Frequency of 'clusters' variable
df_matrix['clusters'].value_counts()

clusters
2 468
1 288
0 1
Name: count, dtype: int64
```

### Hypothesis #1

- The higher the score for Aftertaste, the higher the score for Flavor would be.

```
# Cluster plot - "Aftertaste" and "Flavor"
plt.figure(figsize = (12,8))
ax = sns.scatterplot(x = df_matrix['Aftertaste'], y = df_matrix['Flavor'], hue = kmeans.labels_, s = 100)
ax.grid(False)
plt.xlabel('Aftertaste')
plt.ylabel('Flavor')
plt.show()
```



### Descriptive statistics

```
# Assign corresponding colour to the clusters
df_matrix.loc[df_matrix['clusters'] == 2, 'cluster_color'] = 'dark purple'
df_matrix.loc[df_matrix['clusters'] == 1, 'cluster_color'] = 'purple'
df_matrix.loc[df_matrix['clusters'] == 0, 'cluster_color'] = 'pink'

# Calculate statistics
df_matrix.groupby('cluster_color').agg({'Aroma':['mean', 'median'],
                                         'Flavor':['mean', 'median'],
                                         'Aftertaste':['mean', 'median'],
                                         'Acidity':['mean', 'median'],
                                         'Body':['mean', 'median'],
                                         'Balance':['mean', 'median']})

Aroma   Flavor   Aftertaste   Acidity   Body   Balance
mean   median   mean   median   mean   median   mean   median
cluster_color
dark purple  7.548880  7.58  7.516453  7.50  7.393654  7.42  7.501688  7.50  7.480769  7.50  7.498333  7.50
pink    7.257665  7.25  7.097411  7.17  6.970467  7.00  7.210406  7.25  7.228071  7.25  7.097817  7.17
purple  7.811464  7.75  7.802107  7.75  7.688571  7.67  7.804143  7.75  7.751464  7.75  7.817679  7.75
```

### Group and aggregate

```
# Group by 'Year' and calculate the mean production
df_production_4 = df_production_3.groupby('Year')['Production'].mean().reset_index()

# Rename column
df_production_4.rename(columns = {'Production': 'Average Production'}, inplace = True)

# Display the result
df_production_4
```

### Analysis

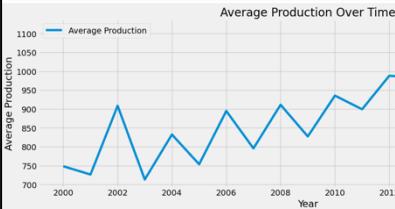
#### Line chart

```
plt.figure(figsize = (15, 5), dpi = 100)

# Specify the column to plot
plt.plot(df_production_5.index, df_production_5['Average Production'], label = 'Average Production')

# Add labels and title
plt.title('Average Production Over Time')
plt.xlabel('Year')
plt.ylabel('Average Production')
plt.legend()

# Show the plot
plt.show()
```



### Dicky-Fuller test - Round 2

```
from statsmodels.tsa.stattools import adfuller

# Define the function
def dickey_fuller(timeseries):
    # Perform the Dickey-Fuller test:
    print('Dickey-Fuller Stationarity test:')
    test = adfuller(timeseries, autolag = 'AIC')
    result = pd.Series(test[0:4], index = ['Test Statistic','p-value','Number of Lags Used','Number of Observations Used'])
    for key,value in test[4].items():
        result['Critical Value (%s)'%key] = value
    print(result)

# Apply the test
dickey_fuller(df_production_7['Average Production'])
```

Dickey-Fuller Stationarity test:  
Test Statistic: -2.157692  
p-value: 0.000000  
Number of Lags Used: 6.000000  
Number of Observations Used: 10.000000  
Critical Value (1%) -4.331573  
Critical Value (5%) -2.861434  
Critical Value (10%) -2.722950  
dtype: float64

#### Comments

- At a significance level of 5%, the test statistic is smaller than the critical value and the null hypothesis can be rejected.
- There is not a unit root in the data, and the data is stationary.

### Indexing

```
# Set 'Date' as datetime index
df_production_5 = df_production_4.copy()

from datetime import datetime

df_production_5['datetime'] = pd.to_datetime(df_production_5['Year'], format='%Y')
df_production_5.set_index('datetime', inplace = True)
df_production_5.drop(['Year'], axis = 1, inplace = True)
```

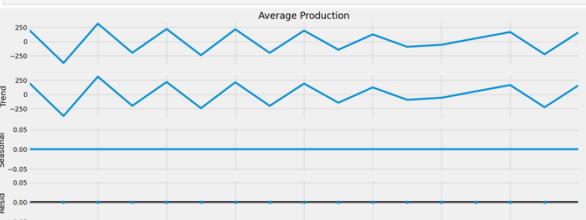
### Decomposition - Round 2

```
from pylab import rcParams

# Decompose time series - additive model
decomposition = sm.tsa.seasonal_decompose(df_production_7['Average Production'], model = 'additive')

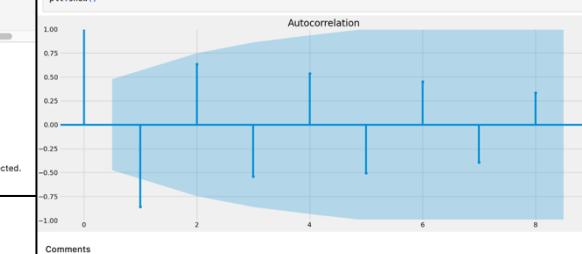
# Set the figure size
rcParams['figure.figsize'] = 18, 7

