

Overview:

This analysis aims to understand the trend of public tastes in music in 3 main East Asia countries: China, Korea, and Japan. A comparison among artists from these countries was done in this notebook.

- **Finding I:** Japanese and Korean artists tends to have higher popularity and influence regionally and world-wide
- **Finding II:** Follower base for Chinese artists are significantly less than that of Japanese and Korean artists
- **Finding III:** Rap is trending in Japan and Korea in the past few years whereas in China it receives less attention
- **Finding IV:** J-pop has been one of the most popular genres in East Asia in the past 20+ years.
- **Finding V:** There is a rising trend of 'Chinese viral music' (songs that gone viral on the Internet, specifically TikTok and will quickly vanish). However, artists who benefit from this types of music tend to be 'short-lived' in terms of their long-term influence and popularity. Thus, this kind of artists need to find ways to retain public attention and actually accumulate follower bases (e.g. increase music quality, be innovative) to avoid being 'nobody' after the trend is gone.

What I hope I can do:

- I hope to collect data on rough demographic of followers for each artists, and analyzed their influence globally. Although I have said Japanese and Korean artists are popular world-wide, that is half data-supported and half observation. Additional data is needed to solidify this.

Note:

1. This analysis assumes that popularity of artists indicates public music tastes
2. It further assumes that top track is representative to the genres artists are involved in
3. Potential deviation in analysis given the fact that there are many alternative music apps that are popular in mainland China i.e., not a lot of people use Spotify there, which may affect calculation of popularity metrics.

In [1]:



```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
```

```
In [2]: 1 dfc = pd.read_csv('chinese_top100_artist.csv', index_col = 0).reset_index()
2 dfe = pd.read_csv('east_asia_top_artists.csv', index_col = 0).reset_index()
3 dfj = pd.read_csv('japanese_top100_artist.csv', index_col = 0).reset_index()
4 dfk = pd.read_csv('korean_top100_artist.csv', index_col = 0).reset_index()
5
6 df_all = {'China': dfc, 'Japan': dfj, 'Korea': dfk, 'East Asia': dfe}
```

```
In [3]: 1 # east_asia_artist has an extra column 'query_genre', not very informative
2 for df in df_all.values():
3     print(df.shape)
```

(100, 13)

(100, 13)

(100, 13)

(700, 14)

```
In [4]: 1 # example head
2 dfj.head(2)
```

Out[4]:

	artist_name	popularity	followers	artist_link	genres
0	YOASOBI	76	5945744	https://open.spotify.com/artist/64tJ2EAv1R6UaZ...	['j-pop', 'japanese teen pop']
1	Vaundy	70	1731531	https://open.spotify.com/artist/2IUl3m1H1EQ7Qf...	['j-pop', 'japanese soul']

```
In [5]: 1 dfe.head(2)
```

Out[5]:

	artist_name	popularity	followers	artist_link	genres	1
0	BTS	88	67507448	https://open.spotify.com/artist/3Nrfpe0tUJi4K4...	['k-pop', 'k-pop boy group', 'pop']	
1	BLACKPINK	82	43385244	https://open.spotify.com/artist/41MozSoPlsD1dJ...	['k-pop', 'k-pop girl group', 'pop']	

Features (dfc, dfj, dfk, dfe)

- artist_name: unique values
- popularity: int, calculated by Spotify, [0, 100]

- followers: int, # of followers
- artist_link: **Drop**
- genres: used list to contain different genres
- top_track: name of the top track for given artist **Drop?**
- top_track_album: name of top track's album **Drop?**
- top_track_popularity: int, [0, 100]
- top_track_release_date: obj, include different format (e.g. idx 100), should be datetime, extract year only
- top_track_duration_ms: int, duration in milisecond consider convert into second/min
- top_track_explicit: bool, 80%+ 'False' in all four datasets **Drop**
- top_track_album_link: **Drop**
- top_track_link: **Drop**

Features (dfe)

- 'query_genre': Not very informative **Drop**

```
In [6]: 1 for df in df_all.values():
2         # extract year
3         df['top_track_release_year'] = df['top_track_release_date'].str.spl
4
5         # add duration minute
6         df['top_track_duration_min'] = (df['top_track_duration_ms']/60000)
7
8         # drop columns
9         df.drop(['artist_link', 'top_track_album_link', 'top_track_link',
10                'top_track', 'top_track_album', 'top_track_release_date'], axi
11
12 dfe.drop('query_genre', axis = 1, inplace = True)
```

PART II. Analysis

I want to explore the following questions based on what are given in datasets:

1. How are popularity of artists in China, Japan, and Korea compared to each other? How are they compared to the East Asia Standard?
2. Assuming popularity of artists indicates public music tastes, are there any similarity in terms of genres?
3. How has the trend changed over year in these countries?
4. Are there any outlier-artists in dataset? (i.e. in any way different than others)
5. Can we segment these artists? How is each group different?

Note: There are some artists in east asia dataset that are also from China, Japan, or Korea. Hence, we shouldn't overinterpret the result but focus on the general trend, similarity, and differences.

