

イントロダクション

コンペティションのタイトルからは分かりにくいのですが、今回のコンペティションの目的は「**広告がクリックされたあとにダウンロードされるかされないか**」を予測することです。

不正クリックかどうかを検知するのに、なぜ広告がクリックされた後アプリがダウンロードされるかを調べる必要があるのでしょうか？

まずドメイン知識として、「なぜ不正クリックをするのか」を調べました。

不正クリックを行って利益を得ようとする人は大きく「競合アプリの開発者」と「広告ネットワーク会社」に分かれます。

1. 競合アプリの開発者

- 広告出稿主は広告費として一日の上限金額を設定しています。1クリックごとに金額が発生するため、一定回数以上クリックされるとその広告は表示されなくなります。これを利用して競合アプリの広告を表示させなくすることで自分たちのアプリの広告を表示させようとしています。

2. 広告ネットワーク会社

- 広告ネットワーク会社は広告出稿主からクリックごとに料金をもらえます。そのため不正クリックでクリックを水増しすることで利益を得ようとします。

[広告業界の裏事情 \(http://www.china-webby.com/%E7%99%BE%E5%BA%A6sem%E5%BA%83%E5%91%8A%E3%81%AE%E8%A3%8F\)](http://www.china-webby.com/%E7%99%BE%E5%BA%A6sem%E5%BA%83%E5%91%8A%E3%81%AE%E8%A3%8F)
[不正業者の実態 \(https://www.gizmodo.jp/2017/06/thai-click-fraud-farm-busted-using-wall-of-iphones.html\)](https://www.gizmodo.jp/2017/06/thai-click-fraud-farm-busted-using-wall-of-iphones.html)

概要でも説明がありましたが、不正業者はなるべく多くのスマートフォンを使って膨大なクリックを発生させたいので、広告をクリックしてもアプリをダウンロードすることは決してありません。

つまり**広告をクリックしたあとアプリをダウンロードされないクリックは不正クリック**とみなします。

このことを念頭においてデータを見ていきます。

全体の俯瞰

train.csvは約2億件のデータがあるのでまずtrain_sample.csvを見ていきます。

In [33]:

```
1 # 必要なモジュールのインポート
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 %matplotlib inline
7
8 import warnings
9 warnings.filterwarnings('ignore')
```

In [34]:

```
1 ls ./competitions/talkingdata-adtracking-fraud-detection/
```

```
sample_submission.csv.zip train.csv.zip
test.csv.zip             train_sample.csv.zip
test_supplement.csv.zip
```

In [35]:

```
1 data_dir_path = './competitions/talkingdata-adtracking-fraud-detection/'
```

In [36]:

```
1 sample_df = pd.read_csv(data_dir_path + 'train_sample.csv.zip')
```

In [37]:

```
1 sample_df.head()
```

Out[37]:

	ip	app	device	os	channel	click_time	attributed_time	is_attributed
0	87540	12	1	13	497	2017-11-07 09:30:38	NaN	0
1	105560	25	1	17	259	2017-11-07 13:40:27	NaN	0
2	101424	12	1	19	212	2017-11-07 18:05:24	NaN	0
3	94584	13	1	13	477	2017-11-07 04:58:08	NaN	0
4	68413	12	1	1	178	2017-11-09 09:00:09	NaN	0

attributed_timeはダウンロードされたときの時間なのでis_attributedが0のものはNaNになっています。

In [38]:

```
1 sample_df.describe()
```

Out[38]:

	ip	app	device	os	channel	is_attribute
count	100000.000000	100000.000000	100000.000000	100000.000000	100000.000000	100000.000000
mean	91255.879670	12.04788	21.771250	22.818280	268.832460	0.00227
std	69835.553661	14.94150	259.667767	55.943136	129.724248	0.04759
min	9.000000	1.00000	0.000000	0.000000	3.000000	0.00000
25%	40552.000000	3.00000	1.000000	13.000000	145.000000	0.00000
50%	79827.000000	12.00000	1.000000	18.000000	258.000000	0.00000
75%	118252.000000	15.00000	1.000000	19.000000	379.000000	0.00000
max	364757.000000	551.00000	3867.000000	866.000000	498.000000	1.00000

train_sample.csvのデータは10万件あり、click_timeとattributed_time以外は数値データで正の数値です。

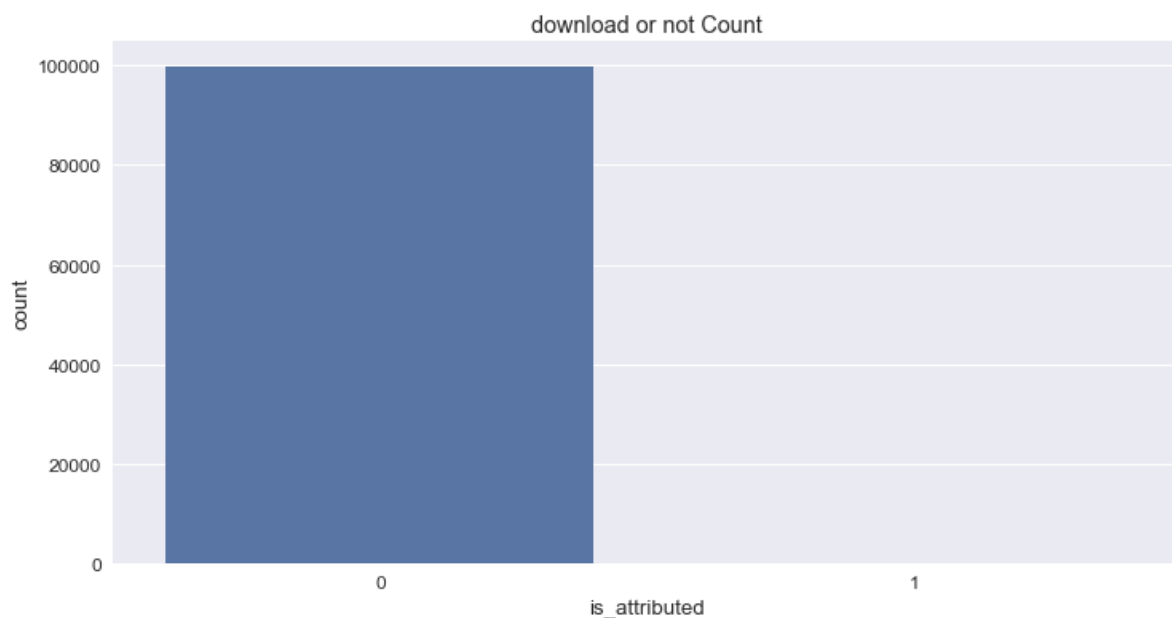
次に広告がクリックされた後ダウンロードされた件数とされていない件数を見比べてみる。

In [39]:

```
1 plt.figure(figsize=(12, 6))
2 click_count = sample_df.groupby('is_attributed', as_index=False)['ip'].count()
3 click_count = click_count.rename(columns={'ip':'count'})
4 sns.barplot(x='is_attributed', y='count', data= click_count);
5 plt.title('download or not Count');
6 click_count.head()
```

Out[39]:

	is_attributed	count
0	0	99773
1	1	227



クリック後にダウンロードされた数はかなり少なくて実際には227件なので約0.23%しかありません。

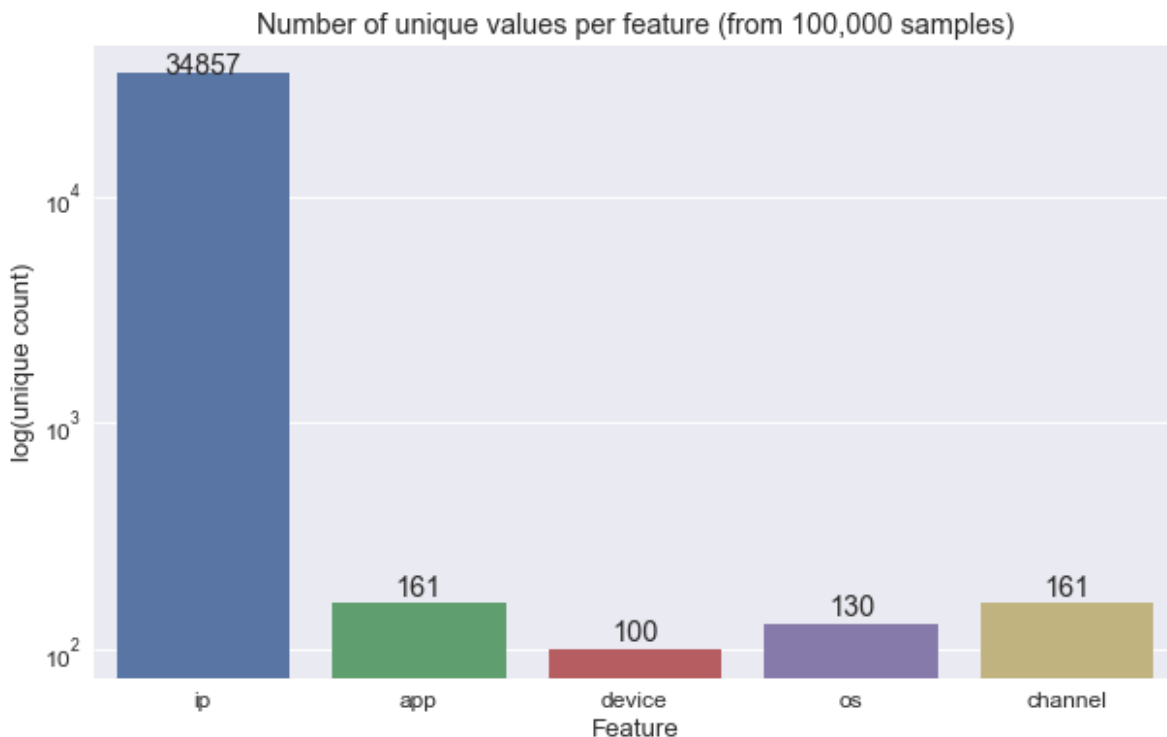
次にそれぞれのユニークな値の件数を見てみましょう。

In [40]:

```

1 plt.figure(figsize=(10, 6))
2 cols = ['ip', 'app', 'device', 'os', 'channel']
3 uniques = [len(sample_df[col].unique()) for col in cols]
4 sns.set(font_scale=1.2)
5 ax = sns.barplot(cols, uniques, log=True)
6 ax.set(xlabel='Feature', ylabel='log(unique count)', title='Number of unique values per feature (from
7 for p, uniq in zip(ax.patches, uniques):
8     height = p.get_height()
9     ax.text(p.get_x()+p.get_width()/2.,
10            height + 10,
11            uniq,
12            ha="center")

