

# PROJECT 02: PHÂN TÍCH DỮ LIỆU

- Lê Minh Hữu - 19120525
- Ninh Duy Huy - 19120533
- Cao Thanh Kiệt - 19120544
- Nguyễn Tuấn Khoa - 19120547
- Trần Tuấn Kiệt - 19120557

## Phân công công việc

- Lê Minh Hữu: Khai thác dữ liệu, đặt câu hỏi 5
- Nguyễn Tuấn Khoa: đặt câu hỏi 1, 2
- Cao Thanh Kiệt: đặt câu hỏi 3, review code
- Ninh Duy Huy: đặt câu hỏi 6
- Trần Tuấn Kiệt: đặt câu hỏi 7, 8

## Cấu trúc thư mục

- data folder: chứa file dữ liệu lấy từ API soundcloud
- cleaned\_data folder: chứa file dữ liệu sau khi 1 vài bước tiền xử lý cơ bản (xóa các hàng hoặc cột khi mất nhiều giá trị hơn ngưỡng cho trước)
- source.ipynb: file source code

## TÌNH XỬ LÝ DỮ LIỆU

```
In [1]: import matplotlib.pyplot as plt
import pandas as pd
import numpy as np

In [2]: playlists_df = pd.read_csv("data/playlists.csv", sep="," )
tracks_df = pd.read_csv("data/tracks.csv", sep="," )
users_df = pd.read_csv("data/users.csv", sep="," )

# configure matplotlib
plt.rcParams["figure.figsize"] = (15, 10)
font = {'family': 'serif',
        'color': 'blue',
        'size': 20}

plt.rcParams['font.size'] = 20

In [3]: ## thông tin về số dòng và cột của dữ liệu ban đầu
print(playlists_df.shape)
print(tracks_df.shape)
print(users_df.shape)

(608, 3)
(6979, 45)
(422, 33)

In [4]: # delete some row with missing values more than the given ratio
preprocess_data(download_ratio):
    """
    Hàm này sẽ tiền xử lý dữ liệu, loại bỏ một vài dòng hoặc cột khi mất mát dữ liệu quá tỉ lệ cho trước
    """
    attributes = list(old_df.columns)
    row = len(old_df.index)
    col = len(attributes)

    a = []
    arr = []
    for i in range(row):
        for j in range(col):
            for k in range(col):
                arr.append(arr[i][j])
            result_arr.append(temp_a)
            temp_a = []

    # xóa dòng
    for value in arr:
        count_missing = 0
        for i in value:
            if i == 'None' or i==1:
                count_missing += 1
            if (count_missing / col) > ratio:
                result_arr.remove(value)

    # xóa cột
    new_row = len(result_arr)
    temp_arr = []
    temp_a = []
    for i in range(new_row):
        for j in range(col):
            temp_a.append(result_arr[i][j])
            temp_arr.append(temp_a)
            temp_a = []

    count_del_col = 0 #kiểm đếm các cột đã bị xóa
    for i in range(col):
        count_missing = 0
        for j in range(new_row):
            if temp_arr[j][i] != temp_arr[j][i] or temp_arr[j][i] == 'None':
                count_missing += 1
            if (count_missing /new_row) > ratio:
                del attributes[attributes.index(i)]
                for r in result_arr:
                    del r[i-count_del_col]
                count_del_col += 1

    # tạo một DataFrame mới để lưu kết quả
    new_data = []
    new_row = len(result_arr)
    for i in range(len(attributes)):
        new_data[attributes[i]] = [result_arr[j][i] for j in range(new_row)]

    new_df = pd.DataFrame(new_data, columns = attributes)
    return new_df

In [5]: """
Tiền xử lý dữ liệu
"""
playlists_new = preprocess_data(playlists_df,0.5)
tracks_new = preprocess_data(tracks_df,0.5)
users_new = preprocess_data(users_df,0.5)

In [6]: # ghi dataframe ra file mới
playlists_new.to_csv('cleaned_data/playlists_new.csv', index=False)
tracks_new.to_csv('cleaned_data/tracks_new.csv', index=False)
users_new.to_csv('cleaned_data/users_new.csv', index=False)

In [7]: """
File dữ liệu sau khi tiền xử lý cơ bản
"""
print(playlists_new.shape)
print(tracks_new.shape)
print(users_new.shape)

(608, 3)
(3485, 38)
(216, 32)

In [8]: My_tracks = tracks_new

Ta cần tiền xử lý một vài dữ liệu bất hợp lý

• Bỏ đi những dòng có dữ liệu bất hợp lý đó là những dòng có downloadable là False nhưng dòng download_count vẫn có giá trị khác 0
• Loại bỏ đi những dòng có giá trị là 0 ở hay NaN trong cột download_count
• Loại bỏ đi những dòng trùng nhau và giữ lại dòng xuất hiện đầu tiên

In [9]: My_tracks = My_tracks[My_tracks['downloadable'] == True]
My_tracks = My_tracks[My_tracks['download_count'] != 0]
My_tracks = My_tracks.drop_duplicates(subset=['followers_count', 'username'], keep='first')

Tiếp theo ta sẽ sắp xếp lại theo thứ tự giảm dần rồi chọn 100 dòng đầu tiên

In [10]: My_df = My_tracks[['download_count', 'title']]
My_df = My_df.drop_duplicates(subset=['download_count', 'title'], keep='first')
My_df = My_df.sort_values(by='download_count', ascending=False)
My_df = My_df.head(100)

Out[10]:
```

	download_count	title
715	94327.0	Better Days
2585	63685.0	ODESSA - My Friends Never Die
860	45030.0	The Process
844	28523.0	Good Morning
771	19074.0	Summertime Love
--	--	---
815	212.0	Pieza 7: Jueces
720	198.0	The Way - Instrumental Music Oud / موسيقى العود
855	189.0	Pieza 40: Mateo
813	176.0	Pieza 23: Isaias
810	173.0	Pieza 17: Nehemias