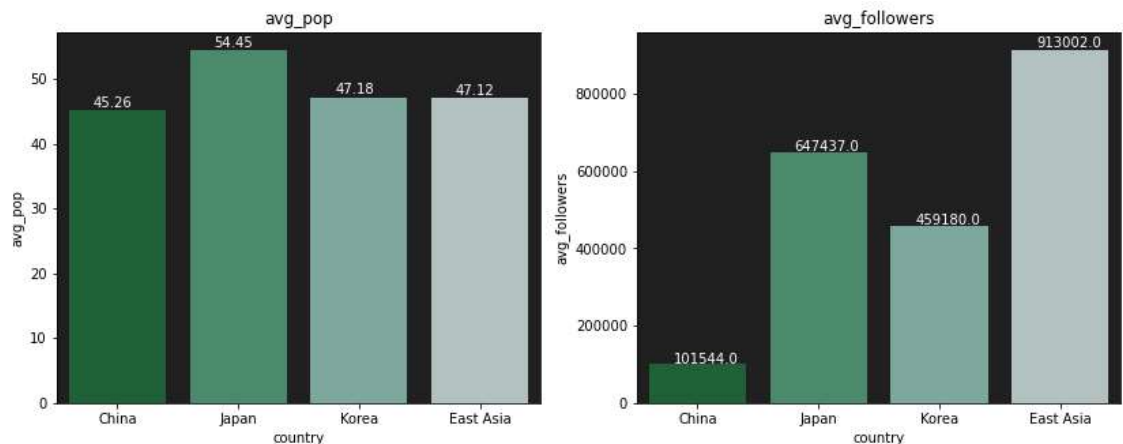
Note2: There might be deviations for Chinese artists as access limitation exists in mainland China.

```
In [7]: ▶ 1 cty_avg = pd.DataFrame({'country': df_all.keys()})
        2 cty_avg['avg_pop'] = [df_all[name]['popularity'].mean().round(2) for name in df_all.keys()]
        3 cty_avg['avg_followers'] = [df_all[name]['followers'].mean().round(2) for name in df_all.keys()]
        4 cty_avg
```

Out[7]:

	country	avg_pop	avg_followers
0	China	45.26	101544.0
1	Japan	54.45	647437.0
2	Korea	47.18	459180.0
3	East Asia	47.12	913002.0

```
In [49]: 1 fig, axes = plt.subplots(1, 2, figsize = (14, 5))
2
3 fig.set_facecolor('white')
4
5 ls = ['avg_pop', 'avg_followers']
6
7 for j in range(len(ls)):
8     p = sns.barplot(data = cty_avg, x = 'country', y = ls[j], ax = axes[j],
9                     palette = 'BuGn_r', alpha = 0.8)
10    p.set(title = ls[j])
11    p.set_facecolor('#202020')
12
13    y = cty_avg[ls[j]].tolist()
14    for i in range(len(y)):
15        p.text(i - 0.2, y[i] + y[i]/100, y[i], color = 'white')
16
17 # sns.barplot(data = cty_avg, x = 'country', y = 'avg_followers', ax = axes[1],
```



1. How are popularity of artists in China, Japan, and Korea compared to each other? How are they compared to the East Asia Standard?

Japanese artists have the highest average popularity and # of followers among the three. Korean and Chinese artists have fairly close popularity (to the east asia standard as well) whereas Japan stands out by 15% higher than the East Asia average.

