

Data science- The first report-Spotify dataset

Sadaf Fatollahy

۶ آبان ۱۴۰۲



نام استاد: دکتر خردپیشه
نام درس : مبانی علوم داده

فهرست مطالب

۳ مقدمه	۱۰۰
۳ درباره دیتاست	۲۰۰
۴ EDA and pre processing	۳۰۰
۶ سوال	۴۰۰

۱.۰ مقدمه

آهنگ های محبوبی که از سال ۱۹۶۰ تا به امروز بر صحنه موسیقی مسلط بوده اند. این مجموعه داده بر اساس رتبه بندی های ARIA (انجمن صنعت ضبط استرالیا) و نمودارهای بیلپورد تنظیم شده است و از نمایش متنوع آهنگ هایی که به موفقیت تجاری و اهمیت فرهنگی بی نظیری دست یافته اند، اطمینان حاصل می کند. مجموعه داده شامل ژانرهای مختلف موسیقی است و تکامل روندهای موسیقی را در طول سال ها به نمایش می گذارد و بینش های ارزشمندی را در مورد چشم انداز همیشه در حال تغییر موسیقی عامه پسند ارائه می دهد. این شامل آهنگ هایی از هنرمندان و گروه های نمادین است که ترکیبی از آثار کلاسیک جاودانه و آثار معاصر را نشان می دهد که تأثیری ماندگار بر دوستداران موسیقی در سراسر جهان گذاشته است. در این تمرین ما قصد داریم آزمون های آماری را بر روی دو دیتاست بررسی کنیم دیتاست دوم مربوط به پیش بینی محبوبیت آهنگ است. در این دیتاست ابتدا مقداری فرایند های EDA و Preprocessing را انجام می دهیم ، بعد از آن تعدادی سوال را خودمان مطرح می کنیم و با آزمون های آماری به آنها پاسخ می دهیم.

۲.۰ درباره دیتاست

هر فیچر نشان دهنده چیست؟

• Popularity: محبوبیت این متغیر هدفی است که می خواهید پیش بینی کنید.