```



ipアドレスが一番多く、それ以外は100~200に収まっています。

時間帯でのクリック数とダウンロード数

①目的であるならば人が寝ている時間に不正クリックが発生しているかもしれません。click_timeを分解して、日、時、分の項目を新しく作成します。ちなみにデータは2017年11月6日から9日の平日で中国で特別なイベントはありません。

In [41]:

```

1 sample_df['click_time_dt'] = pd.to_datetime(sample_df['click_time'])
2 dt = sample_df['click_time_dt'].dt
3 sample_df['day'] = dt.day.astype('uint8')
4 sample_df['hour'] = dt.hour.astype('uint8')
5 sample_df['minute'] = dt.minute.astype('uint8')

```

アプリをダウンロードしたデータと ダウンロードしていないデータで分けてデータを見てみましょう。

In [42]:

```

1 dll_df = sample_df[sample_df['is_attributed'] == 1]
2 not_dll_df = sample_df[sample_df['is_attributed'] == 0]

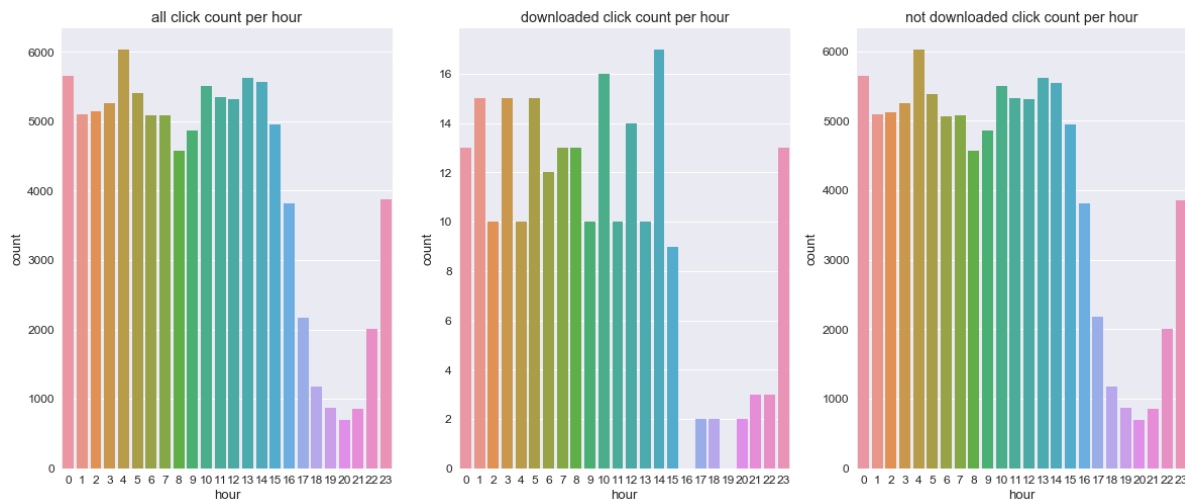
```

In [75]:

```

1 fig ,ax = plt.subplots(1,3, figsize= (20, 8))
2
3
4 sns.countplot(x='hour', data=sample_df, ax=ax[0]);
5 sns.countplot(x='hour', data=dll_df, order=[i for i in range(24)] , ax=ax[1]);
6 sns.countplot(x='hour', data=not_dll_df, ax=ax[2]);
7 ax[0].set_title('all click count per hour');
8 ax[1].set_title('downloaded click count per hour');
9 ax[2].set_title('not downloaded click count per hour');
10

```



時間帯ごとの全体のクリック数、ダウンロードされたクリック数、ダウンロードされなかったクリック数をそれぞれプロットしています。18～22時あたりがクリック数が共通して落ち込んでいますが、3つのグラフで大きな違いは見られません。

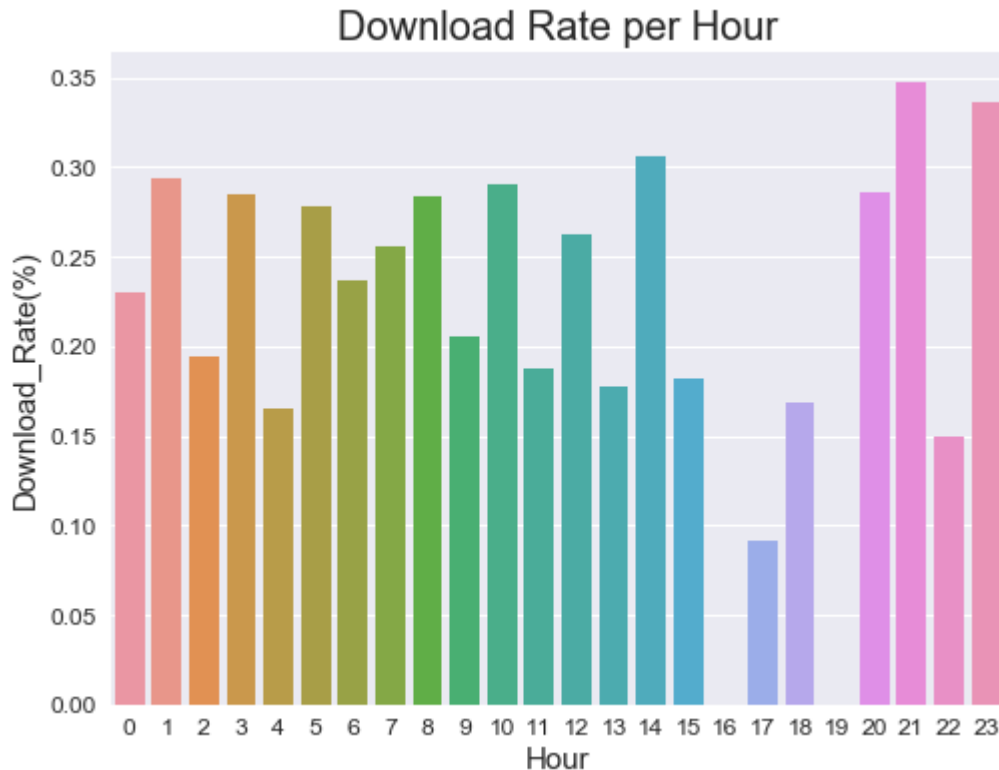
次に時間帯ごとのクリック数に対するダウンロード率をプロットしてみましょう。

In [67]:

```

1 # クリック後のダウンロード率を時間ごとに確認する
2 count_per_hour = sample_df.groupby('hour', as_index = False).count()
3 count_per_hour = count_per_hour[['hour', 'app', 'attributed_time']]
4 count_per_hour = count_per_hour.rename(columns = {'app': 'click_count', 'attributed_time': 'download_count'})
5 count_per_hour['download_rate'] = count_per_hour['download_count'] / count_per_hour['click_count']
6
7 fig, ax = plt.subplots(figsize=(8, 6))
8 sns.barplot(x='hour', y='download_rate', data=count_per_hour);
9 plt.ylabel('Download_Rate(%)', fontsize=15)
10 plt.xlabel('Hour', fontsize=15)
11 ax.set_title('Download Rate per Hour', fontsize = 20);

```



夕方近辺はダウンロード数も落ち込んでいますが、ダウンロード率も落ち込んでいます。これはあまり良いグラフではないね。ダウンロード率が小さすぎるので1件でダウンロード率が大きく変わってしまうからね。

