DAVID EY

# DATA ANALYTICS PORTFOLIO

# PROJECTS

## GameCo

Analyze global video game sales to plan new game development

## Influenza Preparation

Use historical flu & census data to allocate medical staff

## Rockbuster Stealth

Determine movie acquisition strategy for a new streaming service

## Instacart

Initial data & exploratory analysis for insights and better segmentation

## Pig E. Bank

Assess client & transaction risk, build and optimize models to assist a compliance program

## Airbnb Berlin: Pricing Factors & Effects on Local Market

Determine price and rating factors, effects on local markets, impact of recent regulation

GAMECO:
VIDEO
GAMES
SALES
ANALYSIS

# BACKGROUND

## GOAL

Perform a descriptive analysis of a video game sales data set to foster a better understanding of how GameCo's new games might fare in the market.

## TOOLS

Excel
PowerPoint

## DATASET

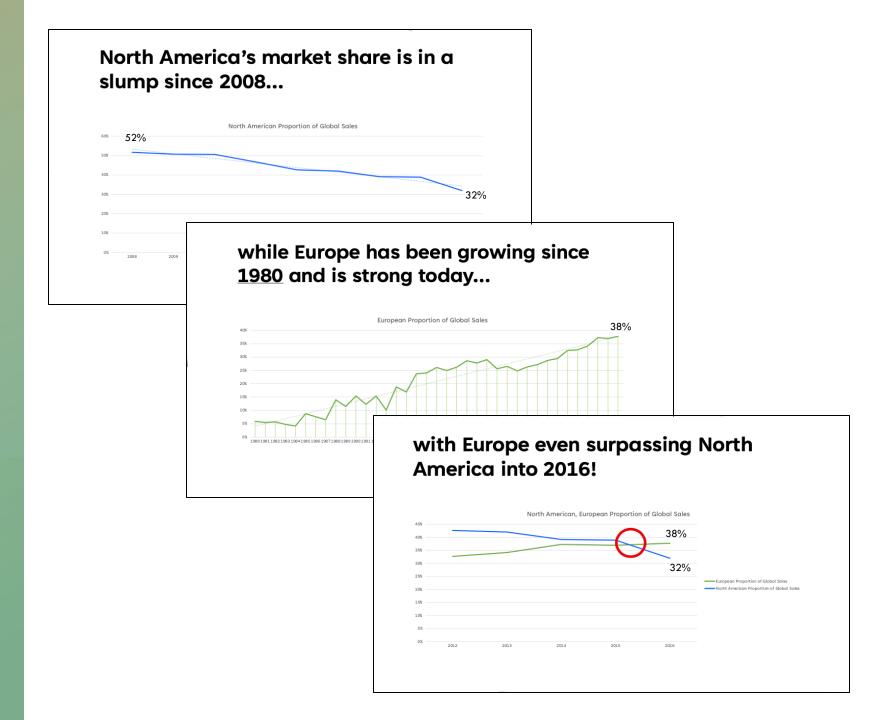VGChartz video game sales data (1980-2016)

## SKILLS USED

Grouping data
Summarizing data
Descriptive analysis
Visualizing results
Presenting results

# GAMECO: SALES SECTION

In market share, Europe is the primary market for GameCo's investments, growing consistently since 1980 and is now the largest (at **38%**). This investment should come at the expense of North America, on a consistent downtrend since 2008 (**52%** then, **32%** now).
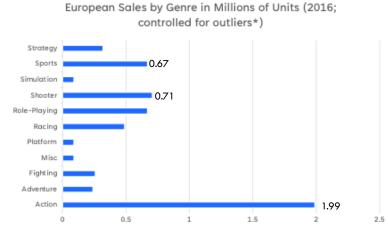
## North America's market share is in a slump since 2008...

North American Proportion of Global Sales

52%

32%

## while Europe has been growing since 1980 and is strong today...

European Proportion of Global Sales

38%

## with Europe even surpassing North America into 2016!

North American, European Proportion of Global Sales

38%

32%

European Proportion of Global Sales
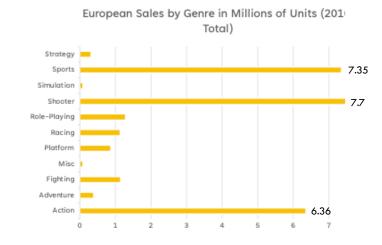North American Proportion of Global Sales

In Europe, as the most important market, publish Action (**1.99m units sold in 2016**), Shooters (**0.71m**) and Sports (**0.67m**) games. These are the strongest categories even when excluding statistical outliers.