• Track URL: لینک ورود به آهنگ

• Track Name: اسم آهنگ

• Artist URL(s): لینک ورود به پروفایل خواننده

• Artist Name(s): نام خواننده

• Album URL: لینک ورود به آلبوم

• Album Name: اسم آلبوم

• Album Artist URI(s): لینک ورود به خواننده آلبوم

• Album Artist Name(s): اسم خواننده آلبوم

• Album Release Date: تاریخ انتشار آلبوم

• Album Image URL: لینک عکس آلبوم

• Disc Number: شماره دیسک

• Track Number: شماره ترک

• Track Duration (ms)1: مدت زمان ترک

• Track Preview URL2: لینک پیشنمای (دمو) ترک

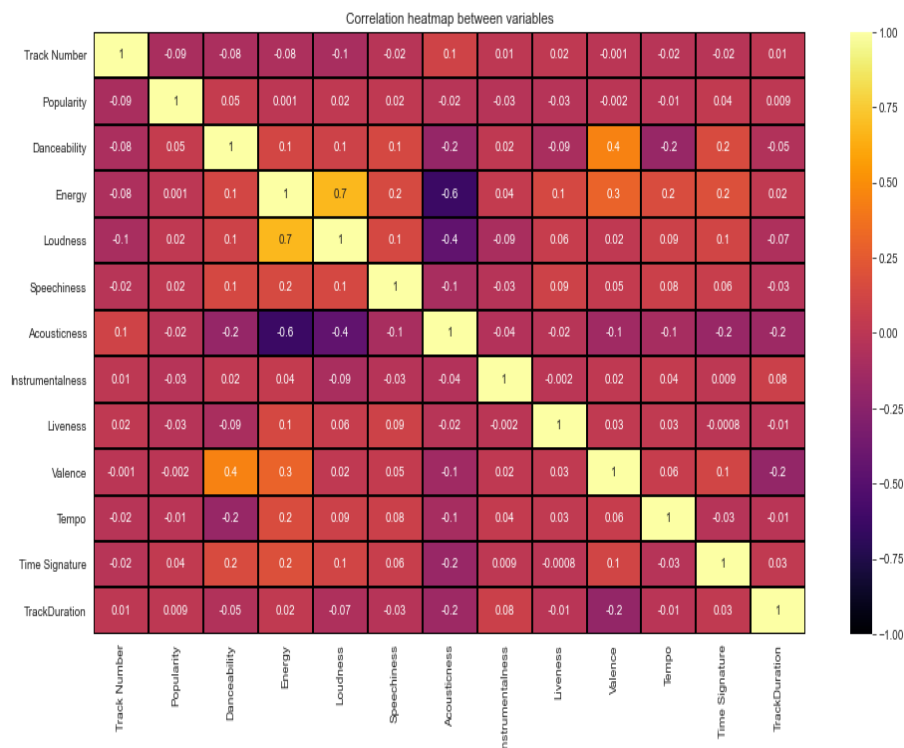
• Explicit: صراحت

- ISRC
- Added By
- Added By : تاریخ و زمان انتشار
- Artist Genres : ژانر خواننده
- Danceability : میزان رقص پذیری
- Energy : میزان انرژی
- Key
- Loudness : بلندی صدا
- Mode
- Speechiness : میزان گفتاری
- Acousticness : میزان اکوستیک بودن
- Instrumentalness : میزان ابزاری بودن
- Liveness : میزان سرزندگی
- Valence : میزان ظرفیت
- Tempo : میزان تمپو
- Time Signature : یک نت موسیقی است که سازماندهی ضربات را در یک قطعه موسیقی نشان می دهد.
- Album Genres : ژانر آلبوم
- Label1 : لیبل آهنگ
- Copyrights

۳.۰ EDA and pre processing

- این دیتاست شامل ۹۹۹۹ سطر و ۳۵ ستون فیچر است. در ابتدا تعدادی از فیچر ها که به آنها نیازی نداریم را حذف میکنیم :
- ["Track URI", "Artist URI(s)", "Album URI", "Album Artist URI(s)", "Album Image URL", "Added By", "Added At", "Key", "Mode", "Explicit", "Album Genres", "ISRC", "Disc Number", "Track Preview URL"]
- حال میخواهیم missing value ها را بررسی کنیم. وقتی تعداد آنها را محاسبه کردیم متوجه شدیم مقدار زیادی از دیتا را تشکیل نمیدهد بنابراین آنها را پاک کردیم به این ترتیب مقدار زیادی دیتا از دست ندادیم.

- فیچر duration را که برحسب میلی ثانیه بود برای راحتی به ثانیه تبدیل کردیم.
- سپس بررسی کردیم که آیا داده تکراری وجود دارد یا خیر که متوجه شدیم ۴۸ مورد دیتا تکراری داریم و آنها را پاک کردیم
- با نگاه کردن به ستون "تاریخ انتشار آلبوم" باید همانطور که از نامش می گوید تاریخ باشد اما در حال حاضر از نوع Object است. به همین دلیل ان را به تاریخ تبدیل کردیم.
- سپس بین فیچر های عددی heatmap رسم کردیم.

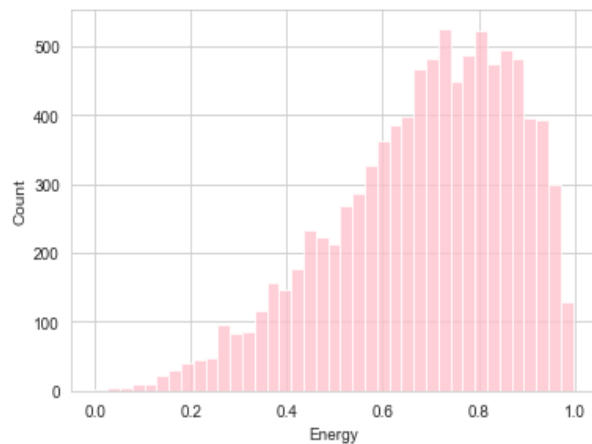


مشاهده میشود بین فیچر های energy و loudness همبستگی بالایی وجود دارد و بین فیچر های energy و Acousticness همبستگی منفی وجود دارد و این بسیار طبیعی است. وقتی صدای آهنگ بلند باشد انرژی بالاتری را منتقل میکند. همبستگی منفی بین آکوستیک و انرژی نشان می دهد که آهنگ های با آکوستیک بالا تمایل به سطوح انرژی پایین تری دارند و آهنگ های با آکوستیک پایین تمایل به سطوح انرژی بالاتری دارند. این نشان می دهد که با آکوستیک تر شدن مسیر (کمتر الکترونیکی) در طبیعت، سطح انرژی آن کاهش می یابد.

