An abstract graphic on the left side of the slide consists of several overlapping and nested geometric shapes. It includes a white circle at the top left, a grey circle below it, a blue triangle containing concentric circles, a red square with a series of nested triangles pointing towards the bottom left corner, and a dark purple triangle at the bottom right.

RUSAMIJAN PERMISON

Data Analytics Portfolio

Email: rsg.design@gmail.com
Portfolio: www.meepermison.com

PROJECTS



Analyze sales using Python Excel and [Tableau Public](#) ([Github](#))



Analyze sales and customer profiling using Python and Excel ([Github](#))



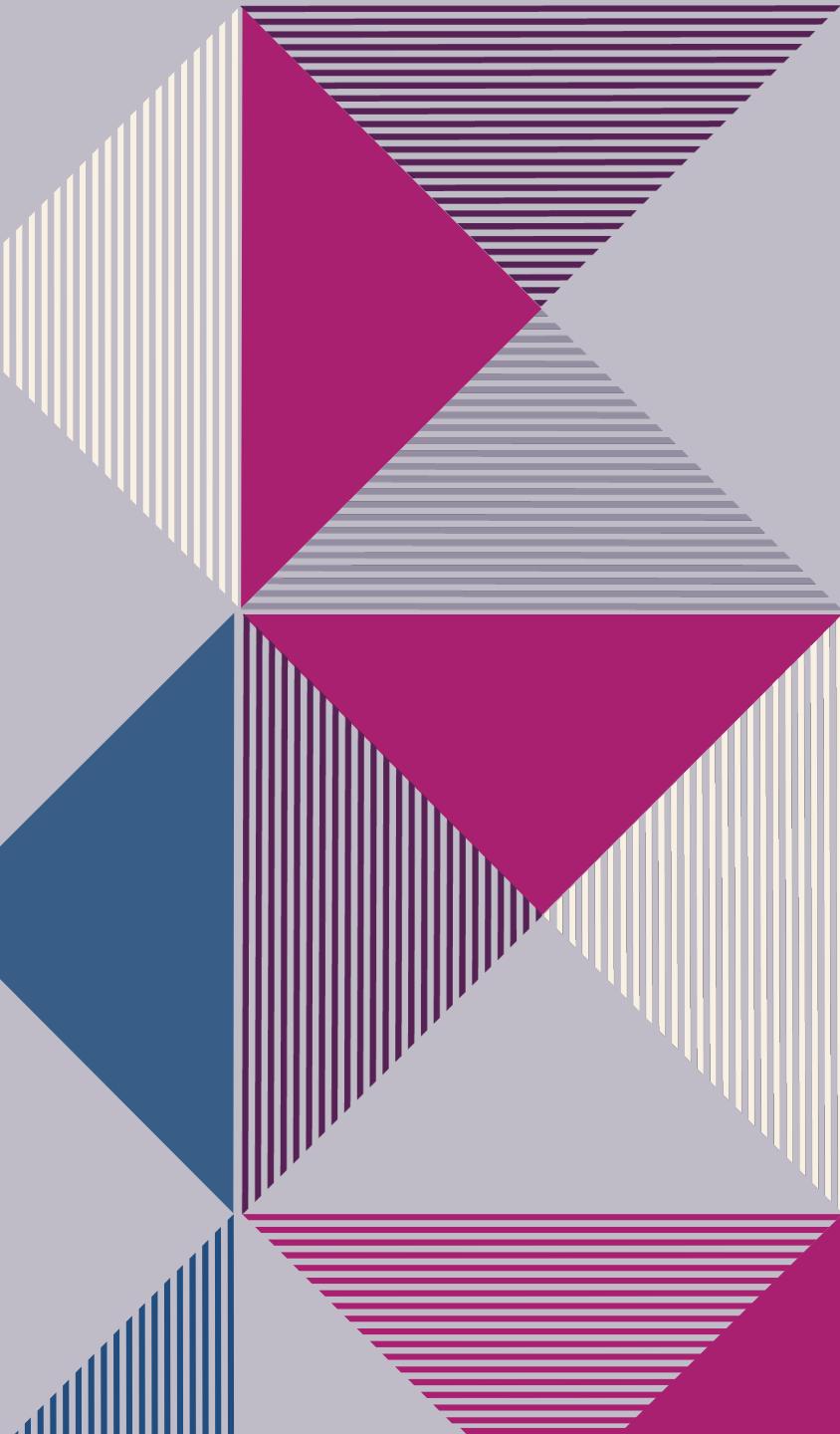
Identify insights and customer analysis by using SQL and [Tableau Public](#) ([Github](#))



Identify flu season and staffing needs in the US by using Excel and [Tableau Public](#)



Analyze global videogame retail sales by using Excel, PowerPoint and [Tableau Public](#)



US ADIDAS SALES DATA ANALYSIS

OBJECTIVE

Develop a current understanding of the adidas total sales in the US. Determine Which factors have the most affect on total sales in the US?

TOOLS & SKILLS

- Data cleaning (wrangle, consistency checks)
- Data manipulation (Derive new variable, grouping, aggregating, subsetting, exporting)
- Advanced analysis(geospatial analysis, linear regression analysis, clustering analysis)

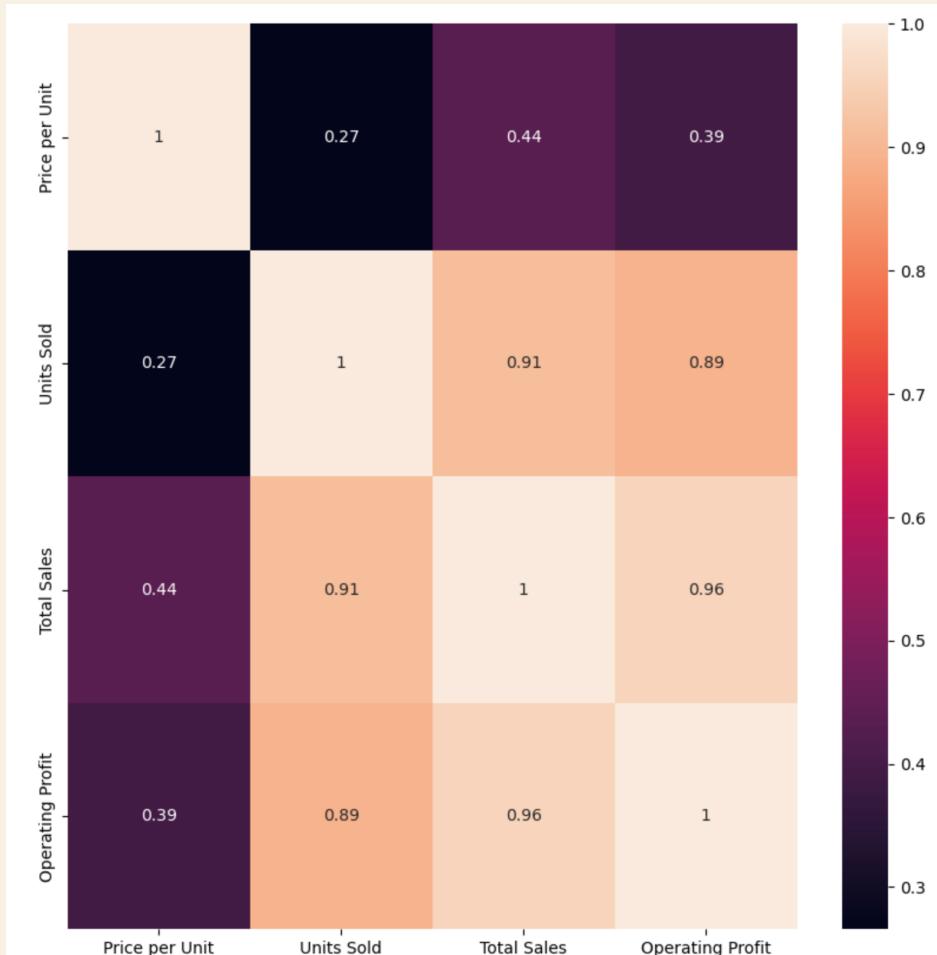
DATA

This data is publicly available open-source data. It was downloaded from Kaggle.com ([Adidas Sales Dataset](#))

LIMITATION

- The data contains incorrect data types for analysis
- The data contains only 2020 to 2021

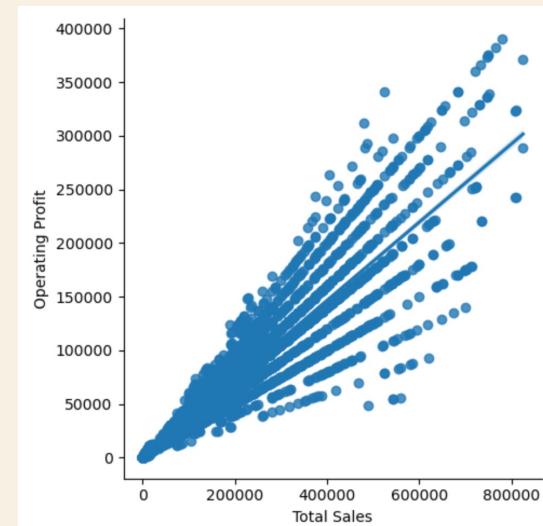
EXPLORING RELATIONSHIPS



```
# Create a subplot with matplotlib
f,ax = plt.subplots(figsize=(10,10))

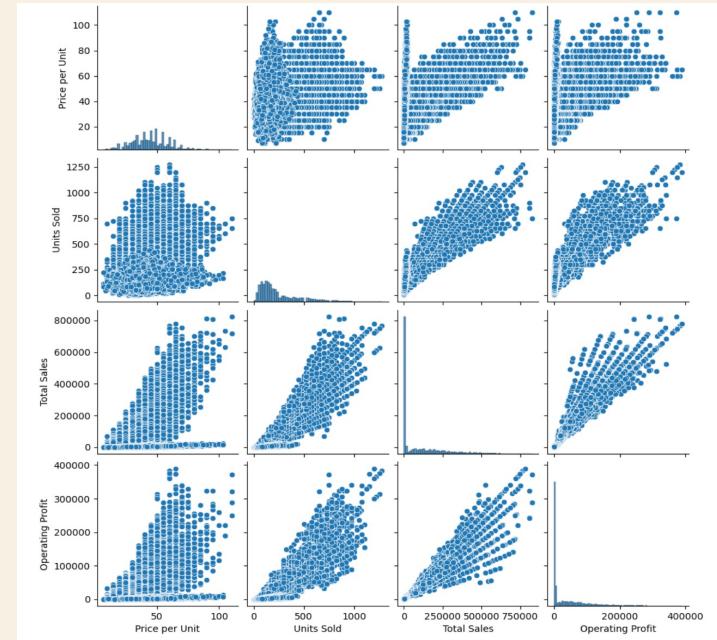
# Create the correlation heatmap in seaborn
corr = sns.heatmap(df_cor.corr(), annot = True, ax = ax)
```

A correlation heatmap helps verify through colors the **strength of relationships** between all variables. The relationship between Total Sales and Operating Profit is 0.96 – it is positive strong relationship and highest.