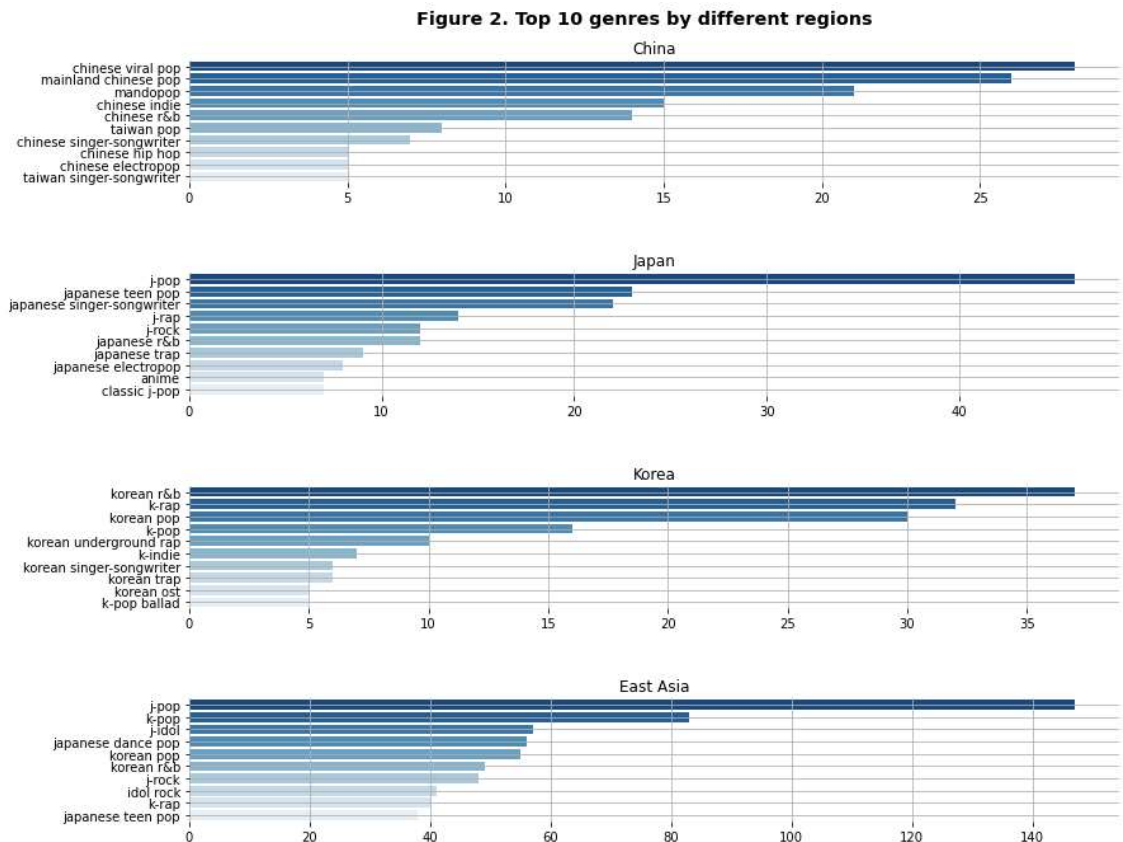
The difference is more significant in average # followers. Chinese artists has the lowest avg # followers (77% and 84% lower than Korea and Japan respectively)

```
In [9]: ▶ 1 # extract and count genres and total popularity from each artists, ret
2 def get_genre_counts(df):
3     df = df.reset_index(drop=True) #make sure filtered df also work
4     genres = dict()
5     popularity = dict()
6
7     for i in range(len(df)):
8         genre_list = df.genres[i].strip("'").split(", ")
9         for g in genre_list:
10             if not g in genres.keys():
11                 genres[g] = 1
12                 popularity[g] = df.top_track_popularity[i]
13             else:
14                 genres[g] += 1
15                 popularity[g] += df.top_track_popularity[i]
16     genres = pd.DataFrame({'genres': genres.keys(),
17                           'total_popularity': popularity.values(),
18                           'count': genres.values()}).sort_values(by =
19
20     return genres
```

```

In [42]: 1 fig, axes = plt.subplots(4, 1, figsize = (12, 10))
2
3 fig.tight_layout(pad = 5)
4 fig.suptitle('Figure 2. Top 10 genres by different regions', weight =
5 fig.set_facecolor('white')
6
7 j = 0
8
9 for i in df_all:
10     genres = get_genre_counts(df_all[i]).head(10)
11     p = sns.barplot(data = genres, x = 'count', y = 'genres', ax = axes[j])
12     p.set(title = f'{i}', xlabel = '', ylabel = '')
13     p.grid(visible = True)
14
15     axes[j].spines['right'].set_color(None)
16     axes[j].spines['top'].set_color(None)
17     axes[j].spines['bottom'].set_color(None)
18     axes[j].spines['left'].set_color(None)
19
20     j += 1

```



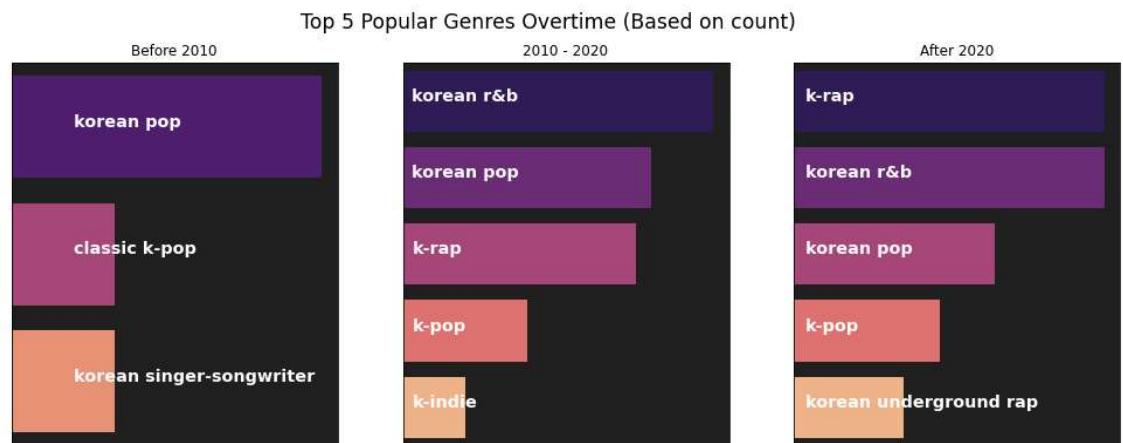
2. Assuming popularity of artists indicates public music tastes, are there any similarity in terms of genres?

