

HAHA GROUP - 22KHDL1 - INTRO2DS

PRESENTED BY  
**CHAU ANH DANG**  
**KIM ANH NGUYEN**  
**DIU HUYEN TRAN**  
**MINH HUY DO**





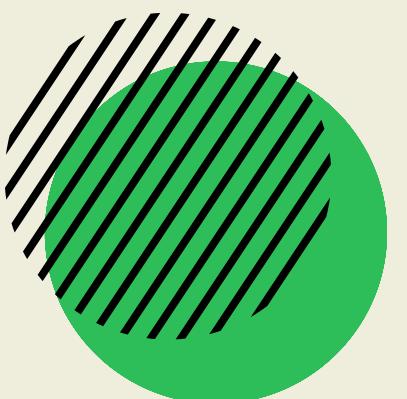
# TABLE OF CONTENTS



PROJECT MILESTONES

CHALLENGES & SOLUTIONS

NEXT STEPS



# TABLE OF CONTENTS

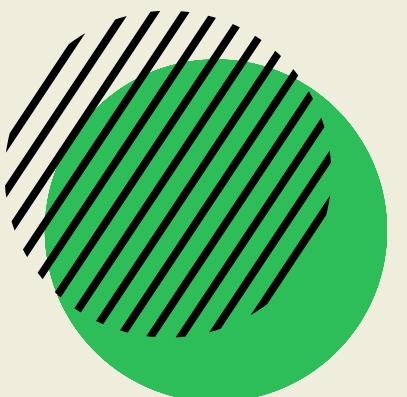


PROJECT MILESTONES



CHALLENGES & SOLUTIONS

NEXT STEPS





# PROJECT MILESTONES

OF OUR PROGRESS

---

COMPLETED TASKS

ON-GOING TASKS

# COMPLETED TASKS

66 %



## DATA COLLECTION

**Determined** project purposes, objectives and data sources.

**Scraped** Spotify's weekly top 200 music charts for Vietnam (Jan–Oct 2024).