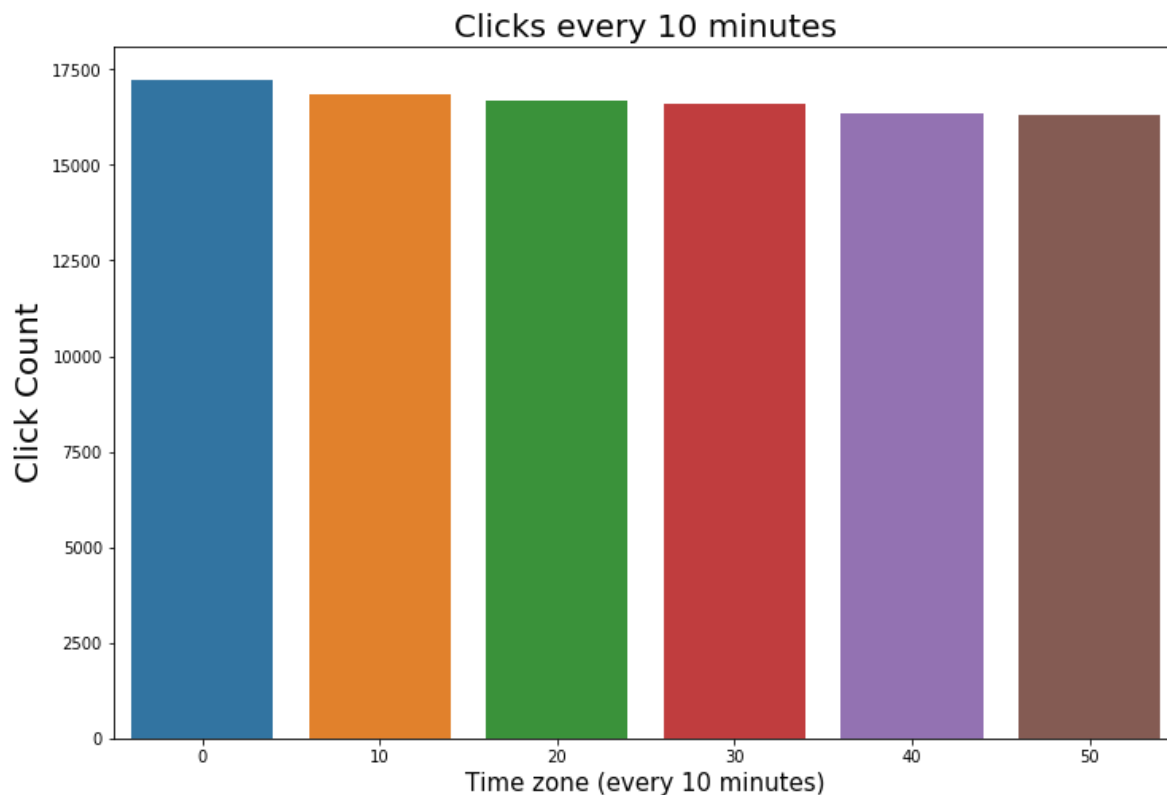
コンピュータで自動クリックしているとしたら10分間隔でクリックが集中するかもしれない。

In [55]:

```

1 sample_df['minute_band'] = pd.cut(sample_df['minute'], 6)
2 # 10分幅でクリック数をカウントして表示する
3 minute_df = sample_df.groupby('minute_band', as_index=False).count().sort_values(by='minute_
4
5 rename_dict = {}
6 minute_band_name = ['0', '10', '20', '30', '40', '50']
7 for minute_band,col in zip(minute_band_name, minute_df['minute_band']):
8     rename_dict[col] = minute_band
9
10 minute_df = minute_df.rename(index=rename_dict)
11
12 def rename_minute_band(s):
13     return rename_dict[s]
14
15 minute_df ['minute_band']= minute_df['minute_band'].map(rename_minute_band)
16
17 plt.figure(figsize=(12,8))
18 sns.barplot(x='minute_band', y='ip', data=minute_df)
19 plt.title("Clicks every 10 minutes", fontsize=20)
20 plt.xlabel('Time zone (every 10 minutes)', fontsize=15)
21 plt.ylabel('Click Count', fontsize=20)
22 plt.show()

```



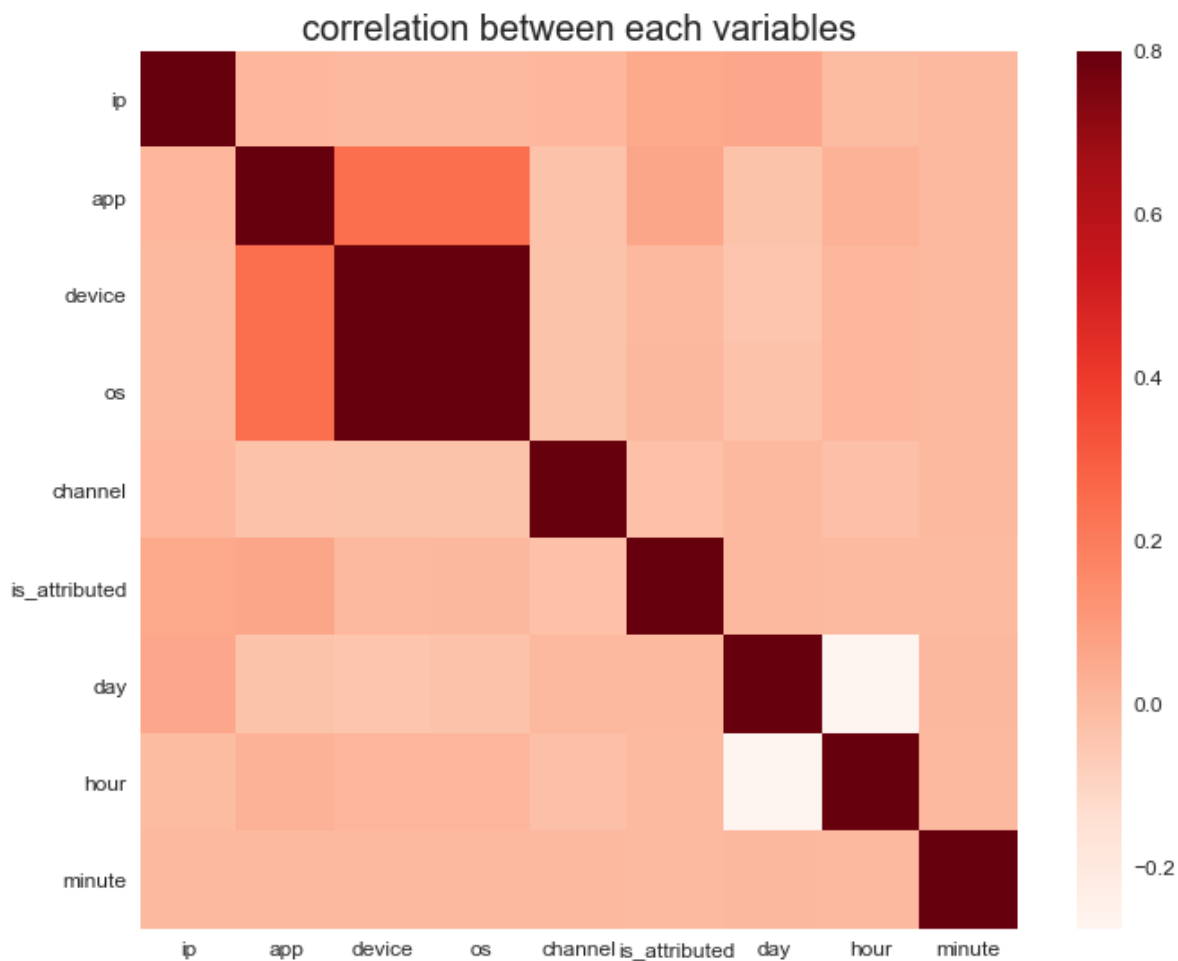
時間と10分間隔で見てきたけど、顕著な違いは現れなかった。思った以上に不正業者は巧妙に不正を行っているようだ。

各変数の相関関係

次に各変数間の相関関係を見てみよう。

In [94]:

```
1 # 相関係数を見る
2 corrmat = sample_df.corr()
3 f, ax = plt.subplots(figsize=(12, 9))
4 sns.heatmap(corrmat, cmap='Reds', vmax=.8, square=True, ax=ax);
5 ax.set_title("correlation between each variables", fontsize= 20);
```



OS、デバイス、アプリに強い相関があるようだけど、目的変数であるis_attributedと相関が強いものが見当たらない。

次に不正業者の実態 (<https://www.gizmodo.jp/2017/06/thai-click-fraud-farm-busted-using-wall-of-iphones.html>)を見て、次の推測は次の4つです。

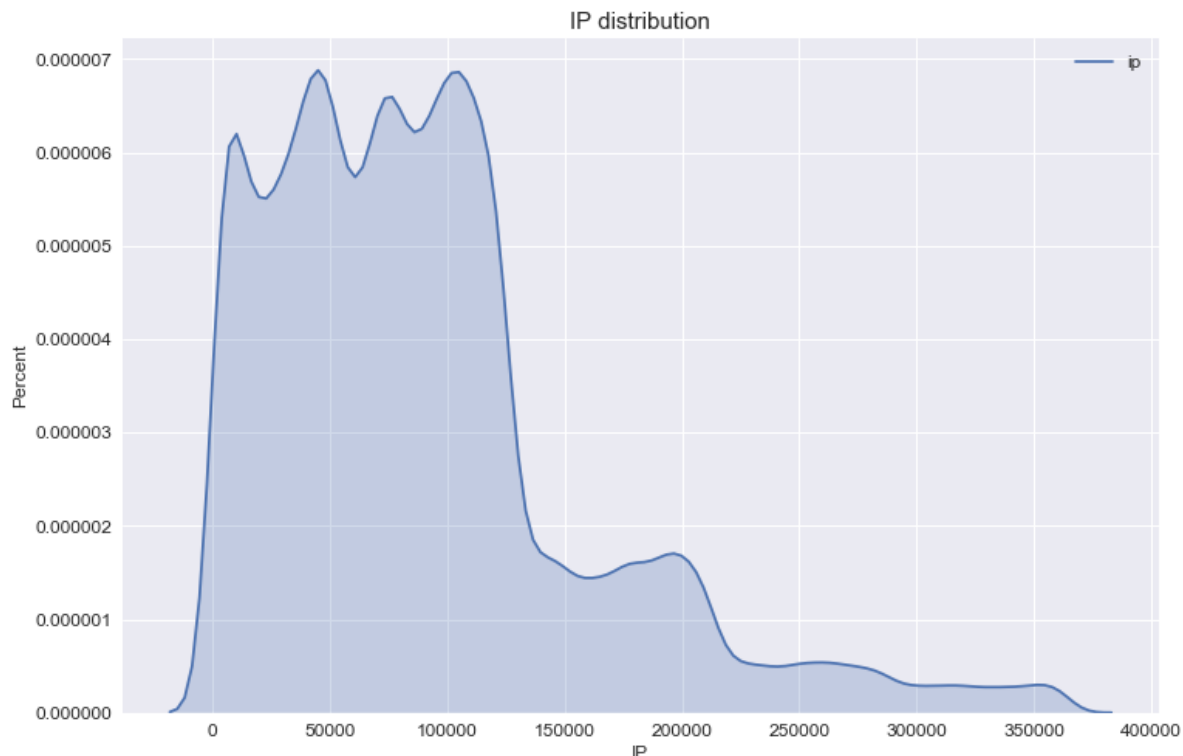
1. 同じipアドレスで大量にクリックしていないか
2. 同じデバイス、OSでクリックしていないか
3. 競合アプリ開発者からの不正で特定のアプリでクリック数が増えていないか
4. 広告パブリッシャーからの不正で特定のチャンネルでクリック数が増えていないか

