

Ukrainian Catholic University

Faculty of Applied Sciences

Business Analytics & Computer Science Programmes

A guide for beginner musicians:

How to become popular in the modern world

Econometrics final project report

Authors:

Anastasiia Kurylets,

Julia Chebotarova,

Viktoriiia Markus



APPLIED
SCIENCES
FACULTY ●

Contents:

Introduction.....	3
Literature review.....	3
Data description.....	3
Data analysis:.....	4
Methods.....	5
Results.....	6
Conclusions.....	6
Appendixes:.....	9

Introduction

Becoming popular in today's music world demands more than talent — it requires smart strategy and constant adaptation.

The aim of this research is to explore how aspiring musicians can achieve popularity in the modern world by analyzing the characteristics of popular music. Using Spotify data we identify possible practical strategies related to musical style and audience engagement.

Literature review

1. Angrist, J. D., & Pischke, J.-S. (2008). *Mostly Harmless Econometrics*.

From this source, we used the concept of omitted variable bias, which underlines the importance of including all relevant predictors in regression analysis. The book also emphasizes that ordinary least squares estimators are only unbiased when all Gauss-Markov assumptions hold, which guided our decision to test model assumptions before interpreting results.

2. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*.

This work provided a structured overview of regression modeling and diagnostics. We used it specifically for understanding model evaluation metrics such as R-squared and RMSE; learning about assumption tests; recognizing the importance of proper model selection and interpretation in statistical learning.

3. Wooldridge, J. M. (2015). *Introductory Econometrics: A Modern Approach*.

This book helped us clarify the correct interpretation of coefficients in linear models and provided practical examples of how multicollinearity and omitted variables can distort inference. It reinforced the need to check variance inflation factors when selecting independent variables.

Data description

For our research, we decided to choose the “[30000 Spotify Songs](#)”. This dataset contains information about 30 000 songs available on Spotify. The data is obtained through the Spotify API and covers a variety of musical characteristics, allowing us to analyze the factors that influence the popularity of tracks.

Variables:

track_id: Unique identifier for the song.

track_name: Name of the song.

track_artist: Artist who performed the song.

track_album_id: Unique identifier for the album.

track_album_name: Name of the album the song is from.

track_album_release_date: Release date of the album.

playlist_name: Name of the playlist containing the song.

playlist_id: Unique identifier for the playlist.

We do not take it into account variables above, because our main goal is to understand what properties of the song itself make it popular.

Selected variables:

track_popularity: Popularity score of the song (0–100), with higher values indicating greater popularity.

playlist_genre: Genre classification of the playlist.

playlist_subgenre: More specific subgenre of the playlist.

danceability: Score (0.0–1.0) indicating how suitable the song is for dancing.

energy: Score (0.0–1.0) reflecting the track's intensity and activity.

key: Musical key of the track, represented as integers (0 = C, 1 = C#/D ♭, etc.).

loudness: Average loudness of the track in decibels (dB).

mode: Indicates if the track is in a major (1) or minor (0) key.

speechiness: Score (0.0–1.0) indicating the presence of spoken words.

acousticness: Confidence score (0.0–1.0) that the track is acoustic.

instrumentalness: Likelihood (0.0–1.0) that the track has no vocals.

liveness: Probability (0.0–1.0) that the track was recorded live.

valence: Score (0.0–1.0) representing musical positivity or mood.

tempo: Estimated tempo in beats per minute (BPM).

duration_ms: Length of the track in milliseconds.

Data analysis:

Missing values:

On the picture 1 (Appendixes, data analysis) you can see that the dataset contains almost no missing values. This allows you to proceed to data analysis and model building without additional cleaning. The completeness of the data ensures higher accuracy of further statistical conclusions.

Distribution of songs:

The plot 2 (Appendixes, data analysis) shows how the popularity of approximately 30,000 songs on Spotify is distributed on a scale from 0 to 100. Most tracks are popular between 40-65 points, where the second highest bar is visible. This suggests a concentration of moderately popular songs in that range.

It is also evident that relatively few songs achieve very high popularity (above 85), while a broad distribution spans from 20 to 80. In summary, most songs on Spotify maintain average to moderately high popularity, with only a few becoming viral hits but many remaining entirely obscure.

