# Project: Detecting fake reviews

- Team members
    - Georgios (Economics and Business)

    - Pedro (Economics and Business)

    - Marius (User Experience Design)

JADS
Jheronimus
Academy
of Data Science

# Objective

Create a ML- pipeline to Detect fraudulent reviews

Given these predictions, they can conduct investigations, decide whether to notify users (i.e flag fake reviews) or take corrective action.

**Main Beneficiary**: e-commerce platforms

**Indirect beneficiaries**: customers, sellers and the regulatory agencies

# Project motivation

- E-commerce on the rise

- Reviews closely determine purchasing decisions

- Dynamic environment: easier than ever to create fake content that seems believable

- Fraudulent reviews damage both the platform, the seller reputation and consumers

# Real world case

UK competition watchdog to probe Google and Amazon over fake reviews  Financial Times

- Competition and Markets Authority says tech groups may not be doing enough to protect consumers.

- The UK competition regulator has opened an investigation into Amazon and Google over fake reviews on their sites that may be duping consumers.

- A thriving industry where potentially hundreds of thousands of reviews are bought and sold for as little as £5 each".

  "Our worry is that millions of online shoppers could be misled by reading fake reviews and then spending their money based on those recommendations,"

https://www.ft.com/content/b7c4b9fc-116e-4681-92cb-f433f2a09aa6

# This is a real business problem

Fake online reviews cost the global economy $152 billion a year. (WEF)

In response, there exist off-the-shelf solutions:

Fraud detection algorithms that combine behavioural analytics and text analysis e.g. Amazon Web Services (AWS):

# Some Statistics (about the e-commerce market)

- Revenue in the eCommerce Market is projected to reach US$3,226.00bn in 2024.

- Revenue is expected to show an annual growth rate (CAGR 2024-2029) of 9.79%, resulting in a projected market volume of US$5,145.00bn by 2029.

- In the eCommerce Market, the number of users is expected to amount to 3.2bn users by 2029.

- Data from STATISTA (world wide in us dolars)

# Some Statistics (about the role of reviews in the buying process)

- 95% of costumers read reviews before making the buying decision (Global Newswire)
- 88% of customers who read an online review say it influenced their buying decision (Zendesk )
- 49% of consumers trust online reviews as much as personal recommendation (Bright Local)
- Positive reviews can increase customer spending by 31%  (Bright Local)
- 86% of people hesitate to do business with a company if it has too many negative customer reviews.
- If consumers found out a platform was censoring reviews, 62% of consumers would stop using it (Trustpilot)
- 52% of a company's market value is attributed to its reputation (PR Week)

# Similar open source projects:

- Mostly focus on NLP ( analysis on text )

| Sr. No. | Model Accuracy (%) | Precision Score | Recall Score | F1 Score |
|---------|--------------------|-----------------|--------------|----------|
| 1 | MultinomialNB | 90.25 | 0.9325 | 0.8601 |
| 2 | Stochastic Gradient Descent (SGD) | 87.75 | 0.8913 | 0.8497 |
| 3 | Logistic Regression | 87.00 | 0.8691 | 0.8601 |
| 4 | Support Vector Machine | 56.25 | 0.525 | 0.9792 |
| 5 | Gaussian Naive Bayes | 63.5 | 0.6424 | 0.6169 |
| 6 | K-Nearest Neighbour | 57.5 | 0.8604 | 0.1840 |
| 7 | Decision tree | 68.5 | 0.6681 | 0.7412 |

Credits:  Salunkhe, Ashish. "Attention-based Bidirectional LSTM for Deceptive Opinion Spam Classification." arXiv preprint arXiv:2112.14789 (2021).

# Possible analyses

1. Analyze the reviewer, not the review
   - Feature engineering: *how old the account is*, *# of reviews made*, *time stamps for behavior*, etc