- From East Asia list, it's very impressive to see top 10 genres are all from Japan and Korean, suggesting their high influence in east Asia (and world-wide)
- Rap appears to be more popular in Japan and Korean than in China (or that they are more popular in these 2 countries)

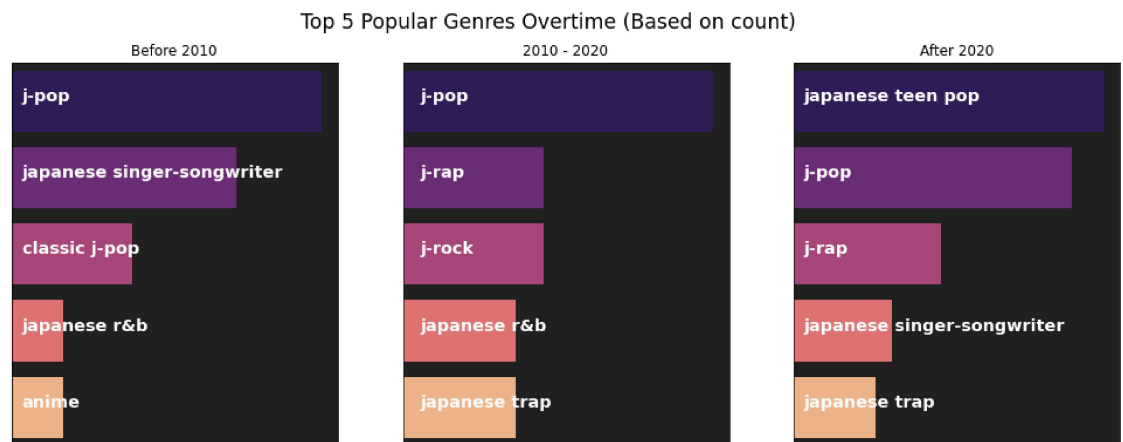
- Common genres in all three genres: pop & r&b
- Special genres in Japan: anime, teen pop
- Special genres in China: viral pop (?), mandopop
- trap which is a type of hip pop also makes its way into top 10 list for Japan and Korea

```
In [45]: ▶ 1 def plot_top_genres_overtime(df, by = 'count'): #another option is 'to
2
3     before_2010 = get_genre_counts(df[df['top_track_release_year'] <=
4     between_2010_and_2020 = get_genre_counts(df[(df['top_track_release
5                                     (df['top_track_release
6     after_2020 = get_genre_counts(df[(df['top_track_release_year'] > '
7
8     period = ['Before 2010', '2010 - 2020', 'After 2020']
9     datasets = [before_2010, between_2010_and_2020, after_2020]
10
11     fig, axes = plt.subplots(1, 3, figsize = (17, 6))
12     fig.set_facecolor('white')
13
14     fig.suptitle(f'Top 5 Popular Genres Overtime (Based on {by})', for
15
16     for i in range(len(datasets)):
17         dataset = datasets[i]
18         data = dataset.head(5)
19         p = sns.barplot(data = data, x = 'count', y = 'genres', orient
20
21         axes[i].xaxis.set_visible(False)
22         axes[i].yaxis.set_visible(False)
23
24         p.set_yticklabels('')
25         p.set_xticklabels('')
26         p.set(title = period[i], ylabel = '', xlabel = '')
27         p.set_facecolor('#202020')
28
29         text = list(data['genres'].values)
30         for j in range(len(text)):
31             axes[i].text(x = 0.6, y = j, s = text[j], fontsize = 'x-la
32
33     return None
```

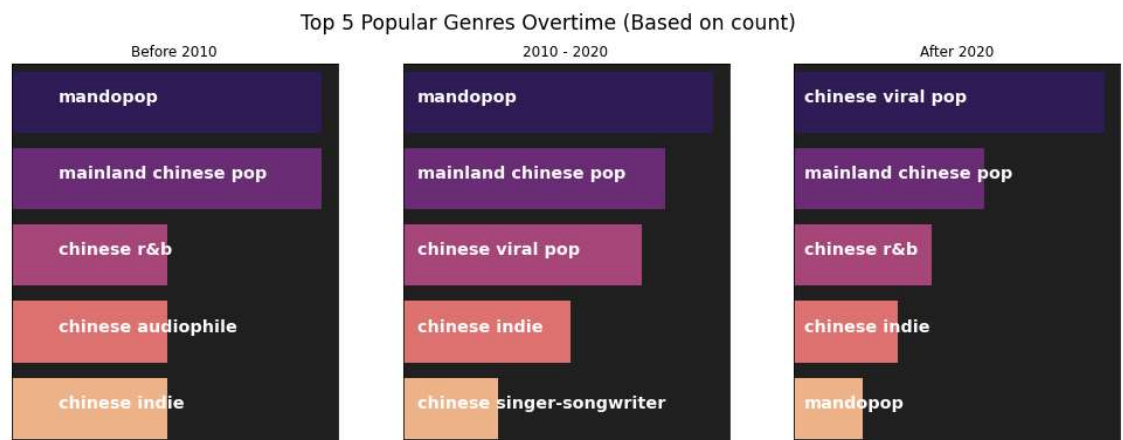

In [46]: 1 plot_top_genres_overtime(dfk, by = 'count')