Question to explore:

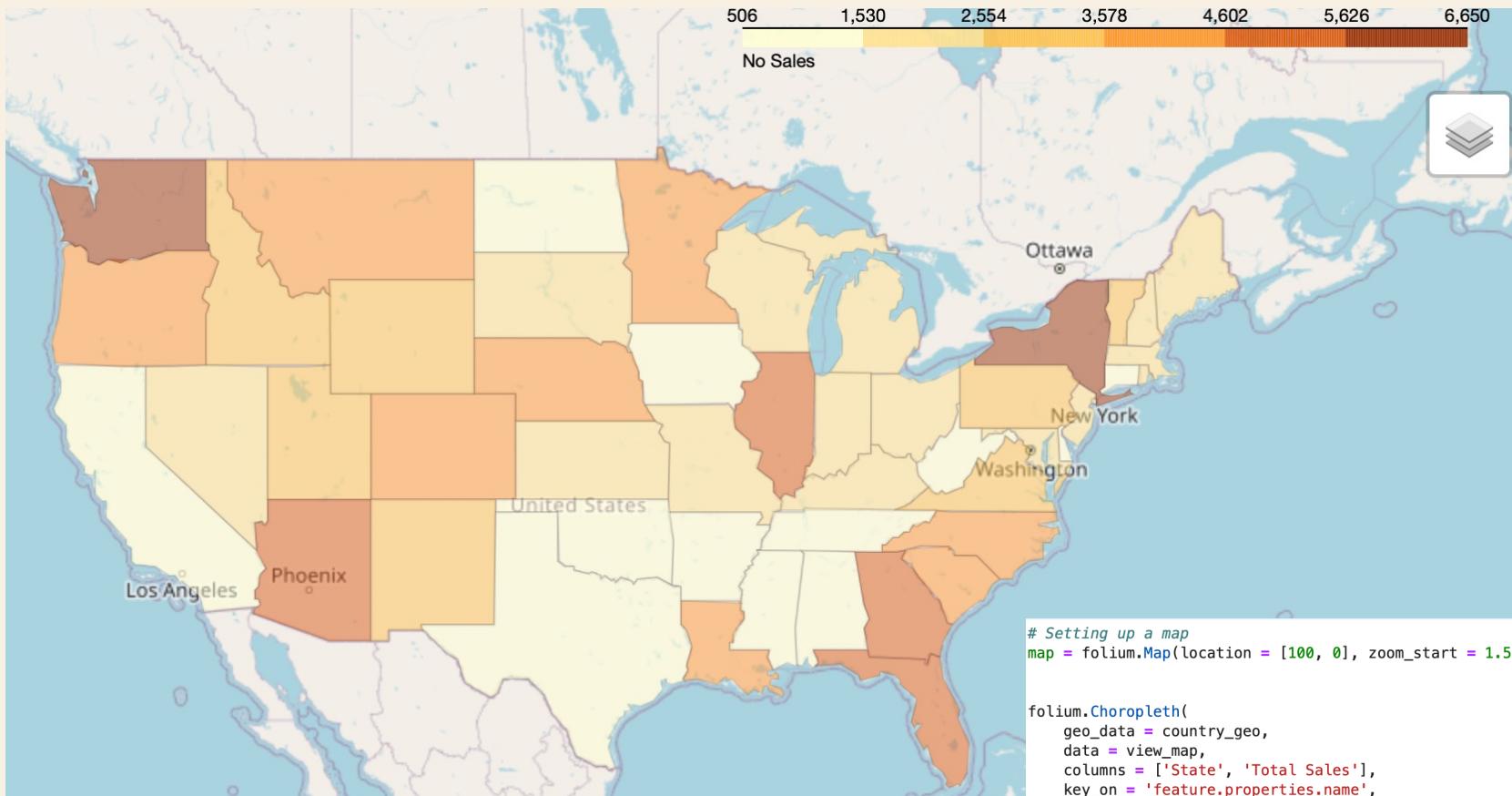
Which factors have the most affect on total sales in the US?



A Pair Plot: At the first and second from the bottom row, the scatterplots look great. There are positive correlations between them.

Scatterplot shows the relationship between Total Sales vs Operating Profit with the line plot. There is strong positive relationship and the slope of the line is steep which is great.

GEOSPATIAL ANALYSIS



The frequency of listing by states
view_map['State'].value_counts(dropna = False)

California	432
Texas	432
New York	360
Florida	360
Mississippi	216
Oregon	216
Louisiana	216
Idaho	216
New Mexico	216
Georgia	216
Arkansas	216
Virginia	216



Questions to explore:
Which region has the most total sales?
From Northeast, Northwest, Southeast and Southwest

Do regions have an impact on the total sales amount? yes, region has an impact on the total sales amount.

ADVANCED TECHNIQUES-REGRESSION



	Actual	Predicted
0	2880.0	2622.825225
1	779.0	342.777563
2	1326.0	867.788538
3	10100.0	15470.593790
4	7392.0	8550.449134
5	7260.0	8385.445685
6	7003.0	4912.873095
7	4884.0	7857.934658
8	3834.0	5060.376178
9	2632.0	2337.819267
10	7000.0	5960.394992

Slope: [[2.50005226]]

Mean squared error: 1796553450.7975018

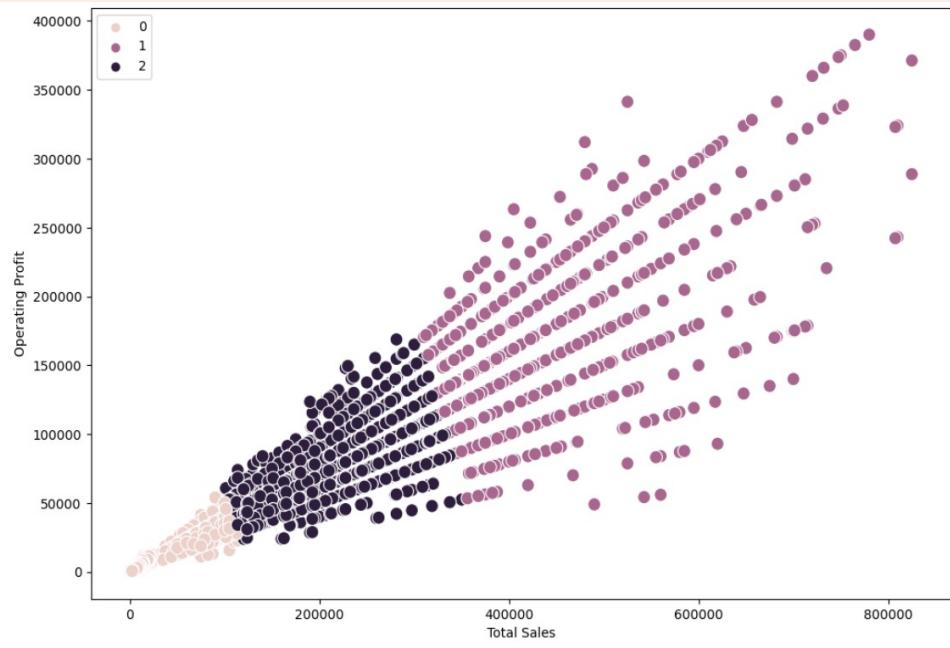
R2 score: 0.9122218446438342

There is a strong positive relationship between variables. The number that represents the total sales increase then the operating profit also increase. The high MSE and high R2 score are good for making predictions

```
# Splitting data into a train set and a test set
X_train_2, X_test_2, y_train_2, y_test_2 = train_test_split(X_2, y_2, test_size=0.3, random_state=0)
```

```
# Creating predictions based on X values from test set
y_predicted_2 = regression.predict(X_test_2)
```

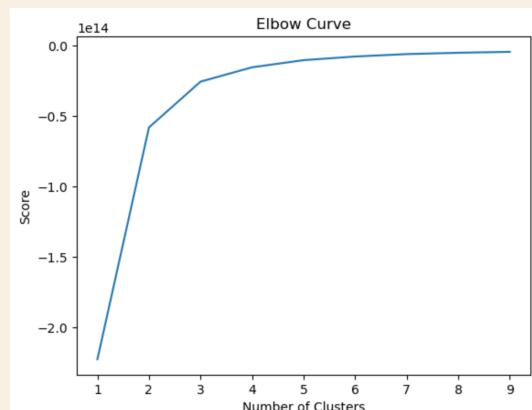
ADVANCED TECHNIQUES-CLUSTERING



Descriptive analysis for clustering

cluster	Total Sales		Price per Unit		Units Sold		Operating Profit	
	mean	median	mean	median	mean	median	mean	median
	dark purple	202416.059281	195000.0	48.542274	50.0	432.896016	425.0	73392.640428
pink	15173.765569	5839.0	42.038849	41.0	143.091785	133.0	5896.094751	2657.5
purple	450352.245863	425000.0	62.458629	60.0	736.347518	725.0	167056.507092	157500.0

The first cluster, in medium purple (coded as “2” in the legend), contains the points with the highest operating profit and the highest total sales. The second cluster, in dark purple (coded as “1” in the legend), is also the most populated cluster. It gathers the data points with high operating profit and relatively high total sales but less than the first cluster. The third cluster, in pink (coded as “0” in the legend), includes points with lowest operating profit and lowest total sales.



Elbow technique. The optimal number of clusters shouldn't be too many (otherwise, there won't be much difference between them), while also not being too few. What the elbow technique does, then, is show you the breaking point which adding more clusters won't help better explain the variances in your data.