2. Use review metadata: helpfulness ranking, verified purchase or not.

3. Analyse the review text
   - Natural Language Processing (NLP)
   - Complex and beyond the scope of this course

# Dataset (correction: huge database)

**Yelp Dataset on Kaggle.**



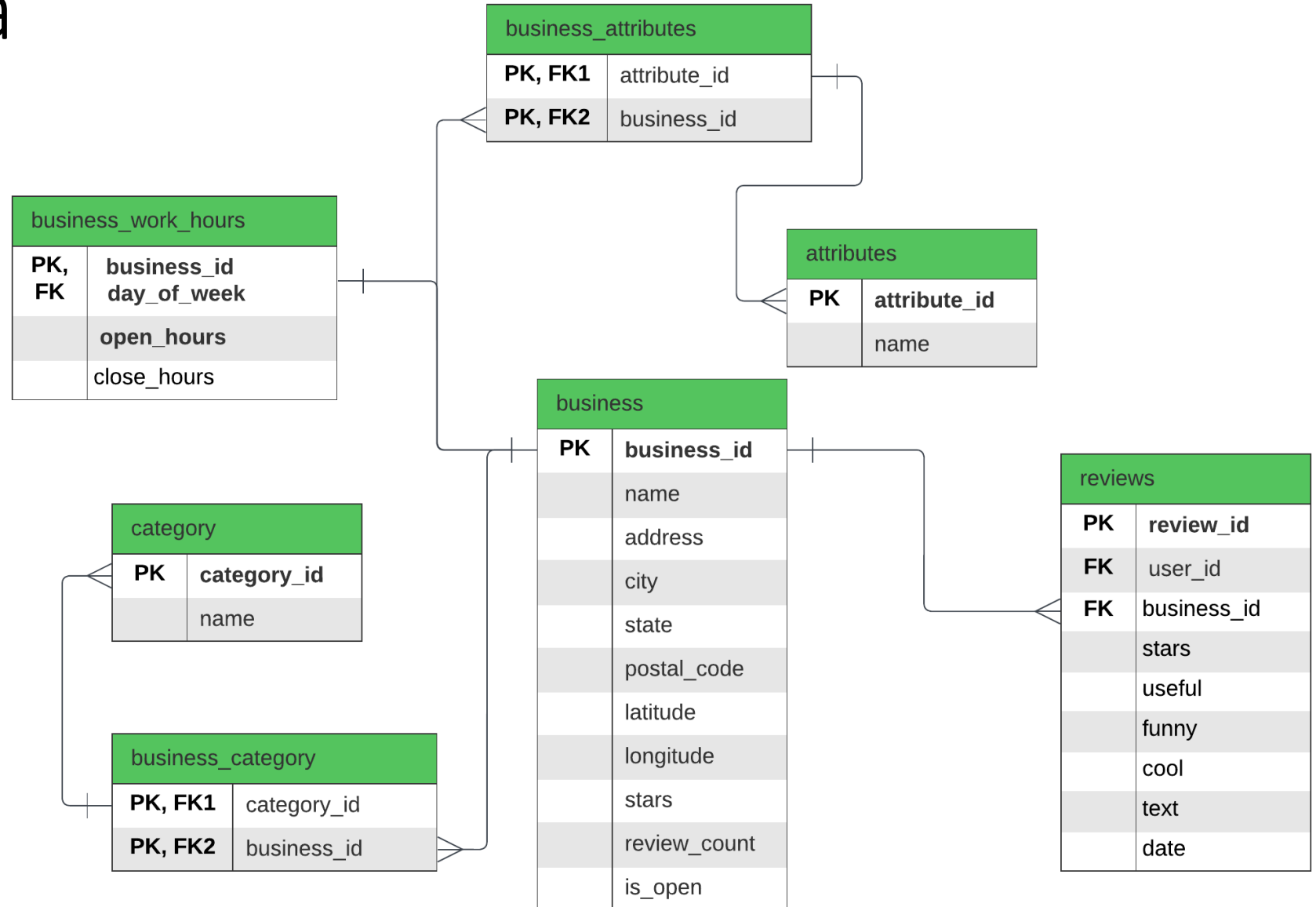kaggle.com/datasets/yelp-dataset/yelp-dataset