Multicollinearity:

The plot 3 (Appendixes, data analysis) presents the correlation matrix heatmap of all numerical independent variables used in the model. It helps to visually detect multicollinearity, i.e. strong linear relationships between predictors. While most variable pairs show low or moderate correlation, a few noteworthy patterns are observed. For instance, energy and loudness are positively correlated ($r = 0.68$), suggesting that more energetic songs tend to be louder. Another significant relationship is the strong negative correlation between energy and acousticness ($r = -0.54$), which reflects that highly energetic songs are less likely to be acoustic. Additionally, danceability and valence show a mild positive correlation ($r = 0.33$).

To formally detect multicollinearity among the independent variables, we computed the Variance Inflation Factor (VIF) for each predictor. The results are shown in the picture 4 (Appendixes, data analysis). All VIF values are below 3, with the highest being for energy (VIF = 2.66) and loudness (VIF = 2.08), which aligns with their moderately strong correlation in the heatmap (plot 3). The intercept (const) naturally shows a very high VIF, which is expected and not a concern.

Methods

1. Ordinary Least Squares (OLS) Regression

To investigate the factors that influence a song's popularity on Spotify, we used Ordinary Least Squares (OLS) regression as our primary estimation method. The dependent variable in our model is `track_popularity`, which ranges from 0 to 100 and reflects how well a song performs on the platform. As independent variables, we included 11 musical characteristics such as danceability, energy, loudness, speechiness, valence, and others.

We chose Ordinary Least Squares (OLS) regression because it provides interpretable estimates of how each musical feature affects song popularity. It allows us to formally test hypotheses using t-tests and F-tests, and under standard assumptions, OLS produces efficient and unbiased results. This method is widely used in econometrics, making our analysis both rigorous and accessible.

2. Logit Regression (Binary Logistic Regression)

We used logistic regression to model the probability that a song becomes highly popular, defining a binary variable where 1 indicates a “hit” (popularity > 80) and 0 otherwise. This

method is suitable for binary outcomes and allows us to estimate how musical features influence the likelihood of success. The model includes the same predictors as the OLS regression, such as danceability, energy, and valence. Logistic regression was chosen because it provides interpretable results in terms of odds and is well-suited for identifying which characteristics increase the chances of a song becoming a hit.

Results

As a result of the study, we found that:

- For the general dataset, both the OLS model and the Logit model have a low R-squared (7.2%) and Pseudo R-squared (5.4%), i.e. the models are weak, but significant due to other indicators. In general, the popularity of songs on Spotify increases with danceability, loudness, valence and tempo, but decreases with energy, speechiness, instrumentalness, liveness, and duration, as confirmed by both regressions, although acousticness was significant only in OLS.
- For Pop, OLS ($R^2 = 5.7\%$) and Logit (Pseudo $R^2 = 4.1\%$) show that danceability, energy, speechiness, and acousticness increase popularity, while instrumentalness, valence, and duration reduce it.
- For Rap, OLS ($R^2 = 4.8\%$) and Logit (Pseudo $R^2 = 5.2\%$) confirm that danceability and acousticness increase popularity, whereas speechiness, valence, instrumentalness, and duration decrease it.
- For EDM, OLS ($R^2 = 11.8\%$) and Logit (Pseudo $R^2 = 9.6\%$) show strong influence of acousticness and valence (positive), and instrumentalness, speechiness, tempo, and duration (negative); danceability is not significant in OLS but has a negative effect in Logit.
- For R&B, OLS ($R^2 = 7.1\%$) and Logit (Pseudo $R^2 = 5.6\%$) reveal that energy and acousticness positively affect popularity, while valence_danceability, instrumentalness, liveness, and duration lower it; speechiness and tempo are not significant.
- For Latin, OLS ($R^2 = 6.5\%$) and Logit (Pseudo $R^2 = 5.6\%$) indicate that danceability, energy_loudness, acousticness, speechiness, and tempo raise popularity, while instrumentalness, liveness, and valence reduce it; duration is significant only in OLS.

You can look at the outputs for each model in Appendixes.

Conclusions

We analyzed the impact of musical characteristics on the popularity of songs on Spotify. The main goal was to identify strategies that aspiring musicians can use to gain popularity in today's music industry. By applying econometric models such as OLS and Logit, we identified several important trends.