- The dark purple cluster has the best stats in almost all categories. Total Sales and Operating Profit are highest in mean and median
- Price per unit doesn't seem to matter much while Units Sold is better both mean and median

SUMMARY

- Relationships: According to analysis, there is a strong positive linear relationship between the total sales and operating profit.
- Regions: West (California) and Northeast (New York) are the highest total sales across all region.

Deliverables:

All analysis and suggestions have been collected in a [Tableau Public](#)

Need extra information?

Please click here and check my GitHub repository

A large, abstract graphic on the left side of the slide features a hexagonal shape composed of several overlapping triangles. The triangles are colored in shades of grey, blue, and magenta. Some triangles have white horizontal stripes, while others are solid colors. The overall effect is a modern, geometric, and slightly abstract design.

INSTACART BASKET DATA ANALYSIS

OBJECTIVE

Provide an analysis of Instacart's sales patterns that will show customer behavior to help develop marketing and sales strategies to increase revenue

TOOLS & SKILLS

- Data wrangling and data frame merging in Python
- Deriving new variables
- Crosstabs and pivot tables in Python
- Visualizations in multiple Python libraries
- Markup and notebook management in Jupyter

DATA

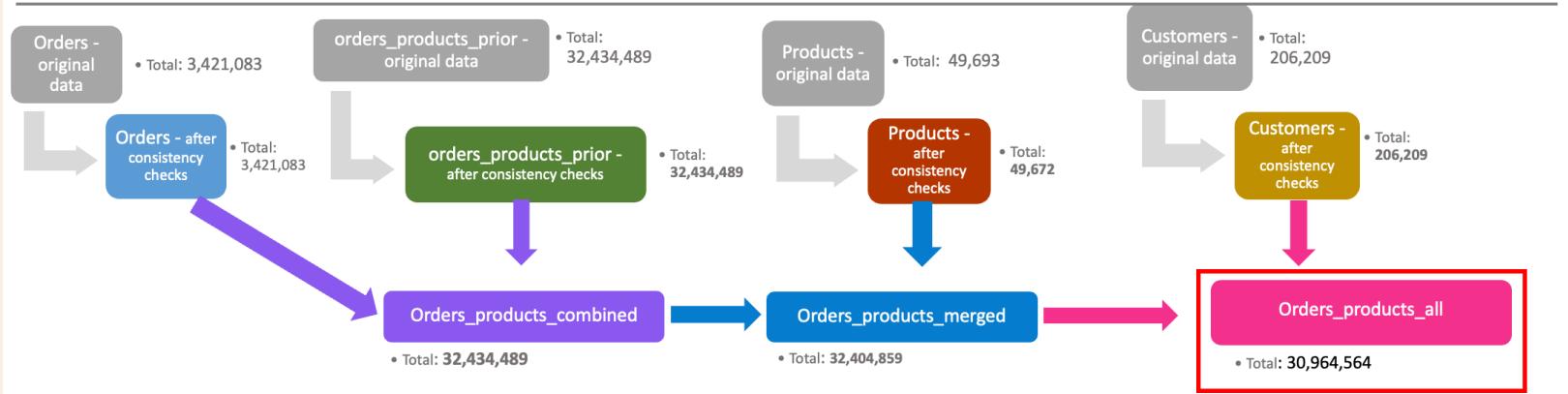
Open source data from Instacart and a customer data set created for the purpose of this project. [Customers Data Set](#)

LIMITATION

- Data only contains records from 2017
- Customer demographics are limited, only including age, family size, income, and marital status

ANALYSIS AND INSIGHT

Population flow



Population Flow gives an overview of all merging phases. Different datasets have been merged to reach the most complete and up to date dataset.

Consistency Checks



Consistency checks

Dataset	
orders	206,209 missing values
products	16 missing values
orders_products_prior	0 missing
customers	0 missing
department	0 missing

Wrangling Steps

Wrangling steps

Dataset	Data set	columns	Column derivations and aggregations
orders	orders	order_id, user_id, eval_set, order_number, order_dow, order_hour_of_day, days_since_prior_order	Dataset: orders_products_merged.pkl (In Task_4.7))
products	Prods	product_id, product_name, aisle_id, department_id, prices	New Column: price_range Column/s It Was Derived From: prices
orders_products_prior	Departments	user_id, First Name, Surname, Gender	New Column: busiest day Column/s It Was Derived From: orders_day_of_week
customers			New Column: busiest_days Column/s It Was Derived From: orders_day_of_week
department			New Column: busiest_period_of_day Column/s It Was Derived From: order_hour_of_day
			Dataset: orders_products_merged_update.pkl (In Task_4.8)
			New Column: loyalty_flag Column/s It Was Derived From: max_order

Consistency Checks:

Checking if values are missing or duplicate and checking for mixed type variables.

Wrangling steps:

Changing columns headers and data types or creating new data frames.

Column derivations and aggregations

Creating new columns/variables and aggregated variables.

ANALYSIS AND INSIGHT

```
[21]: # Creating income_category columns
df_high_act_cust.loc[df_high_act_cust['income'] < 70000, 'income_category'] = 'Low'
df_high_act_cust.loc[(df_high_act_cust['income'] >= 70000) & (df_high_act_cust['income'] < 100000), 'income_category'] = 'Middle-class'
df_high_act_cust.loc[(df_high_act_cust['income'] >= 100000) & (df_high_act_cust['income'] < 130000), 'income_category'] = 'Upper-mid-class'
df_high_act_cust.loc[df_high_act_cust['income'] >= 130000, 'income_category'] = 'High'

[22]: # Checking income_category values
df_high_act_cust['income_category'].value_counts(dropna = False)

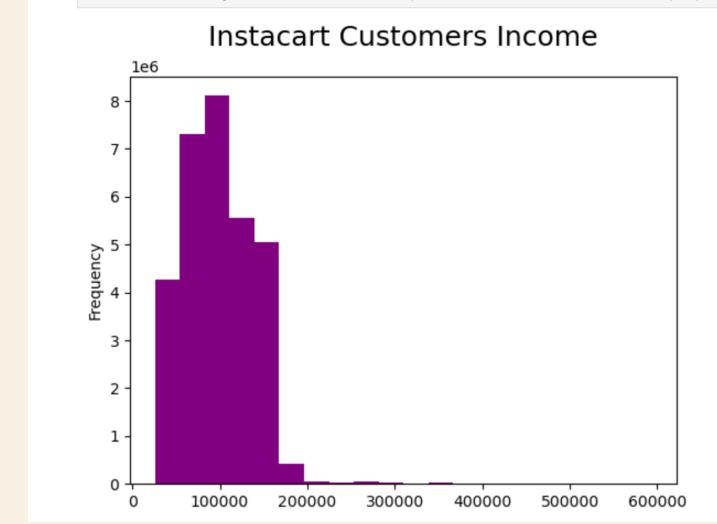
[22]: Low      8520533
Middle-class 8236629
High        7401414
Upper-mid-class 6805988
Name: income_category, dtype: int64

[23]: # Confirming the added column
df_high_act_cust.shape

[23]: (30964564, 36)

[24]: # Create an income histogram
plt.title('Instacart Customers Income', fontsize = 18, pad=20)
hist_inc = df_high_act_cust['income'].plot.hist(bins = 20, color = 'purple')
```

Creating a new column or income_category



Customers with low income has the highest order

```
[16]: # create busiest_of_days_03
busiest_of_days_03 = []

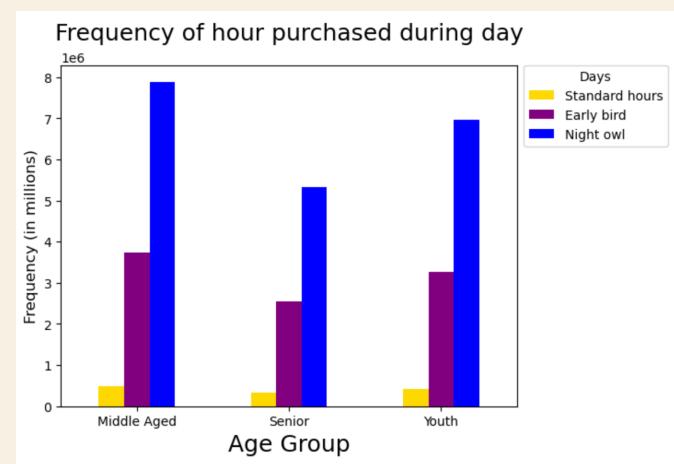
for value in df_high_act_cust["order_hour_of_day"]:
    if value in [10,11,14,15,13,12,16,9]:
        busiest_of_days_03.append("Standard hours")
    elif value in [23,6,0,1,5,2,4,3]:
        busiest_of_days_03.append("Early bird")
    else:
        busiest_of_days_03.append("Night owl")

[17]: df_high_act_cust['busiest_of_days_for_chart']= busiest_of_days_03

[18]: # Customer comparison by region & order hour of day
crosstab_age_hour = pd.crosstab(df_high_act_cust['age_category'], df_high_act_cust['busiest_of_days_for_chart'], d
```

Creating busiest_of_days

Crosstabs were created from the merged data frames to better understand the connections between variables.

