

Musical Carbon Dating

A Statistical Feature Recognition Approach (1960-2020)

敖婧妍, 李赫轩, 田润泽, 王胤博, 肖俊佳

School of Statistics

December 26, 2025

Table of Contents

- 1 Overview & Data
- 2 Methodology Pipeline
- 3 Phases I–III: Baseline & Diagnostics
- 4 Phases IV–VI: Refinement & Results
- 5 Applications

The Research Question

Feature Recognition

"Can we determine the vintage of a musical recording purely from its acoustic properties?"

Objective: To build a regression model that maps audio features to release year, quantifying the "Arrow of Time" in music production.

Hypothesis: Musical eras have distinct, quantifiable acoustic fingerprints (e.g., the "dryness" of 70s rock vs. the "compression" of 2000s pop).

The Data: Spotify 600k Tracks

Data Filtering Strategy:

- **Source:** Spotify 600k Tracks Dataset (Kaggle).
- **Filter 1:** Timeframe $1960 \leq T \leq 2020$ (Modern Era).
- **Filter 2:** popularity > 30 (Focus on culturally significant music).
- **Final Sample:** $N = 250,971$ tracks (1960–2020).

Validation Strategy:

- **Random Split:** 80% Training / 20% Test.
- *Rationale:* We are testing "feature recognition" (interpolating styles), not future forecasting (extrapolating time).

The Feature Set (p=13)

We utilized all 13 available audio features.

Physical Features

- ① Loudness (dB)
- ② Tempo (BPM)
- ③ Duration (ms)

Musical Features

- ④ Key (0-11)
- ⑤ Mode (Major/Minor)
- ⑥ Time Signature

Perceptual Features

- ⑦ Acousticness
- ⑧ Danceability
- ⑨ Energy
- ⑩ Instrumentalness
- ⑪ Liveness
- ⑫ Speechiness
- ⑬ Valence (Positivity)

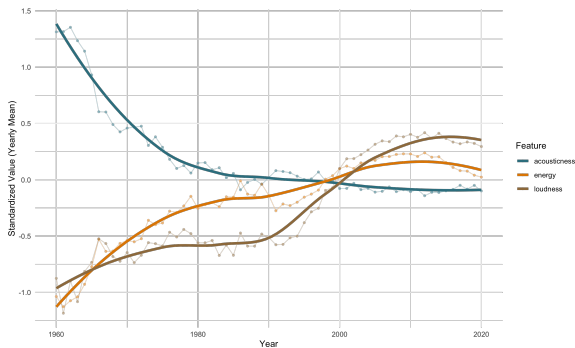
The Regression Pipeline

We followed a rigorous 5-phase statistical workflow:

- ➊ **Phase I: Simple Linear Regression (SLR)** *Hypothesis testing: The "Loudness War".*
- ➋ **Data Standardization** *Z-Score Normalization ($x' = \frac{x - \mu}{\sigma}$) to ensure scale invariance.*
- ➌ **Phase II: Multiple Linear Regression (MLR)** *Baseline model using all $p = 13$ features.*
- ➍ **Phase III: The Diagnostic Audit** (Critical Step) *Testing Linearity, Multicollinearity, Normality, and Homoscedasticity.*
- ➎ **Phase V: Refinement (WLS)** *Weighted Least Squares to correct for Heteroscedasticity.*
- ➏ **Phase VI: Interpretation** *Quantifying the "Arrow of Time" and Cultural Shifts.*

Phase I: The "Loudness War" (SLR)

$$\text{Year}_i = \beta_0 + \beta_{\text{loud}} \cdot \text{Loudness}_i + \varepsilon_i$$



Results:

- t -stat: **183.3** (Highly Significant).
- $\beta_{\text{loud}} \approx 5.44$ years/dB (Standardized).
- $R^2 = 0.1434$.

Conclusion: Tracks have consistently gotten louder, but Loudness alone explains only 14% of the variance.

Phase II: Multiple Linear Regression (Baseline)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- **Algorithm:** Ordinary Least Squares (OLS) with all 13 features.
- **Mean Error (MAE):** 9.72 Years
- **Median Error:** 7.52 Years

Why so low?

