

# What makes a song popular?

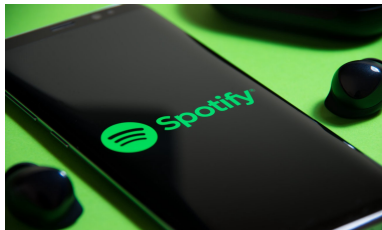
Jiayi CHEN, Marouane FELLOUSSI

French Civil Aviation University

April 30, 2021

# Context

- Spotify is popular and widely used, so it reflects people's opinions or taste in music.
- We'd like to analyse how the acoustic features alone influence the popularity of a song.



# Data Source

The dataset used is from Kaggle.

- Consists of more than 170,000 tracks
- Years 1920 to 2021

# Description

## Dependent Variable: Popularity (POP)

- Concerned demographic: US Spotify users
- Ranges from 0 to 100
- Calculated by algorithm and based on
  - Total number of plays the track has had
  - How recent those plays are

# Description

## Independent Variables - Numerical variables

Name	Variable	Range	Explanation
Acousticness	ACOU	0-1	Usage of acoustic/electric means.
Danceability	DNC	0-1	How suitable a track is for dancing.
Energy	NRG	0-1	Perceptual measure of intensity and activity.
Duration	DUR	3.88-88.97	Length of the track (min).
Instrumentalness	INSTR	0-1	Whether a track contains any or no vocals.
Valence	VAL	0-1	Positiveness of the track.
Tempo	TMP	0-100	Speed in Beat Per Minute (BPM).
Liveness	LIV	0-1	Relative duration sounding as a live show.
Loudness	LOUD	(-60)-0	Relative loudness in decibel (dB).
Speechiness	SPCH	0-1	Relative length containing human voice.
Year	YR	1920-2021	Year of release of track.

# Description

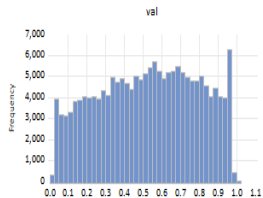
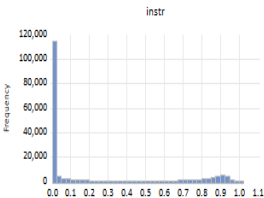
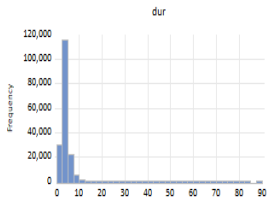
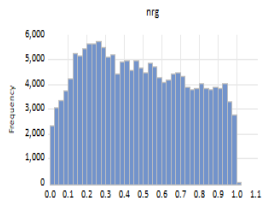
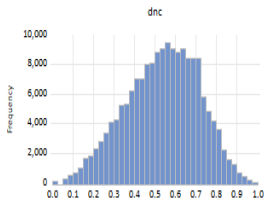
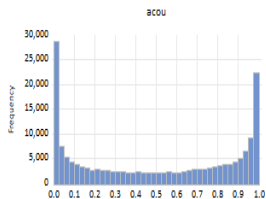
## Independent Variables - Dummy variables

Name	Variable	Explanation
Mode	MODE	Type of scale, 0=Minor, 1=Major
Explicit	XPLC	0=No explicit content, 1=Explicit content

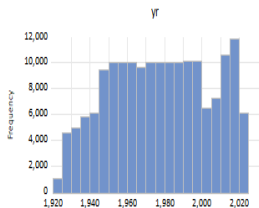
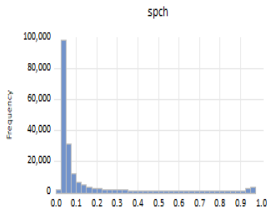
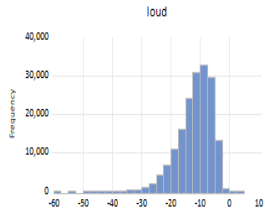
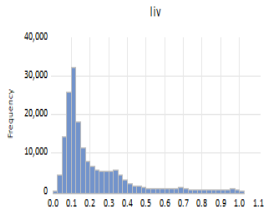
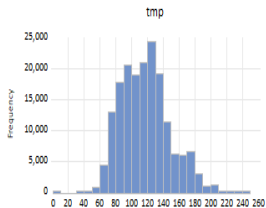
## Independent Variables - Categorical variables

Name	Variable	Explanation
Key	KEY	The primary key of the track encoded as integers in between 0 and 11. E.g. 0 = C, 1 = C $\sharp$ /D $\flat$ , 2 = D, and so on.

# Distribution

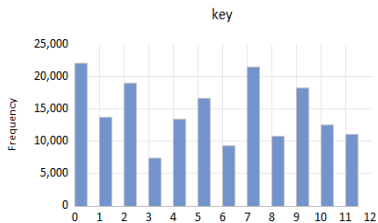
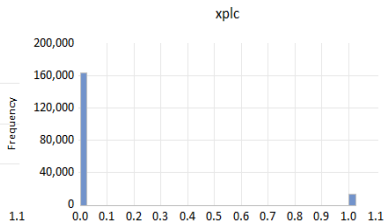
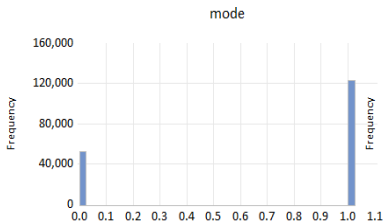