**Utilized** frameworks: *Selenium*, *BeautifulSoup*, *pandas*, and others.

**Saved** and **combined** data into unified CSV files.

**Verified** dataset accuracy to meet project requirements.

# COMPLETED TASKS



## DATA EXPLORATION

**Reviewed** dataset structure, row/column counts, and data types.

**Checked** for missing or duplicated values.

**Generated** summary statistics (mean, min, max).



## DATA PREPARATION

### Data Cleaning:

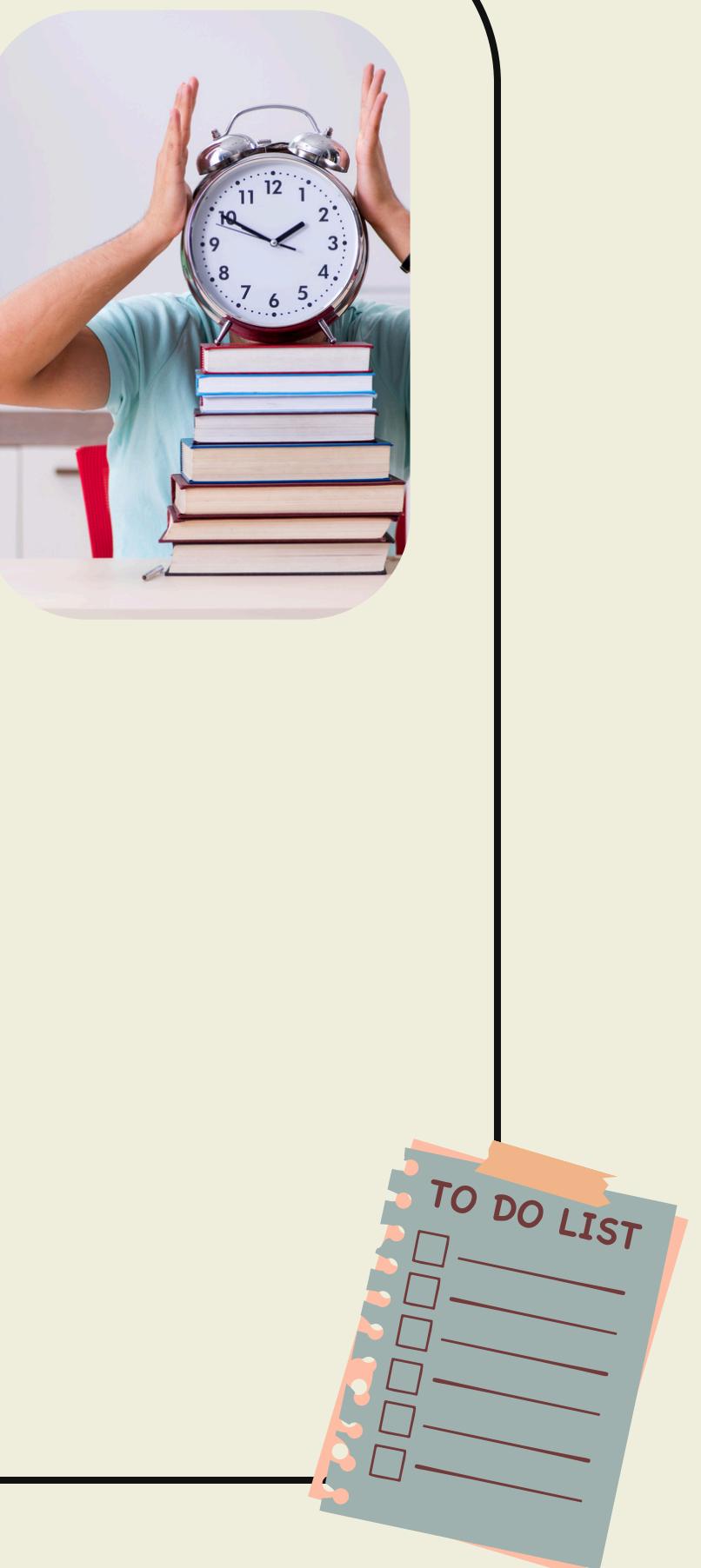
- Remove duplicates.
- Handle missing values.

### Feature Engineering:

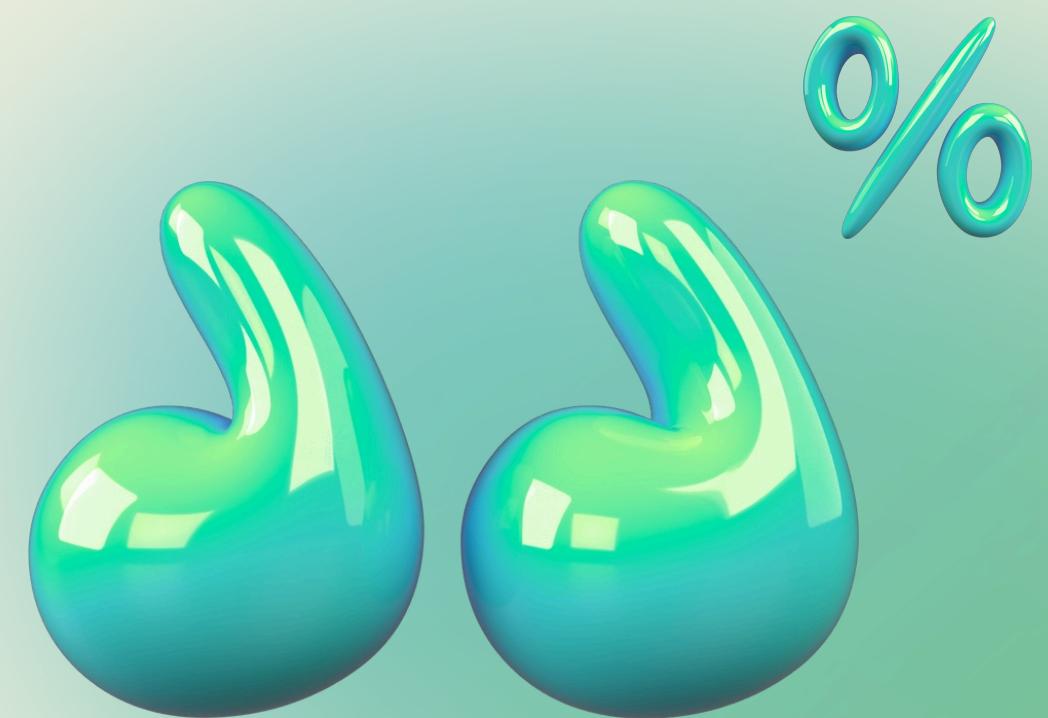
- Create new features such as:
  - Monthly streaming trends.
  - Streams+.
  - Number of peaks in each month.

