

# YouTube Streamer Analysis

## Introduction

This dataset contains valuable information about the Top 1000 YouTube streamers.

The dataset contains 9 columns namely:

1. Rank: Ranking of each account
2. Username: Account name of each Youtuber
3. Categories: Category of content created
4. Subscribers: Number of subscribers
5. Country: Country of the Youtubers
6. Visits: Number of Views
7. Likes: Numbers of Likes
8. Comments: Total numbers of comment
9. Links: link to the account

## Task

The task is to perform a comprehensive analysis to extract insights about the top YouTube content creators.

```
In [1]: # Import the necessary Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

pd.options.display.float_format = '{:.0f}'.format # Suppress scientific notation
```

## Data Exploration:

- Start by exploring the dataset to understand its structure and identify key variables.

- Check for missing data and outliers.

```
In [2]: # Import the dataset  
data = pd.read_csv("youtubers_df.csv")
```

```
# made a copy of the data  
df = data.copy()
```

```
#check information about the data  
df.info()
```

```
# view the data  
df
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 1000 entries, 0 to 999  
Data columns (total 9 columns):  
 #   Column      Non-Null Count  Dtype     
 ---    
 0   Rank        1000 non-null    int64    
 1   Username    1000 non-null    object    
 2   Categories  694 non-null    object    
 3   Suscribers  1000 non-null    float64   
 4   Country     1000 non-null    object    
 5   Visits      1000 non-null    float64   
 6   Likes       1000 non-null    float64   
 7   Comments    1000 non-null    float64   
 8   Links       1000 non-null    object    
dtypes: float64(4), int64(1), object(4)  
memory usage: 70.4+ KB
```

Out[2] :

	Rank	Username	Categories	Suscribers	Country	Visits	Likes	Comments	Links
0	1	tseries	Música y baile	249500000	India	86200	2700	78	<a href="http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...">http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...</a>
1	2	MrBeast	Videojuegos, Humor	183500000	Estados Unidos	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
2	3	CoComelon	Educación	165500000	Unknown	7000000	24700	0	<a href="http://youtube.com/channel/UCbCmjCuTUZos6lnko4...">http://youtube.com/channel/UCbCmjCuTUZos6lnko4...</a>
3	4	SETIndia	NaN	162600000	India	15600	166	9	<a href="http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...">http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...</a>
4	5	KidsDianaShow	Animación, Juguetes	113500000	Unknown	3900000	12400	0	<a href="http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...">http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...</a>
...	...	...	...	...	...	...	...	...	...
995	996	hamzymukbang	NaN	11700000	Estados Unidos	397400	14000	124	<a href="http://youtube.com/channel/UCPKNKldggioffXPkSm...">http://youtube.com/channel/UCPKNKldggioffXPkSm...</a>
996	997	Adaahqueen	NaN	11700000	India	1100000	92500	164	<a href="http://youtube.com/channel/UCk3fFpqI5kDMf_mUP...">http://youtube.com/channel/UCk3fFpqI5kDMf_mUP...</a>
997	998	LittleAngellIndonesia	Música y baile	11700000	Unknown	211400	745	0	<a href="http://youtube.com/channel/UCdrHrQf0o0TO8YDntX...">http://youtube.com/channel/UCdrHrQf0o0TO8YDntX...</a>
998	999	PenMultiplex	NaN	11700000	India	14000	81	1	<a href="http://youtube.com/channel/UCObyBrdrtQ20BU9PxH...">http://youtube.com/channel/UCObyBrdrtQ20BU9PxH...</a>
999	1000	OneindiaHindi	Noticias y Política	11700000	India	2200	31	1	<a href="http://youtube.com/channel/UCOjgc1p2hJ4GZi6pQQ...">http://youtube.com/channel/UCOjgc1p2hJ4GZi6pQQ...</a>

1000 rows × 9 columns

the Categories and Country values are not in english so i translate them to english

i thought of using the googletrans library but i keep getting a timeout error due to the large dataset so i used an alternaive approach: getting the unique data in each columns and using ChatGPT to translate and convert them into dictionary

In [3] :

```
categories_dict = {
    'Música y baile': 'Music and Dance',
    'Videojuegos, Humor': 'Video Games, Humor',
    'Educación': 'Education',
    'Unknown': 'Unknown',
    'Animación, Juguetes': 'Animation, Toys',
    'Películas, Videojuegos': 'Movies, Video Games',
```

```
'Juguetes': 'Toys',
'Videojuegos': 'Video Games',
'Películas, Animación': 'Movies, Animation',
'Películas': 'Movies',
'Noticias y Política': 'News and Politics',
'Animación, Humor': 'Animation, Humor',
'Música y baile, Animación': 'Music and Dance, Animation',
'Música y baile, Películas': 'Music and Dance, Movies',
'Películas, Juguetes': 'Movies, Toys',
'Películas, Humor': 'Movies, Humor',
'Vlogs diarios': 'Daily Vlogs',
'Videojuegos, Juguetes': 'Video Games, Toys',
'Animación, Videojuegos': 'Animation, Video Games',
'Animación': 'Animation',
'Música y baile, Humor': 'Music and Dance, Humor',
'Diseño/arte, DIY y Life Hacks': 'Design/Art, DIY and Life Hacks',
'Ciencia y tecnología': 'Science and Technology',
'Fitness, Salud y autoayuda': 'Fitness, Health and Self-help',
'Belleza, Moda': 'Beauty, Fashion',
'Humor': 'Humor',
'Comida y bebida': 'Food and Drink',
'Deportes': 'Sports',
'Fitness': 'Fitness',
'Viajes, Espectáculos': 'Travel, Entertainment',
'Comida y bebida, Salud y autoayuda': 'Food and Drink, Health and Self-help',
'Diseño/arte': 'Design/Art',
'DIY y Life Hacks, Juguetes': 'DIY and Life Hacks, Toys',
'Educación, Juguetes': 'Education, Toys',
'Juguetes, Coches y vehículos': 'Toys, Cars and Vehicles',
'Música y baile, Juguetes': 'Music and Dance, Toys',
'Animales y mascotas': 'Animals and Pets',
'ASMR': 'ASMR',
'Moda': 'Fashion',
'DIY y Life Hacks': 'DIY and Life Hacks',
'Diseño/arte, Belleza': 'Design/Art, Beauty',
'Coches y vehículos': 'Cars and Vehicles',
'Animación, Humor, Juguetes': 'Animation, Humor, Toys',
'ASMR, Comida y bebida': 'ASMR, Food and Drink',
'Comida y bebida, Juguetes': 'Food and Drink, Toys',
'Juguetes, DIY y Life Hacks': 'Toys, DIY and Life Hacks'
}
```

```
# Function to translate Spanish words to English
```

```
def categories_to_english(word):
    return categories_dict.get(word, word) # Return the translated word if it exists,
# Apply the translation function to the 'Category' column
df['Categories'] = df['Categories'].apply(categories_to_english)
```

```
In [4]: countries_dict = {
    'India': 'India',
    'Estados Unidos': 'United States',
    'Unknown': 'Unknown',
    'Brasil': 'Brazil',
    'México': 'Mexico',
    'Rusia': 'Russia',
    'Pakistán': 'Pakistan',
    'Filipinas': 'Philippines',
    'Indonesia': 'Indonesia',
    'Tailandia': 'Thailand',
    'Francia': 'France',
    'Colombia': 'Colombia',
    'Irak': 'Iraq',
    'Japón': 'Japan',
    'Ecuador': 'Ecuador',
    'Argentina': 'Argentina',
    'Turquía': 'Turkey',
    'Arabia Saudita': 'Saudi Arabia',
    'El Salvador': 'El Salvador',
    'Bangladesh': 'Bangladesh',
    'Reino Unido': 'United Kingdom',
    'Argelia': 'Algeria',
    'España': 'Spain',
    'Perú': 'Peru',
    'Egipto': 'Egypt',
    'Jordania': 'Jordan',
    'Marruecos': 'Morocco',
    'Singapur': 'Singapore',
    'Somalia': 'Somalia'
}

# Function to translate Spanish words to English
def countries_to_english(word):
    return countries_dict.get(word, word) # Return the translated word if it exist
```

```
# Apply the translation function to the 'Category' column  
df['Country'] = df['Country'].apply(countries_to_english)
```

```
In [5]: # View the first 5 rows of the now translated dataset  
df.head()
```

	Rank	Username	Categories	Suscribers	Country	Visits	Likes	Comments	Links
0	1	tseries	Music and Dance	249500000	India	86200	2700	78	http://youtube.com/channel/UCq-Fj5knLsUf-MWSy...
1	2	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...
2	3	CoComelon	Education	165500000	Unknown	7000000	24700	0	http://youtube.com/channel/UCbCmjCuTUZos6Inko4...
3	4	SETIndia		Nan	India	15600	166	9	http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...
4	5	KidsDianaShow	Animation, Toys	113500000	Unknown	3900000	12400	0	http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...

```
In [6]: # checking the numbers of unique values in each columns  
df.nunique()
```

```
Out[6]: Rank      1000  
Username    994  
Categories   45  
Suscribers  274  
Country     29  
Visits      713  
Likes        648  
Comments     389  
Links        994  
dtype: int64
```

the username and links columns has some duplicated data.