# Distribution





# Distribution



# Estimated Equation

A first equation is estimated using least squares, for which the output is given by

	C	ACOU	DNC	DUR	INSTR	KEY	LIV	LOUD	MODE	NRG	SPCH	TMP	VAL	XPLC	YR
pop	-582.25	-7.55	-3.23	-0.04	-16.07	-0.04	-8.22	-0.04	0.71	-3.37	-22.06	-0.01	4.06	9.92	0.31
prob	0.000	0.000	0.000	0.018	0.000	0.001	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000
R <sup>2</sup>	0.36357														
F-stat[prob]	7117.013[0.000]														

Figure: Estimation output for Equation 1

# Estimated equation

Name	Variable	Expectation	Model
Acousticness	ACOU	?	-7.55
Danceability	DNC	?	-3.23
Energy	NRG	+	-3.37
Duration	DUR	-	-0.04
Instrumentalness	INSTR	-	-16.07
Valence	VAL	+	+4.06
Tempo	TMP	+	-0.01
Liveness	LIV	-	-8.22
Loudness	LOUD	+	-0.04
Speechiness	SPCH	?	-22.06
Year	YR	?	+0.31
Mode	MODE	?	+0.71
Explicit	XPLC	?	+9.92

# JB Test

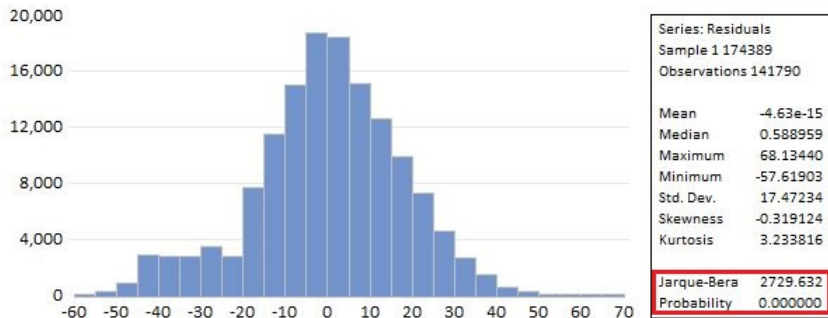


Figure: Jarque-Bera test

# Breusch-Pagan-Godfrey Test

---

HTSK Test: BPG

$H_0$  : HMSK

---

F-stat[prob]    3230.089[0.000]

OBS  $\times$  R<sup>2</sup>    34289.03[0.000]

---

Figure: BPG test : Homoscedasticity

# Variance inflation factors

Variable	ACOU	DNC	DUR	INSTR	KEY	LIV	LOUD	MODE	NRG	SPCH	TMP	VAL	XPLC	YR
Centered VIF	3.03	1.84	1.07	1.25	1.01	1.09	3.20	1.03	5.06	1.55	1.10	1.93	1.26	1.99

$1/(1-R^2) = 1.57$

Figure: Variance Inflation factors

# Stability Diagnostics

---

Ramsey's RESET test

---

	Value	Probability
F-statistic	0.4162	0.5188

---

Figure: Ramsey's RESET test

# Estimated equation

	C	ACOU	DNC	DUR	INSTR	KEY	LIV	LOUD	MODE	NRG	SPCH	TMP	VAL	
pop	-38625.32	9.50	5.23	0.06	-22.52	-0.01	-9.98	1.76	0.15	7.87	-13.49	-0.15	0.82	
prob	0.000	0.000	0.000	0.011	0.000	0.141	0.000	0.000	0.077	0.000	0.000	0.000	0.240	R <sup>2</sup> = 0.4612
	XPLC	YR	ACOU <sup>2</sup>	DNC <sup>2</sup>	DUR <sup>2</sup>	INSTR <sup>2</sup>	LIV <sup>2</sup>	LOUD <sup>2</sup>	NRG <sup>2</sup>	SPCH <sup>2</sup>	TMP <sup>2</sup>	VAL <sup>2</sup>	YR <sup>2</sup>	F-stat = 5973.447
pop	11.63	38.05	-14.65	-0.45	-0.006	15.61	0.53	0.04	-15.07	1.34	0	-2.13	-0.009	
prob	0.000	0.000	0.000	0.689	0.000	0.000	0.514	0.000	0.000	0.174	0.000	0.000	0.000	

Figure: Estimation output for Equation 2

- slight improvement when removing DNC<sup>2</sup> and LIV<sup>2</sup>.



# Further analysis

- Sub-sample based analysis
- Popularity based on the time of the year

Autumn-Winter									Spring-Summer								
NRG	NRG <sup>2</sup>	DNC	VAL	VAL <sup>2</sup>	ACOU	ACOU <sup>2</sup>	LOUD	LOUD <sup>2</sup>	NRG	NRG <sup>2</sup>	DNC	VAL	VAL <sup>2</sup>	ACOU	ACOU <sup>2</sup>	LOUD	LOUD <sup>2</sup>
4.07	-15.81	-2.01	1.19	1.19	10.61	-16.66	2.02	-9.98	21.59	-30.32	10.19	0.31	-4.84	3.34	-4.69	3.13	0.09

Figure: Model based on seasons

- Out-of-sample prediction and best possible combinations
- Socio-cultural factors, language of track, artist familiarity