### Data Transformation:

- Encode categorical data.
- Normalize numerical data.
- Integrate various datasets.



## ON-GOING TASKS



# VISUALIZATION

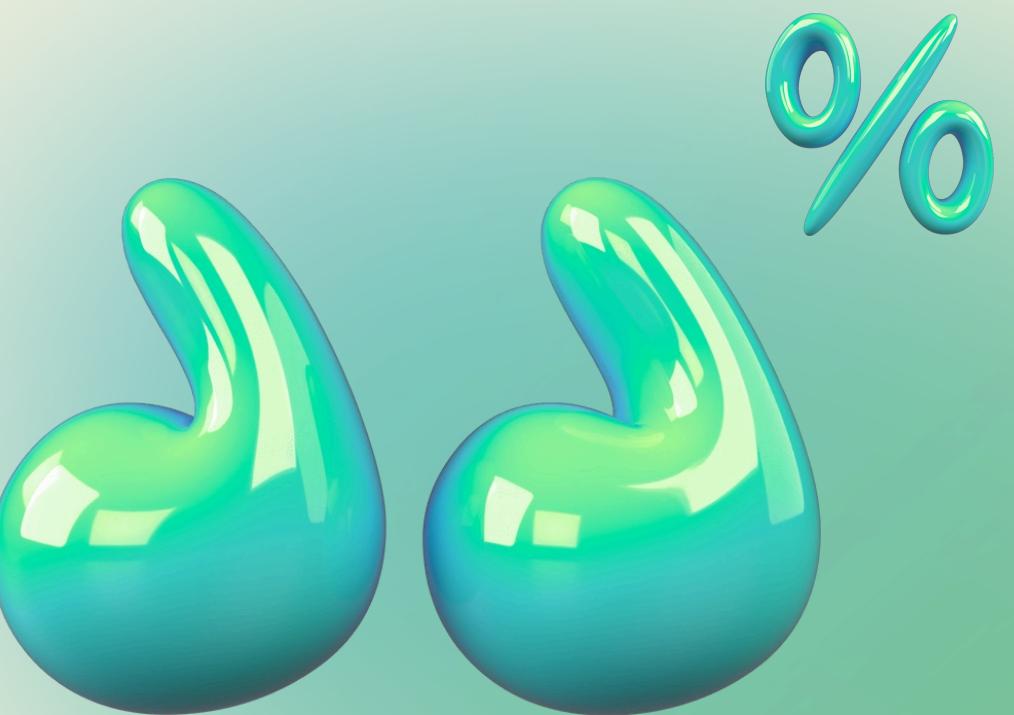
# Calculating and visualizing correlations using NumPy and Matplotlib.