Clean

Popular for analysis

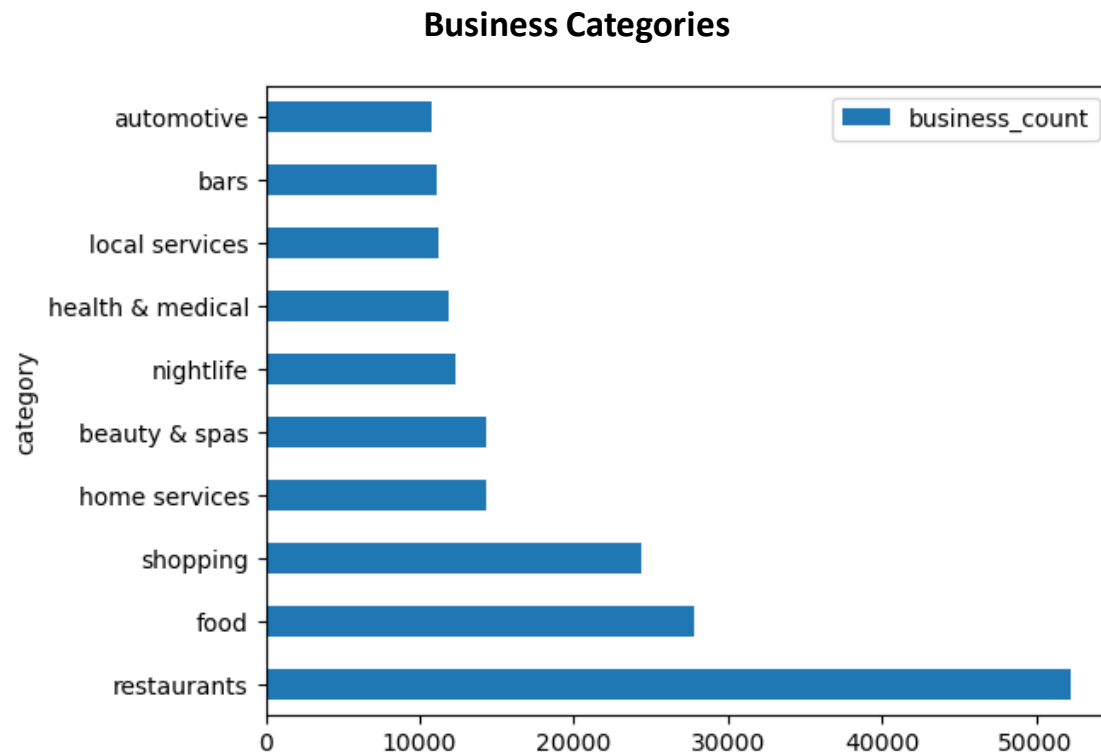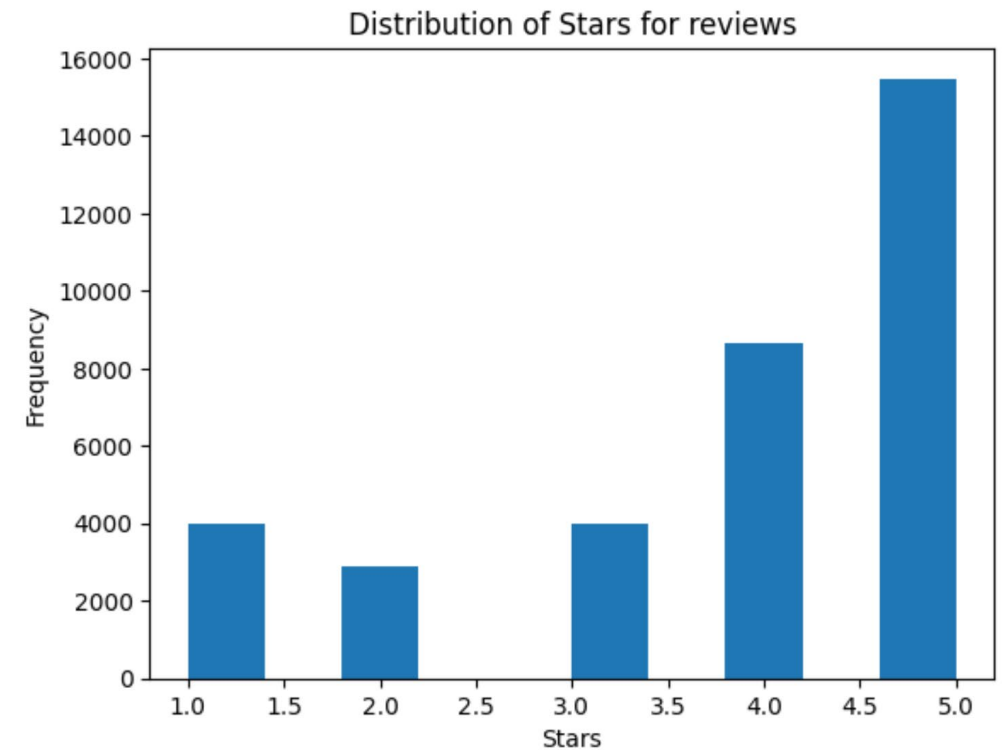JSON (not CSV)