In general, we found that song popularity increases with higher danceability, loudness, positivity (valence), and tempo, while it decreases with higher energy, speechiness, instrumentalness, liveness,

and longer song duration. Although the models had relatively low R-squared values, the results were statistically significant and highlighted key factors influencing success.

When analyzing different genres, we observed specific patterns. In pop music, songs that are danceable, energetic, include spoken elements, and have acoustic qualities tend to be more popular, while those with excessive instrumentalness, high positivity, and long duration tend to be less successful. In rap, danceability and acousticness positively influence popularity, but too much speechiness, positivity, instrumentalness, and longer duration can reduce a song's success. For EDM, acousticness and valence help boost popularity, while too much instrumental focus, high tempo, speechiness, and long duration have a negative effect. In R&B, energy and acousticness support popularity, but tracks that are overly cheerful, dance-heavy, instrumental, or long tend to perform worse. Finally, in Latin music, danceability, energy, loudness, acousticness, speechiness, and a fast tempo are all beneficial, while too much instrumentalness, liveness, and positivity can reduce a song's appeal.

Based on these findings, we suggest that musicians tailor their approach depending on their genre. For pop artists, focusing on catchy, energetic, and moderately acoustic songs is key, avoiding overly long or purely instrumental tracks. Rap musicians should balance rhythm and melody, being careful not to rely too heavily on speech alone and adding acoustic touches to enhance appeal. EDM producers should aim for emotionally uplifting tracks with a clean balance between electronic sounds and vocals, keeping tempo and instrumental sections in moderation. R&B artists are encouraged to create powerful, acoustic-rich songs while avoiding too much cheerfulness or dance emphasis. Latin musicians should prioritize energetic, danceable tracks with strong acoustic and rhythmic elements, while minimizing overuse of instrumental or live performance features.

In conclusion, while talent remains crucial, a thoughtful combination of musical characteristics, adapted to the genre, can significantly enhance an artist's chances of gaining popularity in the modern music scene.

Appendixes:

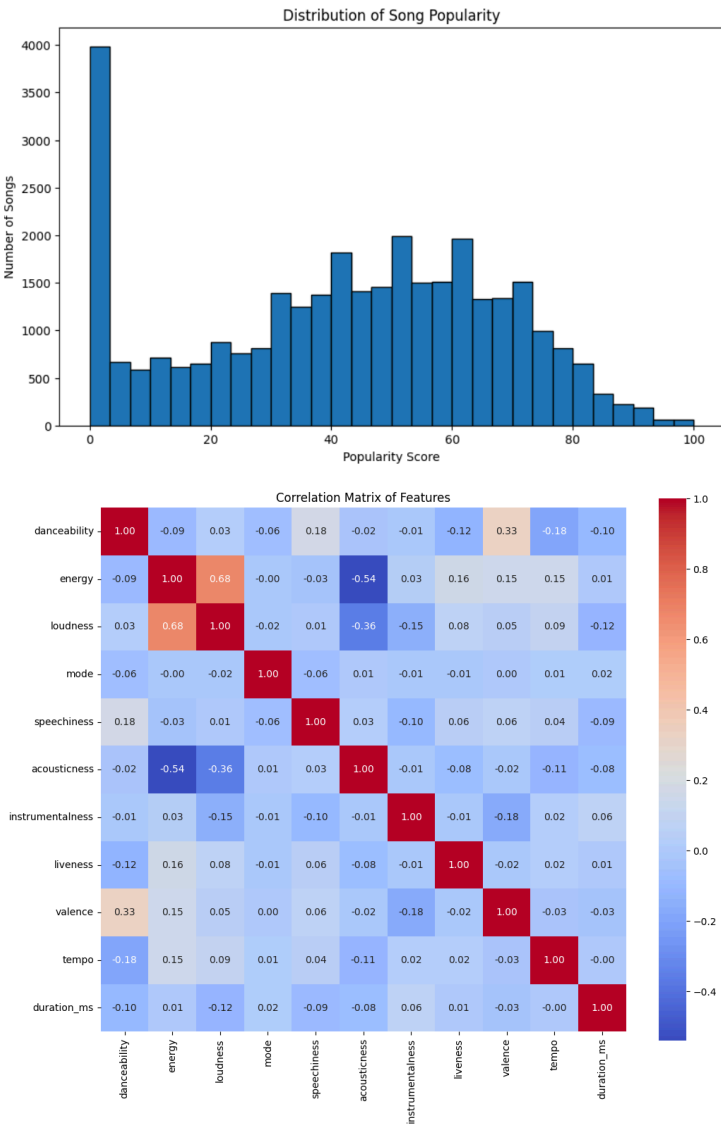
Data Analysis

Picture 1

Missing values per column:	
track_id	0
track_name	5
track_artist	5
track_popularity	0
track_album_id	0
track_album_name	5
track_album_release_date	0
playlist_name	0
playlist_id	0
playlist_genre	0
playlist_subgenre	0
danceability	0
energy	0
key	0
loudness	0
mode	0
speechiness	0
acousticness	0
instrumentalness	0
liveness	0
valence	0
tempo	0
duration_ms	0

	Variable	VIF
0	const	156.463325
1	danceability	1.311575
2	energy	2.661166
3	loudness	2.083386
4	mode	1.007601
5	speechiness	1.069859
6	acousticness	1.464190
7	instrumentalness	1.132005
8	liveness	1.049505
9	valence	1.266500
10	tempo	1.068981
11	duration_ms	1.051185

Picture 2



Picture 3

Picture 4

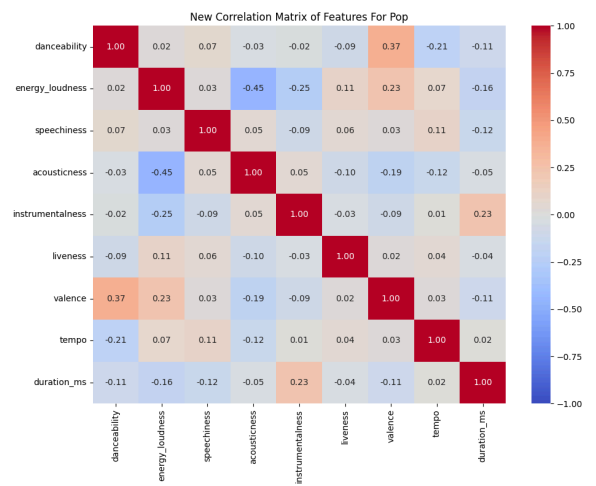
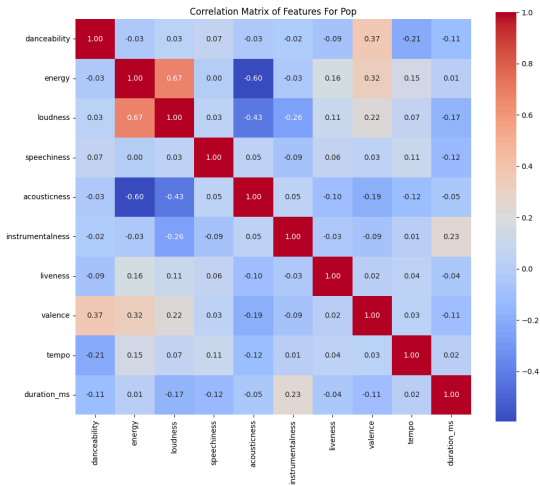
Models:

1. General Data

OLS Regression Results						
Dep. Variable:	track_popularity		R-squared:	0.072		
Model:	OLS		Adj. R-squared:	0.072		
Method:	Least Squares		F-statistic:	232.1		
Date:	Fri, 09 May 2025		Prob (F-statistic):	0.00		
Time:	19:52:20		log-likelihood:	-1.5102e+05		
No. Observations:	32833		AIC:	3.021e+05		
Df Residuals:	32821		BIC:	3.022e+05		
Df Model:	11					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	78.0638	1.662	46.982	0.000	74.807	81.321
danceability	5.0361	1.049	4.803	0.000	2.981	7.091
energy	-29.5015	1.198	-24.629	0.000	-31.849	-27.154
loudness	1.5243	0.064	23.758	0.000	1.399	1.650
mode	0.6635	0.269	2.466	0.014	0.136	1.191
speechiness	-7.2310	1.356	-5.332	0.000	-9.889	-4.573
acousticness	3.2319	0.732	4.416	0.000	1.797	4.666
instrumentalness	-11.9744	0.630	-18.998	0.000	-13.210	-10.739
liveness	-4.3131	0.882	-4.891	0.000	-6.042	-2.585
valence	2.7970	0.641	4.362	0.000	1.540	4.054
tempo	0.0212	0.005	4.156	0.000	0.011	0.031
duration_ms	-4.587e-05	2.28e-06	-20.152	0.000	-5.03e-05	-4.14e-05