GameCo should evaluate its capability to compete with major franchises in these genres

## Action games are solid, but sports and shooters had mega hits in 2016

European Sales by Genre in Millions of Units (2016; controlled for outliers*)

| Genre | Value |
|---|---|
| Strategy | |
| Sports | 0.67 |
| Simulation | |
| Shooter | 0.71 |
| Role-Playing | |
| Racing | |
| Platform | |
| Misc | |
| Fighting | |
| Adventure | |
| Action | 1.99 |

\* Outliers sold more than .125 million units in 2016

European Sales by Genre in Millions of Units (2016 Total)

| Genre | Value |
|---|---|
| Strategy | |
| Sports | 7.35 |
| Simulation | |
| Shooter | 7.7 |
| Role-Playing | |
| Racing | |
| Platform | |
| Misc | |
| Fighting | |
| Adventure | |
| Action | 6.36 |

# PROJECT LINKS

PROJECT BRIEF

FULL PRESENTATION
(PDF)

FULL PRESENTATION
(PPTX)

PROJECT
REFLECTIONS

# INFLUENZA PREPARATION: ALLOCATE MEDICAL STAFF FOR FLU SEASON

# BACKGROUND

## GOAL

Analyze temporal and geographic flu trends to plan for staffing needs across the US and lower the mortality rate

## TOOLS

Excel
Tableau

## DATASET

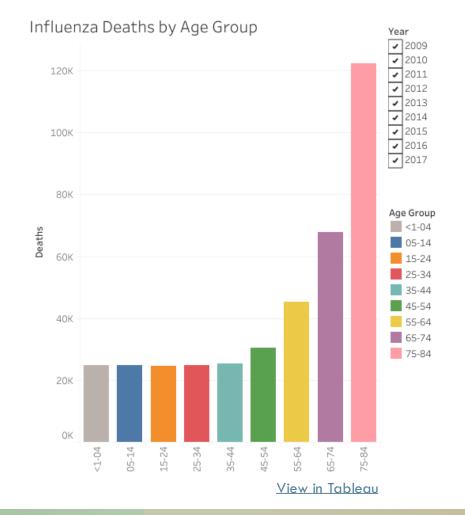US Census Data (2009-2017)

Flu-Related Death Counts (2009-2017)

## SKILLS USED

Translating business requirements
Data cleaning
Data integration
Data transformation
Statistical hypothesis testing
Visual analysis
Forecasting
Storytelling in Tableau
Presenting results to an audience

# FLU SEASON: APPROACH

**Approach:**

Mortality rates from the flu were used to determine where staff would be most needed.

Exploratory analysis revealed that the higher age groups **(especially 65+)** were much more likely to die from influenza than younger groups **(.99 correlation between 65+ age and mortality)**



Influenza Deaths by Age Group

View in Tableau



US 65+ Population by State

US Flu Deaths

Visualized: the 65+ aged population and flu mortality heatmaps are very similar, demonstrating this correlation

View in Tableau

# FLU SEASON: WHEN?

Overall, flu season occurs simultaneously throughout the US, running from September to May.

States stressed by the flu will be stressed simultaneously; staff allocations will be based on most need.

## Influenza Deaths by Month

Month



View in Tableau

# FLU SEASON: WHERE TO SEND STAFF?

Staff ultimately need to be deployed where there is a greater mortality (CA, NY, TX, etc.), which namely tends to be where a larger 65+ and larger general population tends to be.

At the same time some states have a disproportionally large amount of deaths relative to their population (WY, ND, SD, etc.) and would also need a relatively higher amount of staff.