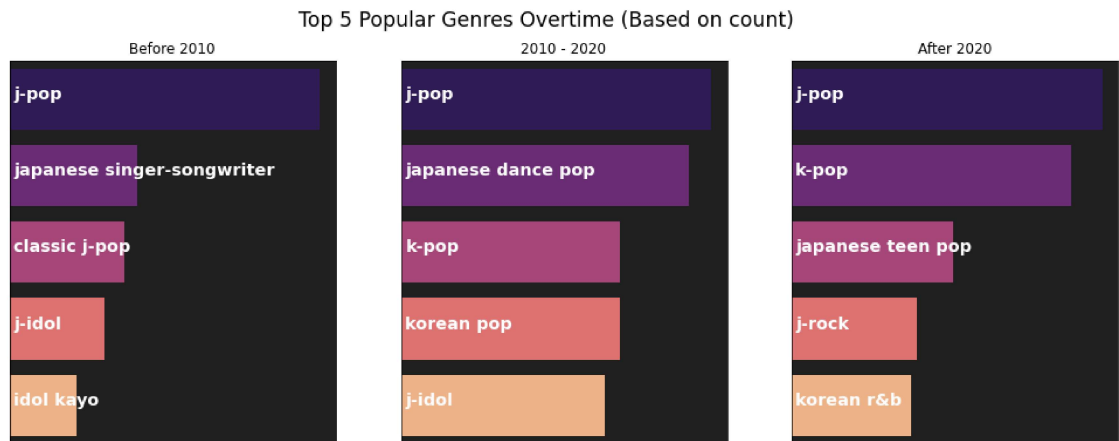
In [13]: 1 plot_top_genres_overtime(dfj, by = 'count')



In [47]: 1 plot_top_genres_overtime(dfk, by = 'count')



```
In [15]: 1 plot_top_genres_overtime(dfe, by = 'count')
```



3. How has the trend changed over years in these countries?

Note: we should be aware that there are some overlapping genres such as 'Korean Pop' and 'k-pop'

Korean

- Korean pop is gradually giving its way to k-rap. There is a rising trend in rap as there are two types of them (rap, underground rap) appearing in the top 5 list after 2023

Japan

- While there is also a rising trend of rap in Japan, j-pop has an unassailable leading advantages. Only that there is changes with in j-pop. A rising of japanese teen pop is seen after 2020.

China

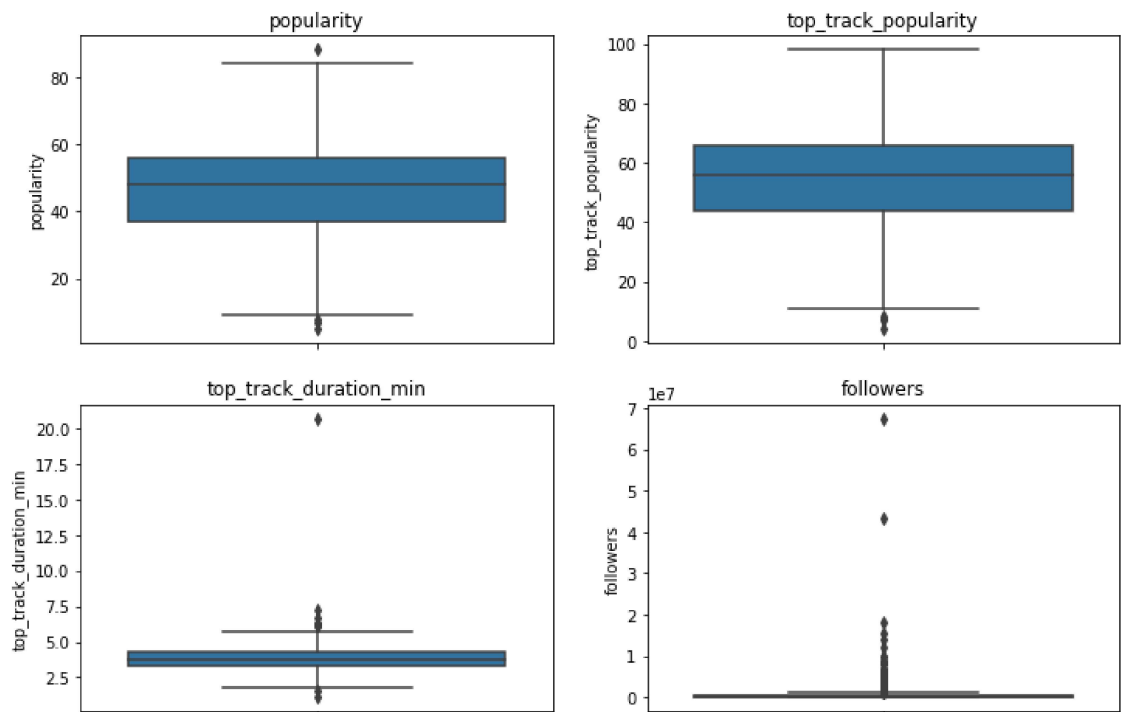
- While the rank may have slightly changed over time, the main categories the public enjoyed listening have no big changes (e.g. pop, r&b, indie).
- It's noticeable that mandopop has dropped to the 5th place after 2020 and chinese viral pop has slowly moved its way up the first place. 这与近几年的音乐流行趋势相符，可能由于抖音等短视频的出现，'网红音乐'变得越来越流行，相反流行华语音乐逐渐式微。 However, due to insufficient data, a firm conclusion can not be drawn.