# Analyzing relationships between different data attributes.

# Extracting insights from visual plots.



# ON-GOING TASKS



# CHALLENGES & SOLUTIONS



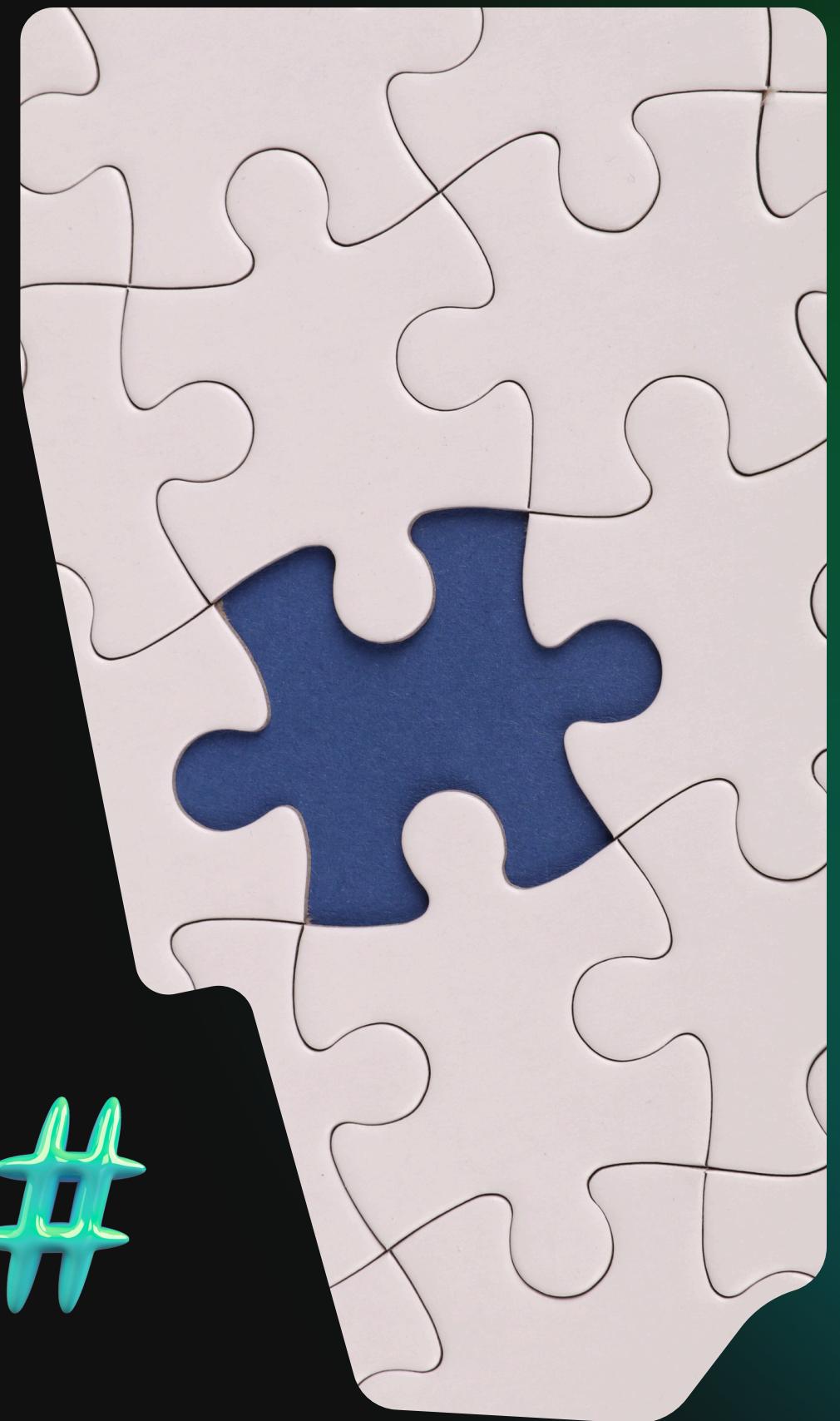
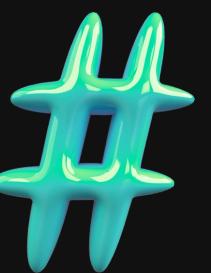
## SCRAPING DATA

Scraping from a dynamic website  
using Selenium

Automating the data downloads

## EXPLORING DATA

Overview of dataset  
Exploratory Data Analysis (EDA)  
Quality check



01

### Dynamic Website Scraping

Spotify charts used dynamic loading and rendering.

Standard HTML parsing was insufficient.

02

### Automating Data Downloads

Weekly chart data had unique URLs based on the date.

Required programmatic construction of URLs.

# SCRAPING DATA (CHALLENGES)



01

### **Dynamic Website Scraping**

Used *Selenium WebDriver* to automate:

- Browser interactions.
- Navigation to the charts page.
- Handling login and cookie pop-ups.