There should be a unique link for each youtubers so we remove all duplicates

```
In [7]: # check for duplicates using the username column  
df[df['Username'].duplicated(keep=False)]
```

Out[7]:	Rank	Username	Categories	Suscribers	Country	Visits	Likes	Comments	Links
	249	250	NickyJamTV	Music and Dance, Movies	23700000	Colombia	15800	1200	58 http://youtube.com/channel/UCpb_iJuhFe8V6rQdbN...
	250	251	NickyJamTV	Music and Dance, Movies	23700000	Colombia	15800	1200	58 http://youtube.com/channel/UCpb_iJuhFe8V6rQdbN...
	447	448	mgcplayhouse	Toys, Cars and Vehicles	17800000	Unknown	56300	96	0 http://youtube.com/channel/UC6zPzUJo8hu-5TzUk8...
	449	450	thexoteam		NaN	17800000	United States	797600	50400 179 http://youtube.com/channel/UCIZAOlfhJQRym39Wo...
	450	451	thexoteam		NaN	17900000	United States	772800	45000 185 http://youtube.com/channel/UCIZAOlfhJQRym39Wo...
	451	452	mgcplayhouse	Toys, Cars and Vehicles	17800000	Unknown	63600	75	0 http://youtube.com/channel/UC6zPzUJo8hu-5TzUk8...
	946	947	Super_Senya_RU	Animation, Toys	12100000	Unknown	47200	192	0 http://youtube.com/channel/UCTn9Vyy-3fzLlr0bqh...
	947	948	HiTechIslamic	Music and Dance	12100000	Pakistan	62200	810	59 http://youtube.com/channel/UCtKKyuORzErSd7TWfk...
	949	950	Family-Box	Movies	12000000	Russia	173600	6600	105 http://youtube.com/channel/UC-jHNWViReG6R_kJ6b...
	952	953	Super_Senya_RU	Animation, Toys	12100000	Unknown	47200	192	0 http://youtube.com/channel/UCTn9Vyy-3fzLlr0bqh...
	953	954	HiTechIslamic	Music and Dance	12100000	Pakistan	62200	810	59 http://youtube.com/channel/UCtKKyuORzErSd7TWfk...
	956	957	Family-Box	Movies	12000000	Russia	177400	6300	86 http://youtube.com/channel/UC-jHNWViReG6R_kJ6b...

```
In [8]: # drop all duplicates
df.drop_duplicates(subset=['Username'], inplace=True, ignore_index=True)

# check again to confirm the dataset is no more having duplicate data
df['Username'].duplicated().sum()
```

Out[8]: 0

```
In [9]: # the rank column is needed for the analysis so i'll drop it
df.drop(columns='Rank', inplace=True)

df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 994 entries, 0 to 993
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Username    994 non-null    object  
 1   Categories  689 non-null    object  
 2   Suscribers  994 non-null    float64 
 3   Country     994 non-null    object  
 4   Visits      994 non-null    float64 
 5   Likes       994 non-null    float64 
 6   Comments    994 non-null    float64 
 7   Links       994 non-null    object  
dtypes: float64(4), object(4)
memory usage: 62.2+ KB
```

we are left with 994rows and 7cols, notice only the categories column is having missing values so i deal with the missing data and also correct the subscribers column name

```
In [10]: # check the total missing values
df['Categories'].isna().sum()
```

```
Out[10]: 305
```

305 missing values, about 30% of the dataset, dropping them will reduce the dataset drastically and this may affect the accuracy of our analysis so i'll fill it with "Unknown"

```
In [11]: df['Categories'].fillna('Unknown', inplace=True)

#check to confirm no more missing values in the dataset
df.isna().sum().sum()
```

```
Out[11]: 0
```

```
In [12]: df.rename(columns={'Suscribers': 'Subscribers'}, inplace=True)

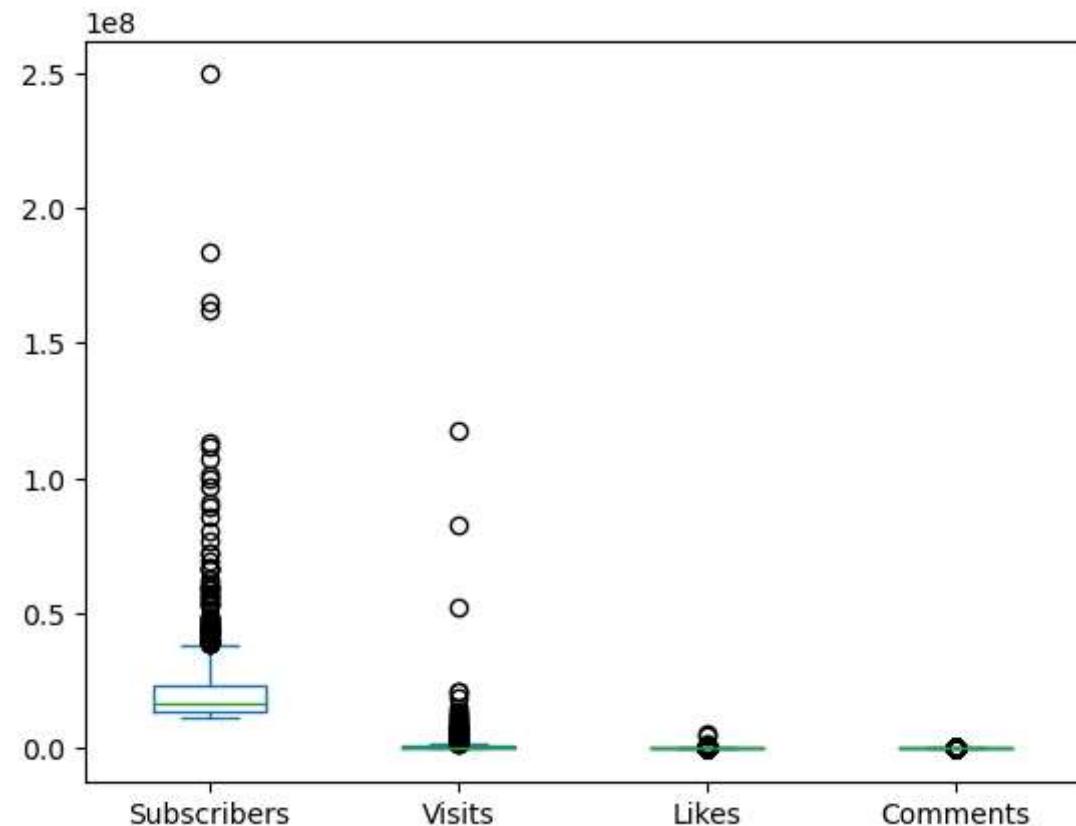
df.head()
```

Out[12]:

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
0	tseries	Music and Dance	249500000	India	86200	2700	78	<a href="http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...">http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...</a>
1	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
2	CoComelon	Education	165500000	Unknown	7000000	24700	0	<a href="http://youtube.com/channel/UCbCmjCuTUZos6Inko4...">http://youtube.com/channel/UCbCmjCuTUZos6Inko4...</a>
3	SETIndia	Unknown	162600000	India	15600	166	9	<a href="http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...">http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...</a>
4	KidsDianaShow	Animation, Toys	113500000	Unknown	3900000	12400	0	<a href="http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...">http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...</a>

In [13]: *#checking for outliers we construct a boxplot*

```
df.plot(kind='box');
```



there are hugh amount of outliers, this is understandable because content creator followers varies based on how their contents gains attraction.