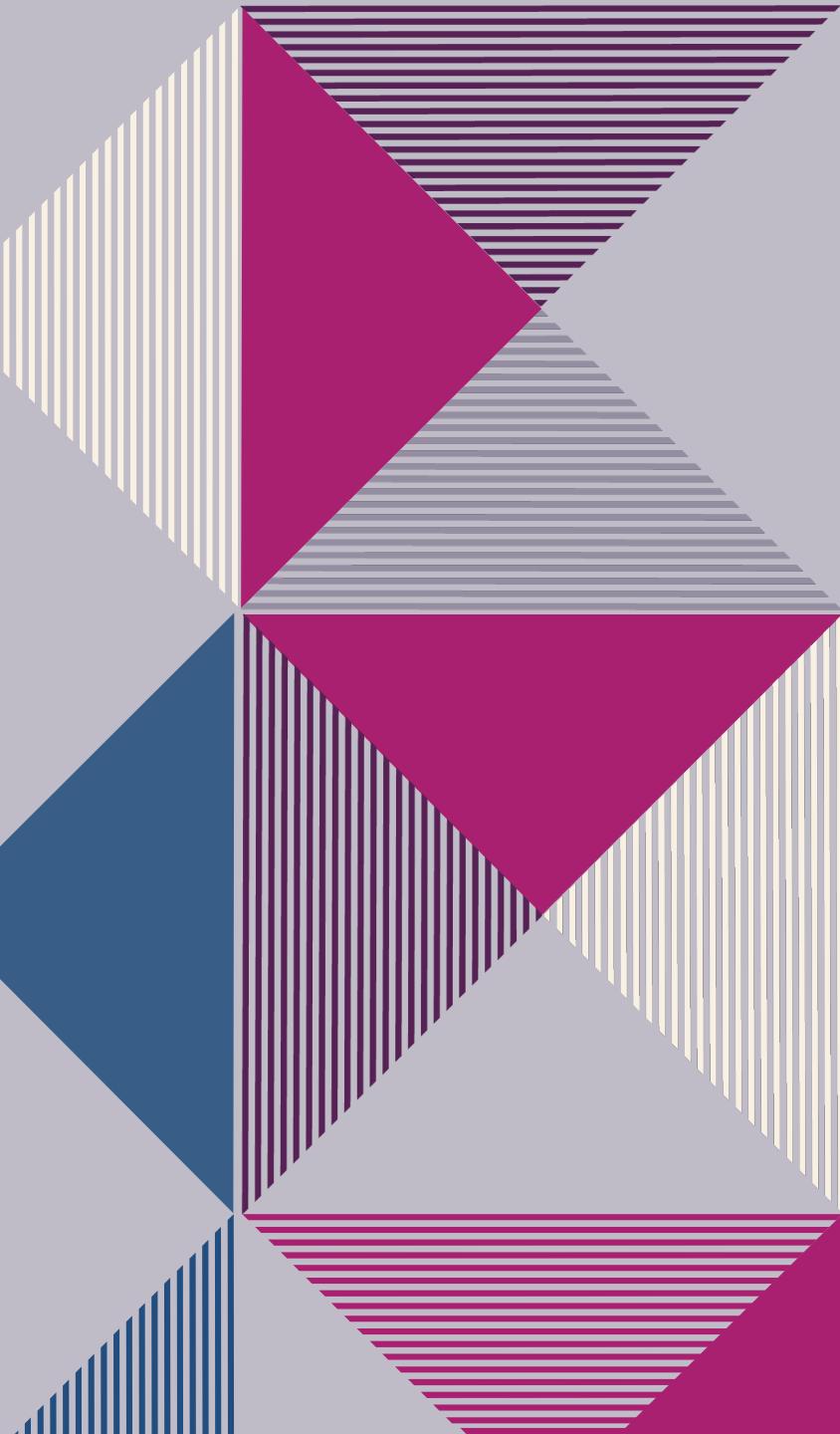


Here is 5 Business key questions python code
[Github](#)

Middle-aged people tend to purchase the most items and at night Or after work.

RECOMMENDATIONS

- **Ads**—Schedule advertisements on the busy weekends so it reaches as many people as possible, specifically between 10am and 2pm.
- **Pricing**—Expand the market of higher priced items to boost their numbers and bring in more revenue.
- **Products**—The most popular products being ordered are those in produce, dairy and eggs, snacks, beverages. Instacart should carry on advertising those product and potentially offering deals to drive sales.
- **Loyalty**—Most customers are new or regular. To ensure that new customers continue to return so Instacart could consider giving percentage discounts for orders to new users to increase uptake.
- **Geography**—The Southern customers tend to be regular customers in terms of ordering-time habits, they also tend to fall into the low-income class. Using this region to test new products would be beneficial and we should focus on growing the customer bases in other regions.

A large, abstract graphic on the left side of the slide features a hexagonal shape composed of several triangles. The triangles are filled with different patterns: some are solid colors (dark blue, light gray, magenta), while others have horizontal or vertical stripes. The overall effect is a modern, geometric, and slightly abstract design.

ROCKBUSTER STEALTH LLC

OBJECTIVE

Rockbuster Stealth is a fictional movie rental company with stores over the world. They are planning to launch an online video rental service in order to stay competitive

DATA

This dataset is provided by PostgreSQL for usage in tutorials. It contains data about film inventory, customers, payments, and associated details. [Rockbuster Data Set](#)

TOOLS & SKILLS

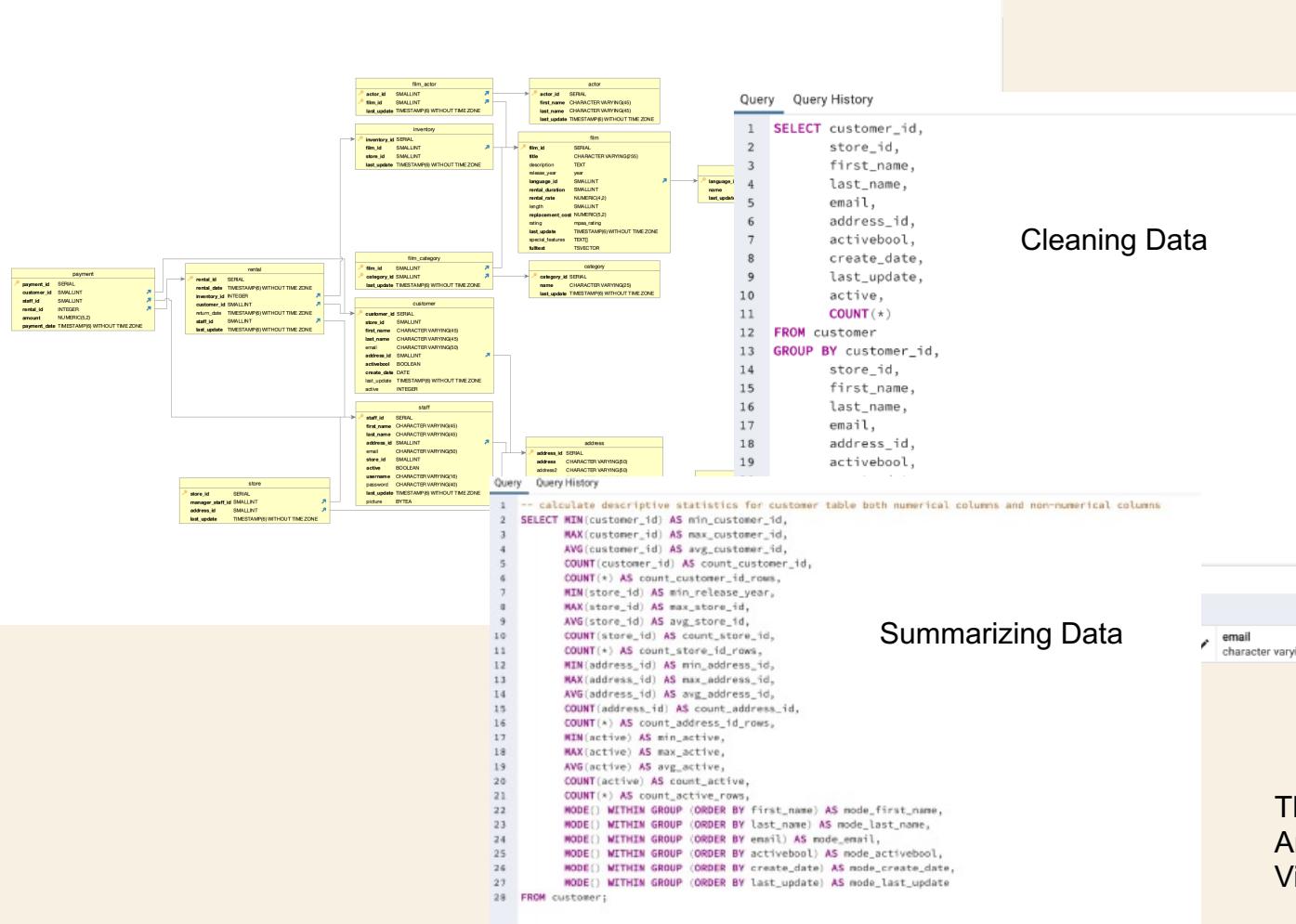
- Relational databases in SQL
- Entity relationship diagram creation and usage
- Data dictionary creation
- Database querying, filtering, and cleaning
- Joining tables in relational database
- Subqueries and common table expressions

LIMITATION

- Only have internal records to work with provided by company

ANALYSIS AND INSIGHT

Entity Diagram



Cleaning Data

```

Query   Query History
1  SELECT customer_id,
2      store_id,
3      first_name,
4      last_name,
5      email,
6      address_id,
7      activebool,
8      create_date,
9      last_update,
10     active,
11     COUNT(*)
12  FROM customer
13 GROUP BY customer_id,
14      store_id,
15      first_name,
16      last_name,
17      email,
18      address_id,
19      activebool,

```

Summarizing Data

```

1 -- calculate descriptive statistics for customer table both numerical columns and non-numerical columns
2 SELECT MIN(customer_id) AS min_customer_id,
3       MAX(customer_id) AS max_customer_id,
4       AVG(customer_id) AS avg_customer_id,
5       COUNT(customer_id) AS count_customer_id,
6       COUNT(*) AS count_customer_id_rows,
7       MIN(store_id) AS min_release_year,
8       MAX(store_id) AS max_store_id,
9       AVG(store_id) AS avg_store_id,
10      COUNT(store_id) AS count_store_id,
11      COUNT(*) AS count_store_id_rows,
12      MIN(address_id) AS min_address_id,
13      MAX(address_id) AS max_address_id,
14      AVG(address_id) AS avg_address_id,
15      COUNT(address_id) AS count_address_id,
16      COUNT(*) AS count_address_id_rows,
17      MIN(active) AS min_active,
18      MAX(active) AS max_active,
19      AVG(active) AS avg_active,
20      COUNT(active) AS count_active,
21      COUNT(*) AS count_active_rows,
22      MODE() WITHIN GROUP (ORDER BY first_name) AS mode_first_name,
23      MODE() WITHIN GROUP (ORDER BY last_name) AS mode_last_name,
24      MODE() WITHIN GROUP (ORDER BY email) AS mode_email,
25      MODE() WITHIN GROUP (ORDER BY activebool) AS mode_activebool,
26      MODE() WITHIN GROUP (ORDER BY create_date) AS mode_create_date,
27      MODE() WITHIN GROUP (ORDER BY last_update) AS mode_last_update
28 FROM customer;