1. 同じipアドレスで大量クリックしているか

In [115]:

```
1 # ipアドレスの分布を確認する。
2 print(sample_df['ip'].describe())
3 plt.figure(figsize=(12, 8))
4 sns.kdeplot(sample_df['ip'], shade=True)
5 plt.title('IP distribution', fontsize = 15)
6 plt.xlabel('IP', fontsize = 12)
7 plt.ylabel('Percent', fontsize = 12)
8 plt.show()
```

```
count    100000.000000
mean      91255.879670
std       69835.553661
min         9.000000
25%      40552.000000
50%      79827.000000
75%     118252.000000
max     364757.000000
Name: ip, dtype: float64
```

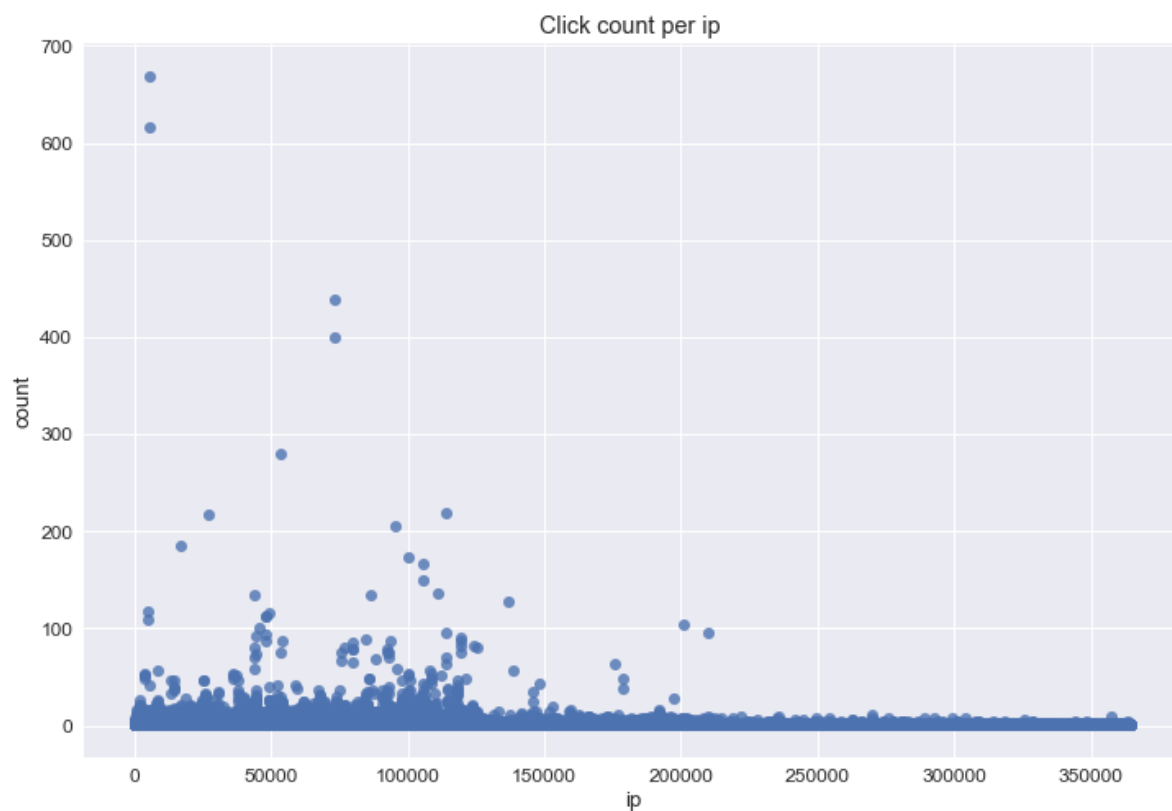


0から10000あたりに集中して分布しています。本物のipアドレス (192.102.10.52)であれば近い場所にあるなどの情報が分かりますが、残念ながらこのデータのipアドレスはエンコードされてしまっているのでipアドレスの数値に意味を見出すことはできません。

次にipアドレスごとの個数をカウントしてみましょう。

In [100]:

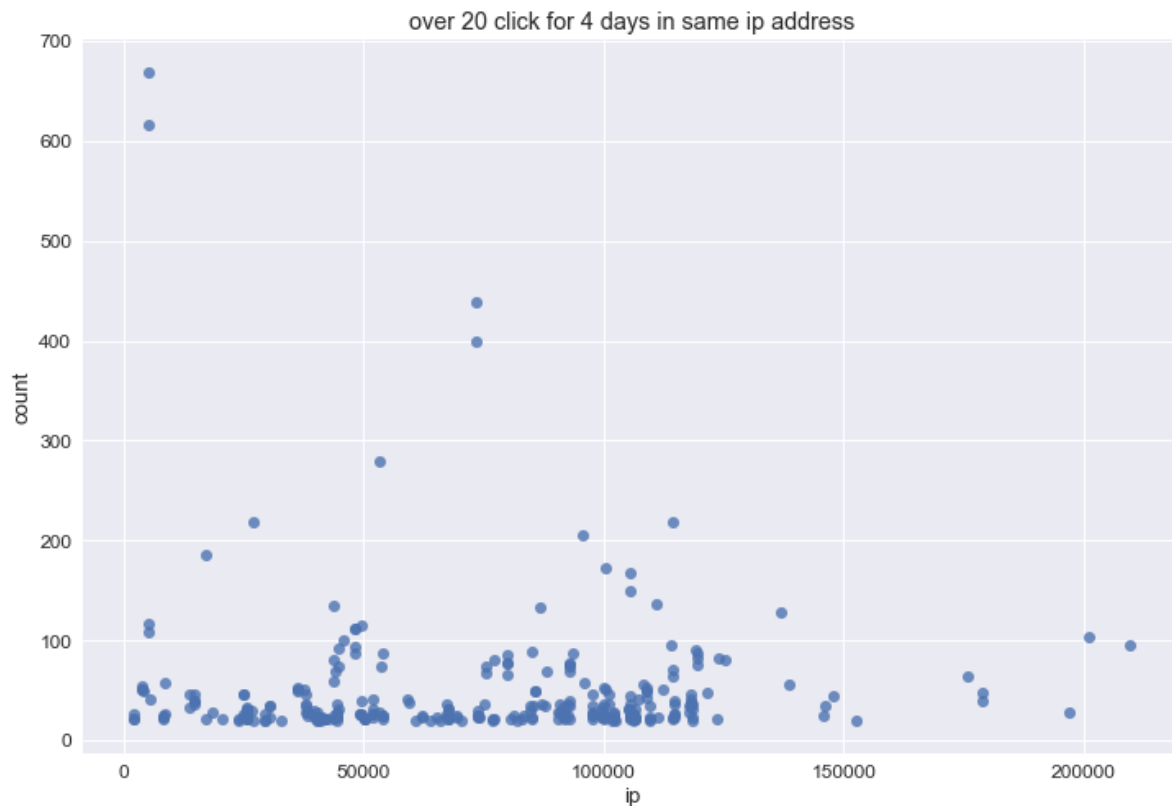
```
1 click_ip_count = sample_df.groupby('ip', as_index=False)['device'].aggregate('count').sort_value
2 click_ip_count = click_ip_count.rename(columns={'device':'count'})
3
4 plt.figure(figsize=(12,8))
5 sns.regplot(x='ip', y='count', fit_reg=False, data=click_ip_count)
6 plt.title("Click count per ip ")
7 plt.show()
```



続いて、クリック数が20以下のipを除いてもう一度表示してみる。

In [99]:

```
1 plt.figure(figsize=(12,8))
2 sns.regplot(x='ip', y='count', fit_reg=False, data=click_ip_count[click_ip_count['count']>19])
3 plt.title('over 20 click for 4 days in same ip address')
4 plt.show()
```



自分がスマートフォンを使っていて4日間で広告を20回以上タップすることはほとんどない。ここにプロットされたものは不正クリックかもしれない。

もう少しipアドレス視点で分析します。同一ipアドレスでのクリック数が多いものを上位10並べてみます。

In [175]:

```
1 ip_click_ranking = sample_df.groupby('ip', as_index = False).count().sort_values(by = 'app', ascen
2 ip_click_ranking = ip_click_ranking[['ip', 'app']]
3 ip_click_ranking = ip_click_ranking.rename(columns={'app':'click_count'})
4 ip_click_ranking.head(10)
```