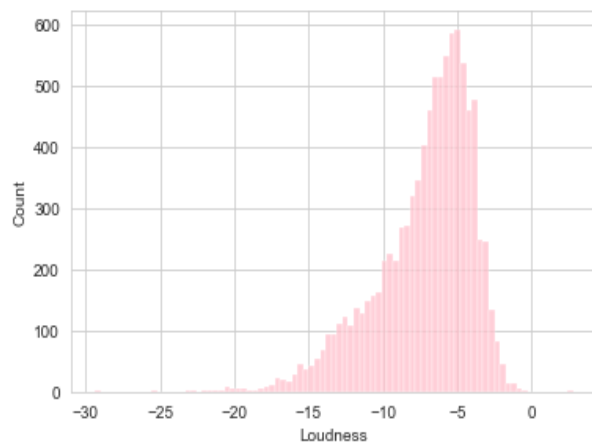
۴.۰ سوال

:

۱. آیا بین انرژی و بلندی صدا رابطه معناداری وجود دارد؟
در ابتدا نگاهی به این توزیع این فیچرها می اندازیم



شکل ۱: توزیع Energy

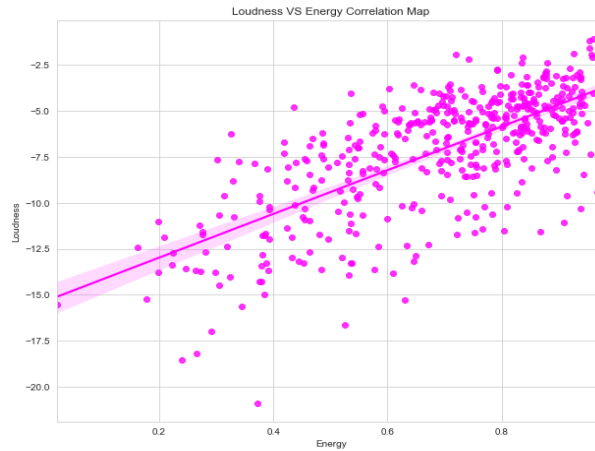


شکل ۲: توزیع Loudness

همانطور که مشخص است این دو فیچر دارای چولگی هستند. برای اطمینان بیشتر از آزمون shapiro استفاده کردیم.

بررسی نرمال بودن دوتا فیچر ذکر شده با استفاده از آزمون shapiro:

- آزمون shapiro
آزمایش می کند که آیا یک نمونه داده توزیع گاوسی دارد یا خیر.
مفروضات:
مشاهدات در هر نمونه مستقل و به طور یکسان توزیع شده است (iid).
تفسیر:
 H_0 : نمونه دارای توزیع گاوسی است.
 H_1 : نمونه توزیع گاوسی ندارد.
نتیجه این آزمون به این صورت بود که هیچ کدام داری توزیع نرمال نیستند.
حال پلات زیر را رسم کردیم تا ببینیم رابطه خطی بین آنها موجود هست یا نه.