```

Data Dictionary

FACT TABLE

Payment

Table	Column	Data Types	Description
	payment_id	integer (32)	a unique identifier for payments
	customer_id	smallint (16)	a unique identifier for customer
	staff_id	smallint (16)	a unique identifier for staff
	rental_id	integer (32)	a unique identifier for rental
	amount	numeric (5,2)	amount paid by the customer
	payment_date	timestamp without time zone (6)	the date and time when a payment is paid

Links to

Table	Join
customer	payment.customer_id = customer.customer_id
staff	payment.staff_id = staff.staff_id
rental	payment.rental_id = rental.rental_id

DIMENSION TABLES

Rental

Table	Column	Data Types	Description
	rental_id	integer (32)	a unique identifier for rental
	rental_date	timestamp without time zone (6)	the date and time when the rental is made
	inventory_id	integer (32)	a unique identifier for inventory

These are the process (Cleaning, Summarizing, Performing Descriptive Analysis) to help to perform query and analyze the data.
View [Data Dictionary](#)

ANALYSIS AND INSIGHT

Subquery

```

SELECT D.country,
       COUNT(DISTINCT A.customer_id) AS all_customer_count,
       COUNT(DISTINCT top_5_customers.customer_id) AS top_customer_count
FROM customer A
INNER JOIN address B ON A.address_id = B.address_id
INNER JOIN city C ON B.city_id = C.city_id
INNER JOIN country D ON C.country_id = D.country_id
LEFT JOIN
    (SELECT A.customer_id, SUM(A.amount) AS total_amount_paid, B.first_name, B.last_name, D.city, E.country
     FROM payment A
     INNER JOIN customer B ON A.customer_id = B.customer_id
     INNER JOIN address C ON B.address_id = C.address_id
     INNER JOIN city D ON C.city_id = D.city_id
     INNER JOIN country E ON D.country_id = E.country_id
     WHERE E.country IN ('India', 'China', 'United States', 'Japan', 'Mexico',
                         'Brazil', 'Russia Federation', 'Philippines', 'Turkey',
                         'Indonesia') AND D.city IN ('Aurora', 'Acua', 'Citrus Heights',
                         'Iwaki', 'Ambattur', 'Shaweei', 'So Leopoldo', 'Teboksary',
                         'Tianjin', 'Cianjur')
    GROUP BY A.customer_id, E.country, D.city, B.first_name, B.last_name
    ORDER BY total_amount_paid DESC
    LIMIT 5)
AS top_5_customers ON D.country = top_5_customers.country
GROUP By D.country
ORDER BY top_customer_count DESC
LIMIT 5;
  
```

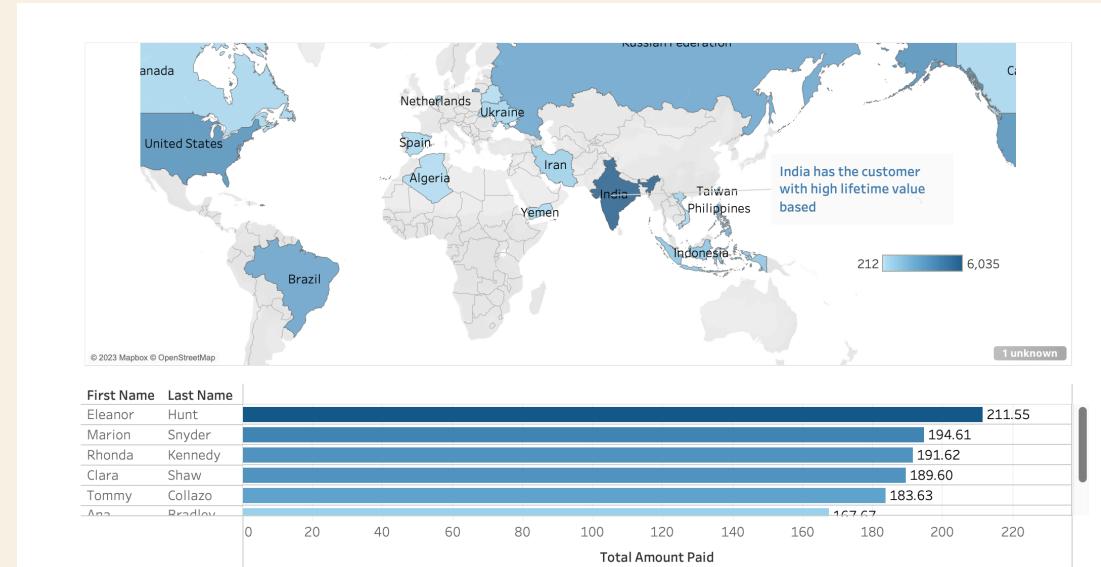
	Data Output	Messages	Notifications
1	country character varying (50)	all_customer_count bigint	top_customer_count bigint
2	1 Japan	31	1
3	2 Mexico	30	1
4	3 China	53	1
5	4 India	60	1
6	5 United States	36	1

```

1   WITH top_ct (ctry, customer_id, top_cnt_customer)
2   AS
3       (SELECT A.customer_id, SUM(A.amount) AS total_amount_paid, B.first_name,
4           B.last_name, D.city, E.country
5           FROM payment A
6           INNER JOIN customer B ON A.customer_id = B.customer_id
7           INNER JOIN address C ON B.address_id = C.address_id
8           INNER JOIN city D ON C.city_id = D.city_id
9           INNER JOIN country E ON D.country_id = E.country_id
10          WHERE E.country IN ('India', 'China', 'United States', 'Japan', 'Mexico',
11                            'Brazil', 'Russia Federation', 'Philippines', 'Turkey',
12                            'Indonesia') AND D.city IN ('Aurora', 'Acua', 'Citrus Heights',
13                            'Iwaki', 'Ambattur', 'Shaweei', 'So Leopoldo', 'Teboksary',
14                            'Tianjin', 'Cianjur')
15          GROUP BY A.customer_id, E.country, D.city, B.first_name, B.last_name
16          ORDER BY total_amount_paid DESC
17          LIMIT 5)
18      SELECT D.country,
19             COUNT(DISTINCT A.customer_id) AS all_customer_count,
20             COUNT(DISTINCT top_ct.top_cnt_customer) AS top_customer_count
21      FROM customer A
22      INNER JOIN address B ON A.address_id = B.address_id
23      INNER JOIN city C ON B.city_id = C.city_id
24      INNER JOIN country D ON C.country_id = D.country_id
25      LEFT JOIN top_ct ON D.country = top_ct.country
26      GROUP By D.country
27      ORDER BY top_customer_count DESC
28      LIMIT 5;
29
30      Data Output Messages Explain X Notifications
31
32      country character varying (50) all_customer_count bigint top_customer_count bigint
33
34      1 Japan 31 1
35      2 Mexico 30 1
36      3 China 53 1
37      4 India 60 1
38      5 United States 36 1
  
```



CTE



To answer the business questions, the right table joins and queries had to be written in SQL. Then the resulting table was exported to a csv file and imported into Tableau. At that point, a visualization showing the answers to the business questions could be created. When performing complex queries, CTE is easier to organize and read. However, readability is not the only consideration when choosing between a CTE and a subquery; performance is also important.

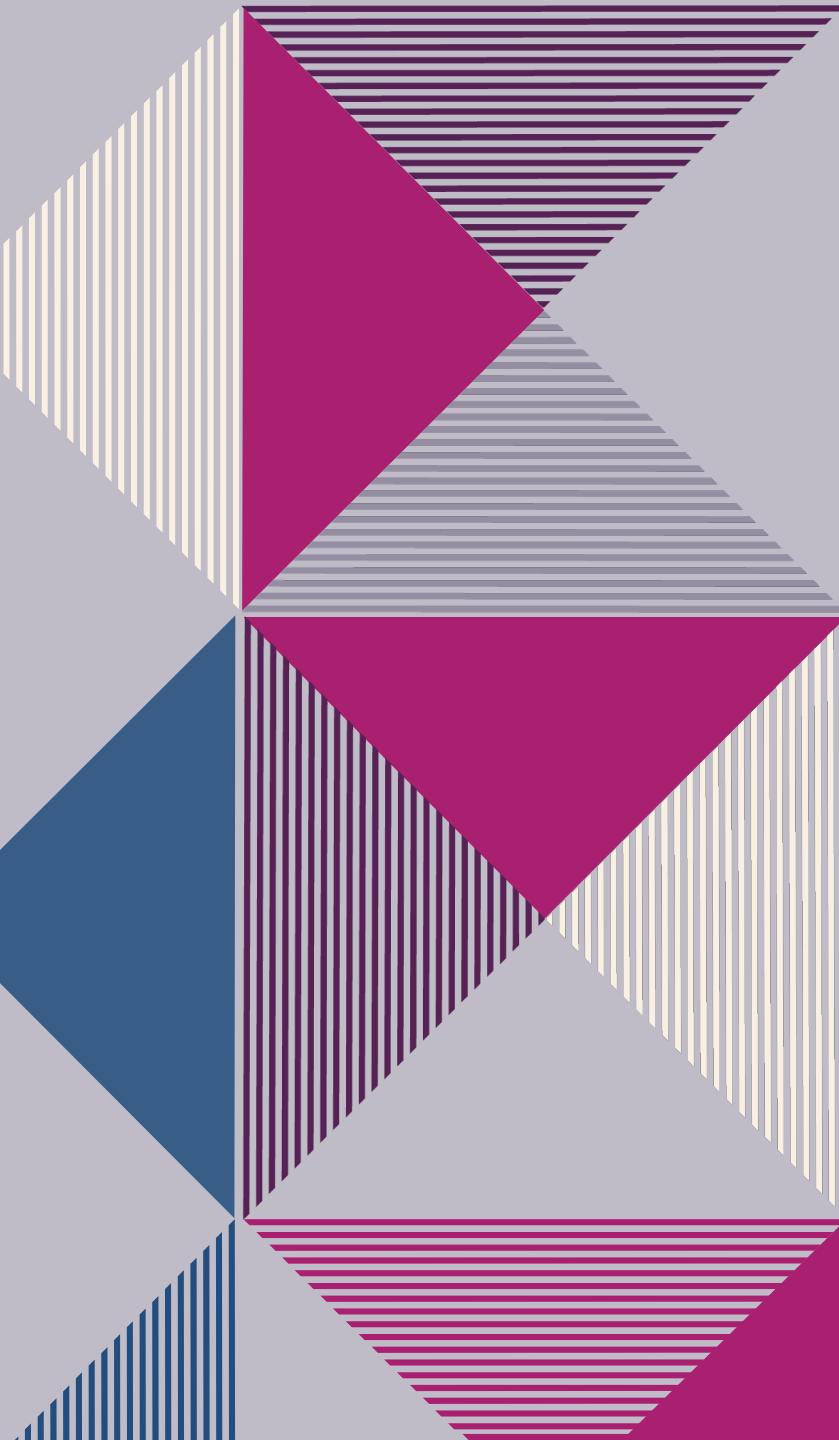
Here is [5 key business questions queries](#) and [Project Visualization](#)

RECOMMENDATIONS

Revenue: To maximize revenue and customer satisfaction, Rockbuster should prioritize launching the top 10 highest movies on their online video rental service.