# Plot the decomposition
decomposition.plot()
plt.show()
```



### Autocorrelation - Round 2

```
# Plot of autocorrelations
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
plot_acf(df_production_7)
plot_pacf(df_production_7)
plt.show()
```

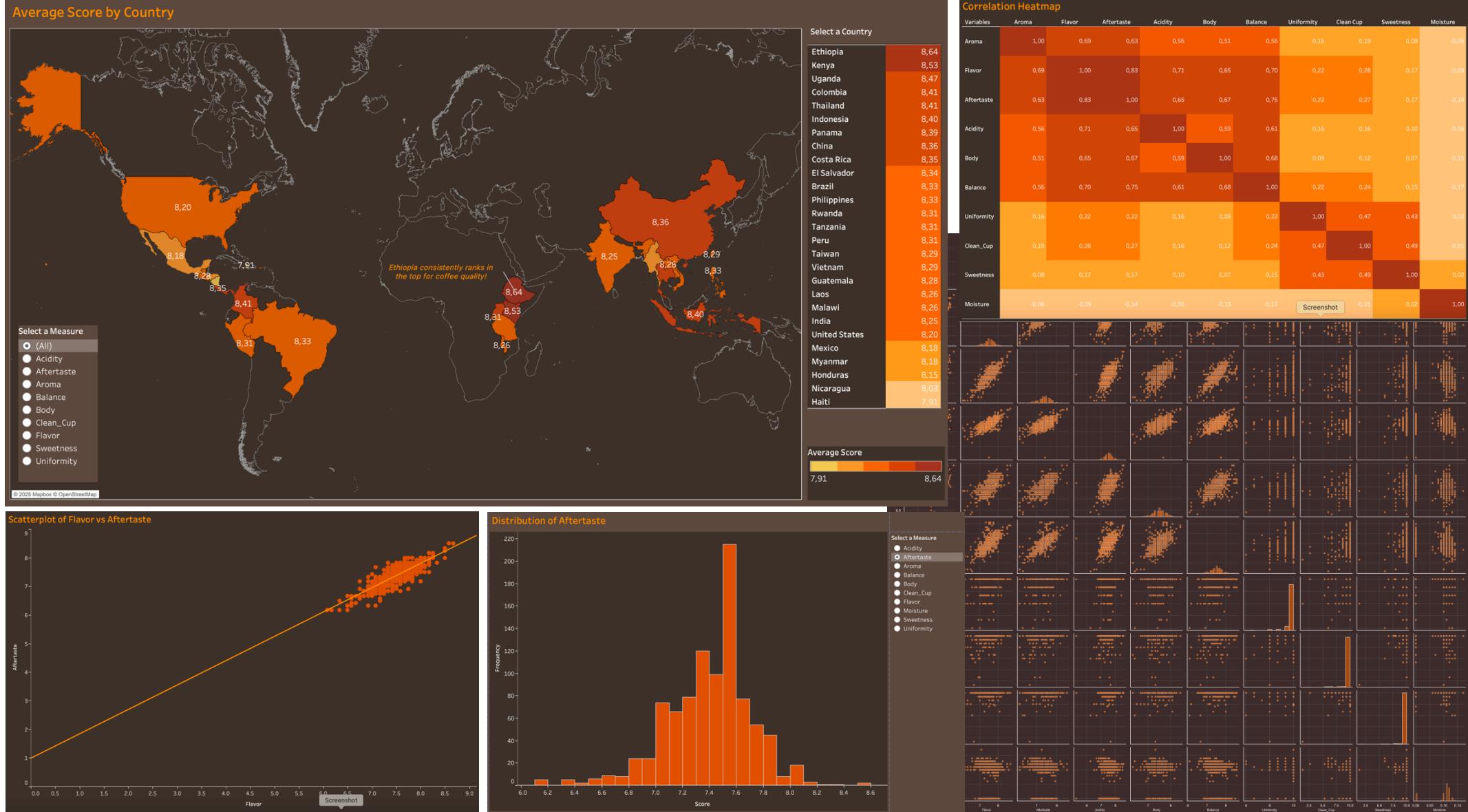


#### Comments

- The plot suggests reduced autocorrelation.

# Analysis & Insights

## Ordering Trends





## Reflection

### Successes

Overall, the project provided meaningful insights into coffee quality metrics and their interactions, highlighting Ethiopia's dominance in producing high-quality coffee and that it is possible to predict scores of certain measures using the scores of others.

### Challenges

It was not possible to obtain a dataset with sufficient granularity to assess whether there is a higher demand for better quality coffee, underscoring the importance of clearly defined data requirements at the outset of a project.

## Recommendations

---

The analysis achieved its goals of identifying top coffee-producing countries and exploring score predictability across quality measures. Ethiopia emerged as a standout, with consistently high coffee ratings, while Central American and Southeast Asian countries scored lower overall.

Strong correlations were found among Flavour, Aftertaste, and Balance, with Aftertaste being a reliable predictor of Flavour, though less so for Balance.

Clustering the data offered no added insights beyond the observed linear relationships.

Due to data limitations, the analysis could not determine whether higher-quality coffee is in greater demand.

### Moving Forward

---

Reach out to the Coffee Quality Institute to inquire about access to a more suitable dataset.

Investigate alternative modelling approaches for predicting the less correlated measures.