



Spotify & YouTube Music

1. Identify and Handle Missing Values:

- **Examine the dataset for any missing values. Which columns contain null values?**
- **How should missing values in the Views and Likes columns be handled? Should they be filled with a default value, removed, or handled in another way? Justify your approach.**

Columns with Missing Values:

- **key, valence, liveness, speechiness, loudness, tempo, danceability, DURATION_MS, instrumentalness, acousticness, ENERGY:** 2 missing values each
- **licensed, OFFICIAL_VIDEO, CHANNEL, youtube_info:** 491 missing values each
- **COMMENTS:** 593 missing values
- **description:** 911 missing values
- **views:** 2,484 missing values
- **LIKES:** 2,685 missing values
- **STREAM:** 610 missing values

Steps :

Step 1: Identify Missing Values

- **Action:** Inspected all columns and found missing values in Views and Likes column

Step 2: Handle Missing Values in Views

- **Action:** Replaced missing values with 0 to reflect no engagement.
- **Reason:** Using 0 aligns with the logic that missing values in Likes likely indicate no likes.

Step 3: Handle Missing Values in Likes

- **Action:** Replaced missing values with 0 to reflect no engagement.
- **Reason:** Using 0 aligns with the logic that missing values in Likes likely indicate no likes.

Step 4: Verify Changes

- **Action:** Rechecked the dataset to confirm no missing values remained.
- **Reason:** Ensures data consistency and completeness for analysis.

2. Fix Irregularities in Merged Columns:

- **The Spotify_Info and Youtube_Info columns contain merged data separated by delimiters. Split these columns back into their original components. What are the original components, and how can you ensure that the split data is clean and accurate?**
- **After splitting, remove any unnecessary delimiters or prefixes/suffixes that do not belong.**

Steps:

Step 1: Split the Spotify_Info and Youtube_Info Columns

Action: Used the Text to Columns feature to split the data.

Components Identified:

Spotify_Info: Spotify Link , Spotify ID.

Youtube_Info: Youtube Link , Youtube video Title.

Reason: Splitting ensures individual components can be analyzed separately.

Step 2: Clean the Split Data

Action: use delimiters (|) pipe symbol for spotify_info and character length for youtube_info to split the data accurately.

Reason: Ensures clean and uniform data.

Step 3: Verify Changes

Action: Reviewed the new columns to ensure data integrity and no residual delimiters.

Reason: Confirmed data is clean and ready for analysis.

3. Correct Case Sensitivity and Naming Conventions:

- **The column names have inconsistent case sensitivity (some are uppercase, others lowercase). Standardize all column names to follow a consistent format (e.g., all lowercase with underscores).**
- **Fix any data entries where case sensitivity might affect consistency (e.g., artist names or track titles). Ensure that the Artist and Track columns are formatted consistently.**

Steps:

Step 1: Renamed columns to lowercase with underscores for uniformity and tool compatibility.

Step 2: Standardized Artist and Track columns to title case for consistency and to avoid duplicates caused by case differences.

Step 3: Verified and removed duplicates after case standardization.

4. Remove or Handle Irrelevant Columns:

- **Identify and remove any irrelevant or randomly generated columns that do not provide useful information for analysis. Which columns should be removed, and why?**
- **If any random data exists in relevant columns, clean or remove those entries.**

Steps:

Step 1 : Inspecting and Identifying Columns:

- Irrelevant Columns Removed:
 - Example: Test_Column (contains randomly generated values not useful for analysis).
 - Reason: Does not contribute to meaningful insights.
- Columns Retained:
 - Example: artist, track, views, likes.

Step 2 : Handling Random Data:

- Checked the views and likes columns:
 - Found random entries like "test", "unknown".
 - Replaced invalid entries in views with the median.
 - Removed rows with invalid data in likes.

Step 3 : Final Verification:

- Confirmed no irrelevant columns remain.
- Ensured all data entries are consistent and usable.

5. Handle Inconsistent Data Types:

- **Some columns that should be numeric (e.g., Danceability, Energy) are stored as text. Convert these columns back to numeric format. What steps would you take to identify and fix any issues that arise during this conversion?**
- **Ensure that all numeric columns are in the correct format and handle any non-numeric values or anomalies.**

Steps:

Step 1: Identified columns (Danceability, Energy) with incorrect data types.

Reason: Numeric values stored as text cannot be used for statistical analysis.

Step 2: Replaced non-numeric entries with the median.

Reason: The median provides a robust replacement unaffected by outliers.

Step 3: Converted columns to numeric format.

Reason: Ensures the columns are ready for analysis and calculations.

Step 4: Verified the conversion.

Reason: Ensures no invalid data remains, and the data type is consistent.

6. Address and Fix Invalid Data Entries:

- **Check the Views column for any entries labeled as "invalid_data" or any other incorrect values. Replace these entries and justify your method.**
- **Ensure that all values in the Album column are correctly labeled and that there are no numeric entries or irrelevant data.**

Steps:

Step 1: Identified invalid entries in the Views column.

Issue: Non-numeric entries like "invalid_data."

Action: Replaced invalid entries with the median of the column.

Justification: The median is less affected by outliers and provides a reliable replacement.

Step 2: Checked the Album column for inconsistencies.

Issue: Numeric entries and irrelevant data.

Action: Replaced invalid entries with "Unknown" and standardized the column to title case.

Justification: Ensures consistency and prevents irrelevant data from affecting analysis.

Step 3: Verified changes.

Action: Rechecked both columns to confirm all values are valid and consistent.

7. Check for and Remove Duplicate Rows:

- **Identify and remove any duplicate rows in the dataset. How can you ensure that the remaining data is unique and accurate**

Steps:

Step 1: Identified duplicate rows.

Action: Inspected the dataset for exact duplicates across all columns.

Justification: Prevents redundant data from skewing analysis.

Step 2: Removed duplicate rows.

Action: Used the "Remove Duplicates" checking key columns (e.g., Artist, Track, Views).

Justification: Ensures concise data without losing relevant information.

Step 3: Validated data uniqueness.

Action: Re-inspected the dataset to confirm no duplicates remain.

Justification: Ensures each row represents unique and accurate information.

8. Reorder and Rename Columns for Clarity:

Reorder the columns in a logical sequence to improve the dataset's readability and usability. What order makes the most sense for this dataset?

Rename columns where necessary to ensure that their names clearly reflect the data they contain.

Steps :

Step 1: Reordered columns for readability.

Action: Rearranged columns to group metadata, metrics, and categorical data logically.

Justification: Improves usability for analysis by organizing related columns together.

Step 2: Renamed ambiguous column names.

Action: Updated names to be descriptive and consistent (e.g., Views → total_views).

Justification: Enhances clarity and aligns with standard naming conventions.

Step 3: Verified the changes.

Action: Reviewed the dataset to confirm the order and names were correct.

Justification: Ensures the dataset is ready for analysis with no confusion about column purposes.