

# Yandex.Music

## Contents

- [Introducción](#)
- [Etapa 1. Descripción de los datos](#)
  - [Conclusiones](#)
- [Etapa 2. Datos preprocessing](#)
  - [2.1 Estilo del encabezado](#)
  - [2.2 Valores ausentes](#)
  - [2.3 Duplicados](#)
  - [2.4 Conclusiones](#)
- [Etapa 3. Prueba de hipótesis](#)
  - [3.1 Hipótesis 1: comparar el comportamiento del usuario en las dos ciudades](#)
  - [3.2 Hipótesis 2: música al principio y al final de la semana](#)
  - [3.3 Hipótesis 3: preferencias de género en Springfield y Shelbyville](#)
  - [Conclusiones](#)

## Introducción

En este proyecto, compararemos las preferencias musicales de las ciudades de Springfield y Shelbyville. Probaremos hipótesis de Yandex.Music para comparar el comportamiento del usuario de esas dos ciudades.

### Objetivo:

Prueba tres hipótesis:

- 1. La actividad de los usuarios difiere según el día de la semana y dependiendo de la ciudad.
- 2. Los lunes por la mañana, los habitantes de Springfield y Shelbyville escuchan diferentes géneros. Lo mismo ocurre con los viernes por la noche.
- 3. Los oyentes de Springfield y Shelbyville tienen preferencias distintas. En Springfield prefieren el pop mientras que en Shelbyville hay más aficionados al rap.

### Etapas

Primero, evaluaremos la calidad de los datos y veremos si los problemas son significativos. Entonces, durante el preprocesamiento de datos, tomaremos en cuenta los problemas más críticos.

Este proyecto consiste en tres etapas:

1. Descripción de los datos
2. Preprocesamiento de datos
3. Prueba de hipótesis

#### Volver a Contenidos

### Etapla 1. Descripción de los datos

```
In [1]: # importando pandas
import pandas as pd

In [2]: # leyendo el archivo y almacenándolo en df
df = pd.read_csv('datasets/music_project_en.csv')
print(df)

   user_id  track  artist  genre  City  time  Day
0  F8692EC  Delayed Because of Accident  The Mass Missle  rock  Springfield  20:28:33  Wednesday
1  5520438  Delayed Because of Accident  Andreas Rönberg  rock  Springfield  14:07:09  Friday
2  28E3C8  Funiculi funiculà  Mario Lanza  pop  Springfield  20:58:07  Wednesday
3  A300K9C3  Dragons in the Sunset  Fire + Ice  funk  Shelbyville  08:37:49  Monday
4  E20C3FAE  Soul People  Space Echo  disco  Springfield  08:34:34  Monday
...
60974  729C8B89  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
60975  D88D4A55  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
60976  C5E3A405  4-Pop!  Karadenziz  Karaoke  Karaoke  Karaoke  Springfield  20:47:49  Wednesday
60977  321D9506  Freight Train  Chas McDevitt  nusrag  Springfield  21:09:41  Friday
60978  3A464F84  Tell Me Sweet Little Lies  Monica Lopez  dance  Springfield  08:34:34  Monday

   genre  City  time  Day
0  rock  Shelbyville  20:28:33  Wednesday
1  rock  Springfield  14:07:09  Friday
2  pop  Shelbyville  20:58:07  Wednesday
3  funk  Shelbyville  08:37:09  Monday
4  dance  Springfield  08:34:34  Monday
...
60974  rnb  Springfield  13:32:28  Wednesday
60975  rnb  Springfield  13:32:28  Wednesday
60976  industrial  Springfield  20:09:26  Friday
60977  rock  Springfield  21:43:59  Friday
60978  country  Springfield  21:09:46  Friday

[65879 rows x 7 columns]

In [3]: # obteniendo las 10 primeras filas de la tabla df
df.head(10)
```

```
Out[3]:   user_id  track  artist  genre  City  time  Day
0  F8692EC  Kungliga To Boots  The Mass Missle  rock  Springfield  20:28:33  Wednesday
1  5520438  Delayed Because of Accident  Andreas Rönberg  rock  Springfield  14:07:09  Friday
2  28E3C8  Funiculi funiculà  Mario Lanza  pop  Springfield  20:58:07  Wednesday
3  A300K9C3  Dragons in the Sunset  Fire + Ice  funk  Shelbyville  08:37:49  Monday
4  E20C3FAE  Soul People  Space Echo  disco  Springfield  08:34:34  Monday
...
60974  729C8B89  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
60975  D88D4A55  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
60976  C5E3A405  4-Pop!  Karadenziz  Karaoke  Karaoke  Karaoke  Springfield  20:47:49  Wednesday
60977  321D9506  Freight Train  Chas McDevitt  nusrag  Springfield  21:09:41  Friday
60978  3A464F84  Tell Me Sweet Little Lies  Monica Lopez  dance  Springfield  08:34:34  Monday
...
65879  rows x 7 columns
```

```
In [4]: # obteniendo información general sobre los datos en df
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 65879 entries, 0 to 65878
Data columns (total 7 columns):
 #   Column  Non-Null Count  Dtype  
---  -
 0   user_id  65879 non-null    object  
 1   track    65879 non-null    object  
 2   artist   65879 non-null    object  
 3   genre    65879 non-null    object  
 4   city     65879 non-null    object  
 5   time     65879 non-null    object  
 6   day      65879 non-null    object  
dtypes: object(7)
memory usage: 3.3+ MB
```