Large: 3GB + 5GB databases

# Data schema

**Covers:**
- Businesses,
- Reviews
- Users

**business_attributes**

| PK, FK1 | attribute_id |
|---------|--------------|
| PK, FK2 | business_id |

**business_work_hours**

| PK, FK | business_id day_of_week |
|--------|-------------------------|
|        | open_hours |
|        | close_hours |

**attributes**

| PK | attribute_id |
|----|--------------|
|    | name |

**business**

| PK | business_id |
|----|-------------|
|    | name |
|    | address |
|    | city |
|    | state |
|    | postal_code |
|    | latitude |
|    | longitude |
|    | stars |
|    | review_count |
|    | is_open |

**category**

| PK | category_id |
|----|-------------|
|    | name |

**business_category**

| PK, FK1 | category_id |
|---------|-------------|
| PK, FK2 | business_id |

**reviews**

| PK | review_id |
|----|-----------|
| FK | user_id |
| FK | business_id |
|    | stars |
|    | useful |
|    | funny |
|    | cool |
|    | text |
|    | date |

JADS

17

# Exploratory data analysis (EDA). Part 1/3
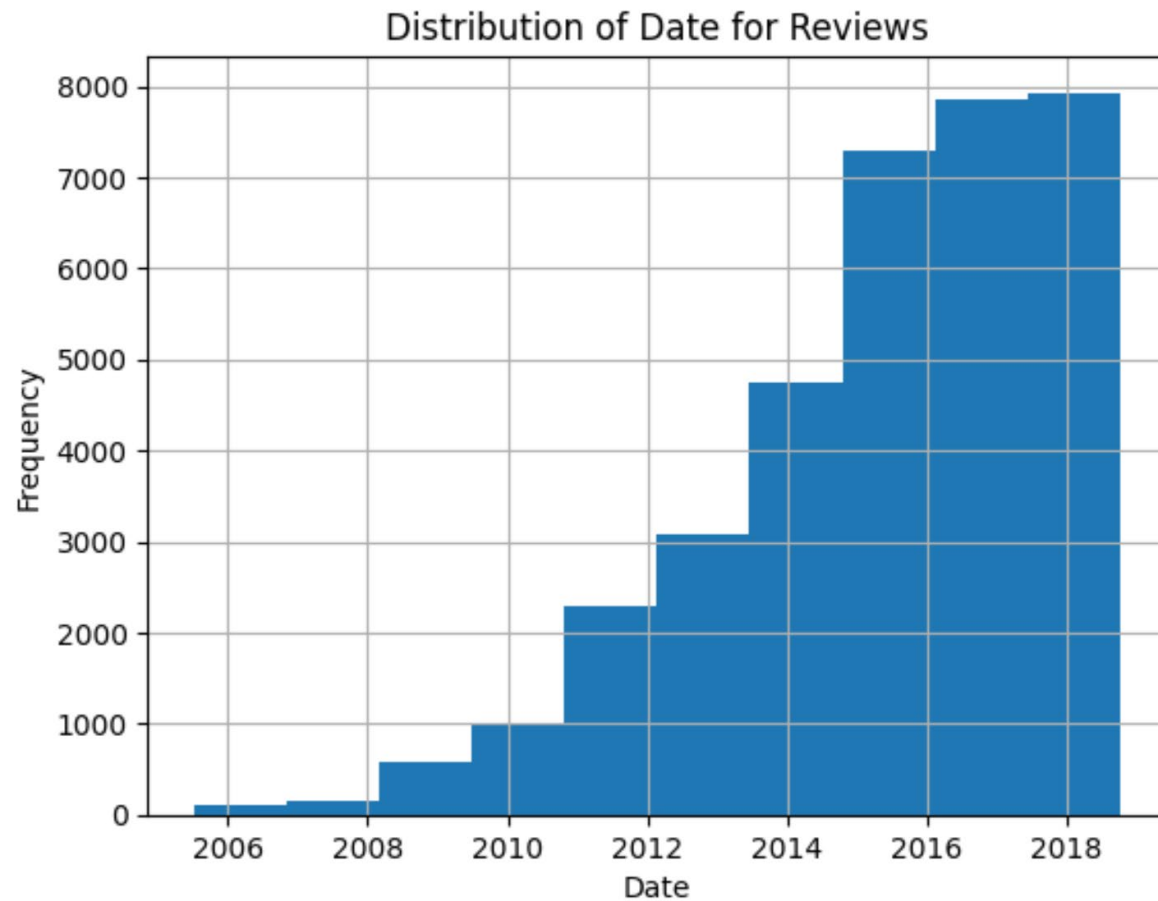
**Business Categories**



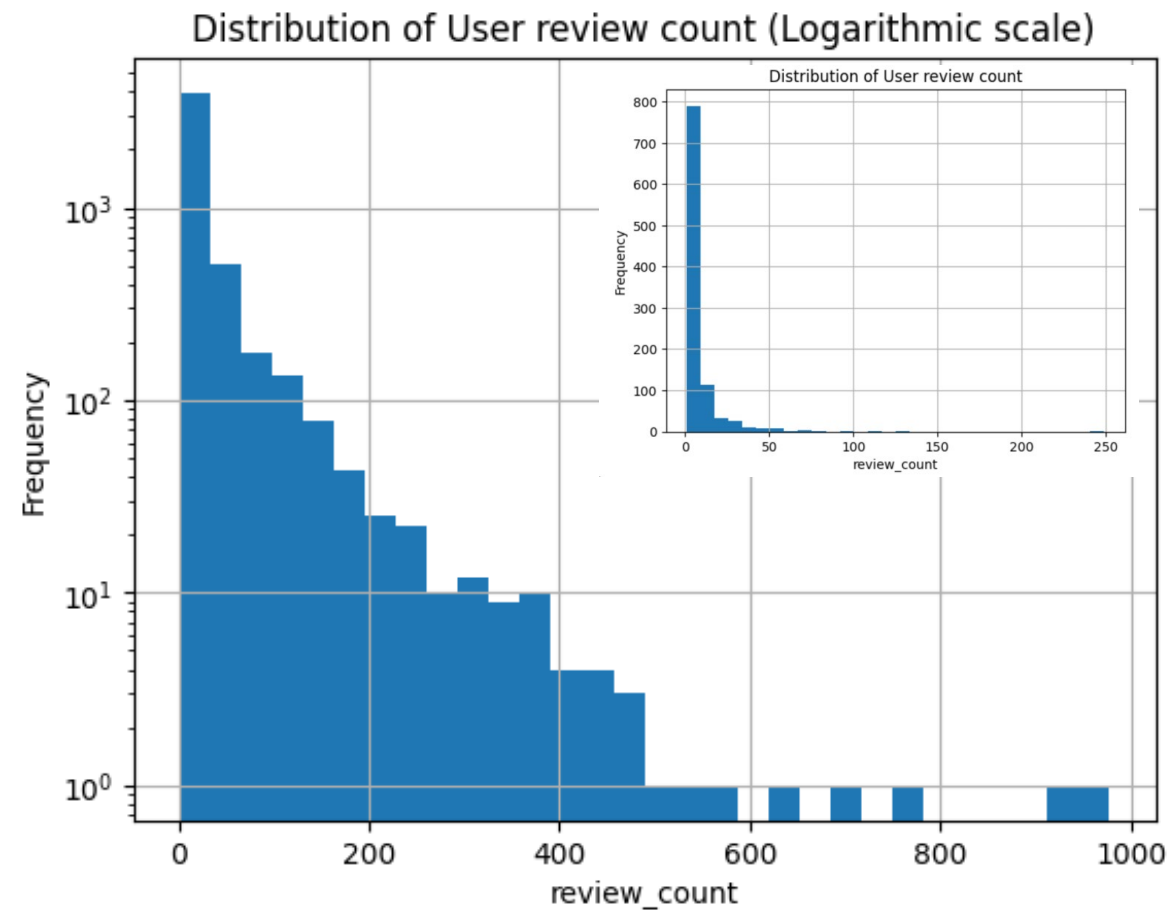Credits: An Hoang Vo on Kaggle, Business Analysis using SQL



