



Spotify & YouTube Music

Introduction to Spotify and YouTube

Spotify and YouTube are top digital media platforms with distinct roles: Spotify is a major music streaming service known for personalized playlists, while YouTube is the leading video-sharing site offering a wide range of videos including music and user content. This document details the data cleaning process done in Power BI, covering how we handled missing values, fixed data issues, adjusted data types, and ensured overall data consistency.

Data Cleaning

1. Reorder and Rename Columns for Clarity

Objective: Improve dataset readability and usability by reorganizing and renaming columns.

Questions to Address:

1. How should the columns be reordered for better clarity?

- **Answer:** Columns should be grouped logically, with Spotify-related columns listed first and YouTube-related columns listed afterward.

2. What is the approach for renaming columns?

- **Answer:** Column names should be standardized to capitalize each word for consistency. For instance, "spotify_info" should be renamed to "Spotify Info."

3. What are the key columns to be retained and organized?

- **Answer:** Key columns for Spotify include Track, Album, Album Type, Key, and others. For YouTube, key columns include Channel, Views, Likes, Comments, and others.

Summary: Columns have been reordered to group related data and renamed for consistency, enhancing clarity and usability.

2. Identify and Handle Missing Values

Objective: Address missing or null values to maintain data integrity.

Questions to Address:

1. Which columns have missing values?

- **Answer:** Columns with missing values include Licensed, Official Video, Views, YouTube Info, Comments, Description, Likes, and Stream.

2. What is the approach to handle missing values in quantitative metrics?

- **Answer:** For columns like Views, Likes, and Stream, rows with missing values are removed to ensure accurate performance analysis.

3. How should missing values in textual data be managed?

- **Answer:** For columns such as Description and Comments, missing or blank entries are retained as they do not significantly affect the analysis.

Summary: Missing values in quantitative columns were handled by removing affected rows, while textual data was retained as is.

3. Fix Irregularities in Merged Columns

Objective: Separate merged columns into their original components for accuracy.

Questions to Address:

1. What are the components of the Spotify_Info and YouTube_Info columns?

- **Answer:** Spotify_Info contains Spotify Link and Spotify Track ID. YouTube_Info contains YouTube Link and YouTube Video Title.

2. What method should be used to separate these components?

- **Answer:** Use delimiters (|)pipe symbol for Spotify_Info and character length for YouTube_Info to split the data accurately.

3. How should the separated components be validated?

- **Answer:** Verify that links are functional and IDs are accurate for Spotify, and ensure video titles are correctly extracted for YouTube.

Summary: Separated Spotify_Info into Spotify Link and Track ID, and YouTube_Info into YouTube Link and Video Title, ensuring accuracy and cleanliness.

4. Correct Case Sensitivity and Naming Conventions

Objective: Standardize column names and data entries for consistent formatting.

Questions to Address:

1. How should column names be standardized?

- **Answer:** Column names should be capitalized in title case (e.g., "Spotify Info" instead of "spotify_info").

2. What approach should be taken for data entry formatting?

- **Answer:** Text entries such as artist names and track titles should be capitalized consistently.

3. What are the benefits of standardizing formatting?

- **Answer:** Enhances clarity, consistency, and readability of the dataset, making it easier to work with.

Summary: Column names and text entries have been standardized to ensure consistency and improve dataset usability.

5. Remove or Handle Irrelevant Columns

Objective: Eliminate columns that do not contribute to the analysis.

Questions to Address:

1. Which columns are considered irrelevant?

- **Answer:** Random Column 1 and Random Column 2.

2. What actions were taken regarding irrelevant columns?

- **Answer:** Removed the irrelevant columns to streamline the dataset.

3. How was random data handled in relevant columns?

- **Answer:** Verified and cleaned any random data in relevant columns to maintain data quality.

Summary: Removed irrelevant columns and ensured no random data remained in the relevant columns.

6. Handle Inconsistent Data Types

Objective: Convert columns with inconsistent data types to their correct numeric format.

Questions to Address:

1. Which columns are affected by incorrect data types?

- **Answer:** Views, Danceability, Energy.

2. How should the columns be converted to the correct numeric format?

- **Answer:** Use Power BI's "Change Data Type" feature to convert columns to decimal format.

3. What issues might arise during conversion?

- **Answer:** Potential non-numeric values or anomalies; use Power BI tools to filter and correct these issues.

4. How can we ensure correct formatting after conversion?

- **Answer:** Verify with a sample of rows to ensure formatting accuracy and validate no text values remain.

Summary: Converted Views, Danceability, and Energy columns to decimal format, correcting any non-numeric values.

7. Address Invalid Data and Ensure Correct Labeling

Objective: Handle invalid entries in the Views column and ensure correct labeling in the Album column.

Questions to Address:

1. What issues are present in the Views column?

- **Answer:** Contains invalid entries like "invalid_data."

2. How should invalid data in the Views column be managed?

- **Answer:** Replace "invalid_data" with null and then convert the column to numeric format.

3. What is the current state of the Album column, and how should it be corrected?

- **Answer:** Verify no numeric or irrelevant entries exist and clean the data to ensure proper labeling.

Summary: Managed invalid data in Views by replacing with null and ensured correct labeling in the Album column.

8. Check for and Remove Duplicate Rows

Objective: Identify and remove duplicate rows to maintain data uniqueness and accuracy.

Questions to Address:

1. How were duplicate rows identified?

- **Answer:** Initial sample check showed no duplicates; full dataset review revealed duplicates.

2. What steps were taken to remove duplicate rows?

- **Answer:** Used the Index Column to identify and remove duplicates.

3. How was data accuracy ensured after removing duplicates?

- **Answer:** Final verification ensured that data remained unique and no critical data was removed.

Summary: Removed duplicates from the dataset, ensuring data accuracy and uniqueness.