East Asia

- Leading position of j-pop even shows up on the East Asia list, suggesting that it's very well received regionally as well.
- Unfortunately, Chinese artists didn't make it to the list.

```
In [16]: 1 # add a column called 'country' before concat
2 for country in df_all:
3     df_all[country]['country'] = country
4
5 dfALL = pd.concat(df_all, axis = 0).reset_index(drop = True).drop_duplicates()
```

```
In [17]: 1 # Univariate Analysis on outliers
2 fig, axes = plt.subplots(2, 2, figsize = (12, 8))
3 y = ['popularity', 'top_track_popularity', 'top_track_duration_min', 'followers']
4
5 for i in range(len(y)):
6     sns.boxplot(data = dfALL, y = y[i], ax = axes[i//2, i%2])
7     axes[i//2, i%2].set_title(y[i])
```

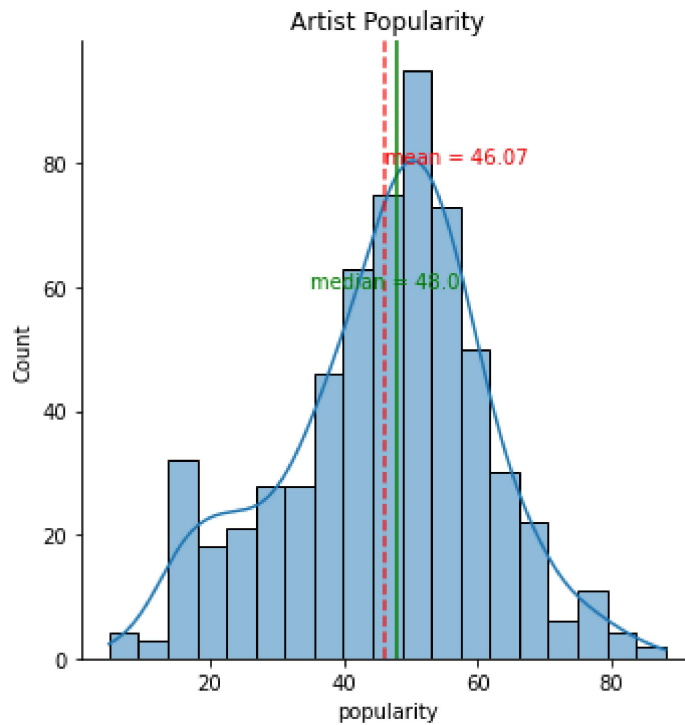


```

In [18]: 1 p = sns.displot(dfALL['popularity'], kde = True)
2 p.set(title = 'Artist Popularity')
3
4 avg = dfALL['popularity'].mean().round(2)
5 plt.axvline(avg, color = 'red', linestyle = '--', alpha = 0.8)
6 plt.text(46, 80, f"mean = {avg}", color = 'red')
7
8 med = dfALL['popularity'].median().round(2)
9 plt.axvline(med, color = 'green')
10 plt.text(35, 60, f"median = {med}", color = 'green')

```

Out[18]: Text(35, 60, 'median = 48.0')

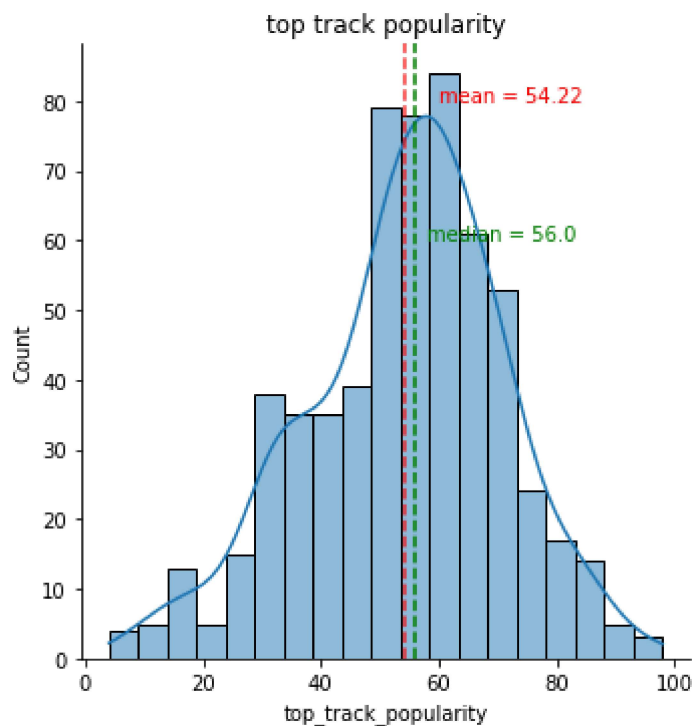


```

In [19]: 1 p = sns.displot(dfALL['top_track_popularity'], kde = True)
2 p.set(title = 'top track popularity')
3
4 avg = dfALL['top_track_popularity'].mean().round(2)
5 plt.axvline(avg, color = 'red', linestyle = '--', alpha = 0.8)
6 plt.text(60, 80, f"mean = {avg}", color = 'red')
7
8 plt.axvline(dfALL['top_track_popularity'].median(), color = 'green', linestyle = '--', alpha = 0.8)
9 plt.text(58, 60, f"median = {dfALL['top_track_popularity'].median().round(0)}", color = 'green')