## Trend Analysis:

- Identify trends among the top YouTube streamers. Which categories are the most popular?
- Is there a correlation between the number of subscribers and the number of likes or comments?

```
In [14]: # we identify trends among the top youtube streamers by the number of subscribers, visit, and likes

# top 10 by numbers of subscribers
display(df.groupby('Categories')[['Subscribers']].sum().sort_values(by='Subscribers', ascending=False).head(10))

# top 10 by numbers of visits
display(df.groupby('Categories')[['Visits']].sum().sort_values(by='Visits', ascending=False).head(10))

# top 10 by numbers of Likes
df.groupby('Categories')[['Likes']].sum().sort_values(by='Likes', ascending=False).head(10)
```

Subscribers	
Categories	
<b>Unknown</b>	6320100000
<b>Music and Dance</b>	4281800000
<b>Movies, Animation</b>	1384300000
<b>Animation, Toys</b>	839800000
<b>Music and Dance, Movies</b>	774800000
<b>News and Politics</b>	676100000
<b>Animation, Video Games</b>	659400000
<b>Daily Vlogs</b>	654900000
<b>Movies, Humor</b>	622100000
<b>Education</b>	600300000

### Visits

#### Categories

<b>Unknown</b>	368427116
<b>Video Games, Humor</b>	174074500
<b>Daily Vlogs</b>	126330500
<b>Animation, Humor</b>	101523400
<b>Music and Dance</b>	59839900
<b>Animation, Video Games</b>	40802000
<b>Movies, Animation</b>	33631100
<b>Food and Drink</b>	32669400
<b>Movies, Humor</b>	31916600
<b>Education</b>	26545000

Out[14]:

### Likes

#### Categories

<b>Unknown</b>	16476915
<b>Video Games, Humor</b>	7148700
<b>Daily Vlogs</b>	6928063
<b>Animation, Humor</b>	3935745
<b>Music and Dance</b>	2784099
<b>Animation, Video Games</b>	2695997
<b>Humor</b>	1699900
<b>Movies, Animation</b>	1565932
<b>Food and Drink</b>	1543977
<b>Movies, Humor</b>	1383277

Since the "Unknown" is the missing values,

the top 5 most popular categories are similar by numbers of Visits and Likes which are:

1. Video Games, Humor
2. Daily Vlogs
3. Animation, Humor
4. Music and Dance
5. Animation, Video Games

but the case is different for the top categories by numbers of subscribers which are:

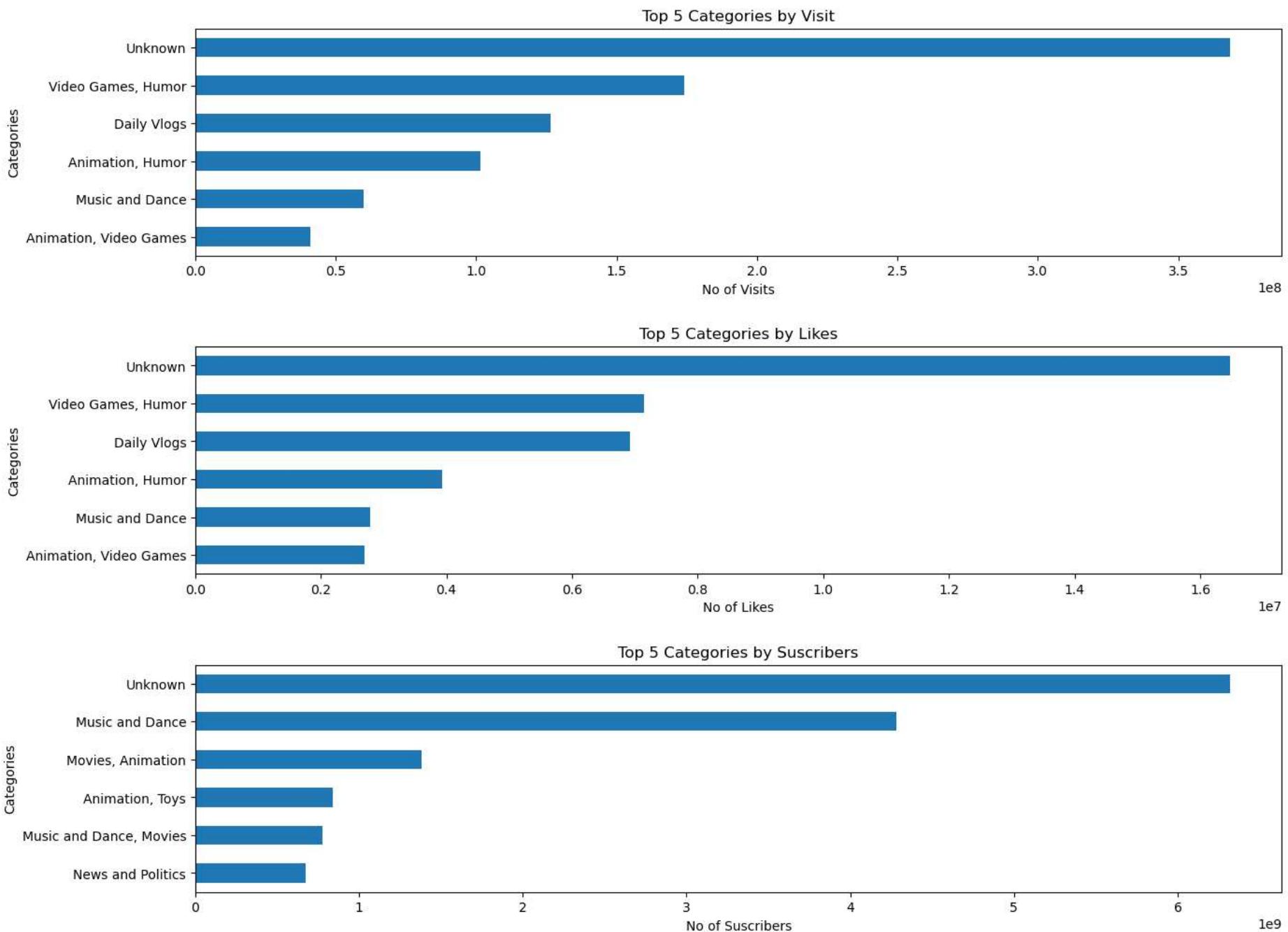
1. Music and Dance
2. Movies, Animation
3. Animation, Toys
4. Music and Dance, Movies
5. News and Politics

```
In [15]: fig, axs = plt.subplots(figsize=(15,12), nrows=3, ncols=1)
plt.subplots_adjust(hspace=0.4)
df.groupby('Categories')['Visits'].sum().sort_values().tail(6).plot(kind='barh', ax=axs[0])
axs[0].set_title("Top 5 Categories by Visit")
axs[0].set_xlabel("No of Visits")

df.groupby('Categories')['Likes'].sum().sort_values().tail(6).plot(kind='barh', ax=axs[1])
axs[1].set_title("Top 5 Categories by Likes")
axs[1].set_xlabel("No of Likes")

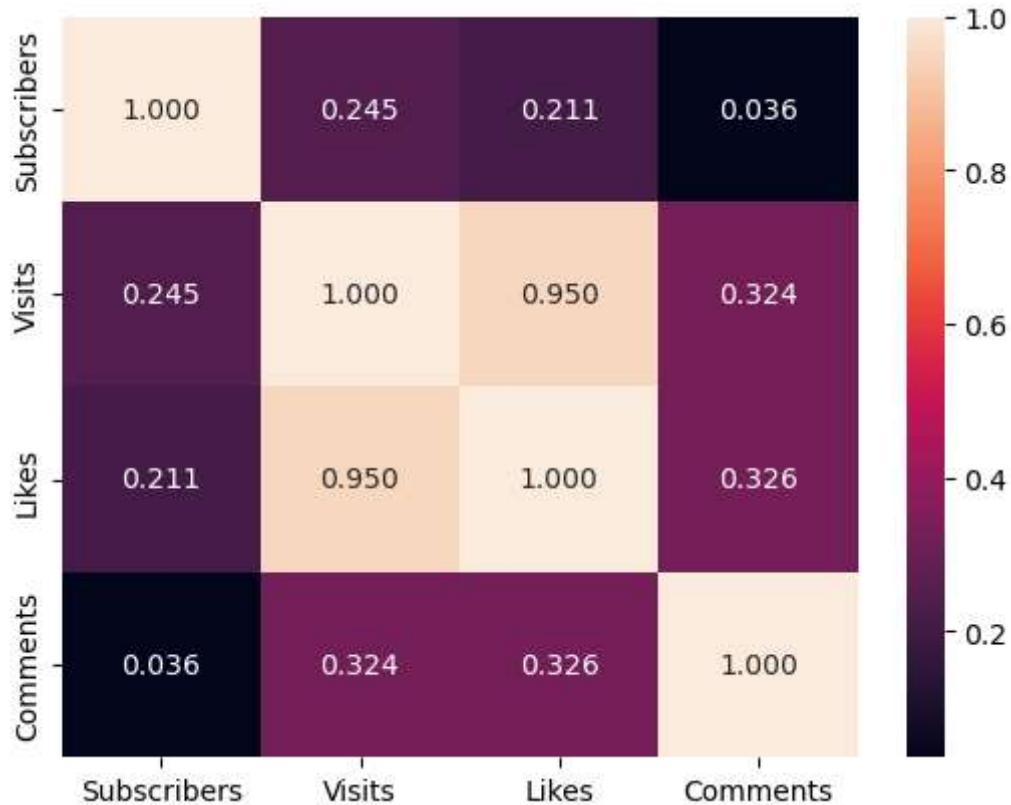
df.groupby('Categories')['Subscribers'].sum().sort_values().tail(6).plot(kind='barh', ax=axs[2])
axs[2].set_title("Top 5 Categories by Suscribers")
axs[2].set_xlabel("No of Suscribers")

plt.show()
```



```
In [16]: # to check for correlation consider the heatmap:
```

```
sns.heatmap(df.select_dtypes('number').corr(), annot=True, fmt='.3f');
```