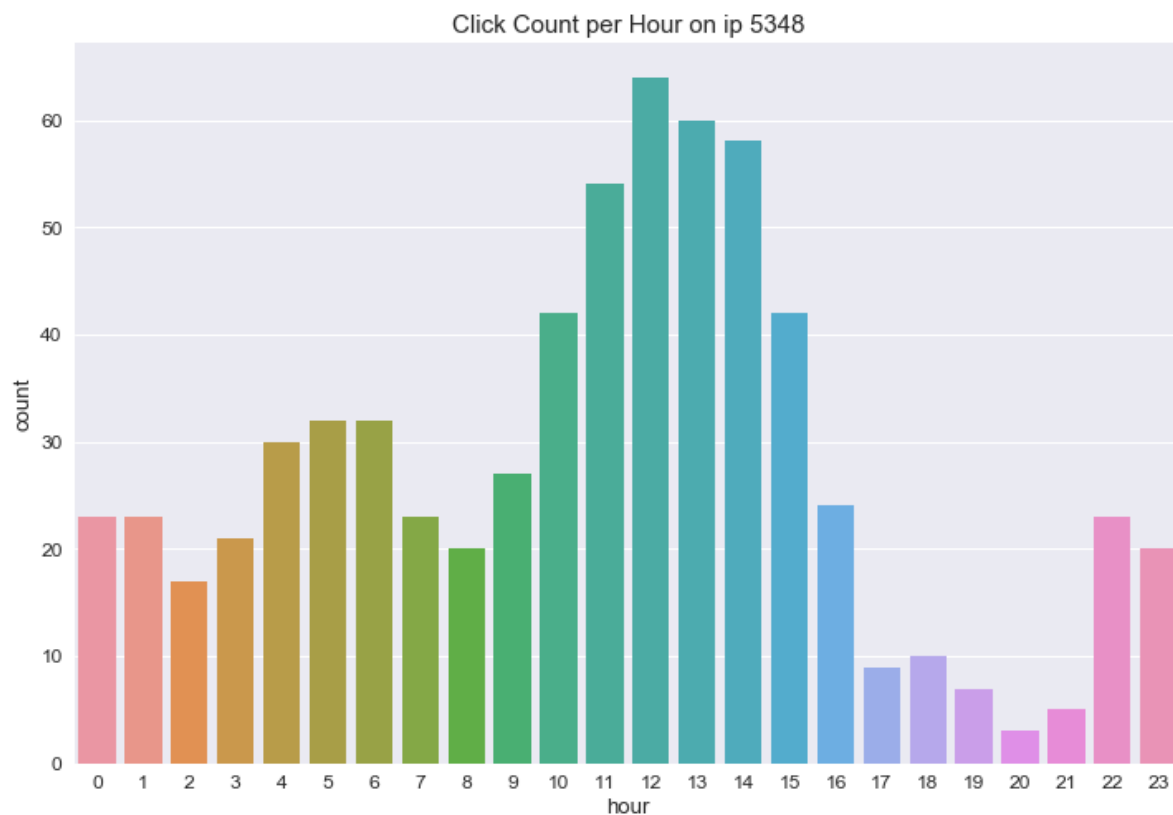
Out[175]:

	ip	click_count
	926	5348
	918	5314
	12833	73487
	12839	73516
	9385	53454
	20021	114276
	4639	26995
	16703	95766
	2990	17149
	17488	100275

試しに一番クリック数の多いip「5348番」のクリック時間を見てみよう。

In [160]:

```
1 sample_df[sample_df['ip'] == 5348].sort_values(by='click_time', ascending=True).head(20)
2
3 plt.figure(figsize=(12,8))
4 sns.countplot(x='hour', data=sample_df[sample_df['ip'] == 5348]);
5 plt.title('Click Count per Hour on ip 5348', fontsize=15)
6 plt.show()
```



寝ずに一日中クリックしていることが分かる。これは推測だが同じIPで時間帯で途絶えなければ不正クリックの匂いがします。

またクリックの多いIPのダウンロード数を見てみましょう。

In [89]:

```

1 ip_click_download = sample_df.groupby('ip', as_index = False).sum().sort_values(by='app', ascend
2 ip_click_download = ip_click_download[['ip', 'app', 'is_attributed']]
3 ip_click_download = ip_click_download.rename(columns={'is_attributed': 'download_count', 'app':
4
5 ip_click_download.head(30)

```

Out[89]:

	ip	click_count	download_count
0	5314	10086	3
1	5348	9478	3
2	73487	4809	0
3	73516	4428	0
4	53454	3548	0
5	114276	3243	0
6	26995	3215	0
7	86767	3003	0
8	17149	2993	0
9	105475	2513	0
10	95766	2323	0
11	114220	2030	0
12	100275	2017	0
13	105560	1833	0
14	111025	1705	1
15	77048	1643	0
16	201182	1630	0
17	5147	1561	0
18	43793	1558	0
19	137052	1476	0
20	119289	1374	0
21	49602	1365	0
22	79827	1322	0
23	84896	1319	0
24	44725	1318	0
25	43827	1311	0
26	48170	1226	0
27	5178	1223	0
28	119369	1202	0
29	45745	1202	0

クリック数が多いipでもダウンロード数が0のipがほとんどです。ipは不正クリックかの判断に有効でしょう。

2.同じデバイス、OSでクリックしていないか

次に同じデバイス、OSでのクリック数を見ていきましょう。

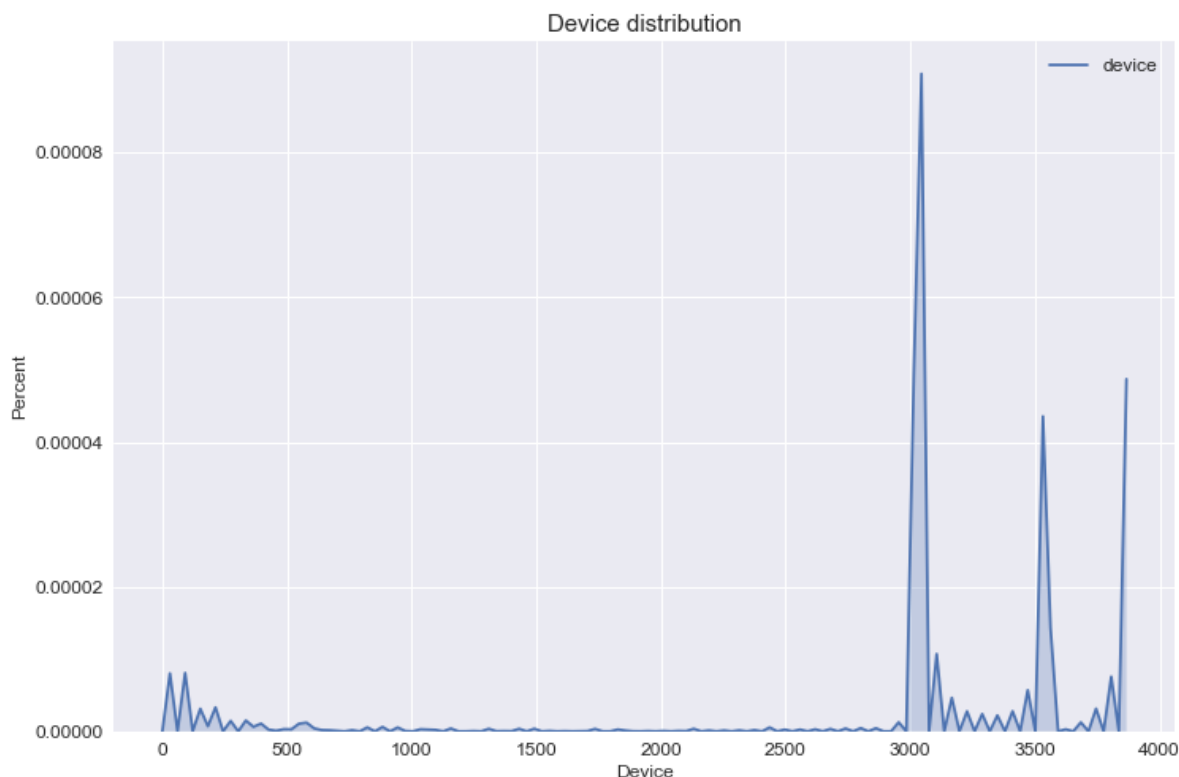
不正業者はなるべく安い端末を大量に持っていて、しかもOSもアップデートせずに使っているかもしれません。

まずデバイスidごとのクリック数を見ていきましょう。

In [132]:

```
1 # デバイスの分布を確認する。
2 print(sample_df['device'].describe())
3 plt.figure(figsize=(12, 8))
4 sns.kdeplot(sample_df['device'], shade=True)
5 plt.title('Device distribution', fontsize = 15)
6 plt.xlabel('Device', fontsize = 12)
7 plt.ylabel('Percent', fontsize = 12)
8 plt.show()
```

```
count    100000.000000
mean       21.771250
std       259.667767
min         0.000000
25%         1.000000
50%         1.000000
75%         1.000000
max      3867.000000
Name: device, dtype: float64
```