Cada fila de la tabla almacena datos de la pista que fue reproducida. Algunas columnas describen la pista en sí: su título, el artista y el género. El resto transmite la información del usuario: la ciudad de la que viene, el tiempo que ha reproducido la pista.

Está claro que los datos son suficientes para probar la hipótesis. Sin embargo, hay valores ausentes.

Para continuar, necesitamos preprocesar los datos.

#### Volver a Contenidos

### Etapla 2. Preprocesamiento de datos

Será corregido el formato en los encabezados de las columnas y los valores ausentes.

#### Estilo del encabezado

```
In [5]: # la lista de los nombres de las columnas en la tabla df
df.columns

Index(['user_id', 'track', 'artist', 'genre', 'City', 'time', 'Day'], dtype='object')

Utilefremos en un criterio de sólo minúsculas y sin espacios los encabezados de las columnas
```

```
In [6]: # renombra las columnas

df.rename(
    columns={
        'user_id':'user_id',
        'track':'track',
        'artist':'artist',
        'genre':'genre',
        'City':'city',
        'time':'time',
        'Day':'day'
    })
```

```
In [7]: # comprobando el resultado: la lista de los nombres de las columnas
df.columns

Out[7]: Index(['user_id', 'track', 'artist', 'genre', 'city', 'time', 'day'], dtype='object')
```

#### Volver a Contenidos

### Valores ausentes

```
In [8]: # calculando valores ausentes
print(df.isna().sum())

user_id      0
track        0
artist       0
genre        0
city         0
time         0
day          0
dtypes: int64 7

No todos los valores ausentes afectan directamente los objetivos del proyecto en el caso de nuestros valores ausentes serán reemplazados con "unknown" dado que el reemplazo de estos no afecta o cambia el resultado final
```

```
In [9]: # recorriendo los nombres de las columnas y reemplazando los valores ausentes con 'unknown'
columns = list(df.columns)

for column in columns:
    df[column].fillna("unknown")
print(column)
```

```
In [10]: # contando valores ausentes
print(df.isna().sum())

user_id      0
track        0
artist       0
genre        0
city         0
time         0
day          0
dtypes: int64 7
```

#### Volver a Contenidos

### Duplicados

Buscaremos el total de duplicados obvios en la tabla.

```
In [11]: # contando duplicados obvios
df.drop_duplicates().sum()

3828

In [12]: # eliminando duplicados obvios
df.drop_duplicates().reset_index(drop=True)
```

```
In [13]: # comprobando duplicados
print(df.drop_duplicates().sum())

0
```