02

### **Automating Data Downloads**

Built the *calculate\_week\_details()* function to:

- Compute start and end dates for each month.
- Construct URLs programmatically.



# **SCRAPING DATA (SOLUTIONS)**

# EXPLORING DATA



- Missing samples compared to expectations.
- Unclear column names or attributes.
- Missing or erroneous values (e.g., negative values for age).
- Duplicated records.
- Categorical attributes with multiple inconsistent values.

## OVERVIEW (CHALLENGES)

# EXPLORING DATA



01

## Dataset Validation

- Use `shape` to check dataset size.
- Use `info()` to inspect column names, missing values, and data types.
- Handle missing values:
  - Mean/median for numerical data.
  - Mode for categorical data.

02

## Duplicate Handling

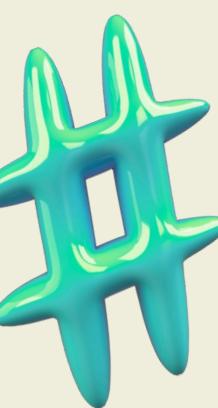
- Use `duplicated()` with `sum()` to identify duplicates.
- Retain only unique records.

03

## Categorical Data

- Use `nunique()` to find unique values.
- Apply encoding methods (one-hot for nominal, label encoding for ordinal).