An R^2 of 0.14 suggests we are missing non-linear patterns or violating OLS assumptions. We initiated a **Diagnostic Audit**.

Phase III: Diagnostic Audit (1/2)

Test 1: Independence of Errors

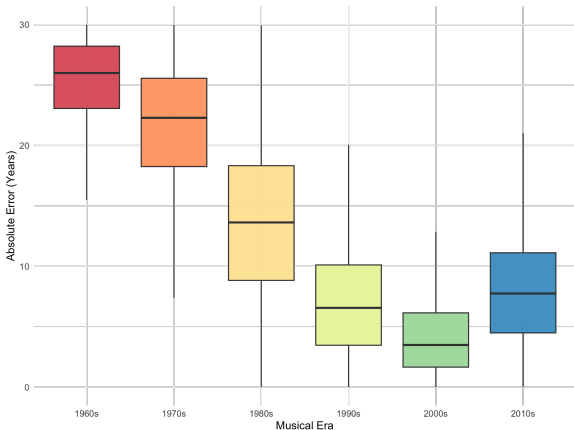
- **Method:** Durbin-Watson Statistic.
- **Result:** $DW \approx 2.001$. (Target = 2.0)
- **Finding:** Assumption Met. No autocorrelation. Each song is statistically unique.

Test 2: Multicollinearity

- **Method:** Variance Inflation Factor (VIF).
- **Concern:** High correlation between Loudness and Energy ($r = 0.74$).
- **Result:** Max VIF (Energy) = **3.58**. All VIFs < 5.0 .
- **Finding:** Assumption Met. No severe multicollinearity.

Phase III: Diagnostic Audit (2/2)

Test 3: Homoscedasticity (Constant Variance)



- **Visual:** Dispersed variance (Right) vs Tight variance (Left).
- **Insight:** "Stylistic Entropy". The definition of specific eras has blurred over time due to technology.
- **Statistic:** $\chi^2 = 20,754$ (Breusch-Pagan).
- **Conclusion:** Variance is expanding. OLS fails because it treats neighboring years inconsistently.

Solution: We must weight the model to trust the "consistent" eras more than the "chaotic" ones.
Enter WLS.

Phase IV: Model Selection

We compared two methods to identify the "True" feature set:

Method	Selected Features	Key Difference
Stepwise (AIC)	12 Features	Dropped Key ($p = 0.94$)
LASSO (L_1)	12 Features	Corroborated AIC selection

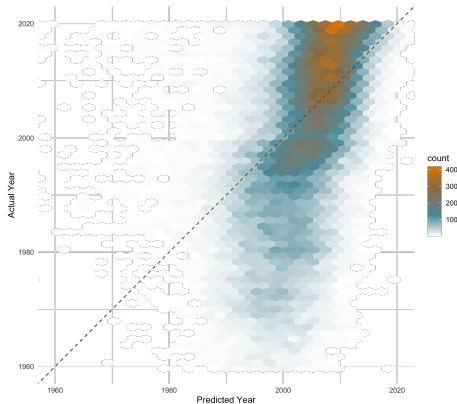
Decision: We utilized the full acoustic feature set.

- Removing variables offered negligible AIC improvement.
- Retaining subtle features (Key, Mode) ensures we capture harmonic evolution.

Phase V: Weighted Least Squares (WLS)

To cure Heteroscedasticity, we implemented WLS:

$$\min_{\beta} \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \beta)^2, \quad \text{where } w_i \propto \frac{1}{\text{Var}(\varepsilon_i)}$$



What WLS Fixes

- **Valid p-values & CIs.**
- **Efficient coefficients.**
- Baseline SLR achieves $R^2 = 0.14$
- Loudness coefficient: $\beta \approx 5.44$ years/dB
- Weighted R^2 : **0.26**
- Chisquare $\approx 20,754$ ($p < 0.001$)
- Test R^2 : ≈ 0.26
- RMSE: 12.83 yrs

Phase VI-A: Technological Drivers ("The Sound of Efficiency")

Regression reveals the impact of technology on composition ($p < 0.001$).