Analysis of a fragment of 35k random (?) entries of 'reviews'
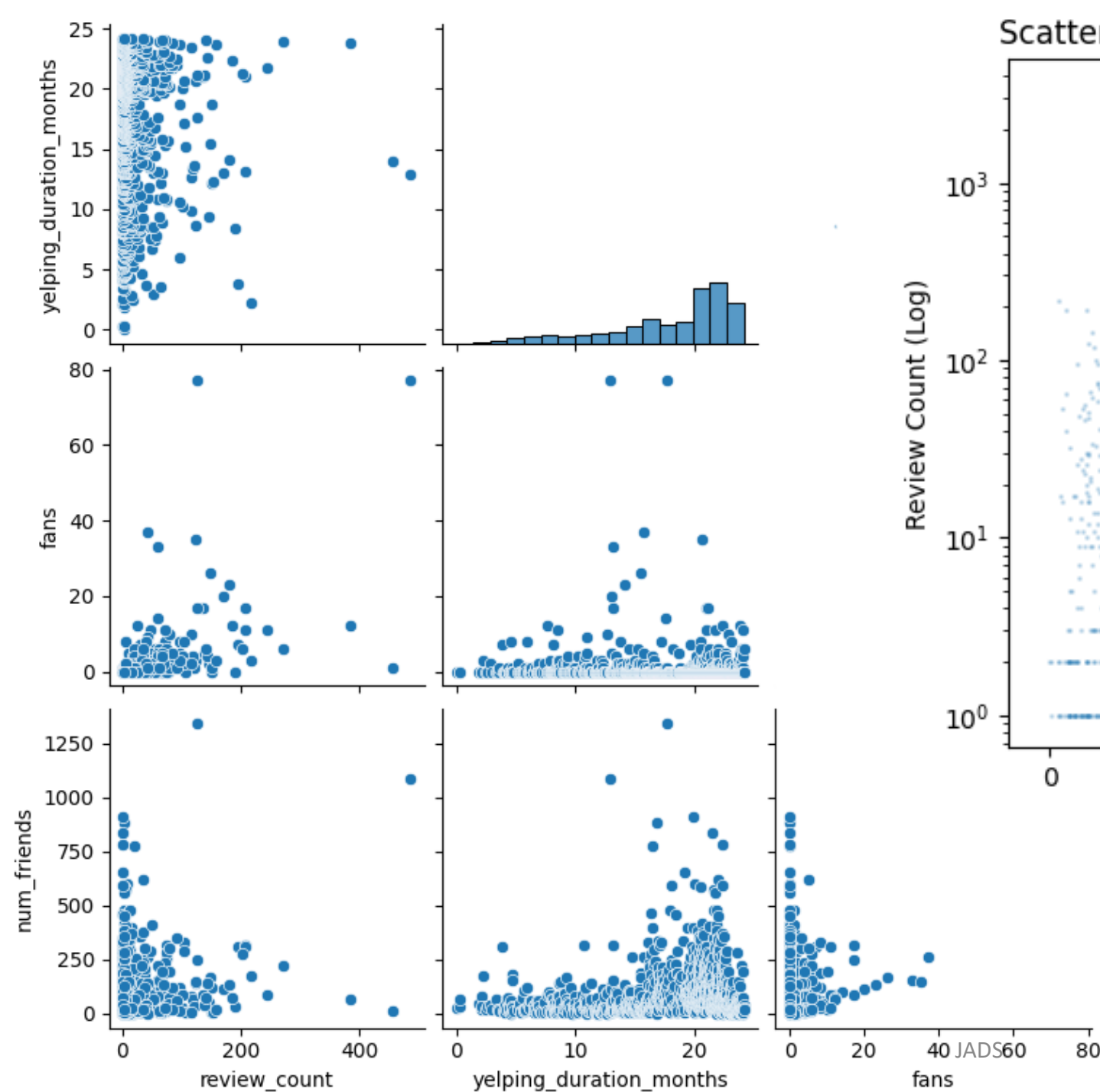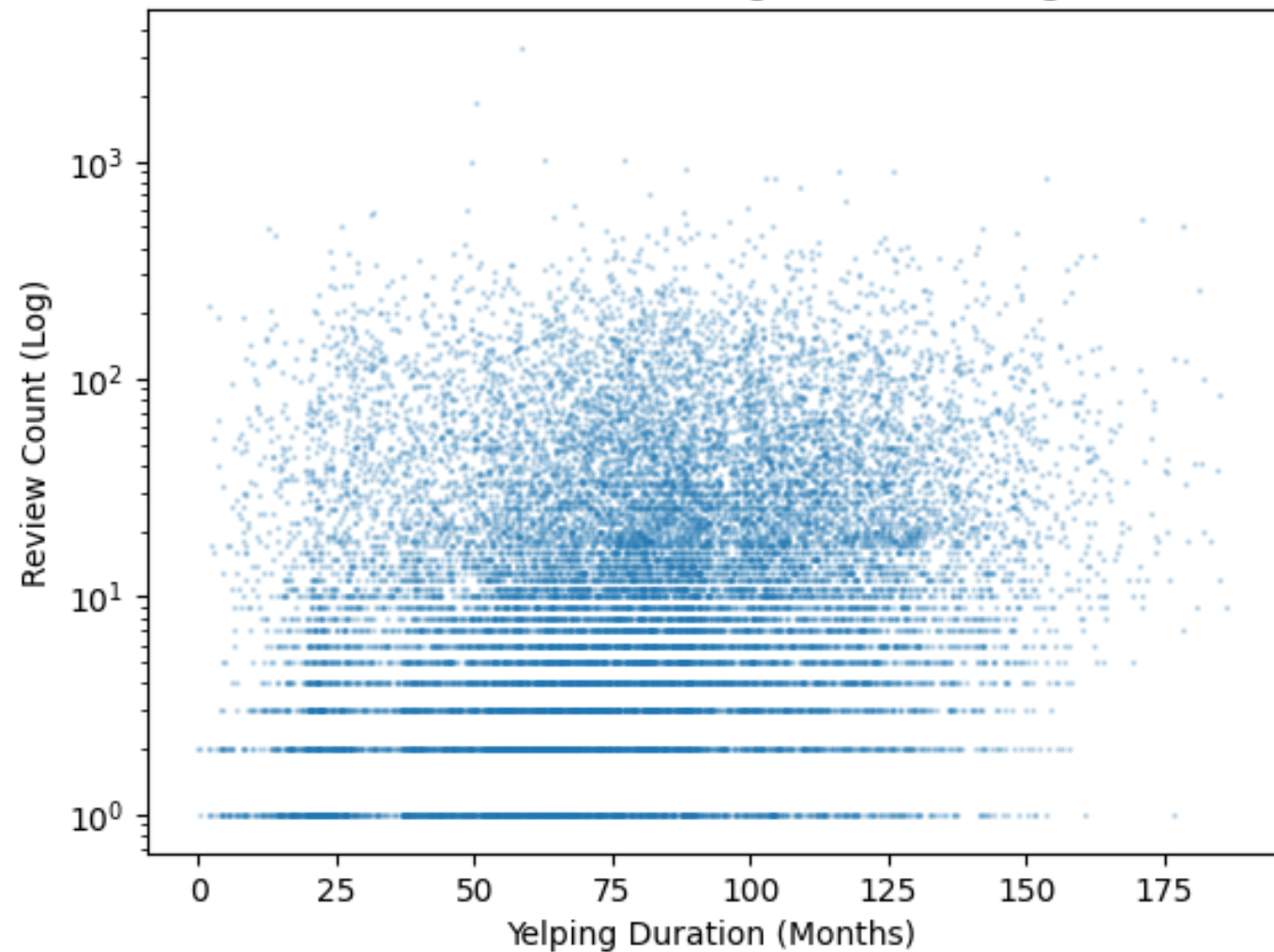
# EDA 2/3