Logit Regression Results						
Dep. Variable:	popular	No. Observations:	32833			
Model:	Logit	Df Residuals:	32822			
Method:	MLE	Df Model:	10			
Date:	Fri, 09 May 2025	Pseudo R-squ.:	0.05358			
Time:	19:52:21	Log-Likelihood:	-18655.			
converged:	True	LL-Null:	-19711.			
Covariance Type:	nonrobust	LLR p-value:	0.000			
	coef	std err	z	P> z	[0.025	0.975]
const	1.9507	0.163	11.996	0.000	1.632	2.269
danceability	0.3084	0.101	3.047	0.002	0.110	0.507
energy	-2.6944	0.118	-22.891	0.000	-2.925	-2.464
loudness	0.1650	0.007	24.591	0.000	0.152	0.178
speechiness	-0.8316	0.129	-6.437	0.000	-1.085	-0.578
acousticness	0.0548	0.069	0.799	0.424	-0.080	0.189
instrumentalness	-1.9380	0.092	-21.053	0.000	-2.118	-1.758
liveness	-0.3088	0.087	-3.534	0.000	-0.480	-0.138
valence	0.5040	0.062	8.119	0.000	0.382	0.626
tempo	0.0018	0.000	3.702	0.000	0.001	0.003
duration_ms	-1.473e-06	2.38e-07	-6.203	0.000	-1.94e-06	-1.01e-06

2. For genre Pop



OLS Regression Results

Dep. Variable:

track_popularity

R-squared:

0.057

Model:

OLS

Adj. R-squared:

0.055

Method:

least Squares

F-statistic:

36.61

Date:

Fri, 09 May 2025

Prob (F-statistic):

1.34e-63

Time:

20:48:08

Log-likelihood:

-25414.

No. Observations:

5507

AIC:

5.085e+04

Df Residuals:

5497

BIC:

5.091e+04

Df Model:

9

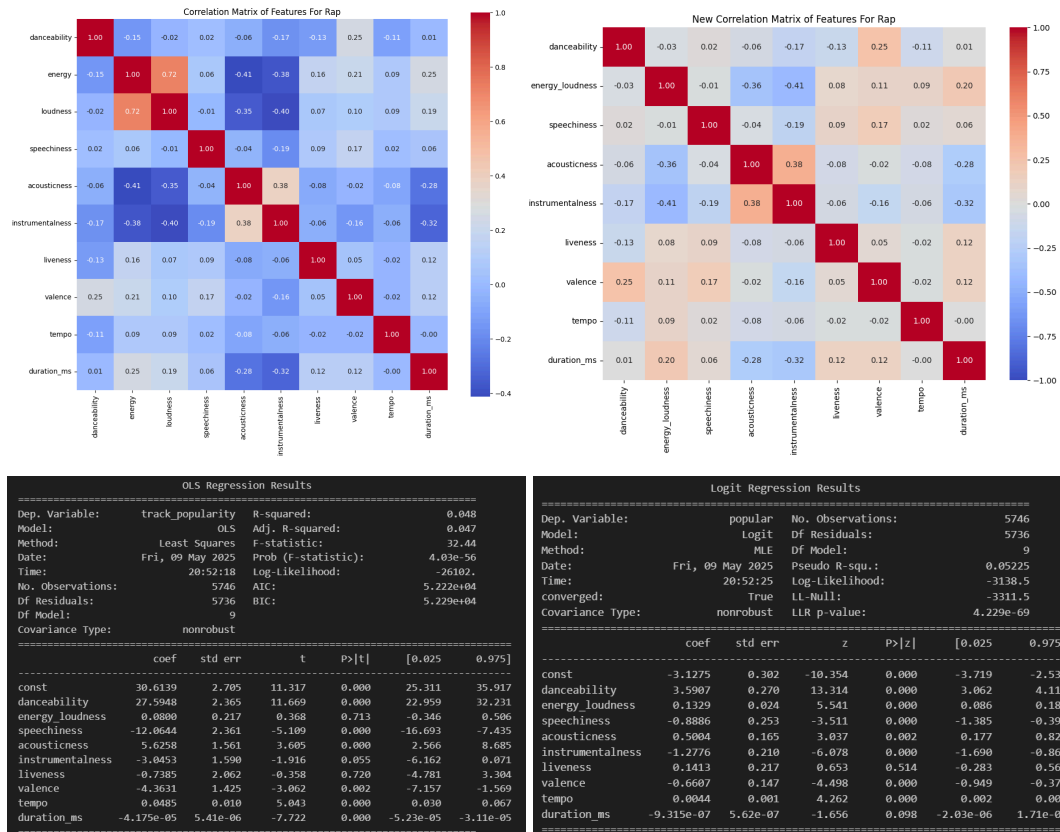
Covariance Type:

nonrobust

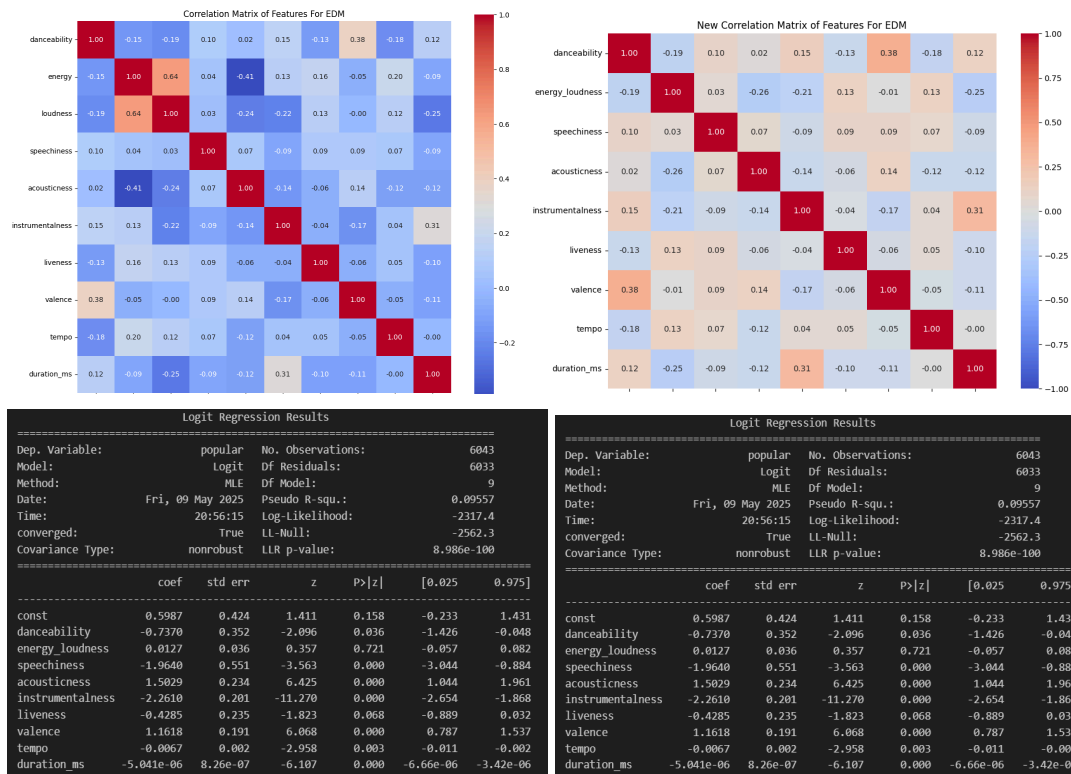
	coef	std err	t	P> t	[0.025	0.975]
const	56.6175	3.363	16.834	0.000	50.024	63.211
danceability	17.9317	2.889	6.207	0.000	12.268	23.595
energy_loudness	2.2392	0.286	7.821	0.000	1.678	2.800
speechiness	20.4675	4.974	4.115	0.000	10.717	30.218
acousticness	7.4135	1.727	4.292	0.000	4.027	10.800
instrumentalness	-13.7755	1.901	-7.247	0.000	-17.502	-10.049
liveness	-2.4541	2.460	-0.997	0.319	-7.278	2.369
valence	-4.8624	1.678	-2.898	0.004	-8.151	-1.574
tempo	-0.0255	0.014	-1.833	0.067	-0.053	0.002
duration_ms	-4.613e-05	7.67e-06	-6.010	0.000	-6.12e-05	-3.11e-05