Ahora eliminaremos los duplicados implícitos en la columna genre. Dado que el nombre de un género se puede escribir de varias formas. Dichos errores también pueden afectar a resultado.

```
In [14]: # inspeccionando los nombres de géneros únicos
unique_genres = df['genre'].unique()
print(sorted(unique_genres))

['acbi', 'acoustic', 'action', 'adult', 'africa', 'afrikaans', 'alternative', 'ambient', 'americana', 'animated', 'anime', 'arabesk', 'arabic', 'arena', 'argentinatang
o', 'art', 'audiobook', 'avantgarde', 'axe', 'baila', 'balkan', 'beats', 'bigroom', 'black', 'bluegrass', 'blues', 'bollywood', 'bossa', 'brazilian', 'breakbeat', 'brea
k', 'broadway', 'cantautor', 'cantopop', 'canzone', 'caribbean', 'caucasian', 'celtic', 'chamber', 'children', 'choral', 'christian', 'christmas', 'classical', 'classica
l', 'club', 'columbian', 'comedy', 'conjaz', 'contemporary', 'country', 'cuban', 'dance', 'dancehall', 'dancepop', 'dark', 'death', 'deep', 'd
e', 'horror', 'house', 'idm', 'independent', 'indian', 'india', 'indipop', 'industrial', 'inspirational', 'instrumental', 'international', 'irish', 'japn', 'japanese', 'jazz', 'j
ewish', 'jpop', 'jungle', 'k-pop', 'karadenziz', 'karaoke', 'kayokyoku', 'korean', 'lalo', 'latin', 'latin', 'leftfield', 'local', 'lounge', 'loungellectronic
c', 'lovers', 'malaysian', 'mandpop', 'marschmusik', 'meditative', 'mediterranean', 'melodic', 'metal', 'metalcore', 'mexican', 'middle', 'minimal', 'miscellaneous', 'othe
r', 'piano', 'pop', 'popelctronic', 'popuoradance', 'post', 'posthardcore', 'postrock', 'power', 'prometal', 'progressive', 'psychdelic', 'punjabi', 'punk', 'quebec', 'quebec
e', 'r', 'roots', 'ruspop', 'rusrap', 'rurrock', 'salsa', 'samba', 'schlager', 'self', 'sertanejo', 'sheezing', 'showtunes', 'singer', 'ska', 'slow', 'smooth', 'soul', 'iso
ster', 'southern', 'specialty', 'speech', 'spiritual', 'sport', 'stonerrock', 'surf', 'swing', 'synthpop', 'sängerportrait', 'tango', 'tanzone
ster', 'tararar', 'tech', 'techno', 'thrash', 'top', 'traditional', 'tradjazz', 'trance', 'tribal', 'trip', 'triphop', 'tropical', 'türk', 'türkce', 'unknown', 'urban',
uzbek', 'variété', 'vsl', 'videogame', 'vocal', 'western', 'world', 'worldbeat', 'yill']
```

Tenemos duplicados implícitos:

- hip
- hip
- hip-hop

```
In [15]: # función para reemplazar duplicados implícitos

def replace_wrong_genres (wrong_genres, correct_genre):
    for wrong_genre in wrong_genres:
        df[genre] = df[genre].replace(wrong_genre, correct_genre)
```

```
In [16]: # eliminando duplicados implícitos
duplicates = ['hip', 'hop', 'hip-hop']
name = 'hiphop'
replace_wrong_genres(duplicates, name)
```

Confirmamos haber eliminado los duplicados implícitos

```
In [17]: # revisando en busca de duplicados implícitos
unique_genres = df['genre'].unique()
print(sorted(unique_genres))

['acbi', 'acoustic', 'action', 'adult', 'africa', 'afrikaans', 'alternative', 'ambient', 'americana', 'animated', 'anime', 'arabesk', 'arabic', 'arena', 'argentinatang
o', 'art', 'audiobook', 'avantgarde', 'axe', 'baila', 'balkan', 'beats', 'bigroom', 'black', 'bluegrass', 'blues', 'bollywood', 'bossa', 'brazilian', 'breakbeat', 'brea
k', 'broadway', 'cantautor', 'cantopop', 'canzone', 'caribbean', 'caucasian', 'celtic', 'chamber', 'children', 'choral', 'christian', 'christmas', 'classical', 'classica
l', 'club', 'columbian', 'comedy', 'conjaz', 'contemporary', 'country', 'cuban', 'dance', 'dancehall', 'dancepop', 'dark', 'death', 'deep', 'd
eutschrock', 'deutschrap', 'dirty', 'disco', 'dub', 'documentary', 'downbeat', 'downtempo', 'drum', 'dub', 'dubstep', 'eastern', 'easy', 'electronic', 'electropop', 'en
g', 'ethno', 'epicmetal', 'estrada', 'ethnic', 'eurofolk', 'european', 'experimental', 'extrememetal', 'fado', 'film', 'fitness', 'flamenco', 'folk', 'folklore', 'fo
k', 'folkrock', 'folktronica', 'forro', 'frankisch', 'französisch', 'french', 'funk', 'future', 'gangsta', 'garage', 'german', 'ghazal', 'glitarre', 'glitch', 'gos
pel', 'gothic', 'grime', 'grunge', 'gyppo', 'hardcup', 'hard'n'heavy', 'hardcore', 'hardstyle', 'hardtechno', 'hip', 'hip-hop', 'hiphop', 'historisch', 'hollyday', 'ho
p', 'horror', 'house', 'idm', 'independent', 'indian', 'india', 'indipop', 'industrial', 'inspirational', 'instrumental', 'international', 'irish', 'japn', 'japanese', 'jazz', 'jewish', 'jpop', 'j
ungle', 'k-pop', 'karadenziz', 'karaoke', 'kayokyoku', 'korean', 'lalo', 'latin', 'latin', 'leftfield', 'local', 'lounge', 'loungellectronic', 'lovers', 'malaysian', 'ma
r', 'marschmusik', 'meditative', 'mediterranean', 'melodic', 'metal', 'metalcore', 'mexican', 'middle', 'minimal', 'miscellaneous', 'modern', 'moody', 'pop', 'p
p', 'pop', 'popelctronic', 'popuoradance', 'post', 'posthardcore', 'postrock', 'power', 'prometal', 'progressive', 'psychdelic', 'punjabi', 'punk', 'quebec', 'quebecis', 'r
ap', 'rave', 'reggae', 'reggaeton', 'regional', 'relax', 'religious', 'retro', 'rhythm', 'rnb', 'rnr', 'rock', 'rockabilly', 'romance', 'roots', 'ruspop', 'rus
r', 'roots', 'ruspop', 'rusrap', 'rurrock', 'salsa', 'samba', 'schlager', 'self', 'sertanejo', 'sheezing', 'showtunes', 'singer', 'ska', 'slow', 'smooth', 'soul', 'sou
l', 'southern', 'specialty', 'speech', 'spiritual', 'sport', 'stonerrock', 'surf', 'swing', 'synthpop', 'sängerportrait', 'tango', 'tanzone
ster', 'tararar', 'tech', 'techno', 'thrash', 'top', 'traditional', 'tradjazz', 'trance', 'tribal', 'trip', 'triphop', 'tropical', 'türk', 'türkce', 'unknown', 'urban',
uzbek', 'variété', 'vsl', 'videogame', 'vocal', 'western', 'world', 'worldbeat', 'yill']
```

#### Volver a Contenidos

### Conclusiones

Detectamos tres problemas con los datos:

- Estilos de encabezados incorrectos
- Valores ausentes
- Duplicados obvios e implícitos

Los encabezados han sido eliminados para conseguir que el procesamiento de la tabla sea más sencillo.

Todos los valores ausentes han sido reemplazados por "Unknown". Pero todavía tenemos que ver si los valores ausentes en "genre" afectan a nuestros cálculos.

#### Volver a Contenidos

### Etapla 3. Prueba de hipótesis

#### Hipótesis 1: comparar el comportamiento del usuario en las dos ciudades

De acuerdo con la primera hipótesis, los usuarios de Springfield y Shelbyville escuchan música de forma distinta. Comprueba esto utilizando los datos de tres días de la semana: lunes, miércoles y viernes.

- Divide a los usuarios en grupos por ciudad.
- Compara cuántas pistas reproducida cada grupo el lunes, el miércoles y el viernes.

```
In [18]: # contando las pistas reproducidas en cada ciudad
df.groupby('city')['track'].count()

Out[18]: city
Shelbyville    1852
Springfield    4711
Name: track, dtype: int64

Springfield ha reproducido más pistas que Shelbyville. Pero eso no implica que los ciudadanos de Springfield escuchan música más a menudo. Esta ciudad es simplemente más grande y hay más usuarios.

Ahora agrupamos los datos por día de la semana y encuentra el número de pistas reproducidas el lunes, miércoles y viernes.
```

```
In [19]: # contando las pistas reproducidas en cada uno de los tres días
df.groupby('day')['track'].count()

Out[19]: day
Friday    21849
Monday    21364
Wednesday 18059
Name: track, dtype: object

El miércoles fue el día más silencioso de todos. Pero si consideramos las dos ciudades por separado podríamos llegar a una conclusión diferente.
```

```
In [20]: # creando la función number_tracks()
# declaramos la función con dos parámetros: days, city=

# deja que la variable track_list almacene las filas df en las que
# el valor en la columna 'day' es igual al parámetro day y, al mismo tiempo,
# el valor de la columna 'city' es igual al parámetro city= (aplica el filtrado consecutivo
# con un índice único).
# deja que la variable track_list_count almacene el número de valores de la columna 'user_id' en track_list
# (encuentra con el método count()).
# genera que la función devuelva un número: el valor de track_list_count.
# la función cuenta las pistas reproducidas en un cierto día y ciudad.
# después recupera las filas del día deseado de la tabla.
# después filtra las filas de la ciudad deseada del resultado,
# entonces, encuentra el número de valores de 'user_id' en la tabla filtrada, b
# devuelve ese número.
# para ver lo que devuelve, envuelve la llamada de la función en print().

def number_tracks(day,city):
    track_list=df[(df['day']==day)&(df['city']==city)]
    track_list_count=track_list['user_id'].count()
    return track_list_count
```

Utilizamos number\_tracks() seis veces, para obtener los datos de las ciudades en cada día estudiado.

```
In [21]: # el número de canciones reproducidas en Springfield el lunes
print(number_tracks('Monday','Springfield'))

15740

In [22]: # el número de canciones reproducidas en Shelbyville el lunes
print(number_tracks('Monday','Shelbyville'))

5614

In [23]: # el número de canciones reproducidas en Springfield el miércoles
print(number_tracks('Wednesday','Springfield'))

11056

In [24]: # el número de canciones reproducidas en Shelbyville el miércoles
print(number_tracks('Wednesday','Shelbyville'))

7903

In [25]: # el número de canciones reproducidas en Springfield el viernes
print(number_tracks('Friday','Springfield'))

15945

In [26]: # el número de canciones reproducidas en Shelbyville el viernes
print(number_tracks('Friday','Shelbyville'))

5895
```

```
In [27]: #Almacenamos los datos recuperados
header=['city','monday','wednesday','friday']
new_data = [['Springfield',15740,11056,15945],
            ['Shelbyville',5614,7903,5895]]

In [28]: # tabla con los resultados
new_table= pd.DataFrame(data=new_data, columns=header)
print(new_table)
```

```
   city  monday  wednesday  friday
0  Springfield    15740    11056    15945
1  Shelbyville     5614     7903     5895

Conclusiones

Los datos revelan las diferencias en el comportamiento de los usuarios:
```

- En Springfield, el número de canciones reproducidas alcanzan el punto máximo los lunes y viernes mientras que los miércoles hay un descenso de la actividad.
- En Shelbyville, al contrario, los usuarios escuchan más música los miércoles. La actividad de los usuarios los lunes y viernes es menor.