- **The Loudness-Energy Paradox (Multicollinearity Insight):**

- Mean Loudness has skyrocketed ($\beta_{\text{loud}} \approx +7.57$).
- Yet, Energy's coefficient is **negative** ($\beta_{\text{energy}} \approx -2.05$).
- **Interpretation:** Loudness comes from compression, not composition.

- **The Attention Economy (Duration):**

- Coefficient $\beta_{\text{duration}} \approx -0.33$ (Negative).
- Songs are getting statistically shorter, likely driven by streaming incentives and skipping behavior.

- **Instrumentalness** ($\beta \approx +0.20$): A shift towards beat-driven (Hip-Hop/EDM) production over vocal-centric ballads.

Phase VI-B: Cultural Evolution ("The Mood of an Era")

Applying Statistical Inference to Cultural Theory.

- **The "Sad Banger" Phenomenon:**

- **Danceability** ($\beta \approx +3.60$): The single strongest predictor. Rhythm is the defining feature of modernity.
- **Valence** ($\beta \approx -3.39$): Optimism has collapsed.
- **Synthesis**: We are dancing more, but feeling less.

- **The Acousticness Paradox (Ceteris Paribus):**

- **Raw Correlation**: Negative ($r \approx -0.12$).
- **WLS Coefficient**: **Positive** ($\beta \approx +0.65$).
- **Discovery**: *Controlling for Loudness*, modern music actually retains significant acoustic elements (Indie, Lo-Fi), hidden by the "Wall of Sound".

Analytical Triumph

By using **Partial Regression Coefficients**, we uncovered trends (like the Acousticness reversal) that simple correlation would have missed.

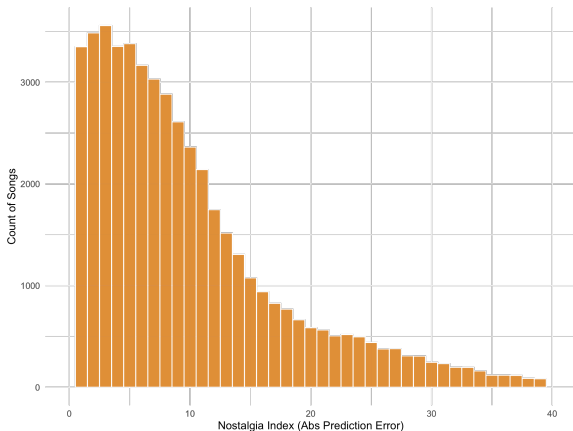
The Nostalgia Index

"One man's error is another man's feature."

We define the **Nostalgia Index** as the model's prediction error:

$$\text{Index} = |\hat{Y}_{\text{predicted}} - Y_{\text{actual}}|$$

High index = A song that "sounds" like it belongs to a different era.



Distribution Stats:

- Mean: **9.72** yrs
- Median: **7.52** yrs
- Max: **85.34** yrs

The index measures the "stylistic distance" between a song's audio and its true era.

Validation: Detecting "Time Travelers"

We validated the index on tracks known for their retro aesthetic.

Song	Year	Predicted	Index	Diagnosis
<i>Uptown Funk</i> (Ronson)	2015	2013.1	1.9	Modern Construction
<i>Physical</i> (Dua Lipa)	2020	2009.0	11.0	Retro Aesthetic
<i>Blinding Lights</i> (Weeknd)	2019	2004.2	14.8	80s Revival

Conclusion: The model correctly identifies these hits as "sounding old", proving it captures aesthetic style rather than just release dates.

Conclusion

- ① **Recognition:** We can date music to within 7.5 years (median) purely from audio.
- ② **Rigor:** Diagnostics proved OLS insufficient; WLS corrected the inference.
- ③ **Insight:** Musical evolution is quantifiable and technology-driven.
- ④ **Value:** The Nostalgia Index provides a metric for "Retro-vibe".

Thank You.