Logit Regression Results						
Dep. Variable:	popular	No. Observations:	5507			
Model:	Logit	Df Residuals:	5497			
Method:	MLE	Df Model:	9			
Date:	Fri, 09 May 2025	Pseudo R-squ.:	0.04092			
Time:	20:49:55	Log-Likelihood:	-3528.8			
converged:	True	LL-Null:	-3679.3			
Covariance Type:	nonrobust	LLR p-value:	1.534e-59			
	coef	std err	z	P> z	[0.025	0.975]
const	-0.2117	0.298	-0.710	0.478	-0.796	0.373
danceability	1.3346	0.255	5.240	0.000	0.835	1.834
energy_loudness	0.2329	0.027	8.698	0.000	0.180	0.285
speechiness	1.4367	0.416	3.452	0.001	0.621	2.252
acousticness	0.4586	0.151	3.029	0.002	0.162	0.755
instrumentalness	-1.9704	0.235	-8.385	0.000	-2.431	-1.510
liveness	-0.1490	0.212	-0.703	0.482	-0.564	0.266
valence	-0.3371	0.147	-2.293	0.022	-0.625	-0.049
tempo	-0.0012	0.001	-0.978	0.328	-0.004	0.001
duration_ms	-9.966e-07	7.06e-07	-1.412	0.158	-2.38e-06	3.87e-07

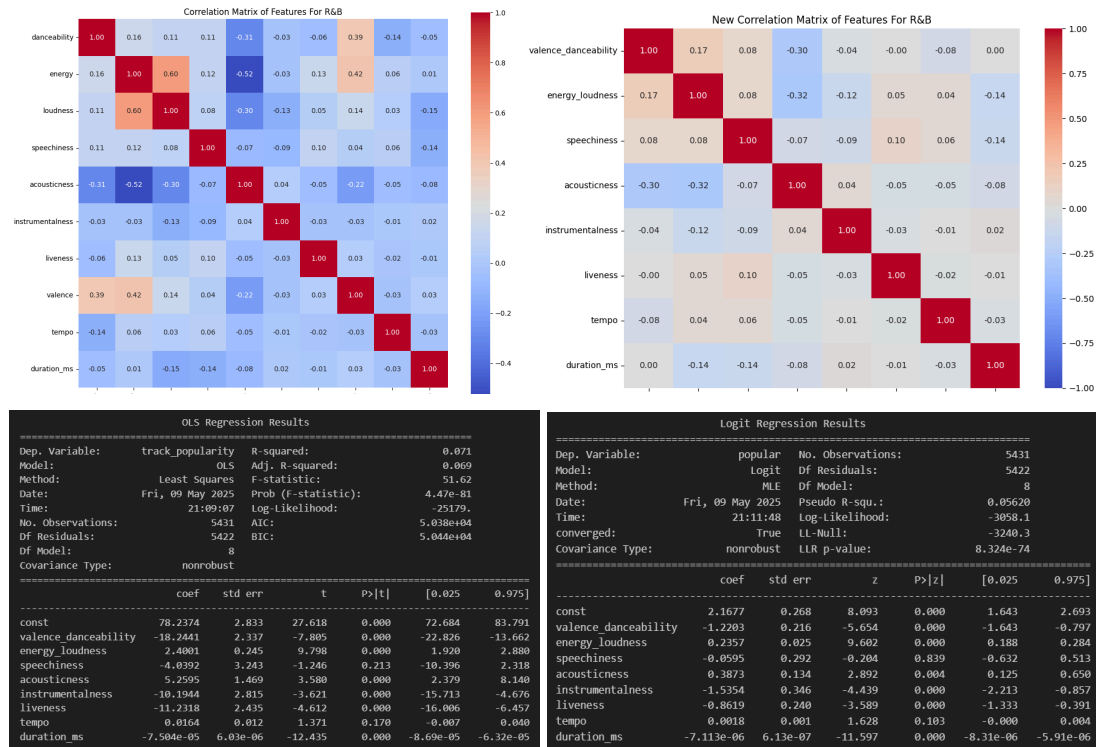
3. For genre Rap



4. For genre EDM (electronic dance music)



5. For genre R&B (rhythm and blues)



6. For genre Latin