```

Out[19]: Text(58, 60, 'median = 56.0')



```

In [20]: 1 dfALL[dfALL['top_track_duration_min'] > 15]

```

Out[20]:

	artist_name	popularity	followers	genres	top_track_popularity	top_track_duration_ms
234	The Quiett	44	109376	['k-rap', 'korean r&b', 'korean trap']	48.0	1240555.0

```
In [21]: 1 dfALL[dfALL['followers'] > 3e7]
```

Out[21]:

	artist_name	popularity	followers	genres	top_track_popularity	top_track_duration_ms
300	BTS	88	67507448	['k-pop', 'k-pop boy group', 'pop']	88.0	154486.0
301	BLACKPINK	82	43385244	['k-pop', 'k-pop girl group', 'pop']	79.0	175889.0

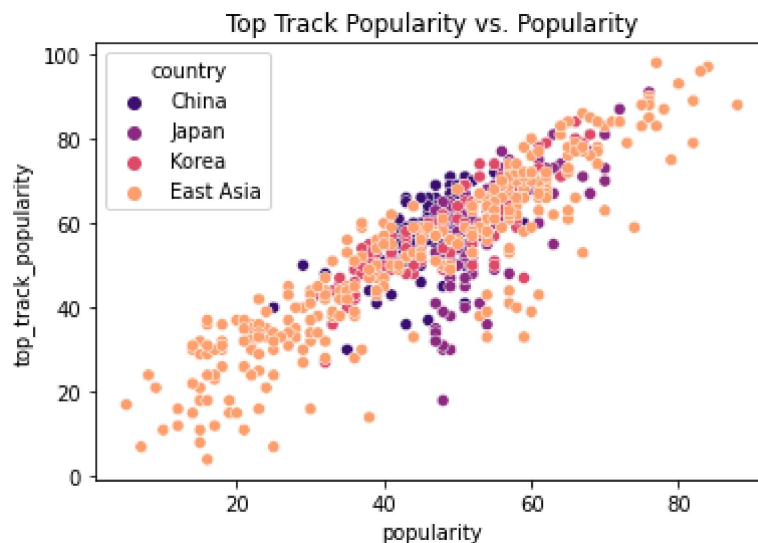
4. Are there any outlier-artists in dataset? (i.e. in any way different than others)

There are 3 noticeable outliers from the plots above.

- Duration of top track by **the Quiett** has more than 20 minutes
- **BTS** and **BLACKPINK** have significantly larger follower bases compared to the rest

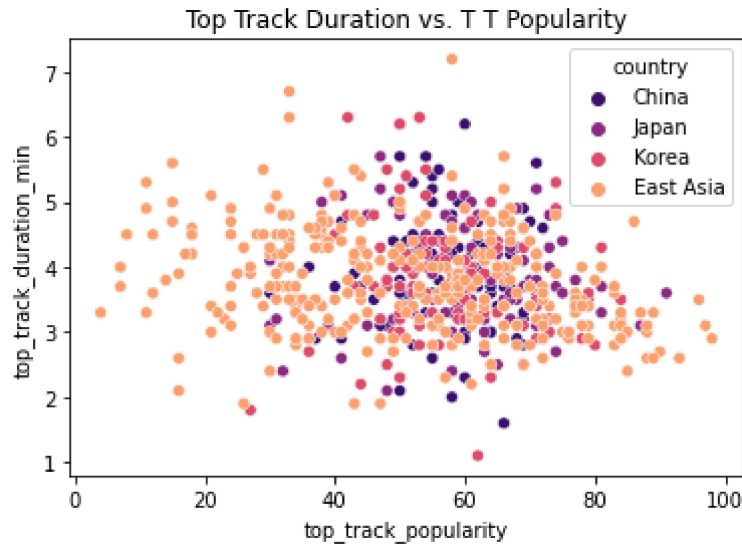
```
In [22]: 1 # Bivariate: Are there any outlier-artists in data set?
2 p = sns.scatterplot(data = dfALL, x = 'popularity', y = 'top_track_popularity',
3                   hue = 'country', palette = 'magma')
4 p.set(title = 'Top Track Popularity vs. Popularity')
```

Out[22]: [Text(0.5, 1.0, 'Top Track Popularity vs. Popularity')]