100 rows × 2 columns

```
In [11]: """
sẽ loại bỏ đi 1 vài giá trị outlier, các giá trị này quá lệch so với các giá trị bình thường dẫn đến việc
thắng xèo đi lên, vì vậy 2 biến like_count và repost_count có sự tương quan với nhau theo hướng dương (positive correlation), nghĩa
là nếu lượt like_count tăng, thì track đó sẽ được yêu thích nhiều hơn dẫn đến lượt repost sẽ cao hơn
"""
My_df = My_df[My_df['download_count'] <= 6000]
plt.hist(My_df['download_count'], bins=100)
plt.show()
```

## Câu 1: top 100 các tracks được download nhiều nhất và vẽ biểu đồ histogram để thể hiện các tracks có số lượng download ở khoảng bao nhiêu chiếm đa số

- Trước khi trả lời câu hỏi này ta sẽ xác định một số ý như sau
  - Các số liệu đã được xác định
  - Ứng dụng câu hỏi là tìm ra các tracks được ưa chuộng nhất
    - Câu hỏi có thể trả lời với dữ liệu có sẵn

```
In [8]: My_tracks = tracks_new

Ta cần tiền xử lý một vài dữ liệu bất hợp lý

• Bỏ đi những dòng có dữ liệu bất hợp lý đó là những dòng có downloadable là False nhưng dòng download_count vẫn có giá trị khác 0
• Loại bỏ đi những dòng có giá trị là 0 ở hay NaN trong cột download_count
• Loại bỏ đi những dòng trùng nhau và giữ lại dòng xuất hiện đầu tiên

In [9]: My_tracks = My_tracks[My_tracks['downloadable'] == True]
My_tracks = My_tracks[My_tracks['download_count'] != 0]
My_tracks = My_tracks.drop_duplicates(subset=['followers_count', 'username'], keep='first')

Tiếp theo ta sẽ sắp xếp lại theo thứ tự giảm dần rồi chọn 100 dòng đầu tiên

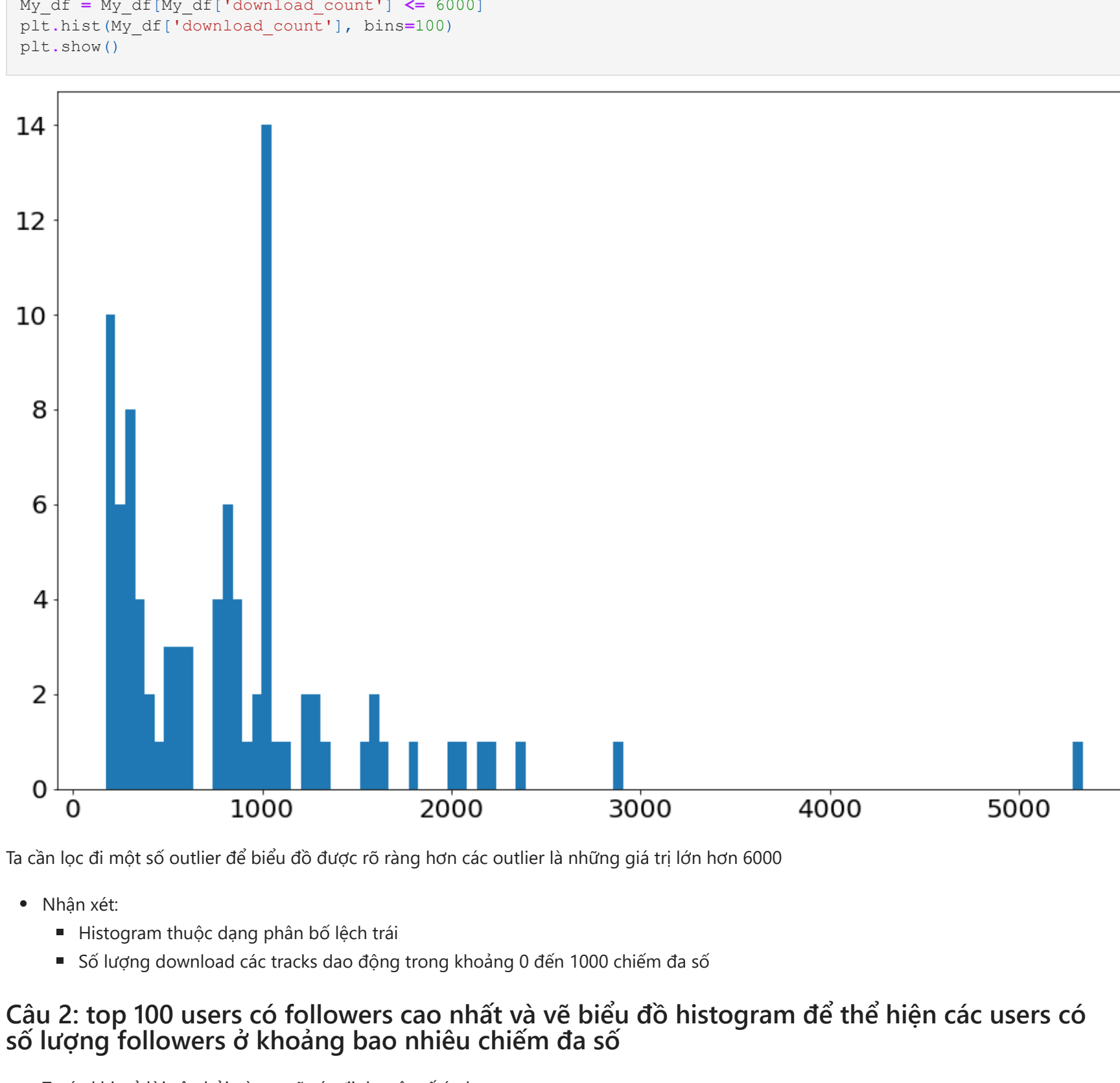
In [10]: My_df = My_tracks[['download_count', 'title']]
My_df = My_df.drop_duplicates(subset=['download_count', 'title'], keep='first')
My_df = My_df.sort_values(by='download_count', ascending=False)
My_df = My_df.head(100)

Out[10]:
```