Así que la primera hipótesis parece ser correcta.

#### Volver a Contenidos

#### Hipótesis 2: música al principio y al final de la semana

De acuerdo con la segunda hipótesis, los lunes por la mañana y los viernes por la noche los ciudadanos de Springfield escuchan géneros que difieren de aquellos que los usuarios de Shelbyville disfrutan.

```
In [29]: # obteniendo la tabla spr_general de las filas de df,
# donde los valores en la columna 'city' es 'Springfield'
spr_general= df[(df['city']=='Springfield')]
print(spr_general)
```

```
   user_id  track  artist  genre  City  time  Day
0  F8692EC  Delayed Because of Accident  Soul People  Andreas Rönberg  Springfield  14:07:09  Friday
1  E20C3FAE  Soul People  Space Echo  disco  Springfield  08:34:34  Monday
2  28E3C8  Funiculi funiculà  Mario Lanza  pop  Springfield  20:58:07  Wednesday
3  A300K9C3  Dragons in the Sunset  Fire + Ice  funk  Springfield  08:37:49  Monday
4  E20C3FAE  Soul People  Space Echo  disco  Springfield  08:34:34  Monday
...
61247  83A4717  I Worship Only What You Ble...  The Black Dahlia Murd...  Springfield  21:08:59  Friday
61248  729C8B89  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
61249  3E3A405  4-Pop!  Karadenziz  Karaoke  Karaoke  Springfield  20:47:49  Wednesday
61251  321D9506  Freight Train  Chas McDevitt  nusrag  Springfield  21:09:41  Friday
61252  3A464F84  Tell Me Sweet Little Lies  Monica Lopez  dance  Springfield  08:34:34  Monday

   genre  City  time  Day
0  rock  Springfield  14:07:09  Friday
1  dance  Springfield  08:34:34  Monday
2  dance  Springfield  20:58:07  Wednesday
3  funk  Springfield  08:37:09  Monday
4  dance  Springfield  08:34:34  Monday
...
61247  extreme  Springfield  21:08:59  Friday
61248  rnb  Springfield  13:32:28  Wednesday
61249  industrial  Springfield  20:09:26  Friday
61250  rock  Springfield  21:43:59  Friday
61251  rnsop  Springfield  21:09:41  Friday
61252  country  Springfield  21:09:46  Friday

[42741 rows x 7 columns]
```

```
In [30]: # obteniendo spr_general de las filas df,
# donde el valor de la columna 'city' es 'Shelbyville'
shel_general= df[(df['city']=='Shelbyville')]
print(shel_general)
```

```
   user_id  track  artist  genre  City  time  Day
0  F8692EC  Kungliga To Boots  The Mass Missle  rock  Springfield  20:28:33  Wednesday
1  5520438  Delayed Because of Accident  Soul People  Andreas Rönberg  Springfield  14:07:09  Friday
2  28E3C8  Funiculi funiculà  Mario Lanza  pop  Springfield  20:58:07  Wednesday
3  A300K9C3  Dragons in the Sunset  Fire + Ice  funk  Springfield  08:37:49  Monday
4  842E29A1  Chains  Obladaet  unknown  Springfield  08:37:49  Monday
...
61239  D94FD5C8  There from the Walking Dead  Proyecto Halloween  Springfield  21:08:59  Friday
61240  B5E85C5F  Red Lips: Rita Bowser Bowser  Bre Petrunko  Springfield  21:08:59  Friday
61241  29E3A405  (Hello) Cloud Mountain  Perunkia Trio  Springfield  21:08:59  Friday
61242  3893D221  (Hello) Cloud Mountain  Perunkia Trio  Springfield  21:08:59  Friday
61249  D88D4A55  Maybe One Day (Feat. Black Spade)  My Name  R'n'B  Springfield  21:08:59  Friday
61252  3A464F84  Tell Me Sweet Little Lies  Monica Lopez  dance  Springfield  08:34:34  Monday

   genre  City  time  Day
0  rock  Springfield  14:07:09  Friday
1  pop  Shelbyville  20:58:07  Wednesday
2  folk  Shelbyville  08:37:09  Monday
3  nusrag  Shelbyville  13:09:41  Friday
4  dance  Shelbyville  21:20:49  Wednesday
...
61239  film  Springfield  21:14:49  Monday
61240  electronic  Springfield  21:08:59  Monday
61241  world  Shelbyville  13:56:09  Monday
61242  postrock  Shelbyville  21:43:59  Monday
61249  hiphop  Springfield  10:09:00  Monday

[18432 rows x 7 columns]
```

```
In [31]: # devolviendo la función genre_weekday() con los parámetros day, time y time2. Debería
# devolver información sobre los géneros más populares de una determinada día y de una determinada hora:

# 1) deja que la variable genre_df almacene las filas que cumplen varias condiciones:
# - el valor de la columna 'day' es igual al valor del argumento day=
# - el valor en la columna 'time' es mayor que el valor del argumento time2=
# - el valor en la columna 'time' es menor que el valor del argumento time2=
# Utiliza un filtrado consecutivo con indexación lógica.

# 2) Agrupa genre_df por la columna 'genre', toma una de sus columnas,
# utiliza el método count() para encontrar el número de entradas por cada uno de
# los géneros representados; almacena los Series resultantes en
# la variable genre_df_sorted

# 3) Ordena genre_df_count en orden descendente de frecuencia y guarda el resultado
# en la variable genre_df_sorted

# 4) Devuelve un objeto Series con los primeros 15 valores de genre_df_sorted ... los 15
# géneros más populares (en un determinado día, en un determinado periodo de tiempo)

# Escribe tu función aquí.
```

```
def genre_weekday(df, day, time2, time22):

    # filtrado consecutivo
    # genre_df solo almacenará aquellas filas df en las que el día sea igual a day=
    genre_df = df[(df['day'] == day)]

    # genre_df solo almacenará aquellas filas df en las que el tiempo sea menor que time2=
    genre_df = genre_df[(genre_df['time'] < time2)]

    # genre_df solo almacenará aquellas filas df en las que el tiempo sea mayor que time22=
    genre_df = genre_df[(genre_df['time'] > time22)]

    # agrupa el DataFrame filtrado por la columna con los nombres de los géneros, toma la columna de género, y encuentra el número de filas por cada género con el método count()
    genre_df_grouped = genre_df.groupby('genre')['genre'].count()

    # ordenamos el resultado en orden descendente (por lo que los géneros más populares aparecerán primero en el objeto Series)
    genre_df_sorted = genre_df_grouped.sort_values(ascending=False)

    # devolvemos el objeto Series que almacena los 15 géneros más populares en un día determinado en un periodo de tiempo determinado
    return genre_df_sorted[0:15]
```