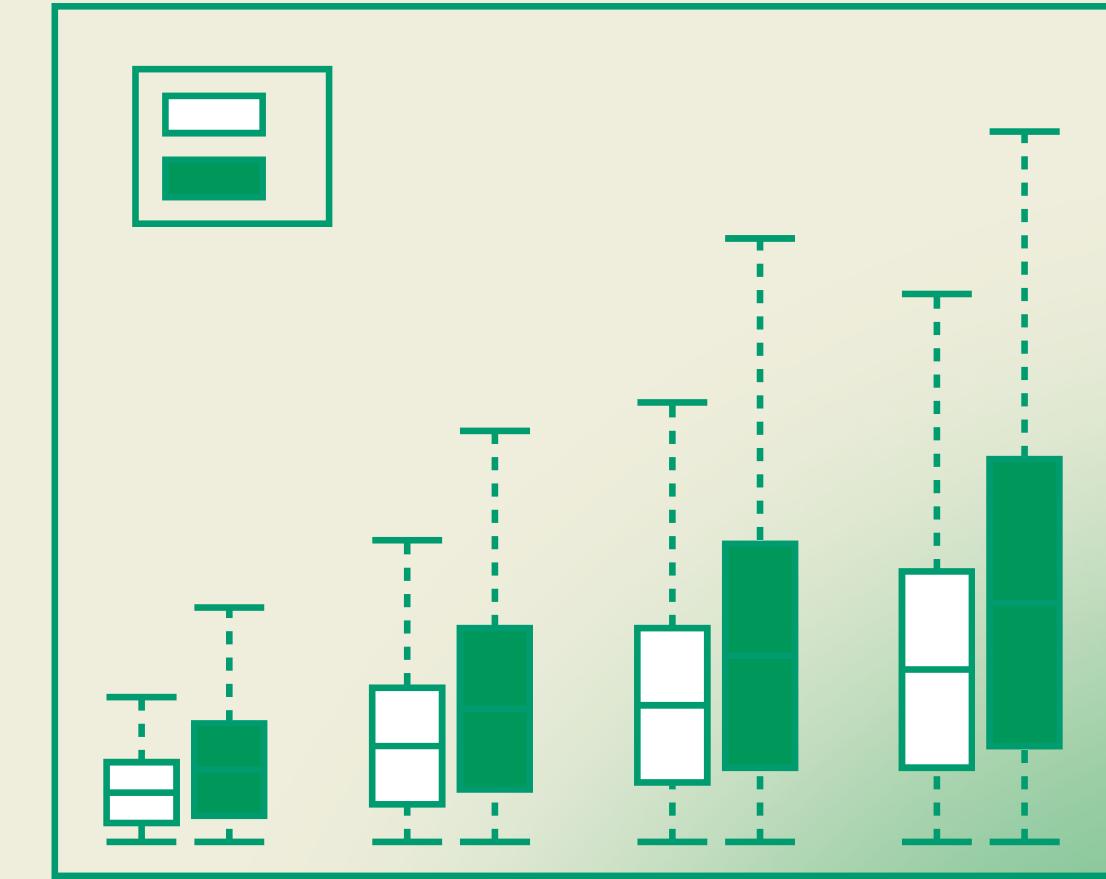
**OVERVIEW  
(SOLUTIONS)**



# EXPLORATORY DATA ANALYSIS (EDA)

Outliers in the data

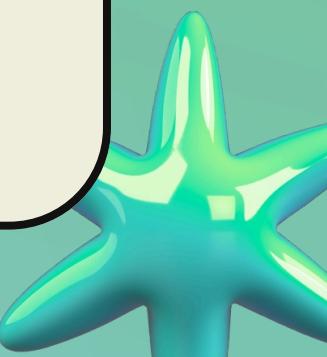
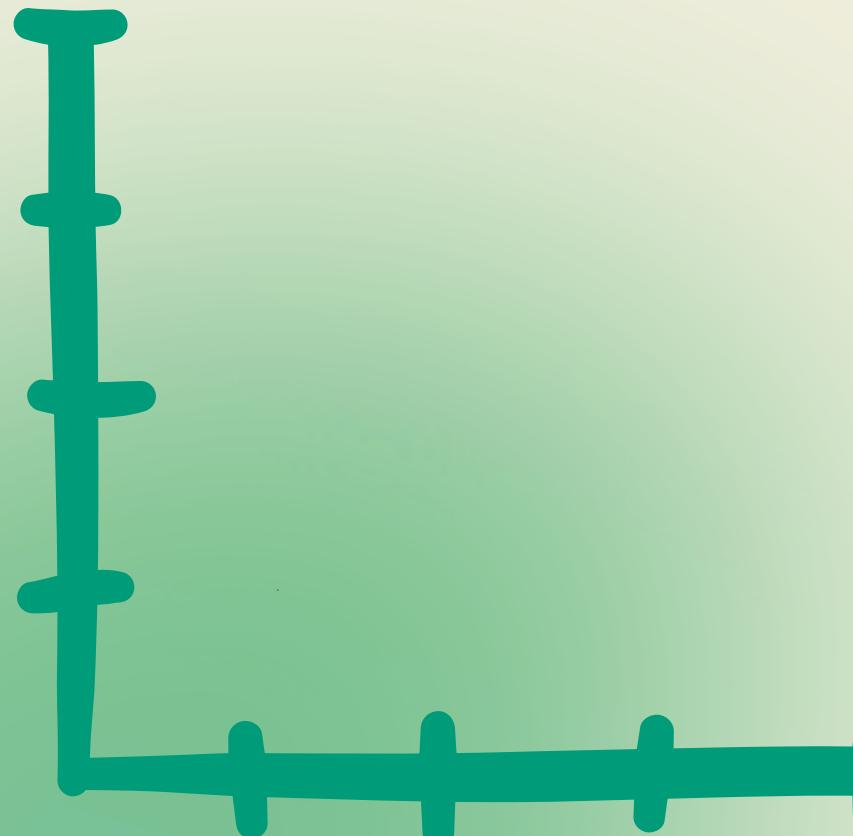
CHALLENGES



SOLUTIONS

Use box plots (IQR) to identify outliers.

- Mitigate outliers with:
  - Z-scores.
  - IQR methods.
  - Log transformations to reduce their effect.





## CHALLENGES

Out-of-range or incorrect values.

## SOLUTIONS

### Rank Validation:

- Ensure rank starts at 1, increments sequentially, and is positive.