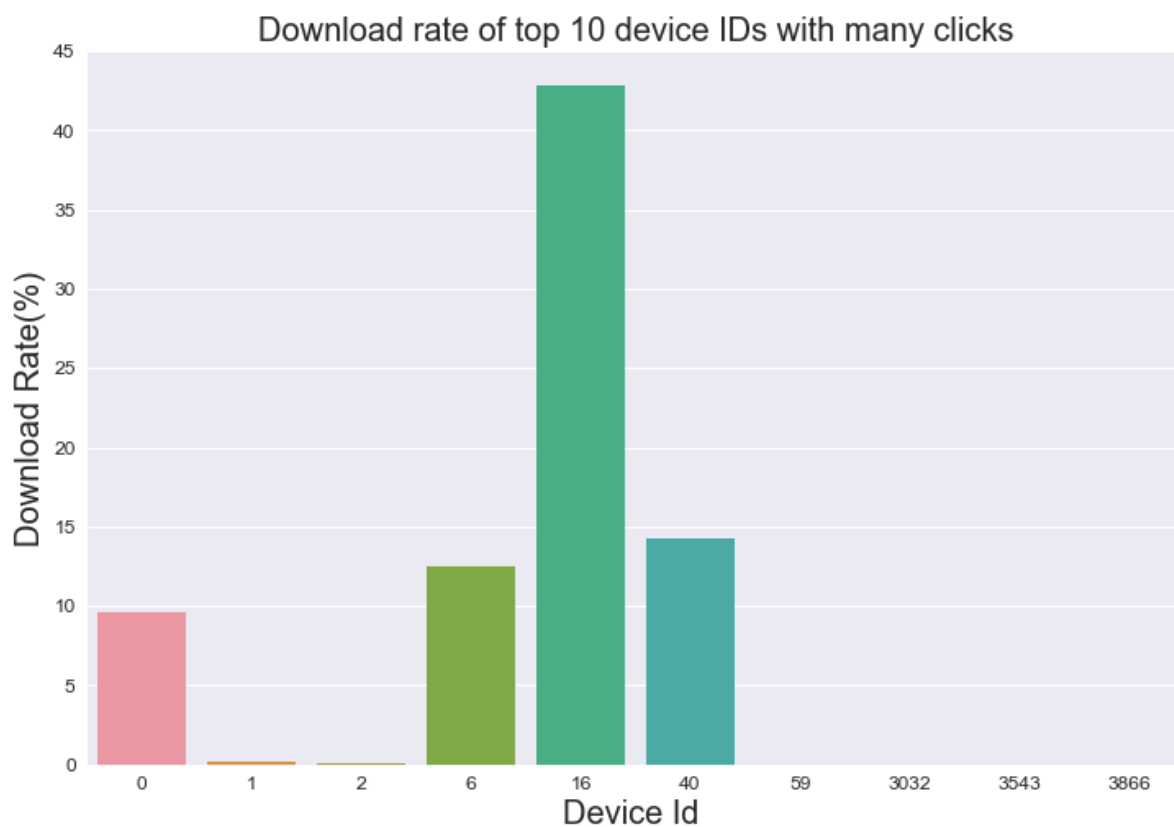
特定のデバイスでのクリック割合が多いことが分かる。クリックの多いデバイスの上位10のダウンロード率を見てみましょう。

In [155]:

```

1 dll_cnt_device = sample_df.groupby('device', as_index=False).sum().sort_values(by='is_attributed')
2 dll_cnt_device = dll_cnt_device[['device', 'is_attributed']]
3 dll_cnt_device = dll_cnt_device.rename(columns={'is_attributed': 'download_count'})
4
5 click_cnt_device = sample_df.groupby('device', as_index=False).count().sort_values(by='app', ascending=False)
6 click_cnt_device = click_cnt_device[['device', 'app']]
7 click_cnt_device = click_cnt_device.rename(columns={'app': 'click_count'})
8
9 device_click_download_df = pd.merge(click_cnt_device, dll_cnt_device, on='device')
10
11 device_click_download_df['download_rate'] = device_click_download_df['download_count'] / device_click_download_df['click_count']
12
13
14
15
16 plt.figure(figsize=(12,8))
17 sns.barplot(x='device', y='download_rate', data=device_click_download_df[:10])
18 plt.xlabel("Device Id ", fontsize=20)
19 plt.ylabel("Download Rate(%)", fontsize=20)
20 plt.title("Download rate of top 10 device IDs with many clicks", fontsize=20)
21 plt.show()
22
23 device_click_download_df.head(10)

```



Out[155]:

	device	click_count	download_count	download_rate
0	1	94338	146	0.154763
1	2	4345	2	0.046030
2	0	541	52	9.611830
3	3032	371	0	0.000000

	device	click_count	download_count	download_rate
4	3543	151	0	0.000000
5	3866	93	0	0.000000
6	59	12	0	0.000000
7	6	8	1	12.500000
8	40	7	1	14.285714
9	16	7	3	42.857143

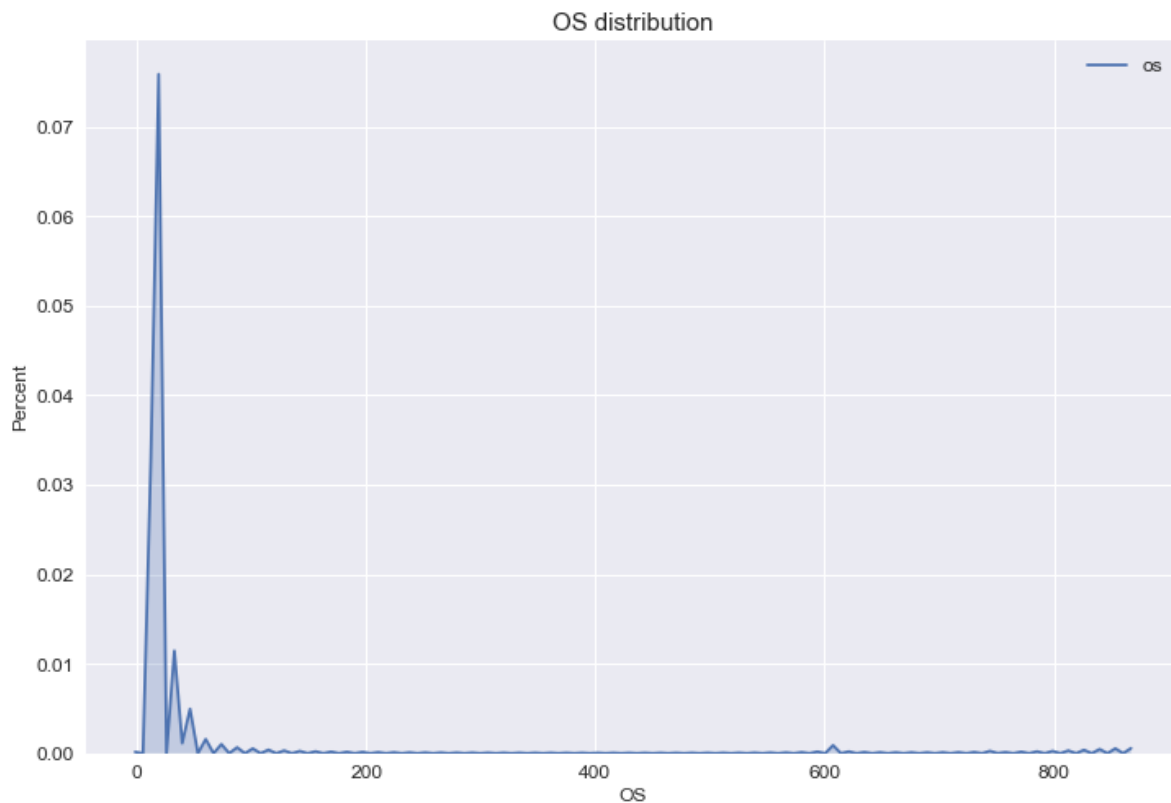
デバイスIDが0,6,16,40以外のIDはほぼダウンロードがない。

同じようにOSごとのクリック数を見ていきましょう。

In [151]:

```
1 # OSの分布を確認する。
2 print(sample_df['os'].describe())
3 plt.figure(figsize=(12, 8))
4 sns.kdeplot(sample_df['os'], shade=True)
5 plt.title('OS distribution', fontsize = 15)
6 plt.xlabel('OS', fontsize = 12)
7 plt.ylabel('Percent', fontsize = 12)
8 plt.show()
```

```
count    100000.000000
mean       22.818280
std        55.943136
min         0.000000
25%        13.000000
50%        18.000000
75%        19.000000
max       866.000000
Name: os, dtype: float64
```