شکل ۳: relation between energy and loudness

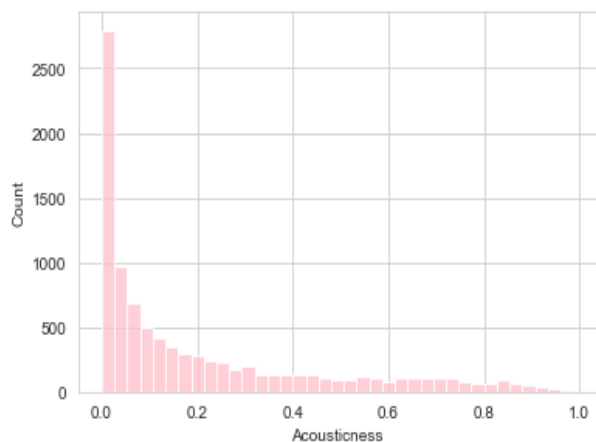
با توجه به شکل این ارتباط وجود دارد اما آیا رابطه معناداری بین آنها از طریق آزمون فرض اثبات میشود؟

به این منظور از آزمون Spearman استفاده میکنیم. این آزمون زمانی استفاده میشود که دو متغیر مورد بررسی عددی باشند و توزیع نرمال نداشته باشند.

- آزمون Spearman
آزمایش می کند که آیا یک نمونه داده به صورت یکپارچه قابل تفکیک است یا خیر.
مفروضات:
الف) مشاهدات در هر نمونه مستقل هستند و به طور یکسان توزیع می شوند.
ب) مشاهدات در هر نمونه رتبه بندی می شوند.
تفسیر:
 H_0 : نمونه ها همبستگی دارند.

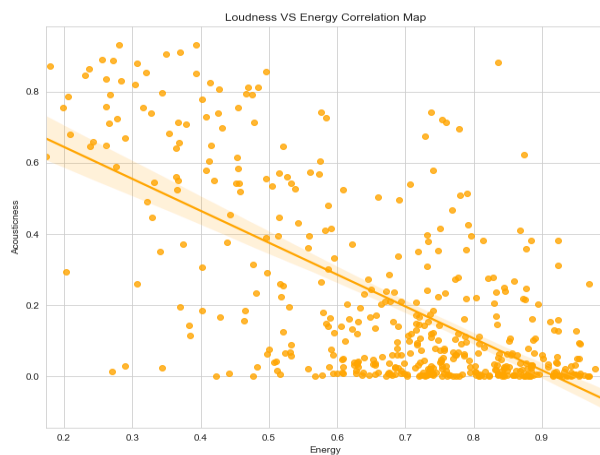
H_1 : نمونه هیچ همبستگی ندارد.
این آزمون مشخص کرد که ارتباط بین این دو فیچر وجود دارد.

۲. آیا بین انرژی و اکوستیک بودن رابطه معناداری وجود دارد؟
به این منظور همانند قبل ابتدا توزیع فیچر Acoustiness را بررسی کردیم که با توجه به شکل نرمال نیست و آزمون فرض هم همین را تایید میکند.



شکل ۴: توزیع Loudness

سپس پلات زیر را رسم کردیم تا ببینیم رابطه خطی بین آنها موجود هست یا نه.



شکل ۵: relation between energy and acoustiness

با توجه به شکل این ارتباط وجود دارد و رابطه عکس وجود دارد همانطور که پیش تر دلیل ان گفته شد و ازمون Spearman نیز ان را اثبات میکند.

۳. آیا بین محبوبیت و نام خواننده رابطه معناداری وجود دارد؟
به دلیل تنوع زیاد نام خواننده ها امکان ترسیم وجود ندارد به همین دلیل به ازمون فرض اکتفا میکنیم.

• ازمون kruskal

این ازمون ارزیابی می کند که آیا بین محبوبیت و نام خواننده رابطه معناداری وجود دارد بدون اینکه توزیع خاصی را فرض کنیم.

مفروضات:

الف) مشاهدات هر نمونه داده مستقل و توزیع می شود.

ب) مشاهدات را می توان رتبه بندی کرد.

تفسیر:

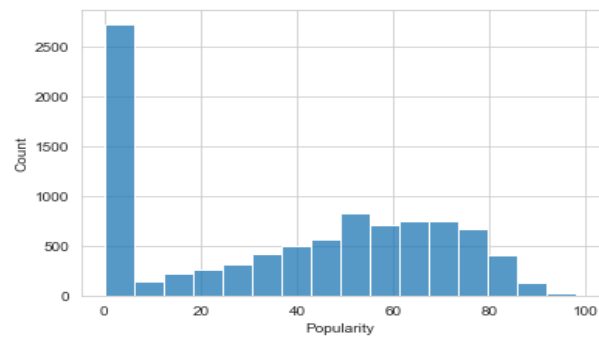
H_0 : مرتبط هستند.