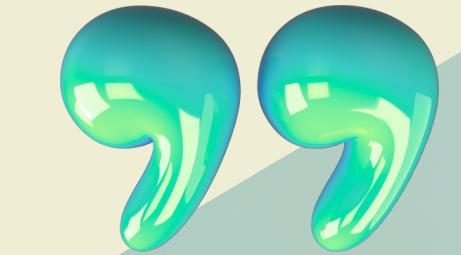
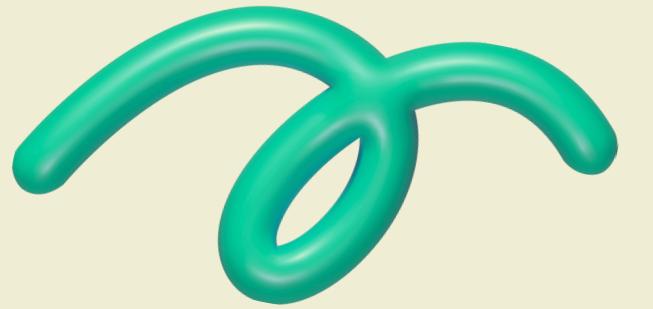
### Stream Validation:

- Check that stream values are non-negative.

### URI Format Validation:

- Verify correct formatting of URI values.

# QUALITY CHECK



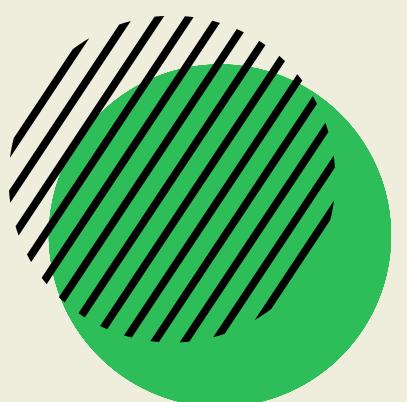
# TABLE OF CONTENTS



PROJECT MILESTONES

CHALLENGES & SOLUTIONS

NEXT STEPS

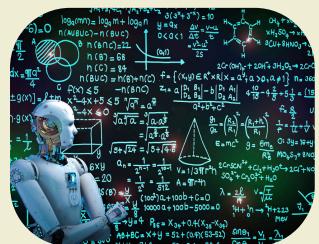
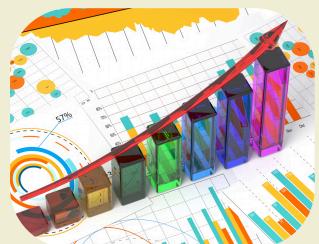


# NEXT STEPS



# DETAILED TASKS FOR NEXT STEPS

PHASE	TASKS	DELIVERABLES
Visualization	<ul style="list-style-type: none"><li>Calculate correlation matrix to uncover relationships.</li><li>Plot data and identify actionable insights.</li><li>Formulate key questions based on findings.</li></ul>	<ul style="list-style-type: none"><li>Correlation matrix</li><li>Insights and relevant questions</li></ul>
Modeling	<ul style="list-style-type: none"><li>Model Selection: Time-series for trends, clustering for artist segmentation.</li><li>Build Models: Analyze genre trends, artist rankings, and peak release periods.</li><li>Model Assessment: Validate and interpret results.</li></ul>	<ul style="list-style-type: none"><li>Model performance report</li><li>Preliminary insights</li></ul>
Evaluation	<ul style="list-style-type: none"><li>Review results to ensure alignment with project objectives.</li><li>Refine models for improved accuracy.</li></ul>	<ul style="list-style-type: none"><li>Evaluation report</li><li>Key insights summary</li></ul>
Deployment	<ul style="list-style-type: none"><li>Summarize findings and recommendations in the final report.</li><li>Present insights through dashboards and presentations.</li></ul>	<ul style="list-style-type: none"><li>Final report</li><li>Stakeholder presentation</li></ul>





**THANK**

**YOU**

FOR YOUR FOLLOWING

A graphic design featuring the word "THANK" in large, bold, black capital letters inside a black-outlined rounded rectangle. A small, shiny green star is positioned at the top right corner of the rectangle. Below "THANK", the word "YOU" is written in large, stylized, glossy green letters that have a thick, rounded, and somewhat melted appearance. At the bottom, the text "FOR YOUR FOLLOWING" is written in a smaller, standard black font. The entire graphic is set against a light beige background.