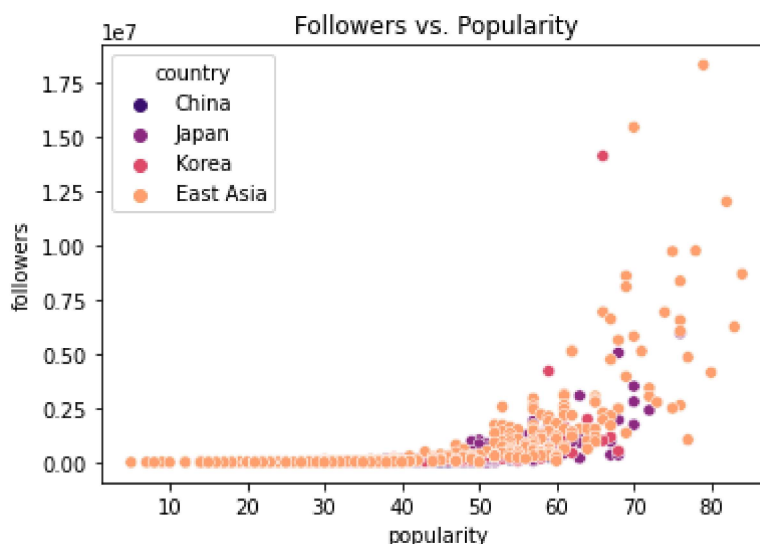
```
In [23]: 1 # removed outlier
2 p = sns.scatterplot(data = dfALL[dfALL['top_track_duration_min'] < 10],
3                    hue = 'country', palette = 'magma')
4 p.set(title = 'Top Track Duration vs. T T Popularity')
```

Out[23]: [Text(0.5, 1.0, 'Top Track Duration vs. T T Popularity')]



```
In [24]: 1 # removed 2 outliers
2 p = sns.scatterplot(data = dfALL[dfALL['followers'] < 3e7], x = 'popu[
3                    hue = 'country', palette = 'magma')
4 p.set(title = 'Followers vs. Popularity')
```

Out[24]: [Text(0.5, 1.0, 'Followers vs. Popularity')]



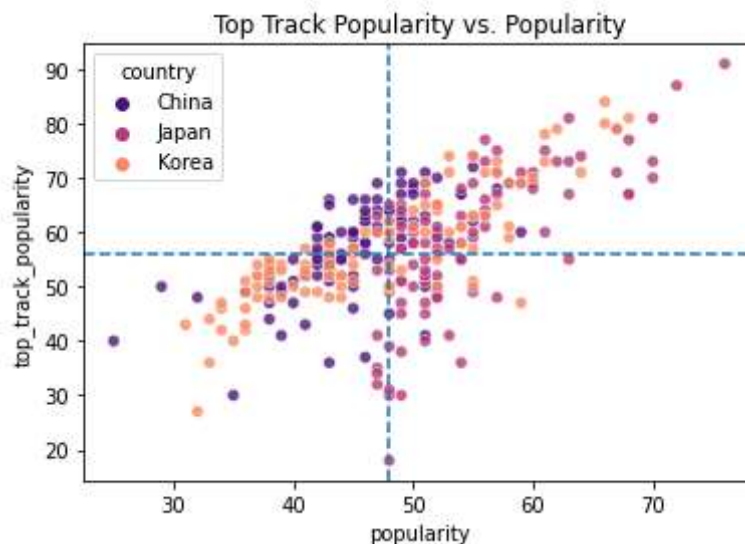
Add-on 4

- By plotting popularity vs. top-track-popularity, we can clearly see the upwards trend. These two variables are definitely positively correlated with each other but causal inference can't be made

- It looks like there is an exponential relationship between followers and popularity, i.e. the higher the popularity, more followers are required to increase popularity score. This could correspond to how Spotify calculate this metrics.
- Another outlier: **V** from Korean has disporportionally high number of followers compared to its popularity score

```
In [48]: 1 # because there are mix of artists from different countries in East Asia
2 fig, ax = plt.subplots()
3
4 fig.set_facecolor('white')
5
6 p = sns.scatterplot(data = dfALL[dfALL['country'] != 'East Asia'],
7                     x = 'popularity', y = 'top_track_popularity',
8                     hue = 'country', palette = 'magma', alpha = 0.8)
9 p.set(title = 'Top Track Popularity vs. Popularity')
10
11 plt.axvline(dfALL['popularity'].median(), linestyle = '--')
12 plt.axhline(dfALL['top_track_popularity'].median(), linestyle = '--')
```

Out[48]: <matplotlib.lines.Line2D at 0x2877e5913c8>



5. Artists Segmentation

For artists in these three countries, I split them into 4 quadrants by median of popularity score and top track popularity score.

- 1st quadrant: 'On Fire'
- 2nd quadrant: 'Flash' artists, their songs or music may become trending at one point but it doesn't last long
- 3rd quadrant: Less popular artists (among all popular ones)
- 4th quadrant: 'High potential'

Something I found interesting:

1. A lot of Chinese artists are on 'Flash' quadrant, which confirm my previous findings. A lot of artists may have one of their songs gone viral on the Internet at a point but soon it will be

replaced by another song. This popularity didn't help them to accumulate their follower/fan bases or click/like (depending on how Spotify calculate popularity score). (People may listen to it because it was trending at that time not because they truly enjoy their music)

2. Almost all top 100 Japanese artists are in the 1st or 4th quadrant, which also align with the fact that Japanese music is well received regionally and even world-wide.