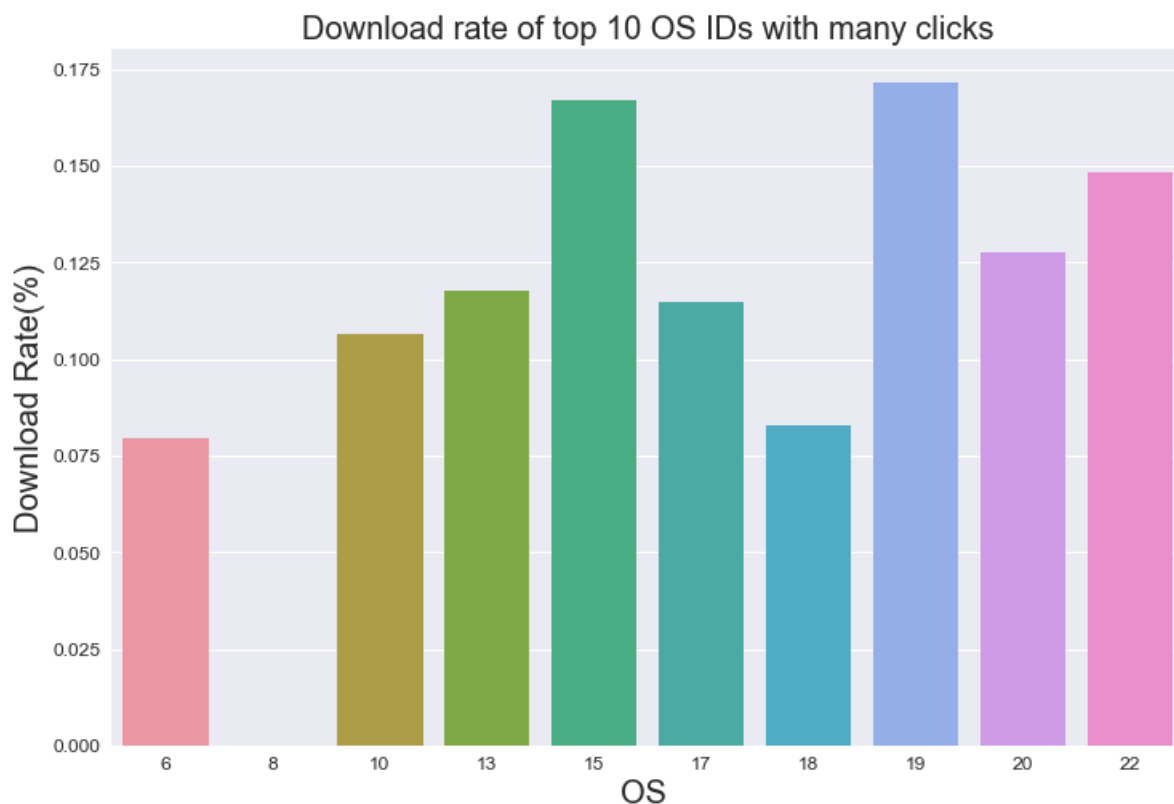


In [156]:

```

1  dll_cnt_os = sample_df.groupby('os', as_index=False).sum().sort_values(by='is_attributed', ascer
2  dll_cnt_os = dll_cnt_os[['os', 'is_attributed']]
3  dll_cnt_os = dll_cnt_os.rename(columns={'is_attributed': 'download_count'})
4
5  click_cnt_os = sample_df.groupby('os', as_index=False).count().sort_values(by='app', ascending=
6  click_cnt_os = click_cnt_os[['os', 'app']]
7  click_cnt_os = click_cnt_os.rename(columns={'app': 'click_count'})
8
9  os_click_download_df = pd.merge(click_cnt_os, dll_cnt_os, on='os')
10
11  os_click_download_df['download_rate'] = os_click_download_df['download_count'] / os_click_do
12
13  os_click_download_df.head(10)
14
15  plt.figure(figsize=(12,8))
16  sns.barplot(x='os', y='download_rate', data=os_click_download_df[:10])
17  plt.xlabel("OS ", fontsize=20)
18  plt.ylabel("Download Rate(%)", fontsize=20)
19  plt.title("Download rate of top 10 OS IDs with many clicks", fontsize=20)
20  plt.show()
21
22  os_click_download_df.head(10)

```



Out[156]:

	os	click_count	download_count	download_rate
0	19	23870	41	0.171764
1	13	21223	25	0.117797
2	17	5232	6	0.114679
3	18	4830	4	0.082816
4	22	4039	6	0.148552

	os	click_count	download_count	download_rate
5	10	2816	3	0.106534
6	8	2775	0	0.000000
7	6	2520	2	0.079365
8	15	2396	4	0.166945
9	20	2347	3	0.127823

こちらは8番OSだけダウンロード数が0以外は目立った違いは見られない。

競合アプリ開発者からの不正で特定のアプリでクリック数が増えているか

次にアプリごとのクリック数を見てみよう。

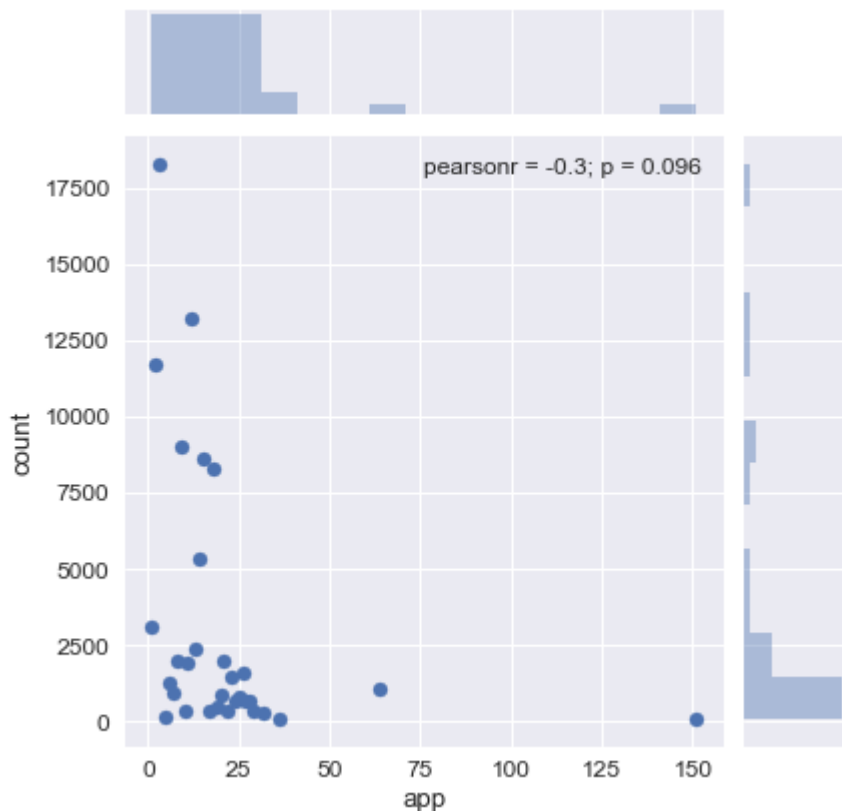
In [79]:

```

1 click_ip_count = sample_df.groupby('app', as_index=False)['device'].aggregate('count').sort_valu
2 click_ip_count = click_ip_count.rename(columns={'device':'count'})
3 click_ip_count = click_ip_count[click_ip_count['count'] > 100]
4
5 plt.figure(figsize=(12,8))
6 sns.jointplot('app', 'count', data=click_ip_count)
7 plt.show()

```

<matplotlib.figure.Figure at 0x1a1aada048>



極端にクリックされてるアプリがありますね。競合アプリ開発者が特定のアプリの広告を表示させないため広告費を使い切らせる目的であるならこの上位のアプリたちはターゲットにされた可能性があります。

4. 広告パブリッシャーからの不正で特定のチャンネルでクリック数が増えているか

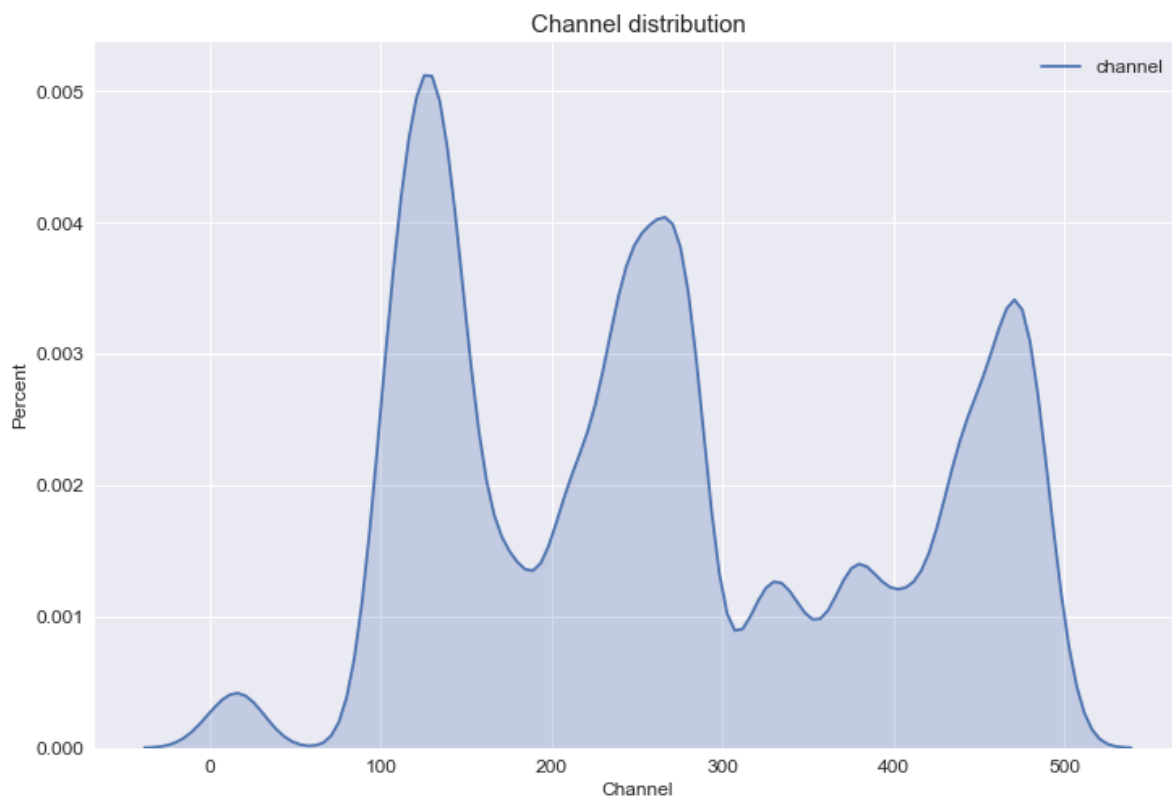
続いてチャンネルも見してみましょう。

広告パブリッシャーが利益を出すために不正クリックをさせているとすると特定のチャンネルのクリック数が多いというデータが出るはずです。

In [162]:

```
1 # チャンネルの分布を確認する。
2 print(sample_df['channel'].describe())
3 plt.figure(figsize=(12, 8))
4 sns.kdeplot(sample_df['channel'], shade=True)
5 plt.title('Channel distribution', fontsize = 15)
6 plt.xlabel('Channel', fontsize = 12)
7 plt.ylabel('Percent', fontsize = 12)
8 plt.show()
```

```
count    100000.000000
mean       268.832460
std        129.724248
min         3.000000
25%        145.000000
50%        258.000000
75%        379.000000
max        498.000000
Name: channel, dtype: float64
```