	download_count	title
715	94327.0	Better Days
2585	63685.0	ODESSA - My Friends Never Die
860	45030.0	The Process
844	28523.0	Good Morning
771	19074.0	Summertime Love
--	--	---
815	212.0	Pieza 7: Jueces
720	198.0	The Way - Instrumental Music Oud / موسيقى العود
855	189.0	Pieza 40: Mateo
813	176.0	Pieza 23: Isaias
810	173.0	Pieza 17: Nehemias

100 rows × 2 columns

```
In [11]: """
sẽ loại bỏ đi 1 vài giá trị outlier, các giá trị này quá lệch so với các giá trị bình thường dẫn đến việc
thắng xèo đi lên, vì vậy 2 biến like_count và repost_count có sự tương quan với nhau theo hướng dương (positive correlation), nghĩa
là nếu lượt like_count tăng, thì track đó sẽ được yêu thích nhiều hơn dẫn đến lượt repost sẽ cao hơn
"""
My_df = My_df[My_df['download_count'] <= 6000]
plt.hist(My_df['download_count'], bins=100)
plt.show()
```



## Câu 2: top 100 users có followers cao nhất và vẽ biểu đồ histogram để thể hiện các users có số lượng followers ở khoảng bao nhiêu chiếm đa số

- Trước khi trả lời câu hỏi này ta sẽ xác định một số ý như sau
  - Các số liệu đã được xác định
  - Ứng dụng câu hỏi là tìm ra users nổi tiếng và có nhiều sản phẩm âm nhạc được ưa thích
    - Câu hỏi có thể trả lời với dữ liệu có sẵn

```
In [12]: My_users = users_new

Ta cần tiền xử lý một vài dữ liệu bất hợp lý

• Loại bỏ đi các dòng trùng nhau và giữ lại dòng xuất hiện đầu tiên

In [13]: My_users = My_users[['followers_count', 'username']]
My_users = My_users.drop_duplicates(subset=['followers_count', 'username'], keep='first')

Tiếp theo ta sẽ sắp xếp lại theo thứ tự giảm dần rồi chọn 100 dòng đầu tiên

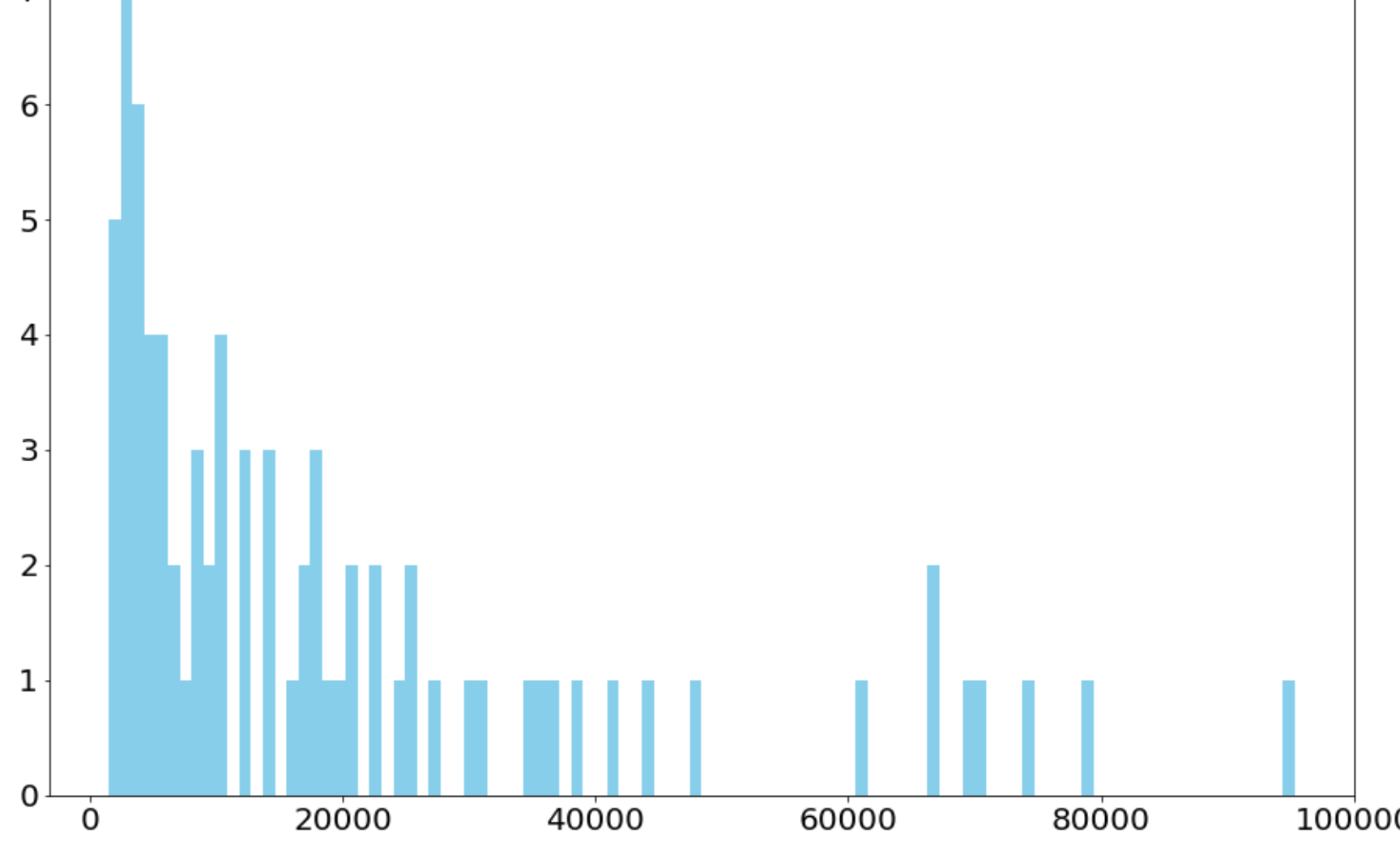
In [14]: My_users = My_users.sort_values(by='followers_count', ascending=False)
My_users = My_users.head(100)
My_users
```

```
Out[14]:
```