Customer: To expand their customer base, Rockbuster should conduct further analysis to understand why other regions are not performing as well as Asia and America.

Location: To get a large and diverse customer base, Rockbuster should prioritize the Asia and American markets as their primary targets,



PREPARING FOR UP COMING FLU SEASON

OBJECTIVE

Identify geographic and seasonal trends for annual influenza outbreaks in the USA.

Provide tools for a medical staffing agency to identify where and when to allocate additional medical support.

TOOLS & SKILLS

- Data research project design
- Data profiling and cleaning
- Data integration and transformation
- Statistical hypothesis testing
- Geographic visualizations and time-series forecasting
- Interactive visualizations and storytelling in Tableau

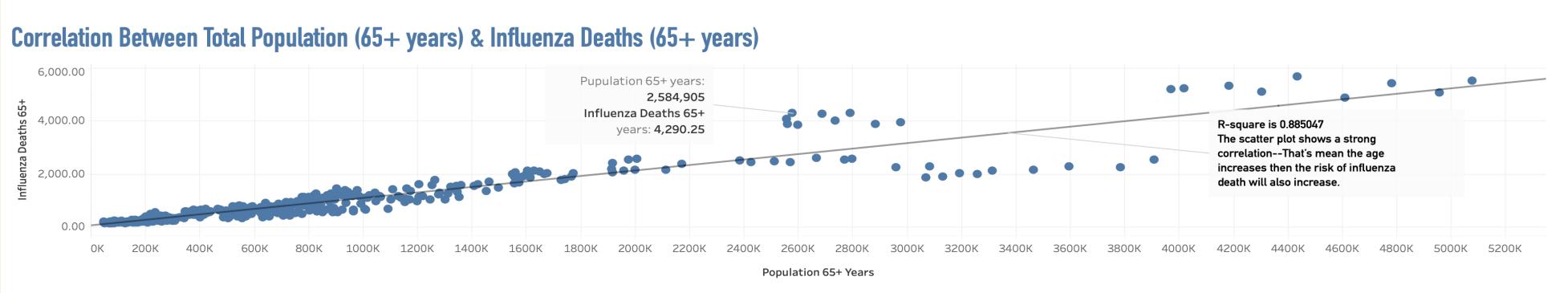
DATA

- [CDC Influenza deaths](#)
- [Population data by geography \(US Census Bureau\)](#)

LIMITATION

- Data is dated from 2009 until 2017
- Influenza death data is 82% suppressed due to confidentiality
- Staffing capacity and sizes of hospitals/clinics is unknown

ANALYSIS AND INSIGHT



Data Spread		
	*Variable 1	*Variable 2
Dataset Name	US Census Data , > 65 years old	Influenza Deaths Data, > 65 years old
Sample or Population?	Sample	Sample
Normal Distribution?	**Yes	**Yes
Variance	7.9E+11	9.5E+05
Standard Deviation (1 Std Dev)	887,017	976
Standard Deviation (2 Std Dev)	1,774,034	1,953
Mean	806,989	889
Outlier Percentage (1 Std Dev)	11%	13%
Outlier Percentage (2 Std Dev)	7%	4%
Upper Limit (1 Std Dev)	1,694,006	1,866
Upper Limit (2 Std Dev)	2,581,023	2,842
Lower Limit (1 Std Dev)	-80,028	-87
Lower Limit (2 Std Dev)	-967,045	-1,064
Total Outliers (Upper Limit, 1 Std Dev)	49	58
Total Outliers (Upper Limit, 2 Std Dev)	30	18
Total Outliers (Lower Limit, 1 Std Dev)	N/A (***)	N/A
Total Outliers (Lower Limit, 2 Std Dev)	N/A	N/A
Correlation		
	Variable 1	Variable 2
Variables:	US Census Data, > 65 years old	Influenza Deaths Data, > 65 years old
Proposed Relationship:	It's Positive relationship. if the person is older than 65 years, then the risk for influenza death is	
Correlation Coefficient	0.9	
Strength of Correlation	It is a strong correlation. The age increases then the risk of influenza death will also increase	
Usefulness / Interpretation	This is useful and it supports my hypothesis.	

Research hypothesis

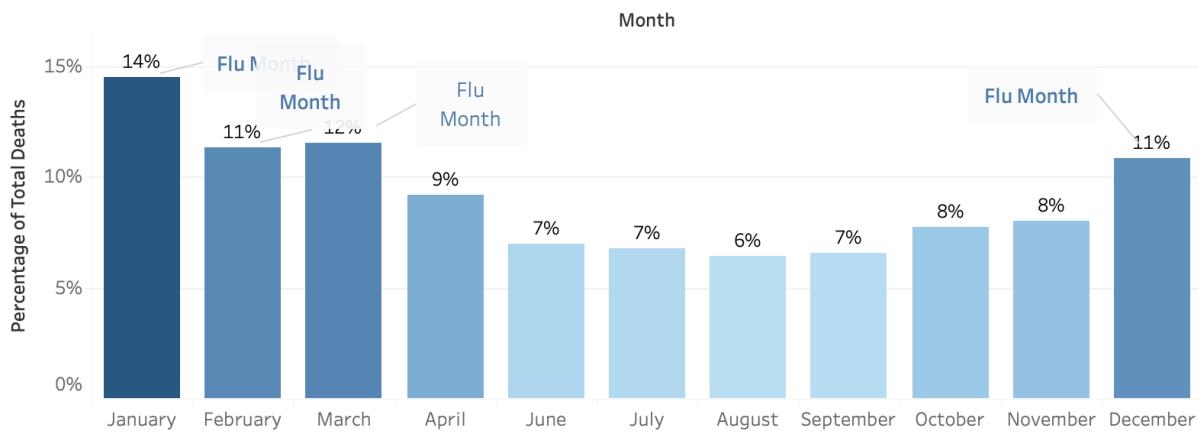
If persons who over the aged of 65 years, they are a higher risk for influenza death

Descriptive Analysis

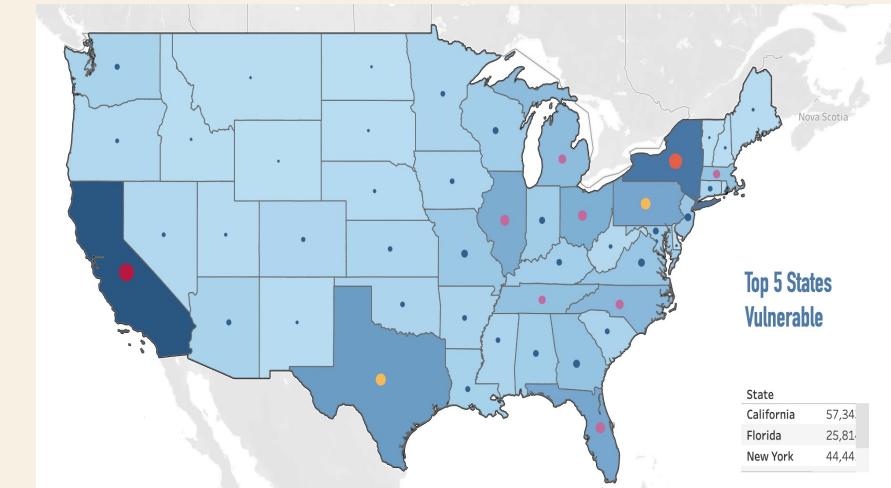
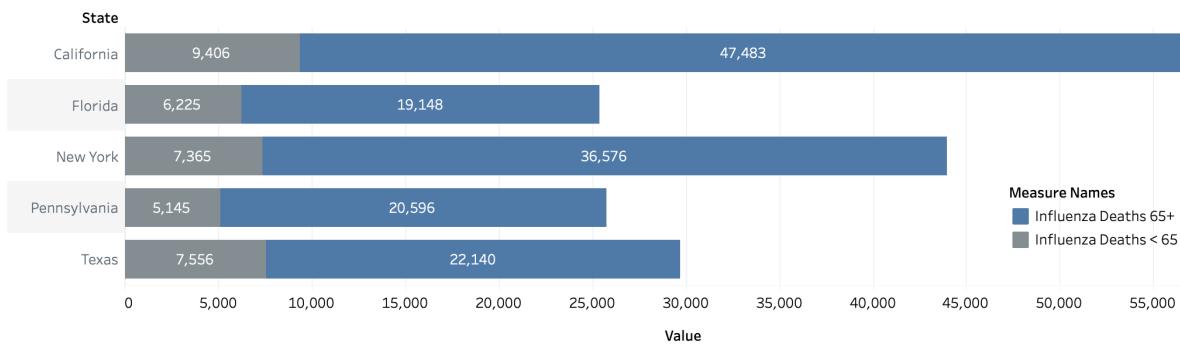
According to our hypothesis, mortality rate increases or is higher with increased age. The correlation study suggests strong correlation between age and mortality rate. The statistics for the same are summarized in table

ANALYSIS AND INSIGHT

Monthly Influenza Deaths



Number of Deaths due to Influenza in the top 5 states



Visualization Project

Influenza Seasonal Peak:

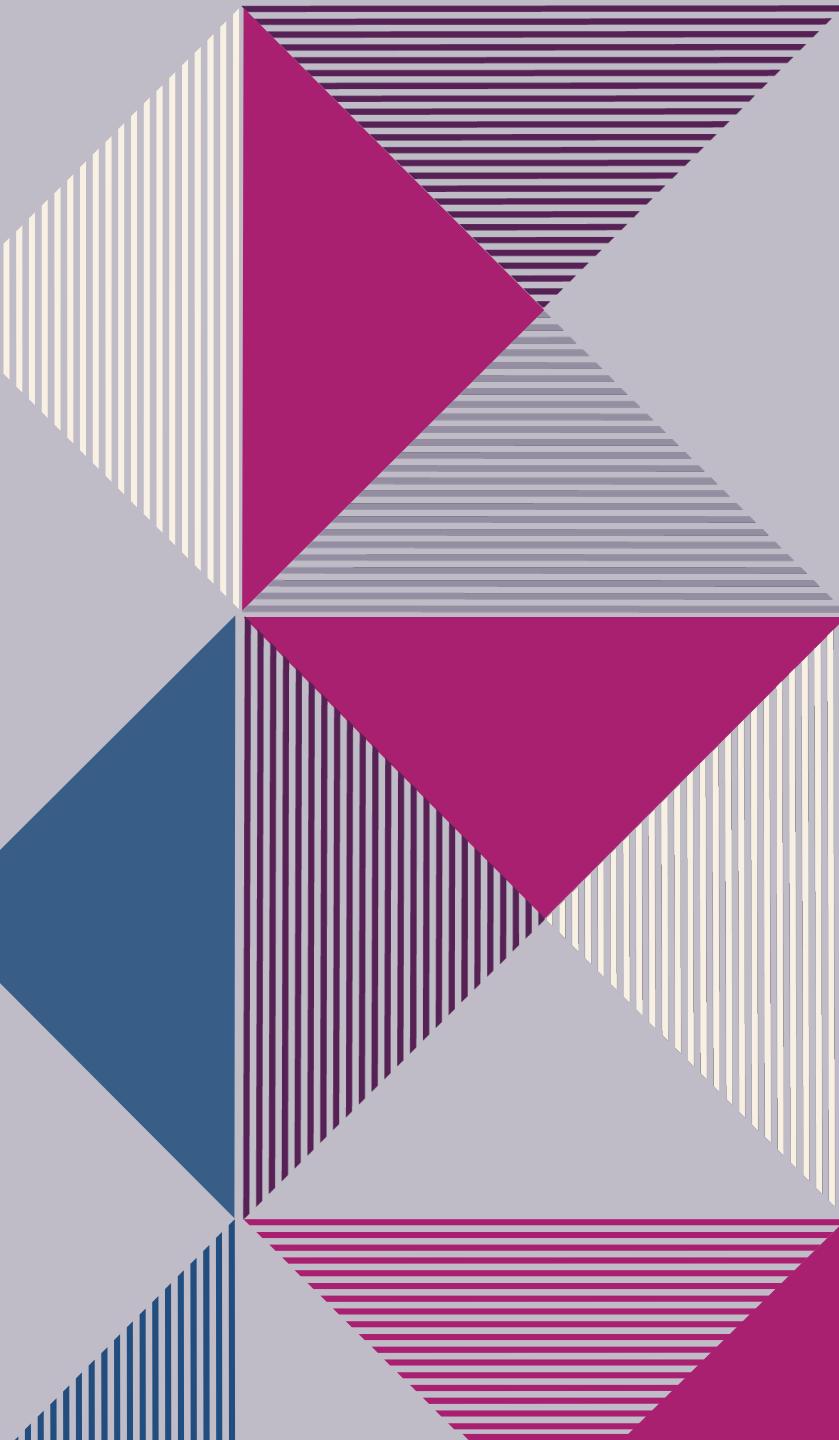
December, January, February and March

The Most Populous States that have the Highest Number of Deaths:

California, New York, Texas, Pennsylvania, Florida. Ages of vulnerable population are 65 years and older

RECOMMENDATIONS

- **Staff:** Since the peak of flu is usually between December to January. More staff should be made available in this peak period. Staff deployments should also consider and prioritize the states that have the most flu deaths plus the prioritize of the vulnerable population (65 years and older)
- **Flu Shot:** Vaccines should be ready mostly in states with the most vulnerable population
- **Survey Evaluation:** To make sure your process is effective and adjust if needed.
Recommend to monitor the effectiveness by using KPIs on Jan, Feb, March and April.

A large, abstract graphic on the left side of the slide features a hexagonal shape composed of several triangles. The triangles are filled with different patterns: some are solid dark red, others have horizontal white stripes, and some have vertical grey stripes. The overall effect is a dynamic, modern, and geometric design.

GAMECO MARKETING DATA ANALYSIS

OBJECTIVE

Develop a current understanding of the global retail videogame sales market, to inform GameCo's efforts to increase market share.

TOOLS & SKILLS

- Data quality, integrity, and consistency checks
- Data cleaning
- Pivot tables (data grouping & summarizing)
- Descriptive analysis
- Excel visualizations
- Reporting in PowerPoint

DATA

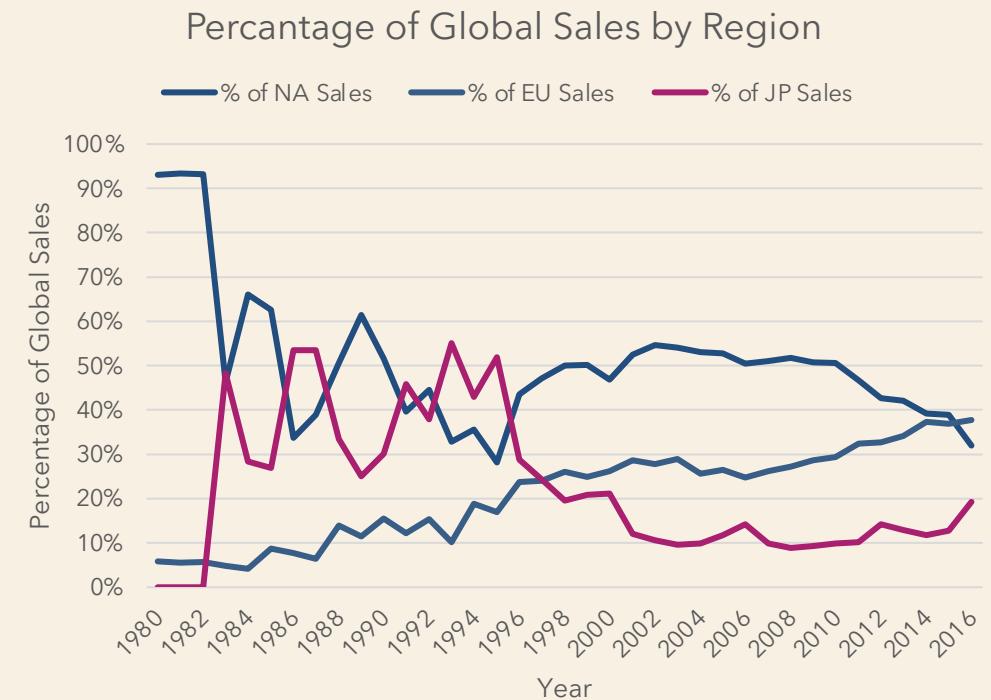
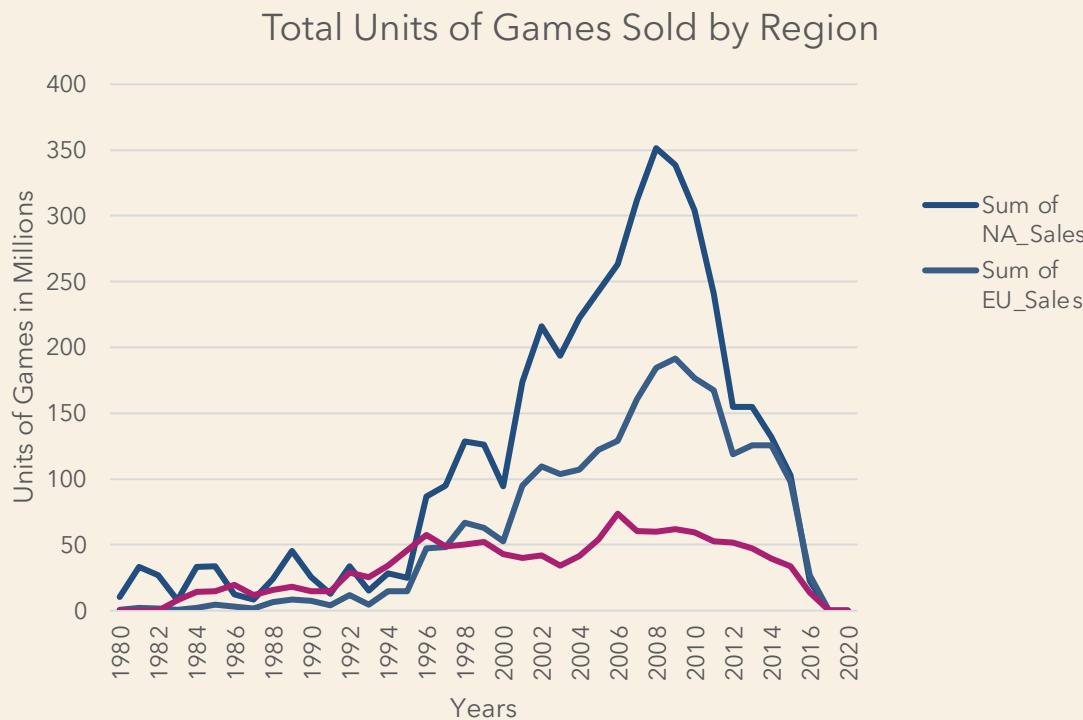
The data is made publicly available by [VGChartz](#). It covers historical retail sales of videogames for games that sold more than 100,000 copies, until 2016.