Comparamos los resultados de la función genre\_weekday() para Springfield y Shelbyville el lunes por la mañana (de 7 a 11) y el viernes por la tarde (de 17:00 a 23:00):

```
In [32]: # llamando a la función para el lunes por la mañana en Springfield (utilizando spr_general en vez de la tabla df)
genre_weekday(spr_general, 'Monday', '07:00', '11:00')
```

```
Out[32]: genre
pop      781
dance     549
rock      480
electronic 480
hiphop     474
world      286
respop     185
world      151
rusrap     175
alternative 164
unknown    161
classical   157
metal       120
jazz        98
folk        97
soundtrack  95
Name: genre, dtype: int64
```

```
In [33]: # llamando a la función para el lunes por la mañana en Shelbyville (utilizando shel_general en vez de la tabla df)
genre_weekday(shel_general, 'Monday', '17:00', '23:00')
```

```
Out[33]: Series([], Name: genre, dtype: int64)

In [34]: # llamando a la función para el viernes por la tarde en Springfield
genre_weekday(spr_general, 'Friday', '17:00', '23:00')
```

```
Out[34]: genre
pop      713
rock     517
dance     495
electronic 482
hiphop     479
world      268
respop     163
alternative 163
jazz       111
unknown    110
soundtrack  89
rnb        80
metal       80
Name: genre, dtype: int64

Conclusion

Habiendo comparado los 15 géneros más populares del lunes por la mañana podemos concluir lo siguiente:
```