Our analysis of a fragment of 35k random (?) entries of 'Reviews'

Our analysis of a fragment of 1000 random entries of 'Users'

Scatter Plot of Review Count (Log) vs Account age in months

Both analyses performed on a random (?) sample of 20k Users

EDA 3/3

# Questions we need help with:

1. Dealing with a large dataset: in practice vs for this course.
Is 'split' acceptable ?
How can we see how randomly distributed it is now ?

2. Should we continue working with a large dataset that requires significant JOINs between different datasets ?

3. How can we evaluate the model if real world data does not contain a label: fake/real. Prediction tested by what ?

# Next actions

- Closer look at how others evaluate their result given data not labelled on authenticity

- Replicating another project to see what we learn

- Comparing 3 key algorithms and decide on one

OR

- coaching session w/ teacher + Plan B for a more manageable objective.

# Follow our progress

github.com/ursumarius/review-analysis-intro-ml-jads/

kaggle.com/code/mariusursu/review-analysis-intro-ml-jads

Slides:

- Thanks for your time.
  - Marius
  - Pedro
  - Georgios

# Project: Detecting fake reviews

- Team members
  - Georgios (Economics and Business)

  - Pedro (Economics and Business)

  - Marius (User Experience Design)

JADS
Jheronimus
Academy
of Data Science