

Collins Tse Jacky Ke Rebecca Bachtra Christy Yau

Thursday 18th December, 2025



```
import os
import pandas as pd
import numpy as np

def load_and_clean_spotify(path="data/spotify.csv"):
    # Load CSV and skip malformed rows
    df = pd.read_csv(path, on_bad_lines="skip", low_memory=False)

    # Drop unnamed index column
    df = df.loc[:, ~df.columns.str.contains("^\u0331named")]

    # Strip whitespace in all strings
    for col in df.select_dtypes(include=["object"]):
        df[col] = df[col].astype(str).str.strip()

    # Fix booleans
    if "explicit" in df.columns:
        bool_map = {"True": True, "False": False, "1": True, "0": False}
        df["explicit"] = df["explicit"].map(bool_map)

    # Convert numeric columns
    numeric_cols = [
        "popularity", "duration_ms", "danceability", "energy", "key",
        "loudness", "speechiness", "acousticness", "instrumentalness",
        "liveness", "valence", "tempo"
    ]
    for col in numeric_cols:
        if col in df.columns:
            df[col] = pd.to_numeric(df[col], errors="coerce")

    # Remove duplicate track IDs
    if "track_id" in df.columns:
        df = df.drop_duplicates(subset=["track_id"])

    # Add new engineered features
```

```
if "duration_ms" in df.columns:
    df["duration_min"] = df["duration_ms"] / 60000

if "year" in df.columns:
    df["decade"] = (df["year"] // 10) * 10

# Clean track_genre text
if "track_genre" in df.columns:
    df["track_genre"] = (
        df["track_genre"]
        .str.replace(",.*", "", regex=True)
        .str.lower()
        .str.strip()
    )

# Remove rows with missing popularity (required for modeling)
if "popularity" in df.columns:
    df = df.dropna(subset=["popularity"])

return df

df_clean = load_and_clean_spotify("data/spotify.csv")
df_clean.head()
```

	track_id	title	artists	popularity	mode	danceability	energy	loudness	tempo	duration_ms	time_signature	min
0	5SuOikCQGdyvBQfMNgSW4610	... Hoshino	... Hoshino	0	0.1430320003637187.917	acoustic	0.814433					
1	4qPBDGNG13513110001201660	... WoofAcous- wartic)Acous- tic	... WoofAcous- wartic)Acous- tic	1	0.07692100000002677.489	acoustic	0.8143500					
2	1iJBsg7EjYX5M8H02005b83590	... MicBeBe- songZAYN Again	... MicBeBe- songZAYN Again	1	0.055210000000026.382	acoustic	0.813767					
3	6lfxK30CG4C7Tg20198002060596	... GraRichHelp nisAsiaFalling (OrIg- i- Love nal Mo- tion Pic- ture Sou...	... GraRichHelp nisAsiaFalling (OrIg- i- Love nal Mo- tion Pic- ture Sou...	1	0.03690500007204813740	acoustic	0.81365550					
4	5vjJCSHdHP162G5W3K6184430	... OveOn On street	... OveOn On street	1	0.05246000008296719449	acoustic	0.814217					

5 rows × 21 columns

```
os.makedirs("data/processed", exist_ok=True)

outpath = "data/processed/spotify_clean.csv"
df_clean.to_csv(outpath, index=False)

print("Saved cleaned dataset to:", outpath)

Saved cleaned dataset to: data/processed/spotify_clean.csv
```