US Flu Deaths by Population

Year
✓ 2009
✓ 2010
✓ 2011
✓ 2012
✓ 2013
✓ 2014
✓ 2015
✓ 2016
✓ 2017

Deaths (State)

5,310          56,803

Relative Deaths to Population
·  0.110%
○  0.400%
○  0.600%
○  0.800%
○  1.010%

© 2023 Mapbox © OpenStreetMap

View in Tableau

# PROJECT LINKS

PROJECT BRIEF

TABLEAU LINK

PRESENTATION (PPTX)

PRESENTATION (PDF)

VIDEO PRESENTATION

# ROCKBUSTER STEALTH: MOVIE ACQUISITION STRATEGY

# BACKGROUND

## GOAL

Utilize existing data to inform the strategic approach of a video rental store franchise in introducing a new online service, with a focus on revenue, customer base, and geography.

## TOOLS

PostgreSQL
Tableau
DbVisualizer

## DATASET

Fictional company Rockbuster's inventory, customer, and payment data

## SKILLS USED

Relational databases
SQL
Database querying
Filtering
Cleaning and summarizing
Joining tables
Subqueries
Common table expressions

# ROCKBUSTER: CUSTOMER COUNTS AND REVENUE

| Country | Customers | Revenue |
|---|---|---|
| **India** | **60** | **$ 2,367.29** |
| **China** | **53** | **$ 2,089.80** |
| United States | 36 | $ 1,344.31 |
| Japan | 31 | $ 1,185.61 |
| Mexico | 30 | $ 1,143.73 |
| Brazil | 28 | $ 1,144.78 |
| Russian Federation | 28 | $ 1,044.89 |
| Philippines | 20 | $ 812.40 |
| Turkey | 15 | $ 581.86 |
| Indonesia | 14 | $ 548.92 |

India and China are the largest countries by far in terms of customers and revenue.

Customers and Revenue by Country



© 2023 Mapbox © OpenStreetMap

Revenue by Region

| Asia-Pacific $9,094.84 | North America $3,000.02 | South America $2,881.28 |
| Europe $4,387.39 | Middle East $2,105.86 | Africa $2,081.83 |

# ROCKBUSTER: REGIONS AND RECOMMENDATIONS

- Focus especially on Asia-Pacific as this has been Rockbuster Stealth's strongest consumer base, but…

- Keep internet accessibility in mind! Some of our customers may be better served by physical locations due to poor internet access

# PROJECT LINKS

PROJECT BRIEF

PRESENTATION
(PPTX)

PRESENTATION
(PDF)

DATA
DICTIONARY

# INSTACART: INITIAL DATA & EXPLORATORY ANALYSIS

# BACKGROUND

## GOAL

Perform an initial data and exploratory analysis on sales information in order to derive insights and suggest strategies for better segmentation

## TOOLS

Python
Jupyter
Pandas
NumPy
Matplotlib
Seaborn

## DATASET

The Instacart Online Grocery Shopping Dataset 2017

## SKILLS USED

Data wrangling
Data merging
Deriving variables
Grouping data
Aggregating data
Reporting in Excel
Population flows

Income Group Prevalence by Profile

# INSTACART: CUSTOMER ORDER INSIGHTS

Instacart's orders are primarily from those w/dependents earning >$60k, **61% of all orders.**

**75%** of Instacart's orders come from those w/dependents.

# INSTACART: PURCHASING BEHAVIOR

Dairy eggs **(5.9m)** and produce **(9.1m)** are by far the two largest categories purchased in orders while the meat/seafood is pretty tiny **(679k)**, which could suggest that these customers are more likely to be vegetarian.

However, the canned goods **(662k)**, snacks **(2.8m)**, and frozen **(2.1m)** departments are also sizable and would best be broken down into further categories in order to make the best judgment there.

# INSTACART: PRODUCT CLASSES



Purchases of Products by Class

Most Instacart orders (67.5%) are for mid-range products ($5-15), a newly derived category.