H_1 : مرتبط نیستند..

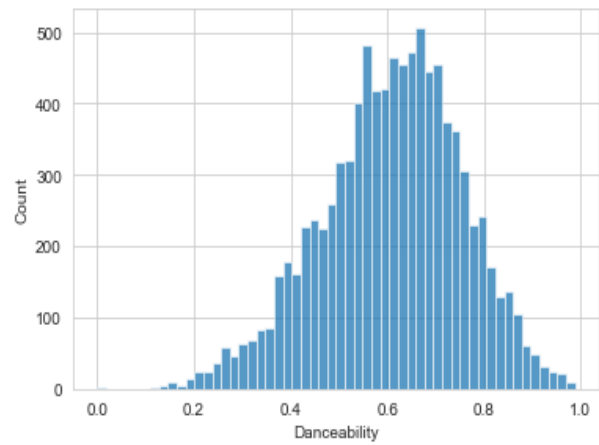
این تست اثبات کرد که بین این دو فیچر ارتباط وجود دارد.

۴. آیا بین سرزندگی اهنگ و نام خواننده رابطه معناداری وجود دارد؟
انتظار داریم که این رابطه وجود داشته باشد زیرا تعداد از خوانندگان فقط اهنگ های شاد و بعضی بیشتر اهنگ های غمگین میخوانند . از ازمون kruskal استفاده کردیم و پاسخ ما را اثبات کرد.

۵. آیا بین محبوبیت و و قابل رقص بودن اهنگ رابطه معناداری وجود دارد؟
در ابتدا نگاهی به این توزیع این فیچرها می اندازیم
از شکل مشخص است که توزیع نرمال نیست . حال بررسی میکنیم ایا رابطه خطی بین انها وجود دارد یا خیر

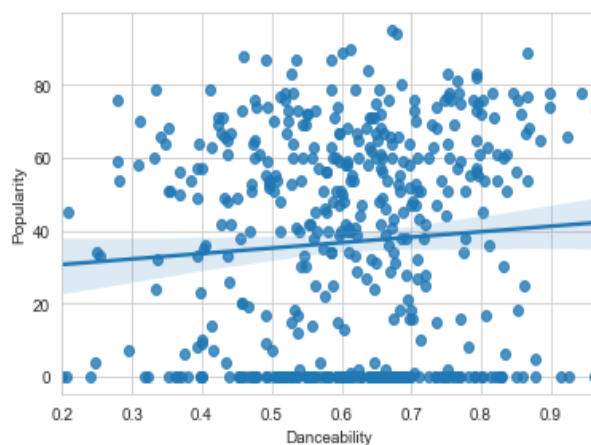


شکل ۶: توزیع popularity



شکل ۷: توزیع danceability

با توجه به شکل رابطه خیلی قوی مشخص نیست . حال با آزمون Spearman انرا بررسی میکنیم .
این آزمون بیان کرد که رابطه بین آنها وجود دارد اما همانطور که مشخص است قوی نیست.



شکل ۸: relation between danceability and popularity

۶. آیا بین نام خواننده و و ژانر اهنگ رابطه معناداری وجود دارد؟
 باتوجه به اینکه دو فیچر کتگوریکال هستند باید از آزمون کای دو استفاده کنیم و به دلیل تنوع اهنگ ها و ژانرها نمیتوانیم ان ها را رسم کنیم.

• آزمون Chi-Squared

این آزمون ارزیابی می کند که آیا بین نام خواننده و و ژانر اهنگ رابطه معناداری وجود دارد.

مفروضات:

الف) مشاهدات استفاده شده در جدول مستقل هستند.

ب) بیش از ۲۵ نمونه در جدول وجود دارد. تفسیر:

H_0 : مرتبط هستند.

H_1 : مرتبط نیستند..

این تست اثبات کرد که بین این دو فیچر ارتباط وجود دارد.

در نهایت شکل زیر نشان میدهد که در سال های اخیر مدت زمان اهنگ ها کم شده است.