	followers_count	username
113	6498535.0	Skrillex
110	3649072.0	Seven Lions
205	3467276.0	Damian Marley
189	2785568.0	dubmatix
190	2768063.0	VP RECORDS
--	--	---
215	2349.0	Black Roots
120	2114.0	Tracey Chattaway
48	2011.0	Copeland Network
132	1814.0	Francesco Berta
183	1480.0	BOST & BM

100 rows × 2 columns

```
In [15]: My_users = My_users[My_users['followers_count'] <= 100000]
plt.hist(My_users['followers_count'], bins=100, color = "skyblue")
plt.show()
```



## Câu 3: Về biểu đồ histogram để xem lượt like\_counts của mỗi track sẽ phân bố như thế nào ?

```
In [16]: his_df = tracks_new[['likes_count']]
his_df = tracks_new[tracks_new['likes_count'].notnull()]
his_df = his_df[his_df['likes_count'] <= 80000] ## Ta sẽ loại bỏ đi các outliers để có thể dễ trực quan hơn.
x = his_df['likes_count'].to_list()
x = np.asarray(x)

In [17]: # calculate bins by Freedman-Diaconis method and plot the chart
q25, q75 = np.percentile(x, [25, 75])
bin_width = 2 * (q75 - q25) * len(x) ** (-1/3)
bins = round((x.max() - x.min()) / bin_width)
print("Freedman-Diaconis number of bins:", bins)

# plot distribution
plt.hist(x, bins=bins)
plt.xlabel('Like_counts')
plt.ylabel('Probability')
```

```
Out[17]:
```

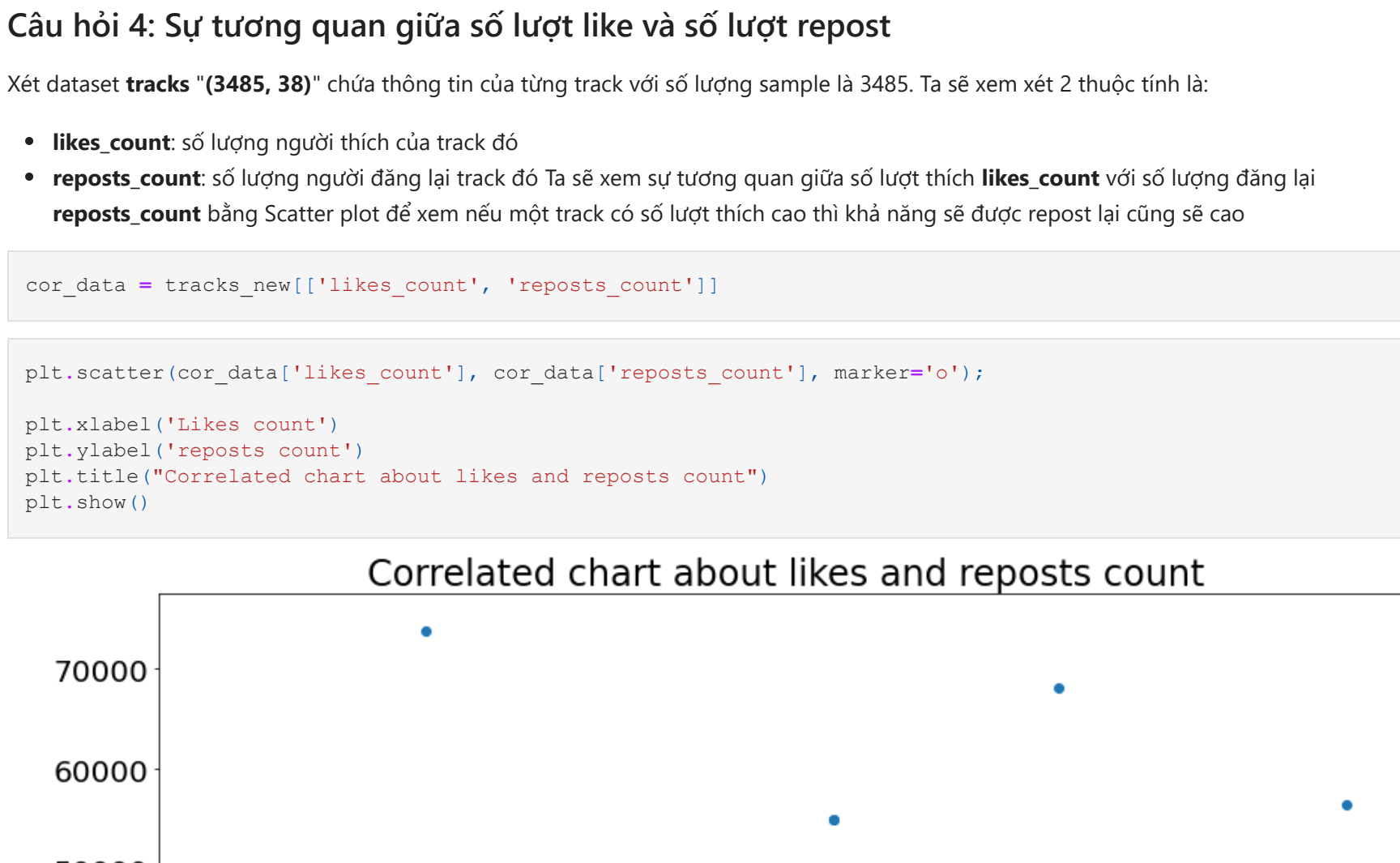
## Câu 4: Sự tương quan giữa số lượt like và số lượt repost

Xét dataset tracks ['3485,38'] chứa thông tin của từng track với số lượng sample là 3485. Ta sẽ xem xét 2 thuộc tính là:

- likes\_count: số lượng người thích của track đó
- reposts\_count: số lượng người đăng lại track đó Ta sẽ xem sự tương quan giữa số lượt thích likes\_count và số lượt đăng lại reposts\_count bằng Scatter plot để xem nếu một track có số lượt thích cao thì khả năng sẽ được repost lại cũng sẽ cao hơn

```
In [18]: cor_data = tracks_new[['likes_count', 'reposts_count']]

In [19]: plt.scatter(cor_data['likes_count', 'reposts_count'], marker='o');
plt.xlabel('likes count')
plt.ylabel('reposts count')
plt.title("Correlated chart about likes and reposts count")
plt.show()
```