# PIG E. BANK: ASSESSING RISK

# BACKGROUND

## GOAL

Identify factors to
predict the likelihood
of customers leaving
the bank

## TOOLS

Excel
GitHub

## SKILLS USED

Big data
Data ethics
Data mining
Predictive analysis
Time series analysis
Forecasting
Using GitHub

# PIG E. BANK: LEAVER CHART

## Pig E Bank Leaver Decision Chart

**Active Status?**
- Yes → **Female?**
  - Yes → **Between 39 and 51?**
    - Yes → **Balance > $93,147?**
      - Yes → Most likely to leave
      - No
    - No → **Balance > $93,147?**
      - Yes
      - No
  - No → **Between 39 and 51?**
    - Yes → **Balance > $93,147?**
      - Yes
      - No
    - No → **Balance > $93,147?**
      - Yes
      - No
- No → **Female?**
  - Yes → **Between 39 and 51?**
    - Yes
    - No
  - No → **Between 39 and 51?**
    - Yes
    - No → **Balance > $93,147?**
      - Yes
      - No → Least likely to leave

### Decision Factors

1. **Lacking active status: 70%** of leavers do not have active status

2. B**eing a woman:** about **60%** of leavers are women

3. **Being between age 39 and 51:** about **55%** of leavers are of this age

4. **Have a higher balance:** leavers' median and average balances (**$93,147**) are **20%** higher than those who remain (**$74,830**).

# AIRBNB BERLIN: PRICING FACTORS & EFFECTS ON LOCAL MARKET

# BACKGROUND

## GOAL

Conduct time series, geospatial, and clustering analysis on data from Airbnb to determine price and rating factors, judge effects on local markets, determine effect of recent regulation

## TOOLS

Python
Jupyter
Pandas
NumPy
Matplotlib
Seaborn
SciPy
Tableau
Excel

## DATASET

Data assembled from
Insideairbnb.com

## SKILLS USED

Sourcing open data
Data wrangling & cleaning
Geospatial analysis
Dickey-Fuller Test
Autocorrelations
Regression analysis
Clustering
Time series data sourcing and analysis
Supervised and unsupervised machine learning

# PRICE & RATING CORRELATION

From correlation, the only quantitative factors affecting price were accommodation capacity and related comfort factors (bed counts, bedrooms, bathrooms, etc.), and at weak levels (**0.5** and lower)

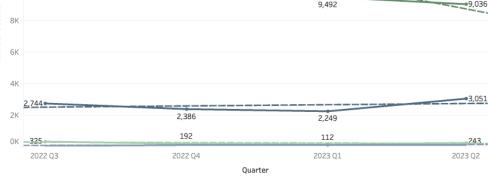Ratings only seemed to be correlated with each other

# CLUSTERING

Through **unsupervised machine learning,** we derived four different clusters primarily based on price and listing counts, helping us determine that:

- **Small-time Locals (79%** of listings) are leaving the platform (**33%** between Sep 2022–Jun 2023), likely driven by a March 2023 licensing requirement

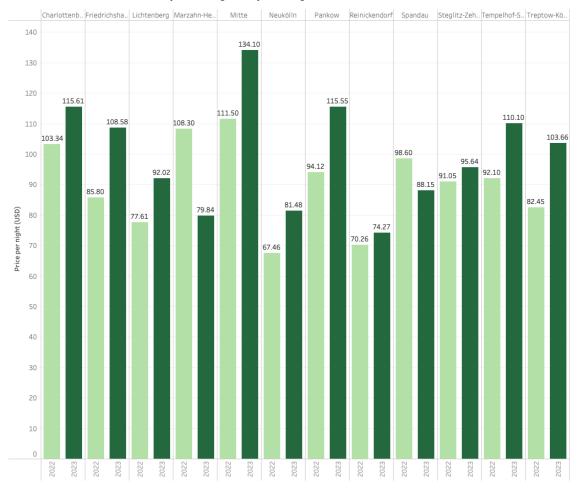- Other, more professionalized clusters tended to stay with the platform

# PRICE EFFECTS

- Airbnb prices stayed flat 2022-2023, but **Small-time Locals** charged the lowest prices (**$63 per night**)

- As **33%** left the platform, the average cost of Airbnb rose **11.6%**, exceeding local rental price increases by **8%**, German national inflation by **9.6%**

- Airbnb presence is known to drive up local rental prices and this will continue, but prices will stabilize once the Small-Time Local exodus is complete
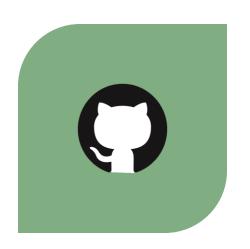


Price per Night by Neighborhood and Year

# PROJECT LINKS

PROJECT BRIEF

TABLEAU LINK

GITHUB

THANK YOU