### Observations:

There is no relationship between the Subscribers, Likes and Comment but there is a strong relationship between Visits and Likes

### Audience Study:

- Analyze the distribution of streamers' audiences by country. Are there regional preferences for specific content categories?

```
In [17]: df.head()
```

Out[17]:

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
0	tseries	Music and Dance	249500000	India	86200	2700	78	<a href="http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...">http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...</a>
1	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
2	CoComelon	Education	165500000	Unknown	7000000	24700	0	<a href="http://youtube.com/channel/UCbCmjCuTUZos6Inko4...">http://youtube.com/channel/UCbCmjCuTUZos6Inko4...</a>
3	SETIndia	Unknown	162600000	India	15600	166	9	<a href="http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...">http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...</a>
4	KidsDianaShow	Animation, Toys	113500000	Unknown	3900000	12400	0	<a href="http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...">http://youtube.com/channel/UCk8GzjMOrta8yxDcKf...</a>

Notice "Unknown" Values in both the Categories and Countries columns i filter them out

In [18]:

```
# removing the rows with "Unknown" values in both the Categories and Country column
df2 = df[(df['Categories'].str.contains("Unknown") != True) & (df['Country'].str.contains("Unknown") != True)]

df2.head()
```

Out[18]:

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
0	tseries	Music and Dance	249500000	India	86200	2700	78	<a href="http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...">http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy...</a>
1	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
5	PewDiePie	Movies, Video Games	111500000	United States	2400000	197300	4900	<a href="http://youtube.com/channel/UC-IHJZR3Gqxm24_Vd_...">http://youtube.com/channel/UC-IHJZR3Gqxm24_Vd_...</a>
8	zeemusiccompany	Music and Dance	99700000	India	74300	2600	32	<a href="http://youtube.com/channel/UCFFbwnve3yF62-tVXk...">http://youtube.com/channel/UCFFbwnve3yF62-tVXk...</a>
9	WWE	Video Games	97200000	United States	184500	6300	214	<a href="http://youtube.com/channel/UCJ5v_MCY6GNUBTO8-D...">http://youtube.com/channel/UCJ5v_MCY6GNUBTO8-D...</a>

In [19]:

```
# region preference by numbers of subscribers
display(df2.groupby(['Country', 'Categories'])[['Subscribers']].sum().sort_values(by=['Country', 'Subscribers'], ascending=False))

# region preference by numbers of visits
display(df2.groupby(['Country', 'Categories'])[['Visits']].sum().sort_values(by=['Country', 'Visits'], ascending=False))

# region preference by numbers of Likes
display(df2.groupby(['Country', 'Categories'])[['Likes']].sum().sort_values(by=['Country', 'Likes'], ascending=False))
```

### Subscribers

Country	Categories	Subscribers
United States	<b>Music and Dance</b>	1425800000
	<b>Video Games, Humor</b>	437300000
	<b>Animation, Video Games</b>	393200000
	<b>Animation, Humor</b>	334000000
	<b>Video Games</b>	302100000
	...	...
Argentina	<b>Movies, Animation</b>	34400000
	<b>Animation</b>	25700000
	<b>Movies, Humor</b>	18800000
	<b>Movies</b>	14200000
Algeria	<b>Education</b>	12200000

151 rows × 1 columns

**Visits**

Country	Categories	Visits
United States	<b>Video Games, Humor</b>	159509700
	<b>Daily Vlogs</b>	86443400
	<b>Animation, Humor</b>	75702600
	<b>Music and Dance</b>	31250600
	<b>Food and Drink</b>	30352400
	...	...
Argentina	<b>Animation</b>	5700000
	<b>Movies, Humor</b>	5600000
	<b>Movies, Animation</b>	594800
	<b>Movies</b>	76600
Algeria	<b>Education</b>	333500

151 rows × 1 columns

Out[19]:

Country	Categories	Likes		
		Subscribers	Visits	Comments
United States	Video Games, Humor	6450500		
	Daily Vlogs	5126400		
	Animation, Humor	2584706		
	Animation, Video Games	1625202		
	Food and Drink	1449177		
	...	...	...	
Argentina	Movies, Humor	401700		
	Animation	208400		
	Movies, Animation	28100		
	Movies	1200		
Algeria	Education	16000		

151 rows × 1 columns

There are regional preferences for specific content categories and it's different based on the Numbers of Subscribers, Visits and Likes

## Performance Metrics:

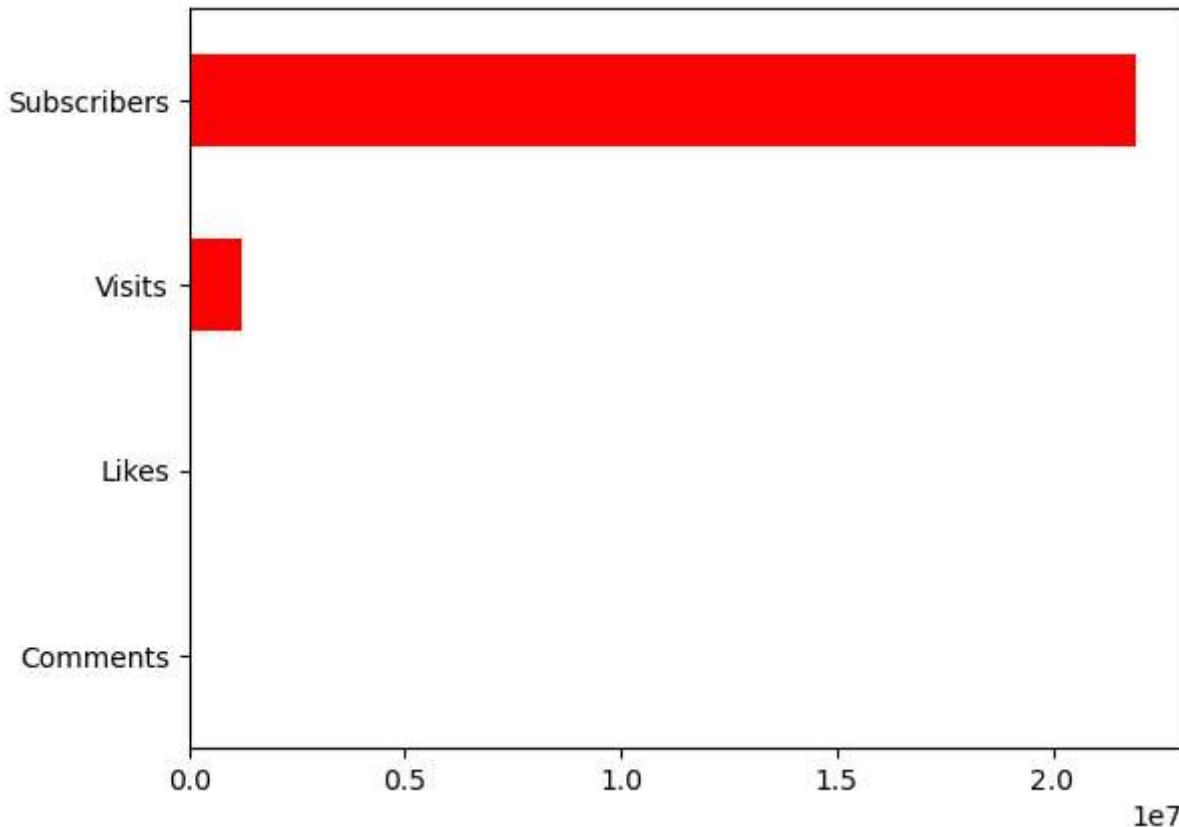
- Calculate and visualize the average number of subscribers, visits, likes, and comments.
- Are there patterns or anomalies in these metrics?