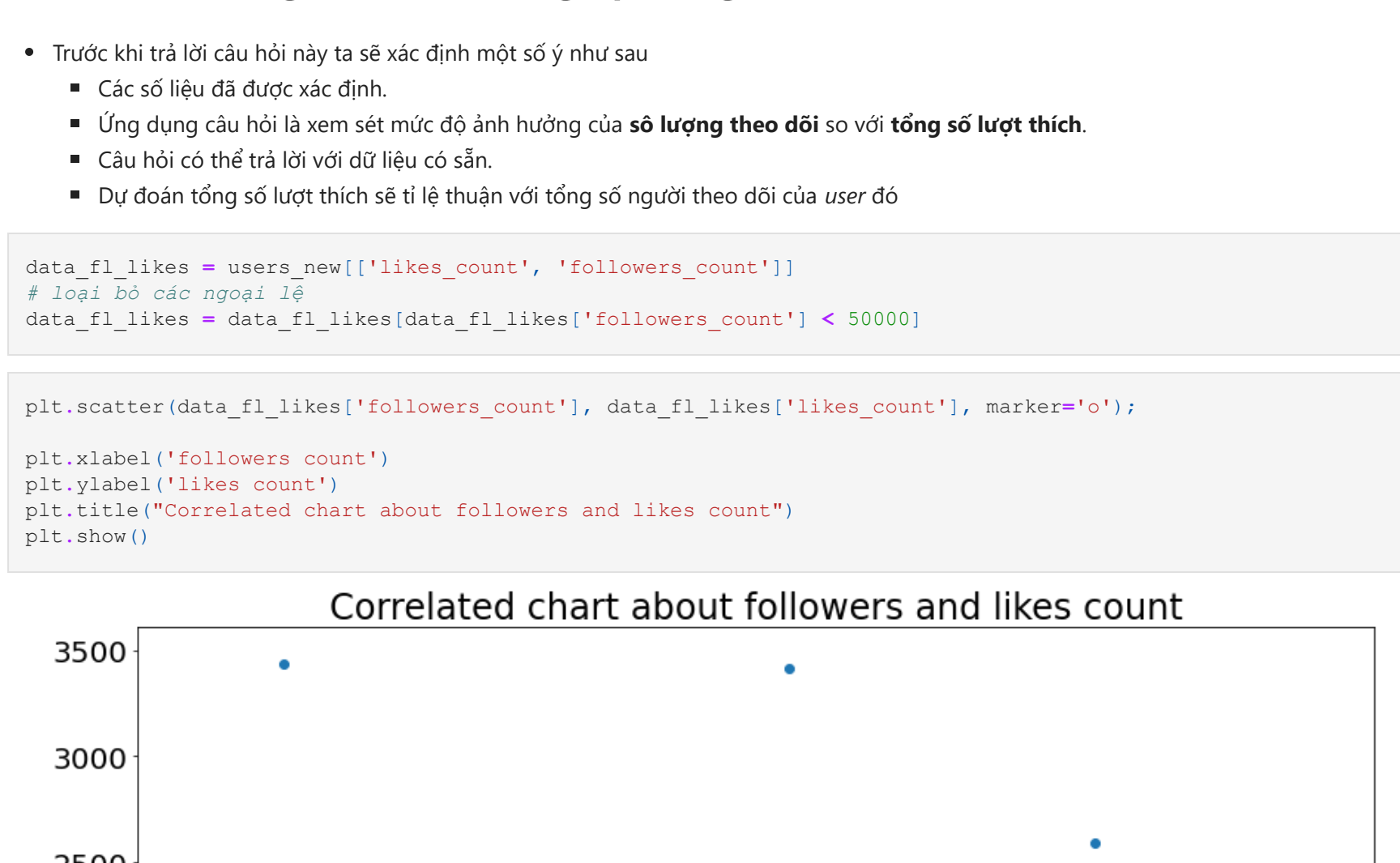


## Câu 5: Đánh giá mối tương quan giữa followers và likes\_count

- Trước khi trả lời câu hỏi này ta sẽ xác định một số ý như sau
  - Các số liệu đã được xác định
  - Ứng dụng câu hỏi là xem xét mức độ ảnh hưởng của số lượng theo dõi so với tổng số lượt thích
    - Câu hỏi có thể trả lời với dữ liệu có sẵn
  - Dự đoán tổng số lượt thích sẽ tỉ lệ thuận với tổng số người theo dõi của user đó

```
In [20]: data_fi_likes = users_new[['likes_count', 'followers_count']]
# loại bỏ các ngoại lệ
data_fi_likes = data_fi_likes[data_fi_likes['followers_count'] < 50000]

In [21]: plt.scatter(data_fi_likes['followers_count', 'likes_count'], marker='o');
plt.xlabel('followers count')
plt.ylabel('likes count')
plt.title("Correlated chart about followers and likes count")
plt.show()
```



## Nhận xét qua biểu đồ:

- Lượng theo dõi không ảnh hưởng quá nhiều đến tổng số lượt thích của user đó.
- Ngữ cảnh biết được xu hướng âm nhạc trong năm 2021, cho thấy cho ra những sản phẩm mang top 10 thể loại phổ biến, có thể giúp cho bài hát được nhiều người biết tới.
- Mức độ đánh giá cho câu trả lời không cao (dữ liệu có thêm thuộc tính lượt view thì kết quả sẽ tốt hơn).

- nhìn vào biểu đồ ta có thể thấy thể loại "Religion & Spirituality" là thể loại có nhiều bài hát nhất năm 2021 ta có thể đánh giá nó là thể loại đang được hưởng tới từ những người nổi sĩ.

## Câu 6: Top 10 thể loại nhạc phổ biến nhất năm 2021

```
In [22]: df = pd.read_csv('cleaned_data/tracks_new.csv')
# lấy 2 thuộc tính là 'genre' và 'created_at' và bỏ qua dòng có giá trị null.
dt_df = df[['genre', 'created_at']].dropna()

# đổi kiểu dữ liệu thuộc tính 'created_at' thành kiểu datetime
dt_df['created_at'] = pd.to_datetime(dt_df['created_at'])

# lấy top 10 thể loại phổ biến nhất năm 2021.
dt_df['year'] = dt_df['created_at'].apply(lambda x: x.year)
top_10_popular = dt_df.query('year == 2021')['genre'].value_counts()[0:10]

In [23]: # visualization
fig, ax = plt.subplots(figsize=(14, 9))
ax.barh(top_10_popular[0:10].index, top_10_popular[0:10])
ax.set_title('Top 10 thể loại phổ biến nhất năm 2021', fontdict=font, fontsize=32)
ax.set_xlabel('Số lượng bài hát', fontdict=font, fontsize=20)
ax.set_ylabel('Thể loại', fontdict=font, fontsize=20)
ax.tick_params(labelsize=14)
# fig.autofmt_xdate()

plt.show()
```

Top 10 thể loại phổ biến năm 2021

- Số liệu trước khi bắt đầu cột thuộc tính "genre": tên thể loại và "created\_at": ngày tạo.
- Ngữ cảnh biết được xu hướng âm nhạc trong năm 2021, cho thấy cho ra những sản phẩm mang top 10 thể loại phổ biến, có thể giúp cho bài hát được nhiều người biết tới.
- Mức độ đánh giá cho câu trả lời không cao (dữ liệu có thêm thuộc tính lượt view thì kết quả sẽ tốt hơn).

- nhìn vào biểu đồ ta có thể thấy thể loại "Religion & Spirituality" là thể loại có nhiều bài hát nhất năm 2021 ta có thể đánh giá nó là thể loại đang được hưởng tới từ những người nổi sĩ.

## Câu 7: Top 10 thể loại nhạc phổ biến dựa trên likes\_count và playback\_count

Trong file "tracks\_new.csv":

- Tim tổng số lượt likes\_count, playback\_count của các 'genre' và lưu vào dataframe
- Sắp xếp dataframe này theo thứ tự tăng dần likes\_count
- Vẽ 2 biểu đồ để hiển thị
- Mức độ: tìm được thể loại có tổng số lượt likes\_count cao nhất và playback\_count cao nhất

```
In [24]: tracks_df = pd.read_csv('cleaned_data/tracks_new.csv', index_col='Unnamed: 0')
users_df = pd.read_csv('cleaned_data/users_new.csv', index_col='Unnamed: 0')

In [25]: genres = tracks_df['genre'].unique()
# create genre df
cc = users_df.groupby(columns = ['sum_likes_count', 'sum_playback_count'], index = genres)

for g in genres:
    genre_df.loc[g, 'sum_likes_count'] = tracks_df[tracks_df['genre'] == g]['likes_count'].sum()
    genre_df.loc[g, 'sum_playback_count'] = tracks_df[tracks_df['genre'] == g]['playback_count'].sum()

genre_df.sort_values(by='sum_likes_count', ascending=False, inplace=True)

Xem kết quả của 10 dòng đầu

In [26]: genre_df.head(10)

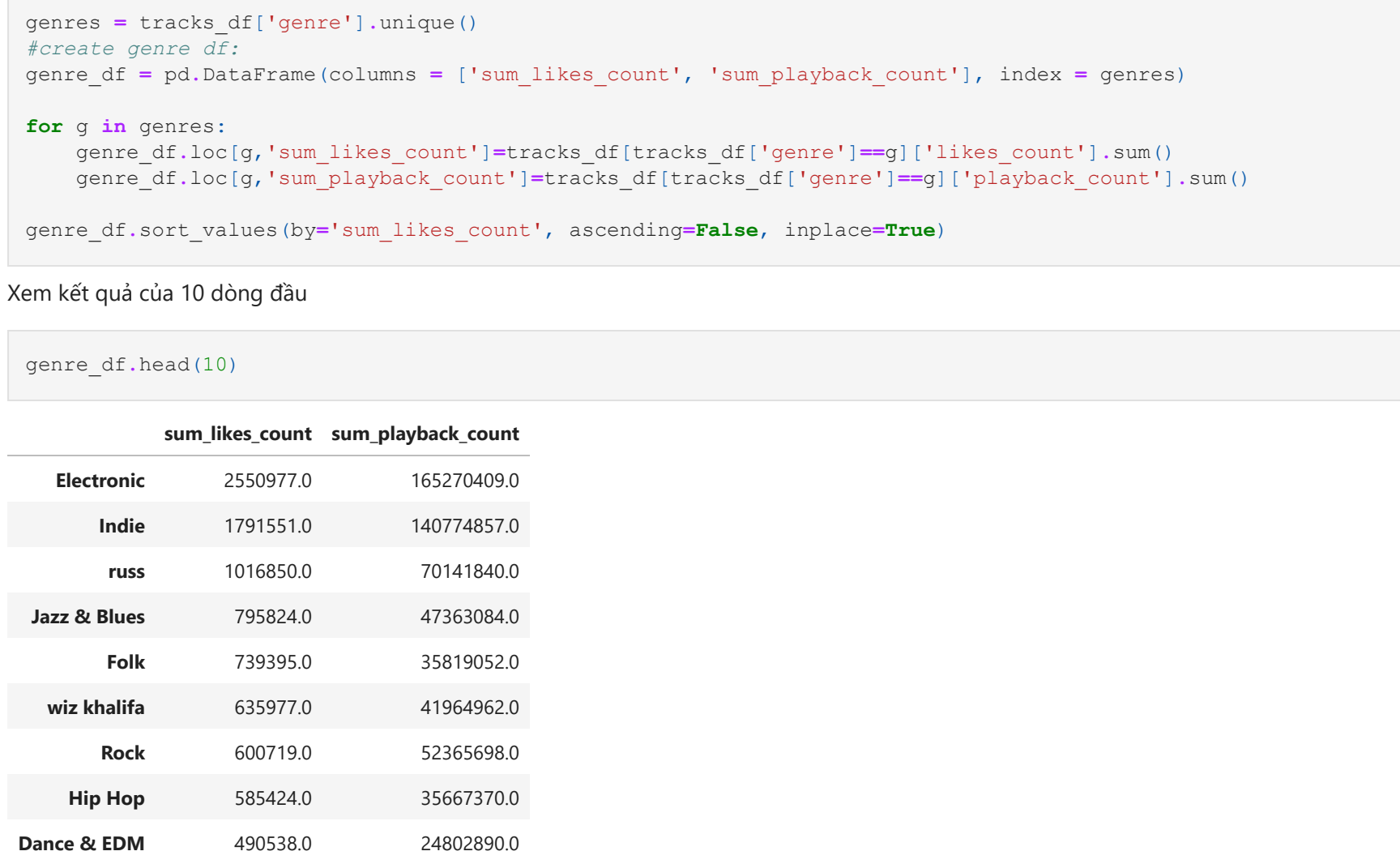
Out[26]:
```

	sum_likes_count	sum_playback_count
Electronic	2550977.0	165270409.0
Indie	1791551.0	140774857.0
russ	1016850.0	70141840.0
Jazz & Blues	793824.0	47363084.0
Folk	733935.0	35819052.0
wiz khalfia	635977.0	41964962.0
Rock	600719.0	52365698.0
Hip Hop	585424.0	35667370.0
Dance & EDM	490538.0	24802890.0
Alternative	481036.0	3187256.0

Vẽ biểu đồ của 10 dòng đầu có sum\_likes\_count cao nhất:

```
In [27]: genre_df['sum_likes_count'].head(10).plot.barh()

Out[27]:
```



## Câu 8: Users đến từ đâu là nhiều nhất

Trong "users\_new.csv":

- Cột "country\_code" thể hiện user này đến từ nước nào (VD: US: Hoa Kỳ, VN: Việt Nam)
- Tạo ra một DataFrame để lưu lại các thông tin sau: Index = "country\_code", 1 cột number\_of\_users lưu số lượng user.
- Không xử lý ô trống
- Tim hiểu users trong file này đến từ nước nào nhiều nhất, sắp xếp theo thứ tự giảm dần

```
In [29]: #country_code viết tắt là cc
cc = users_df['country_code'].notna()
# create df
cc_df = pd.DataFrame(columns = ['number_of_users'], index=cc)

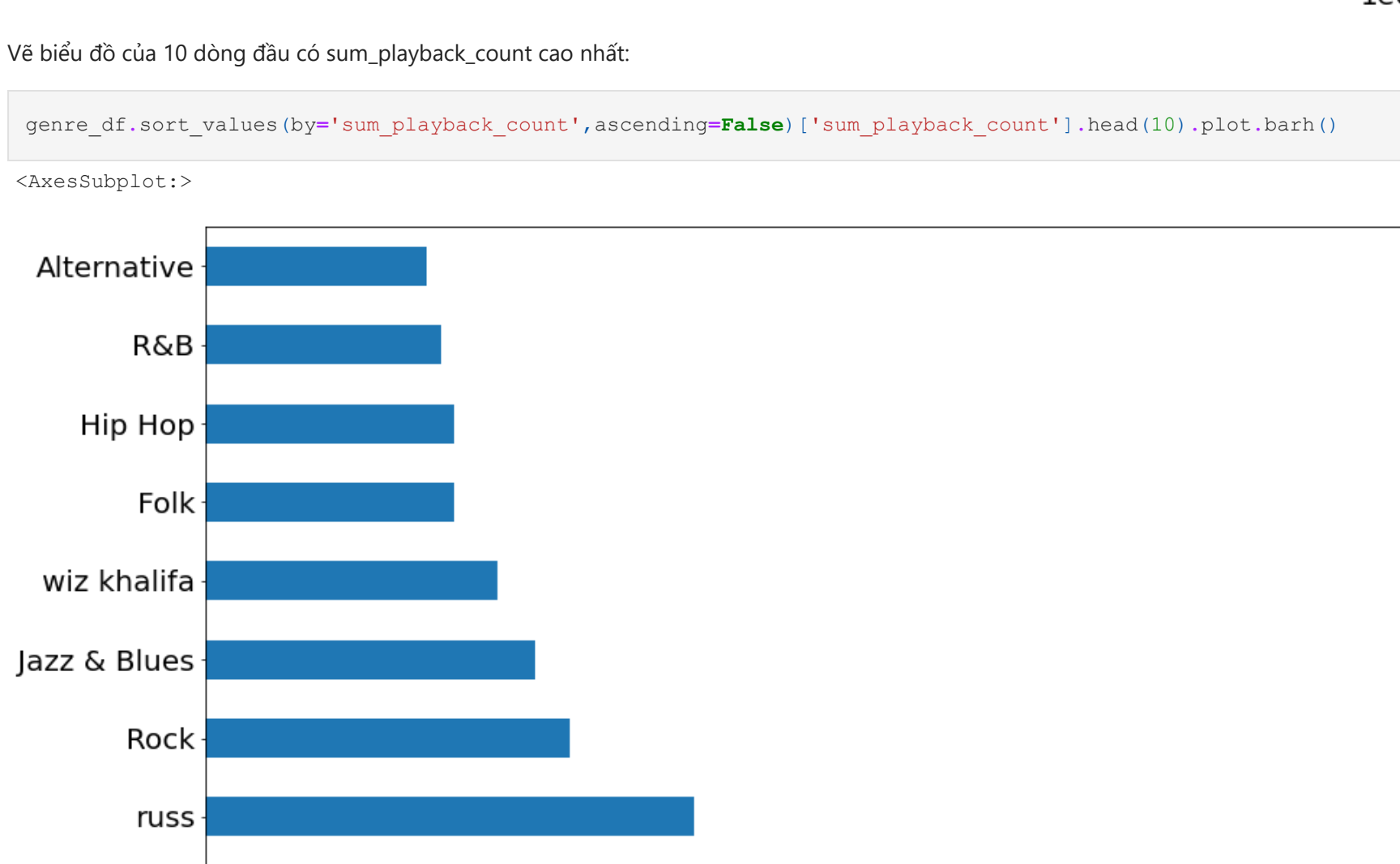
for c in cc:
    cc_df.loc[c, 'number_of_users'] = users_df['country_code'] == c).sum()

cc_df.sort_values(by='number_of_users', ascending=False, inplace=True)

Xem kết quả:

In [30]: cc_df.head(20).plot.barh()

Out[30]:
```



## Câu 9: Users đến từ đâu là nhiều nhất

```
In [31]: cc_df.head(20).plot.barh()

Out[31]:
```