特定のチャンネルのクリックが多いがこれだけではただ広告の表示割合が大きいだけかもしれない。

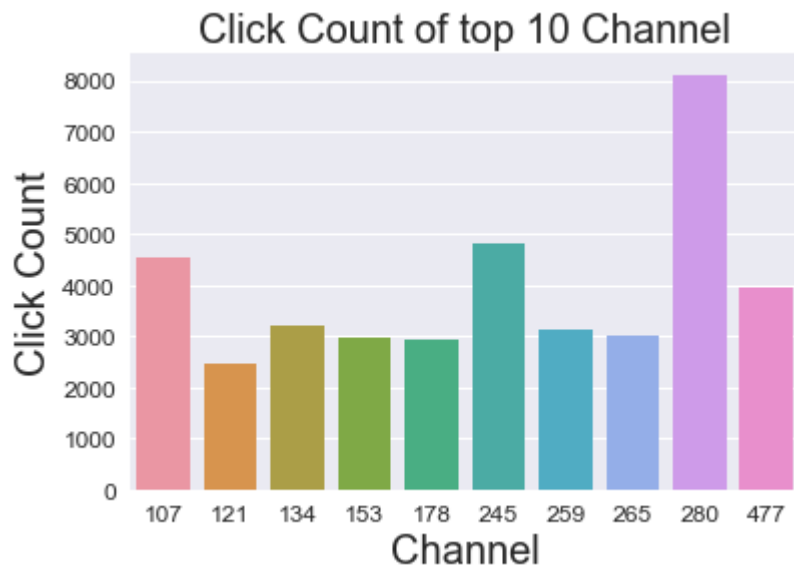
クリック数の多い上位10のチャンネルのクリック数を見てみよう。

In [170]:

```

1 click_cnt_channel = sample_df.groupby('channel', as_index=False).count().sort_values(by='app', a
2 click_cnt_channel = click_cnt_channel[['channel', 'app']]
3 click_cnt_channel = click_cnt_channel.rename(columns={'app': 'click_count'})
4
5 sns.barplot(x='channel', y='click_count', data=click_cnt_channel.sort_values(by='click_count', asc
6 plt.xlabel("Channel", fontsize=20)
7 plt.ylabel("Click Count", fontsize=20)
8 plt.title("Click Count of top 10 Channel", fontsize=20)
9 plt.show()
10
11 click_cnt_channel.head()
12

```



Out[170]:

	channel	click_count
77	280	8114
64	245	4802
14	107	4543
146	477	3960
33	134	3224

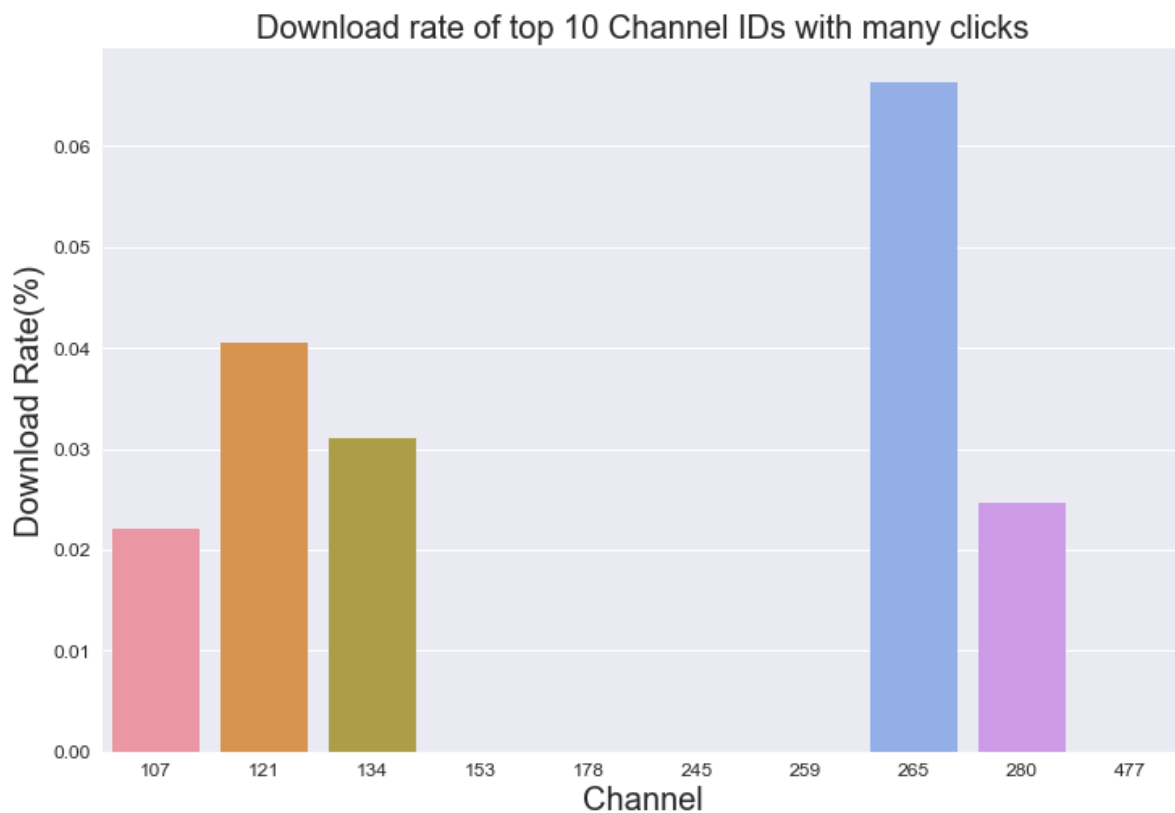
280番が多いが、それ以外はほぼ横ばいです。ここからこの上位10チャンネルのダウンロード率も見てみましょう。

In [164]:

```

1 dll_cnt_channel = sample_df.groupby('channel', as_index=False).sum().sort_values(by='is_attributed')
2 dll_cnt_channel = dll_cnt_channel[['channel', 'is_attributed']]
3 dll_cnt_channel = dll_cnt_channel.rename(columns={'is_attributed': 'download_count'})
4
5
6
7 channel_click_download_df = pd.merge(click_cnt_channel, dll_cnt_channel, on='channel')
8
9 channel_click_download_df['download_rate'] = channel_click_download_df['download_count'] / click_cnt_channel['click_count']
10
11 plt.figure(figsize=(12,8))
12 sns.barplot(x='channel', y='download_rate', data=channel_click_download_df[:10])
13 plt.xlabel("Channel", fontsize=20)
14 plt.ylabel("Download Rate(%)", fontsize=20)
15 plt.title("Download rate of top 10 Channel IDs with many clicks", fontsize=20)
16 plt.show()
17
18 channel_click_download_df.head(10)

```



Out[164]:

	channel	click_count	download_count	download_rate
0	280	8114	2	0.024649
1	245	4802	0	0.000000
2	107	4543	1	0.022012
3	477	3960	0	0.000000
4	134	3224	1	0.031017
5	259	3130	0	0.000000
6	265	3013	2	0.066379

	channel	click_count	download_count	download_rate
7	153	2954	0	0.000000
8	178	2936	0	0.000000
9	121	2472	1	0.040453

クリック数よりも極端で、まったくダウンロードのないチャネルもある。もしかしたらダウンロード率の低いチャネルは不正業者に不正を依頼しているかもしれません。

まとめ

- ipアドレスと時間で見ると定期的にクリックしていることが分かる。明らかに普通の人の使い方ではない。
- クリックが多いがダウンロード率が低いデバイスがある。
- 不正も目論む人たち（広告パブリッシャー、競合アプリ開発者）の目的から推測されたチャネルやアプリIDにも不正クリックとの関連性がありそう。

反省点

- とにかく時間がかかった。
- 何をプロットするか、そして何をプロットするか決めたあとに実際にプロットをコーディングするのにも時間がかかった。

改善点

- EDAは結論を急ぎすぎず、データをいろいろな方法で眺めるつもりでプロットしていく。有用な情報ばかりでなく、「これはあまり関係がなさそうだ」といった情報も1つの情報となる。
- プロットのパターンは決まっているので自分がプロットしたパターンをストックして次に活かす。また他の人のEDAをたくさん読む。

質問やコメントは大歓迎です！