In [20]: `df.select_dtypes('number').mean()`

Out[20]:

Subscribers	21930382
Visits	1215601
Likes	53902
Comments	1296
	dtype: float64

```
In [21]: df.select_dtypes('number').mean().sort_values().plot(kind='barh', color='r');
```



There is an anomalies, the average number of Subscribers is greater than the average number of Visits and Likes.

## Content Categories:

- Explore the distribution of content categories. Which categories have the highest number of streamers?
- Are there specific categories with exceptional performance metrics?

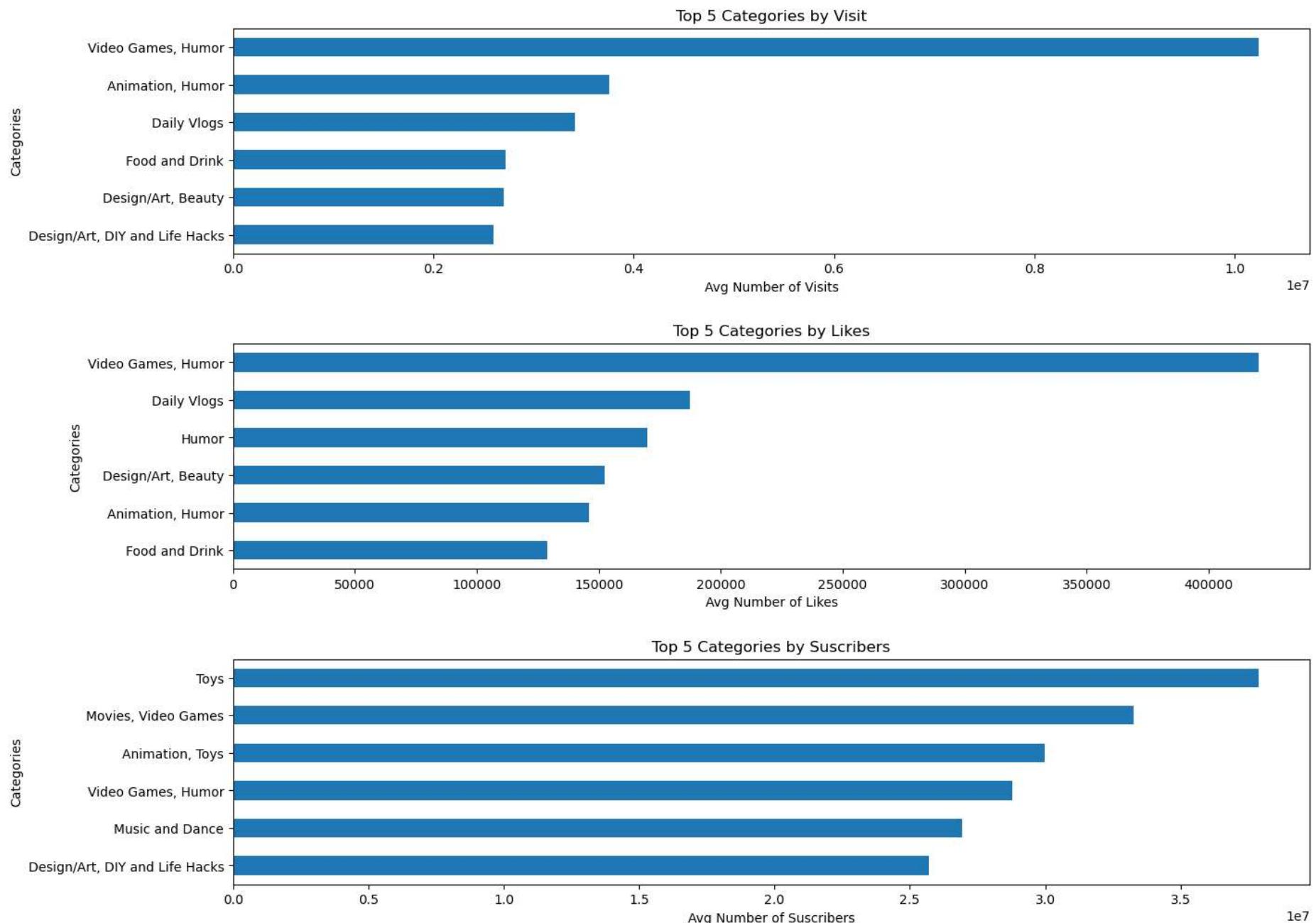
```
In [22]: fig, axs = plt.subplots(figsize=(15,12), nrows=3, ncols=1)
plt.subplots_adjust(hspace=0.4)
df.groupby('Categories')['Visits'].mean().sort_values().tail(6).plot(kind='barh', ax=axs[0])
axs[0].set_title("Top 5 Categories by Visit")
```

```
axs[0].set_xlabel("Avg Number of Visits")

df.groupby('Categories')['Likes'].mean().sort_values().tail(6).plot(kind='barh', ax=axs[1])
axs[1].set_title("Top 5 Categories by Likes")
axs[1].set_xlabel("Avg Number of Likes")

df.groupby('Categories')['Subscribers'].mean().sort_values().tail(6).plot(kind='barh', ax=axs[2])
axs[2].set_title("Top 5 Categories by Subscribers")
axs[2].set_xlabel("Avg Number of Subscribers")

plt.show()
```



Notice how the Top 6 Categories varies by number of Visits, Likes and Subscribers. The "Video Games, Humor" is consistent at the top by number of Visits and Likes.