1. Los usuarios de Springfield y Shelbyville escuchan música similar. Los cinco géneros más populares son los mismos, solo rock y electrónica han intercambiado posiciones.
2. En Springfield el número de valores ausentes resultaron ser tan altos que el valor "unknown". Llegó al décimo. Esto significa que los valores ausentes forman una parte considerable de los datos, lo que podría ser la base de la cuestión sobre la fiabilidad de nuestras conclusiones.

Para el viernes por la tarde, la situación es similar. Los géneros individuales varían algo pero, en general, los 15 más populares son parecidos en las dos ciudades.

De esta forma, la segunda hipótesis ha sido parcialmente demostrada:

- Los usuarios escuchan música similar al principio y al final de la semana.
- No hay una gran diferencia entre Springfield y Shelbyville. En ambas ciudades, el pop es el género más popular.

Si embargo, el número de valores ausentes hace este resultado un tanto cuestionable. En Springfield, hay tantos que afectan a nuestros 15 más populares. De no faltarlos esos valores, las cosas podrían parecer diferentes.

#### Volver a Contenidos

#### Hipótesis 3: preferencias de género en Springfield y Shelbyville

Hipótesis: Springfield ama la música rap. A los ciudadanos de Springfield les gusta más el pop.

```
In [35]: # en una línea, agrupa la tabla spr_general por la columna 'genre',
# cuenta los valores 'genre' con count() en el agrupamiento con count().
# ordena el Series resultante en orden descendente, y almacenado en spr_genres

spr_genres=spr_general.groupby('genre')['genre'].count().sort_values(ascending=False)
spr_genres
```

```
Out[35]: genre
pop      5892
dance     4435
rock      3965
electronic 3795
hiphop     2896
classical   1016
alternative  645
jazz        511
unknown     419
rusrap      364
soundtrack  350
rnb         348
metal       340
Name: genre, dtype: int64
```

```
In [36]: # en una línea, agrupa la tabla shel_general por la columna 'genre',
# cuenta los valores 'genre' con count() en el agrupamiento con count().
# ordena el Series resultante en orden descendente y guárdalo en shel_genres

shel_genres=shel_general.groupby('genre')['genre'].count().sort_values(ascending=False)
shel_genres
```

```
Out[36]: genre
pop      2421
dance     1879
rock      1879
electronic 1758
hiphop     968
classical   646
rusrap      564
soundtrack  515
rnb         515
metal       515
Name: genre, dtype: int64

Conclusion
```

La hipótesis ha sido parcialmente demostrada:

- La música pop es el género más popular en Springfield, tal como se esperaba.
- Sin embargo, la música pop ha resultado ser igual de popular en Springfield que en Shelbyville y el rap no estaba entre los 5 más populares en ninguna de las ciudades.

#### Volver a Contenidos

### Conclusiones

Hemos probado las siguientes tres hipótesis:

1. La actividad de los usuarios difiere dependiendo del día de la semana y de las distintas ciudades.
2. Los lunes por la mañana los residentes de Springfield y Shelbyville escuchan géneros distintos. Lo mismo ocurre con los viernes por la noche.
3. Los oyentes de Springfield y Shelbyville tienen distintas preferencias. En ambas ciudades, Springfield y Shelbyville, se