LIMITATION

- The data only goes until 2016
- No revenue data, just units sold

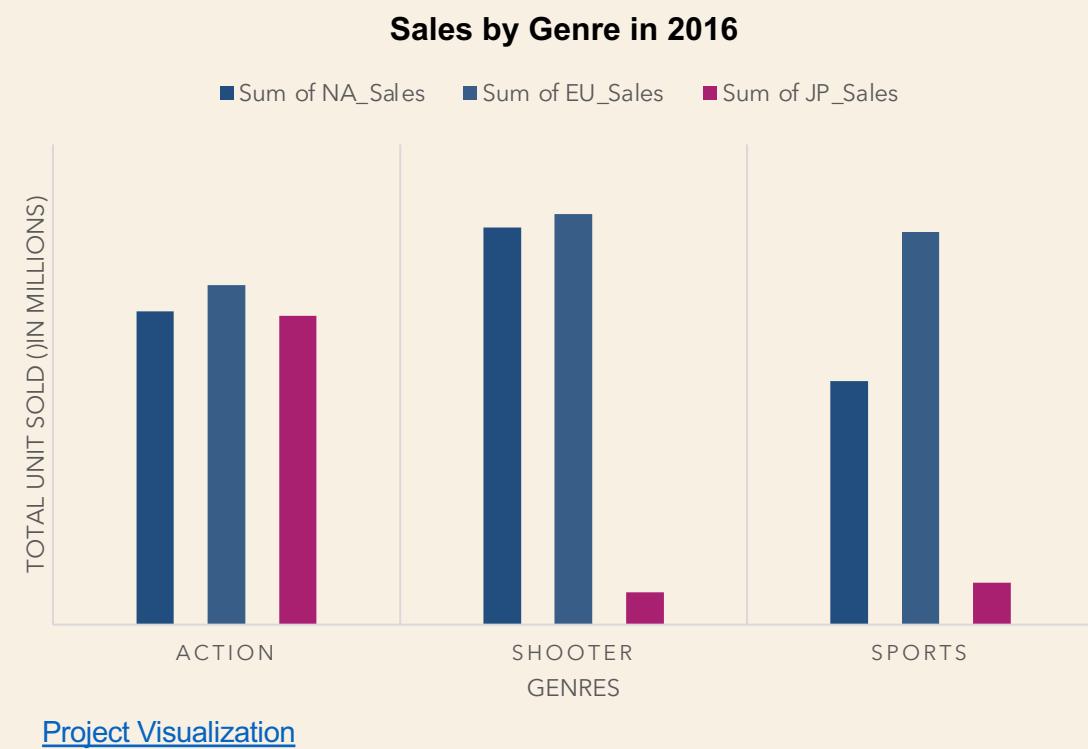
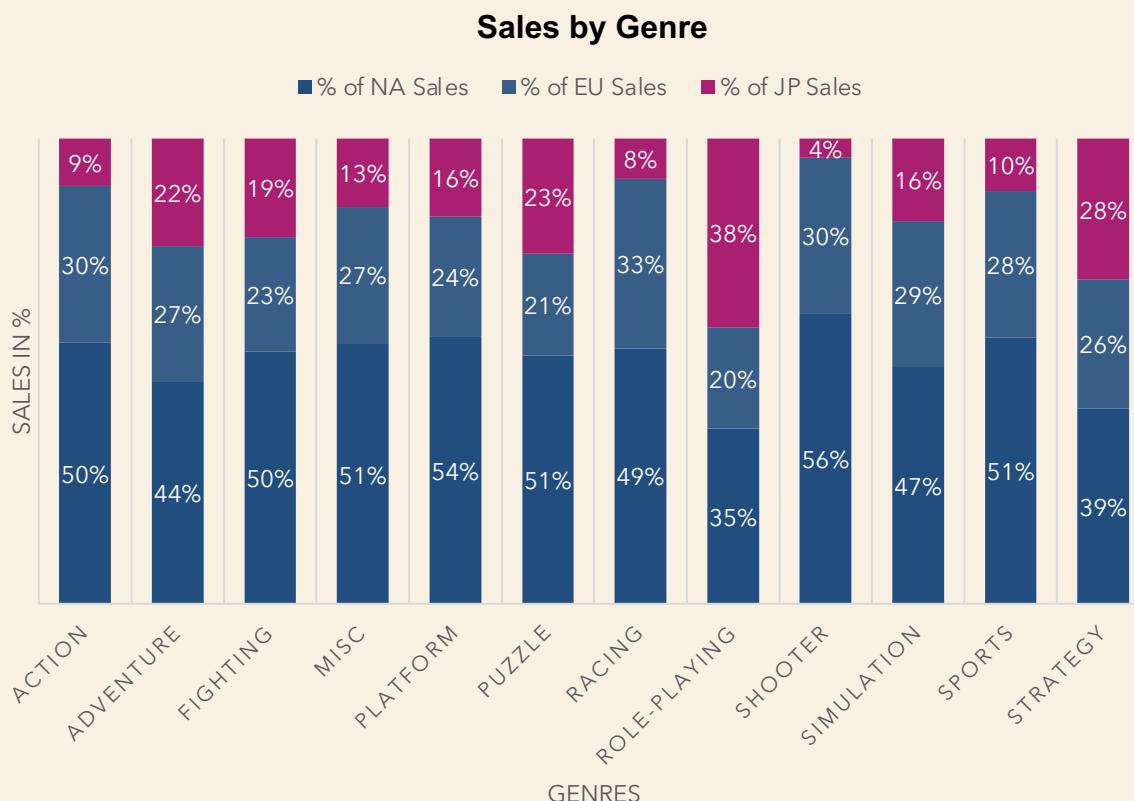
ANALYSIS AND INSIGHT

- Using line chart that represent the proportion of global sales for North America, Europe and Japan by years. There are significant changes between them and should focus deeper insight.
- After seeing the data behaves in 2016, GemeCo marketing team should investigate more in each region for seeing the deeper insights and see which region they should spend more marketing budget in 2017



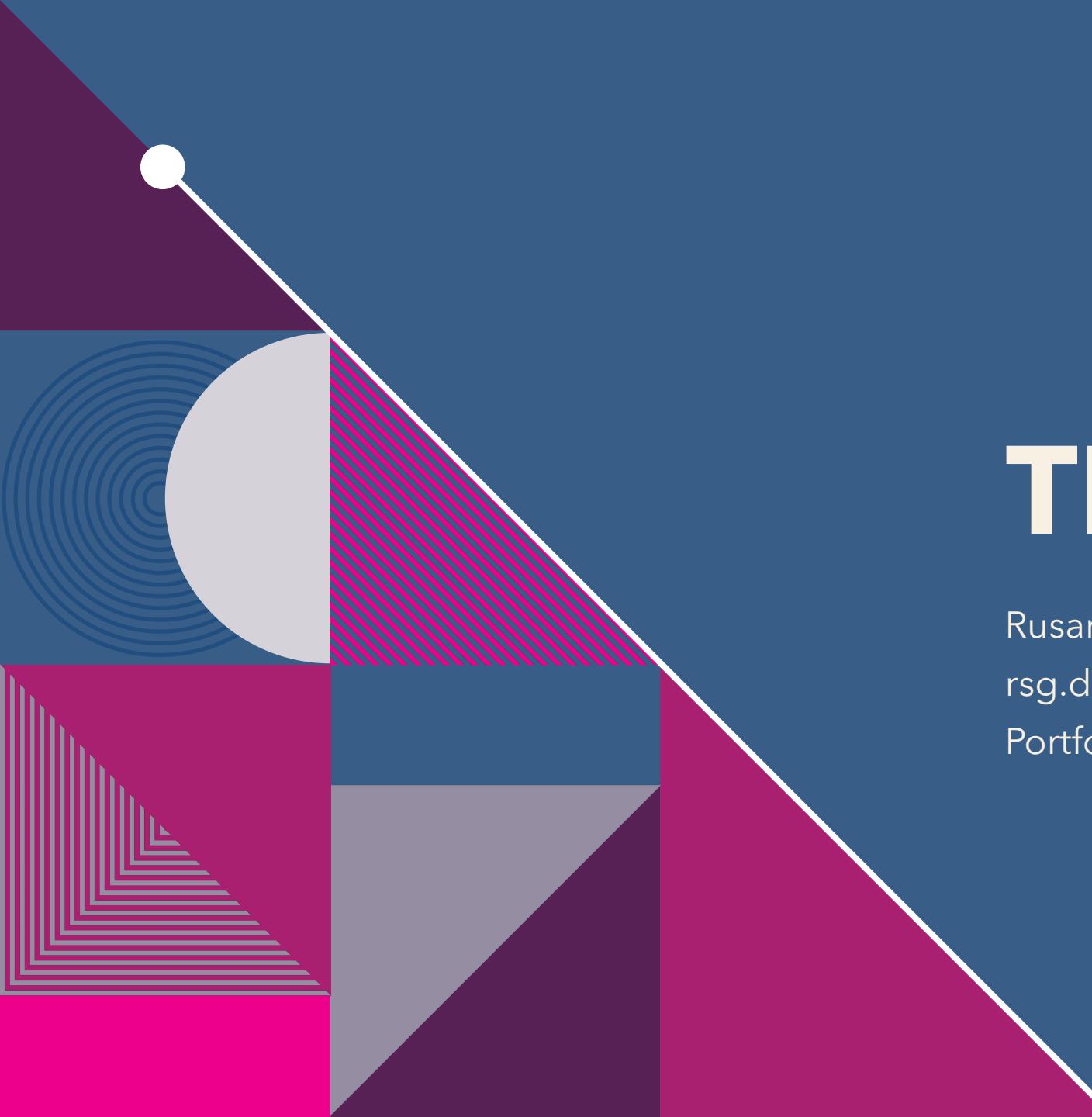
ANALYSIS AND INSIGHT

- Shooter game is the highest sales in genres for North America both overall and in the past year.
- Shooter game is the highest sales in genres for Europe over year however Action game also made the highest sales in genre when it represents the proportion of global sales over years
- Role-Playing game is the highest proportion for Japan sales over years however Action game has been in the top spot in 2016.



RECOMMENDATIONS

- **Budget:** Break it proportionally by region for sales numbers. Focus for the top selling genres within each region.
- **Marketing:** Put more money into the growing markets to increase revenue where demand is higher.
- **Growth:** North America is a huge market and sales there have been declining. Update and create a lot of new games (genres, title, publisher, platform) to help increase revenue there.



THANK YOU

Rusamijan Permison

rsg.design@gmail.com

Portfolio: www.meepermison.com