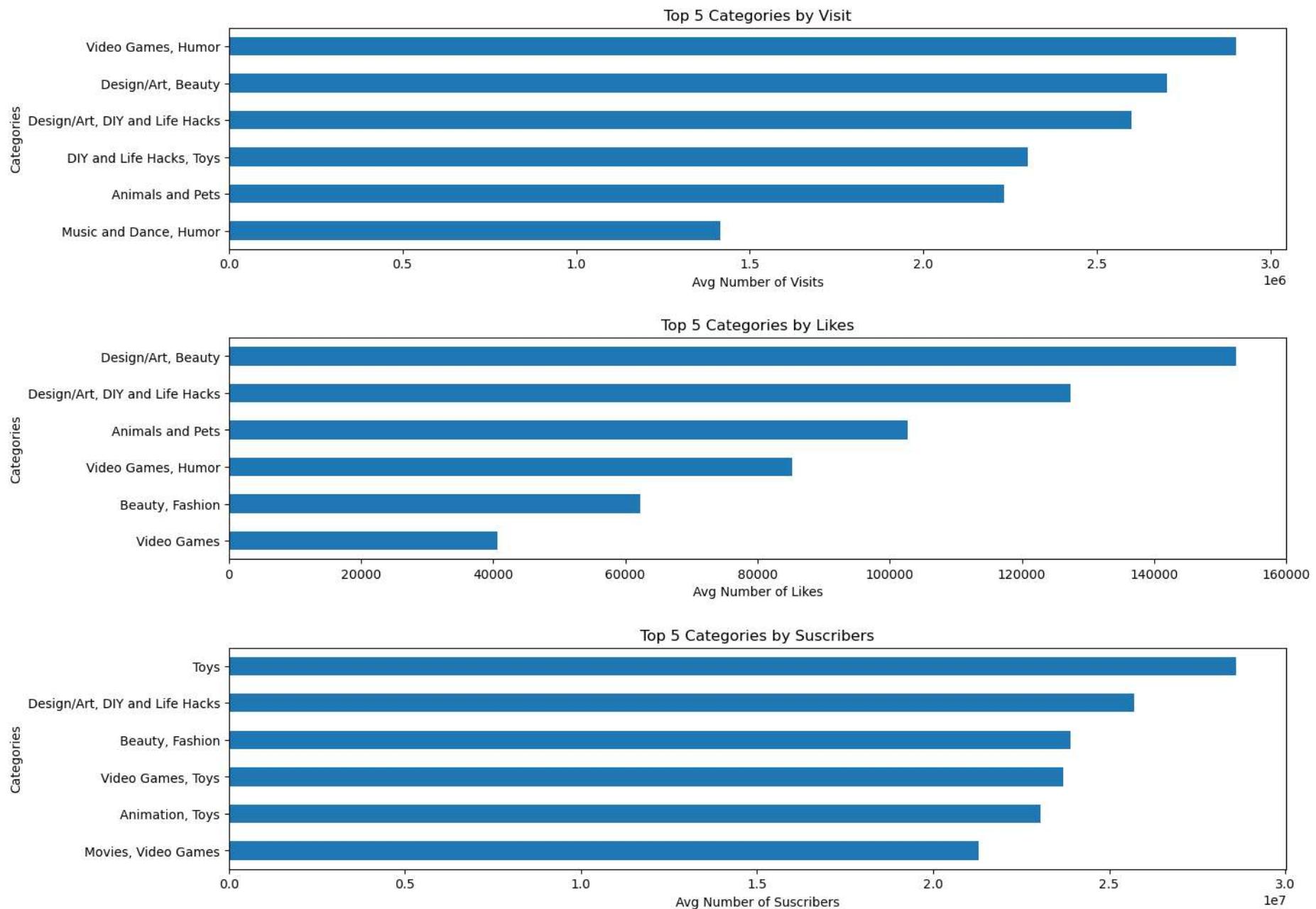
But from our data exploration there are outliers in the data. So i'll visualize the Top 6 using the "Median"

```
In [23]: fig, axs = plt.subplots(figsize=(15,12), nrows=3, ncols=1)
plt.subplots_adjust(hspace=0.4)
df.groupby('Categories')['Visits'].median().sort_values().tail(6).plot(kind='barh', ax=axs[0])
axs[0].set_title("Top 5 Categories by Visit")
axs[0].set_xlabel("Avg Number of Visits")

df.groupby('Categories')['Likes'].median().sort_values().tail(6).plot(kind='barh', ax=axs[1])
axs[1].set_title("Top 5 Categories by Likes")
axs[1].set_xlabel("Avg Number of Likes")

df.groupby('Categories')['Subscribers'].median().sort_values().tail(6).plot(kind='barh', ax=axs[2])
axs[2].set_title("Top 5 Categories by Subscribers")
axs[2].set_xlabel("Avg Number of Subscribers")

plt.show()
```



using the Median I discover that the Top Categories varies also but taking a closer look, it all contains the same Categories but ranking differently except for the numbers of Subscribers where the Toys appears for the first time among the Top 6, eliminating: Animal, Pets and Humor from the

list.

```
In [24]: # To know the specific categories with exceptional performance metrics i obtain the top 10% across all metrics  
df2.select_dtypes('number').quantile(0.90)
```

```
Out[24]: Subscribers    34200000  
Visits        2800000  
Likes         122820  
Comments       3920  
Name: 0.9, dtype: float64
```

```
In [25]: subscribers, visit, likes, comments = round(df2.select_dtypes('number').quantile(0.90))  
  
df_top_cat = df2[(df2['Subscribers'] > subscribers) & (df2['Visits'] > visit) & (df2['Likes'] > likes) & (df2['Comments'] > comments)]  
  
display(df_top_cat)  
  
print("==" * 38)  
print(f"Specific Categories with exceptional performance metrics are: \n\n {df_top_cat['Categories'].unique()}")  
print("==" * 38)
```

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
1	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
26	dudeperfect	Video Games	59700000	United States	5300000	156500	4200	<a href="http://youtube.com/channel/UCRijo3ddMTh_tIHyNS...">http://youtube.com/channel/UCRijo3ddMTh_tIHyNS...</a>
34	TaylorSwift	Music and Dance	54100000	United States	4300000	300400	15000	<a href="http://youtube.com/channel/UCqECaJ8Gagnn7YCbPE...">http://youtube.com/channel/UCqECaJ8Gagnn7YCbPE...</a>
43	A4a4a4a4	Animation, Humor	47300000	Russia	9700000	330400	22000	<a href="http://youtube.com/channel/UC2tsySbe9TNrl-xh2l...">http://youtube.com/channel/UC2tsySbe9TNrl-xh2l...</a>
62	KimberlyLoaiza	Music and Dance	42100000	Mexico	5300000	271300	16000	<a href="http://youtube.com/channel/UCQZfFRohQ7UX-0CdXI...">http://youtube.com/channel/UCQZfFRohQ7UX-0CdXI...</a>

=====

Specific Categories with exceptional performance metrics are:

```
[ 'Video Games, Humor' 'Video Games' 'Music and Dance' 'Animation, Humor' ]
```

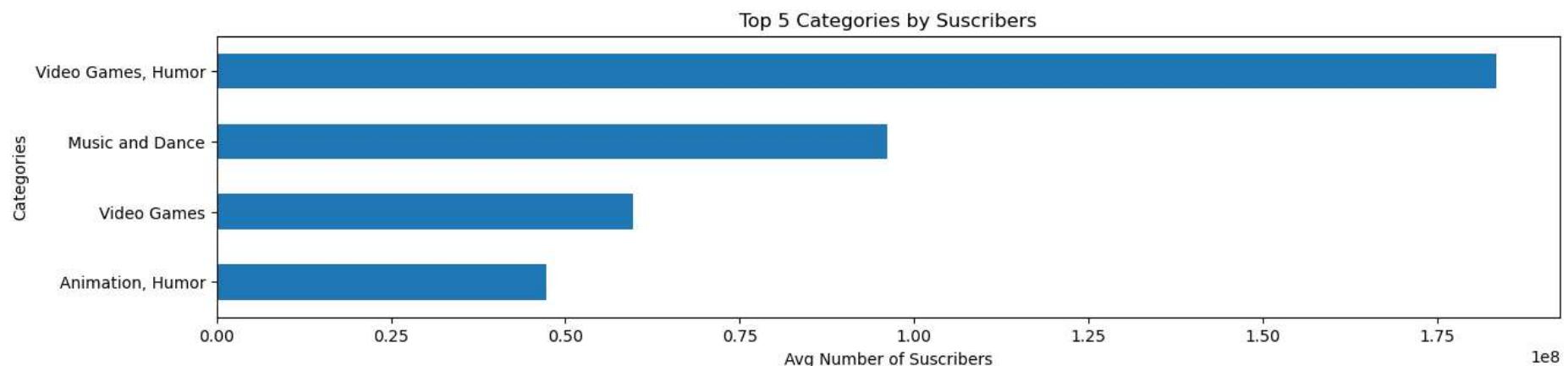
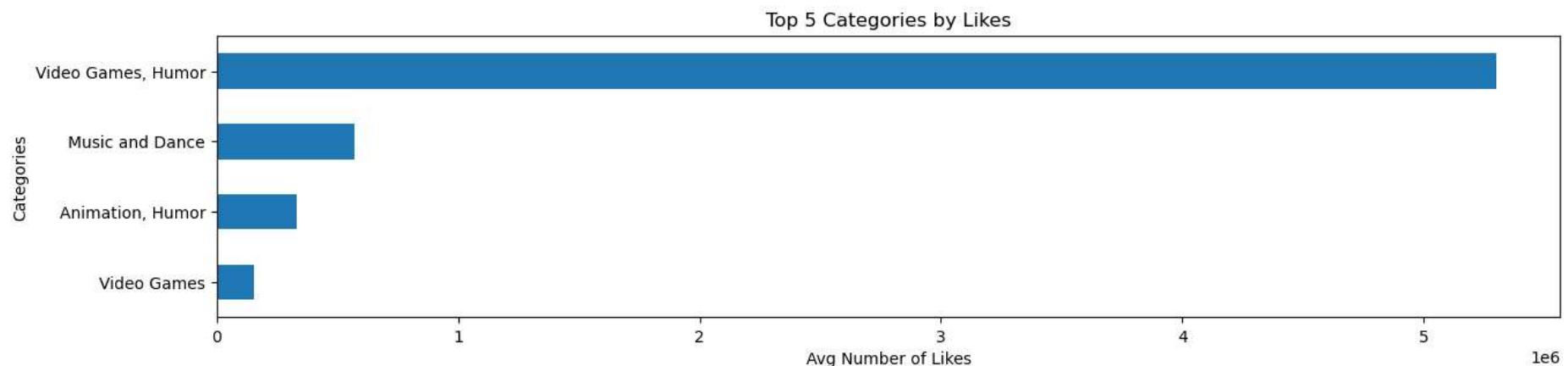
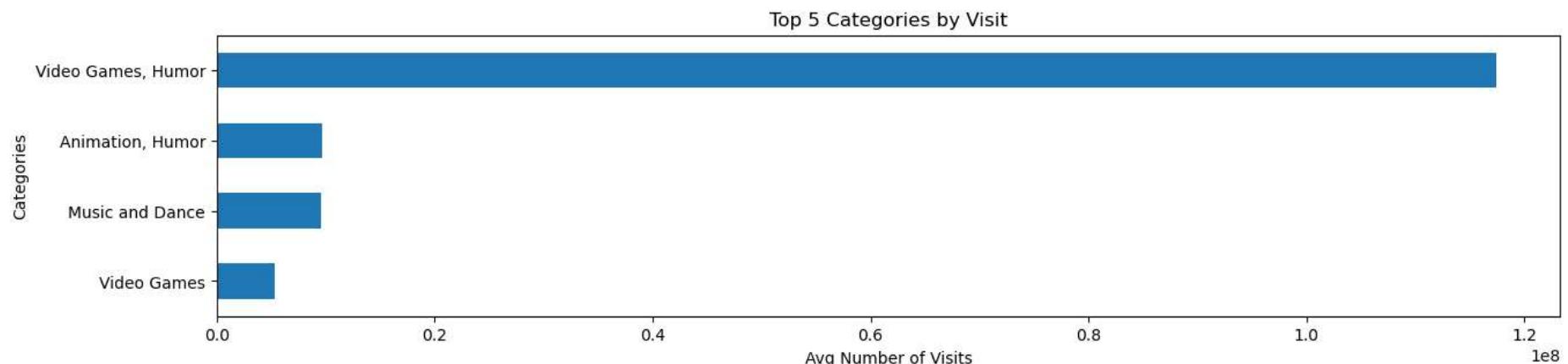
=====

```
In [26]: fig, axs = plt.subplots(figsize=(15,12), nrows=3, ncols=1)
plt.subplots_adjust(hspace=0.4)
df_top_cat.groupby('Categories')['Visits'].sum().sort_values().plot(kind='barh', ax=axs[0])
axs[0].set_title("Top 5 Categories by Visit")
axs[0].set_xlabel("Avg Number of Visits")

df_top_cat.groupby('Categories')['Likes'].sum().sort_values().plot(kind='barh', ax=axs[1])
axs[1].set_title("Top 5 Categories by Likes")
axs[1].set_xlabel("Avg Number of Likes")

df_top_cat.groupby('Categories')['Subscribers'].sum().sort_values().plot(kind='barh', ax=axs[2])
axs[2].set_title("Top 5 Categories by Subscribers")
axs[2].set_xlabel("Avg Number of Subscribers")

plt.show()
```



## Benchmarking:

- Identify streamers with above-average performance in terms of subscribers, visits, likes, and comments.
- Who are the top-performing content creators?

```
In [27]: df.select_dtypes('number').mean()
```

```
Out[27]: Subscribers    21930382
Visits        1215601
Likes         53902
Comments      1296
dtype: float64
```

```
In [28]: sub_mean, visit_mean, likes_mean, comments_mean = round(df.select_dtypes('number').mean())

df_top = df[(df['Subscribers'] > sub_mean) & (df['Visits'] > visit_mean) & (df['Likes'] > likes_mean) & (df['Comments'] > comments_mean)]

# Top performing content creators by Subscribers
df_top_sub = df_top[['Username', 'Subscribers']].nlargest(6, 'Subscribers')
display(df_top_sub)
print("==" * 38)
print(f"Top 6 content creators by numbers of subscribers are: \n {df_top_sub['Username'].unique().tolist()}", end=' '))

# Top performing content creators by Visits
df_top_visit = df_top[['Username', 'Visits']].nlargest(6, 'Visits')
display(df_top_visit)
print("==" * 38)
print(f"Top 6 content creators are by numbers of visit are: \n {df_top_visit['Username'].unique().tolist()}", end=' '))

# Top performing content creators by Likes
df_top_likes = df_top[['Username', 'Likes']].nlargest(6, 'Likes')
display(df_top_likes)
print("==" * 38)
print(f"Top 6 content creators are by numbers of likes are: \n {df_top_likes['Username'].unique().tolist()}", end=' '))
```

**Username Subscribers**

<b>1</b>	MrBeast	183500000
<b>5</b>	PewDiePie	111500000
<b>26</b>	dudeperfect	59700000
<b>34</b>	TaylorSwift	54100000
<b>39</b>	JuegaGerman	48600000
<b>43</b>	A4a4a4a4	47300000

=====

Top 6 content creators by numbers of subscribers are:

[ 'MrBeast', 'PewDiePie', 'dudeperfect', 'TaylorSwift', 'JuegaGerman', 'A4a4a4a4' ]

**Username Visits**

<b>1</b>	MrBeast	117400000
<b>136</b>	MrBeast2	83100000
<b>153</b>	DaFuqBoom	52700000
<b>287</b>	VillageCookingChannel	21500000
<b>277</b>	StokesTwins	11700000
<b>43</b>	A4a4a4a4	9700000

=====

Top 6 content creators are by numbers of visit are:

[ 'MrBeast', 'MrBeast2', 'DaFuqBoom', 'VillageCookingChannel', 'StokesTwins', 'A4a4a4a4' ]

	Username	Likes
1	MrBeast	5300000
136	MrBeast2	5000000
153	DaFuqBoom	1700000
123	MRINDIANHACKER	617400
238	alanbecker	582600
131	fedevigevani	412200

=====

Top 6 content creators are by numbers of likes are:

```
[ 'MrBeast', 'MrBeast2', 'DaFuqBoom', 'MRINDIANHACKER', 'alanbecker', 'fedevigevani' ]
```

In [29]: df\_top

Out[29]:

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
1	MrBeast	Video Games, Humor	183500000	United States	117400000	5300000	18500	<a href="http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...">http://youtube.com/channel/UCX6OQ3DkcsbYNE6H8u...</a>
5	PewDiePie	Movies, Video Games	111500000	United States	2400000	197300	4900	<a href="http://youtube.com/channel/UC-IHJZR3Gqxm24_Vd...">http://youtube.com/channel/UC-IHJZR3Gqxm24_Vd...</a>
26	dudeperfect	Video Games	59700000	United States	5300000	156500	4200	<a href="http://youtube.com/channel/UCRijo3ddMTh_tIHyNS...">http://youtube.com/channel/UCRijo3ddMTh_tIHyNS...</a>
34	TaylorSwift	Music and Dance	54100000	United States	4300000	300400	15000	<a href="http://youtube.com/channel/UCqECaJ8Gagnn7YCbPE...">http://youtube.com/channel/UCqECaJ8Gagnn7YCbPE...</a>
39	JuegaGerman	Movies, Animation	48600000	Mexico	2000000	117100	3000	<a href="http://youtube.com/channel/UCYiGq8XF7YQD00x7wA...">http://youtube.com/channel/UCYiGq8XF7YQD00x7wA...</a>
43	A4a4a4a4	Animation, Humor	47300000	Russia	9700000	330400	22000	<a href="http://youtube.com/channel/UC2tsySbe9TNrl-xh2l...">http://youtube.com/channel/UC2tsySbe9TNrl-xh2l...</a>
58	Mikecrack	Movies, Animation	43400000	Mexico	2200000	183400	1800	<a href="http://youtube.com/channel/UCqJ5zFEED1hWs0KNQC...">http://youtube.com/channel/UCqJ5zFEED1hWs0KNQC...</a>
62	KimberlyLoaiza	Music and Dance	42100000	Mexico	5300000	271300	16000	<a href="http://youtube.com/channel/UCQZfFRohQ7UX-0CdXI...">http://youtube.com/channel/UCQZfFRohQ7UX-0CdXI...</a>
64	luisitocomunica	Unknown	41100000	Mexico	2500000	128900	1800	<a href="http://youtube.com/channel/UCECJDeK0MNapZbpaoZ...">http://youtube.com/channel/UCECJDeK0MNapZbpaoZ...</a>
70	JessNoLimit	Movies, Animation	39600000	Indonesia	1300000	73500	1600	<a href="http://youtube.com/channel/UCvh1at6xpV1ytYOAzx...">http://youtube.com/channel/UCvh1at6xpV1ytYOAzx...</a>
96	TotalGaming093	Movies, Video Games	36300000	India	1500000	129400	4900	<a href="http://youtube.com/channel/UC5c9VIYTSvBSCaoMu...">http://youtube.com/channel/UC5c9VIYTSvBSCaoMu...</a>
98	TechnoGamerzOfficial	Unknown	35600000	India	6200000	341800	16500	<a href="http://youtube.com/channel/UCX8pnu3DYUnx8qy8V...">http://youtube.com/channel/UCX8pnu3DYUnx8qy8V...</a>
100	markiplier	Animation, Video Games	35500000	United States	2100000	126500	3800	<a href="http://youtube.com/channel/UC7_YxT-KID8kRbqZo7...">http://youtube.com/channel/UC7_YxT-KID8kRbqZo7...</a>
122	AboFlah	Animation, Video Games	32700000	Iraq	3300000	382000	11400	<a href="http://youtube.com/channel/UCqq5n-Oe-r1EEHI3yv...">http://youtube.com/channel/UCqq5n-Oe-r1EEHI3yv...</a>

	Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
123	MRINDIANHACKER	Unknown	32600000	India	6500000	617400	26000	<a href="http://youtube.com/channel/UCSiDGb0MnHFGjs4E2W...">http://youtube.com/channel/UCSiDGb0MnHFGjs4E2W...</a>
131	fedevigevani	Animation, Humor	32000000	Mexico	7700000	412200	17000	<a href="http://youtube.com/channel/UCoQm-PeHC-cbJclKJY...">http://youtube.com/channel/UCoQm-PeHC-cbJclKJY...</a>
132	dream	Animation, Video Games	31900000	United States	3300000	309200	19000	<a href="http://youtube.com/channel/UCTkXRDQl0luXxVQrRQ...">http://youtube.com/channel/UCTkXRDQl0luXxVQrRQ...</a>
136	MrBeast2	Daily Vlogs	31300000	United States	83100000	5000000	11600	<a href="http://youtube.com/channel/UC4-79UOIP48-QNGgCk...">http://youtube.com/channel/UC4-79UOIP48-QNGgCk...</a>
145	jacksepticeye	Animation, Humor	30400000	United States	1600000	83400	2300	<a href="http://youtube.com/channel/UCYzPXprvl5Y-Sf0g4v...">http://youtube.com/channel/UCYzPXprvl5Y-Sf0g4v...</a>
153	DaFuqBoom	Animation, Humor	29800000	United States	52700000	1700000	82800	<a href="http://youtube.com/channel/UCsSsgPaZ2GSmO6il8C...">http://youtube.com/channel/UCsSsgPaZ2GSmO6il8C...</a>
176	CrazyXYZ	Unknown	27800000	India	4200000	284100	8600	<a href="http://youtube.com/channel/UCebC4x5l2-PQxg46Uc...">http://youtube.com/channel/UCebC4x5l2-PQxg46Uc...</a>
177	DanTDM	Animation, Video Games	27800000	United States	3500000	285000	52500	<a href="http://youtube.com/channel/UCS5Oz6CHmeoF7vSad0...">http://youtube.com/channel/UCS5Oz6CHmeoF7vSad0...</a>
179	brentrivera	Video Games, Humor	27600000	United States	6400000	154100	5000	<a href="http://youtube.com/channel/UC56D-IHcUvLVFTX_8N...">http://youtube.com/channel/UC56D-IHcUvLVFTX_8N...</a>
180	NichLmao	Daily Vlogs	27500000	United States	1500000	85800	1600	<a href="http://youtube.com/channel/UC6VAIqNQBc7ggiqvsO...">http://youtube.com/channel/UC6VAIqNQBc7ggiqvsO...</a>
195	nickiminaj	Music and Dance	26100000	United States	1600000	98300	7600	<a href="http://youtube.com/channel/UC3jOd7GUMhpgJRBhiL...">http://youtube.com/channel/UC3jOd7GUMhpgJRBhiL...</a>
206	Alejolgoa	Animation	25700000	Argentina	5700000	208400	1700	<a href="http://youtube.com/channel/UCZs0WwC0Dn_noiQE2B...">http://youtube.com/channel/UCZs0WwC0Dn_noiQE2B...</a>
207	ZHCYT	Design/Art, DIY and Life Hacks	25700000	United States	2600000	127300	2200	<a href="http://youtube.com/channel/UCIQubH2NeMmGLTLgNd...">http://youtube.com/channel/UCIQubH2NeMmGLTLgNd...</a>
234	rug	Video Games, Humor	24300000	United States	3200000	85300	5100	<a href="http://youtube.com/channel/UCilwZiBBfl9X6yiZRz...">http://youtube.com/channel/UCilwZiBBfl9X6yiZRz...</a>

		Username	Categories	Subscribers	Country	Visits	Likes	Comments	Links
238		alanbecker	Animation, Video Games	24300000	United States	7600000	582600	5900	<a href="http://youtube.com/channel/UCbKWv2x9t6u8yZoB3K...">http://youtube.com/channel/UCbKWv2x9t6u8yZoB3K...</a>
241		juandediospantojaa	Music and Dance, Movies	24000000	Mexico	3000000	133200	3600	<a href="http://youtube.com/channel/UCzoUWqjCbcfWFdOMvo...">http://youtube.com/channel/UCzoUWqjCbcfWFdOMvo...</a>
265		DrossRotzank	Unknown	23100000	Mexico	1700000	105900	3900	<a href="http://youtube.com/channel/UCNYW2vfGrUE6R5mIJY...">http://youtube.com/channel/UCNYW2vfGrUE6R5mIJY...</a>
271		AmiRodrigueZZ	Animation, Humor	22900000	Colombia	4300000	294400	1300	<a href="http://youtube.com/channel/UCwohKbRK2QnojtLh3p...">http://youtube.com/channel/UCwohKbRK2QnojtLh3p...</a>
277		StokesTwins	Video Games, Humor	22700000	United States	11700000	235000	10000	<a href="http://youtube.com/channel/UCbp9MyKCTEww4CxEzc...">http://youtube.com/channel/UCbp9MyKCTEww4CxEzc...</a>
280		SSundee	Animation, Video Games	22700000	United States	1700000	59800	1800	<a href="http://youtube.com/channel/UCke6I9N4KfC968-yRc...">http://youtube.com/channel/UCke6I9N4KfC968-yRc...</a>
281		souravjoshivlogs7028	Daily Vlogs	22700000	India	5600000	382300	8900	<a href="http://youtube.com/channel/UCjvgGbPPn-FgYeguc5...">http://youtube.com/channel/UCjvgGbPPn-FgYeguc5...</a>
287		VillageCookingChannel	Unknown	22500000	India	21500000	321500	5900	<a href="http://youtube.com/channel/Uck3JZr7eS3pg5AGEvB...">http://youtube.com/channel/Uck3JZr7eS3pg5AGEvB...</a>

In [ ]: