

Profitable App Profiles for the App Store and Google Play Markets

The goal of the project is to find the type of mobile apps that are profitable on the App Store and Google Play markets for an app development company. They only build free apps and their revenues come from in-app ads. This means that the number of users is a big factor in revenue. More users equates to more people who engage with the ads and thus more revenue.

The goal of this analysis is to identify what type of apps are more likely to attract users and help the developers understand this.

```
In [2]: # App Store Data Set #
from csv import reader
opened_file = open('AppStore.csv')
read_file = reader(opened_file)
ios = list(read_file)
ios_header = ios[0]
ios = ios[1:]

# Google Play Store Data Set #
from csv import reader
opened_file = open('googleplaystore.csv')
read_file = reader(opened_file)
android = list(read_file)
android_header = android[0]
android = android[1:]

To make it easier to explore the two data sets, we'll first write a function named explore_data() that we can use repeatedly to explore rows in a more readable way. We'll also add an option for our function to show the number of rows and columns for any data set.
```

```
In [3]: def explore_data(dataset, start, end, rows_and_columns=False):
dataset_slice = dataset[start:end]
for row in dataset_slice:
    print(row)
    print('\n') # adds a new (empty) line after each row

    if rows_and_columns:
        print('Number of rows:', len(dataset))
        print('Number of columns:', len(dataset[0]))

    print(android_header)
    print('\n')
    explore_data(android, 0, 3, True)

['App', 'Category', 'Rating', 'Reviews', 'Size', 'Installs', 'Type', 'Price', 'Content Rating', 'Genres', 'Last Updated', 'Current Ver', 'Android Ver']

['Photo Editor & Candy Camera & Grid & ScrapBook', 'ART_AND_DESIGN', '4.1', '159', '19M', '10,000+', 'Free', '0', 'Everyone', 'Art & Design', 'January 7, 2018', '1.0.0', '4.0.3 and up']

['Coloring book moana', 'ART_AND_DESIGN', '3.9', '967', '14M', '500,000+', 'Free', '0', 'Everyone', 'Art & Design;Pretend Play', 'January 15, 2018', '4.0.0', '4.0.3 and up']

['U Launcher Lite - FREE Live Cool Themes, Hide Apps', 'ART_AND_DESIGN', '4.7', '87510', '8.7M', '5,000,000+', 'Free', '0', 'Everyone', 'Art & Design', 'August 1, 2018', '1.2.4', '4.0.3 and up']

Number of rows: 10841
Number of columns: 13

In [4]: print(ios_header)
print('\n')
explore_data(ios, 0, 3, True)

['id', 'track_name', 'size_bytes', 'currency', 'price', 'rating_count_tot', 'rating_count_ver', 'user_rating', 'user_rating_ver', 'ver', 'content_rating', 'prime_genre', 'sup_devices.num', 'ipadSc_urls.num', 'lang.num', 'vpp_lic']

['284882215', 'Facebook', '389879808', 'USD', '0.0', '2974676', '212', '3.5', '3.5', '95.0', '4+', 'Social Networking', '37', '1', '29', '1']

['389801252', 'Instagram', '113954816', 'USD', '0.0', '2161558', '1289', '4.5', '4.0', '10.23', '12+', 'Photo & Video', '37', '0', '29', '1']

['529479190', 'Clash of Clans', '116476928', 'USD', '0.0', '2130805', '579', '4.5', '4.5', '9.24.12', '9+', 'Games', '38', '5', '18', '1']

Number of rows: 7197
Number of columns: 16

Deleting Wrong Data
Row 10472 has wrong data in it, so we will be deleting it
```

```
In [5]: del android[10472]

Removing Duplicate Entries

There are duplicate rows in both datasets as seen below.
```

```
In [6]: for app in android:
    name = app[0]
    if name == 'Instagram':
        print(app)

['Instagram', 'SOCIAL', '4.5', '66577313', 'Varies with device', '1,000,000,000+', 'Free', '0', 'Teen', 'Social', 'July 31, 2018', 'Varies with device', 'Varies with device']
['Instagram', 'SOCIAL', '4.5', '66577446', 'Varies with device', '1,000,000,000+', 'Free', '0', 'Teen', 'Social', 'July 31, 2018', 'Varies with device', 'Varies with device']
['Instagram', 'SOCIAL', '4.5', '66577313', 'Varies with device', '1,000,000,000+', 'Free', '0', 'Teen', 'Social', 'July 31, 2018', 'Varies with device', 'Varies with device']
['Instagram', 'SOCIAL', '4.5', '6659917', 'Varies with device', '1,000,000,000+', 'Free', '0', 'Teen', 'Social', 'July 31, 2018', 'Varies with device', 'Varies with device']

In [7]: duplicate_apps = []
unique_apps = []

for app in android:
    name = app[0]
    if name in unique_apps:
        duplicate_apps.append(name)
    else:
        unique_apps.append(name)

print('Number of duplicate apps:', len(duplicate_apps))
print('\n')
print('Examples of duplicate apps:', duplicate_apps[:15])

Number of duplicate apps: 1181

Examples of duplicate apps: ['Quick PDF Scanner x OCR FREE', 'Box', 'Google My Business', 'Zoom Cloud Meetings', 'join.me - Simple Meetings', 'Box', 'Zenefits', '5005501', 'Google My Business', 'Slack', 'FreshBooks Classic', 'Insightly CRM', 'QuickBooks Accounting: Invoicing & Expenses', 'HipChat - Chat Built for Teams', 'Xero Accounting Software']
```

```
In [8]: reviews_max = {}
for app in android:
    name = app[0]
    n_reviews = float(app[3])

    if name in reviews_max and reviews_max[name] < n_reviews:
        reviews_max[name] = n_reviews
    elif name not in reviews_max:
        reviews_max[name] = n_reviews

In [11]: print(len(reviews_max))

9659
```

```
In [12]: android_clean = []
already_added = []

for app in android:
    name = app[0]
    n_reviews = float(app[3])

    if (reviews_max[name] == n_reviews) and (name not in already_added):
        android_clean.append(app)
        already_added.append(name)
```

```
In [13]: explore_data(android_clean, 0, 3, True)

['Photo Editor & Candy Camera & Grid & ScrapBook', 'ART_AND_DESIGN', '4.1', '159', '19M', '10,000+', 'Free', '0', 'Everyone', 'Art & Design', 'January 7, 2018', '1.0.0', '4.0.3 and up']

['U Launcher Lite - FREE Live Cool Themes, Hide Apps', 'ART_AND_DESIGN', '4.7', '87510', '8.7M', '5,000,000+', 'Free', '0', 'Everyone', 'Art & Design', 'August 1, 2018', '1.2.4', '4.0.3 and up']

['Sketch - Draw & Paint', 'ART_AND_DESIGN', '4.5', '215644', '25M', '50,000,000+', 'Free', '0', 'Teen', 'Art & Design', 'June 8, 2018', 'Varies with device', '4.2 and up']

Number of rows: 9659
Number of columns: 13

Remove Non-English Apps
```

```
In [14]: def is_english(string):
    non_ascii = 0

    for character in string:
        if ord(character) > 127:
            non_ascii += 1

    if non_ascii > 3:
        return False
    else:
        return True

Filtering the non-English apps from both datasets:
```

```
In [15]: android_english = []
ios_english = []

for app in android_clean:
    name = app[0]
    if is_english(name):
        android_english.append(app)

for app in ios:
    name = app[1]
    if is_english(name):
        ios_english.append(app)

explore_data(android_english, 0, 3, True)
print('\n')
explore_data(ios_english, 0, 3, True)

['Photo Editor & Candy Camera & Grid & ScrapBook', 'ART_AND_DESIGN', '4.1', '159', '19M', '10,000+', 'Free', '0', 'Everyone', 'Art & Design', 'January 7, 2018', '1.0.0', '4.0.3 and up']

['U Launcher Lite - FREE Live Cool Themes, Hide Apps', 'ART_AND_DESIGN', '4.7', '87510', '8.7M', '5,000,000+', 'Free', '0', 'Everyone', 'Art & Design', 'August 1, 2018', '1.2.4', '4.0.3 and up']

['Sketch - Draw & Paint', 'ART_AND_DESIGN', '4.5', '215644', '25M', '50,000,000+', 'Free', '0', 'Teen', 'Art & Design', 'June 8, 2018', 'Varies with device', '4.2 and up']

Number of rows: 9614
Number of columns: 13

['284882215', 'Facebook', '389879808', 'USD', '0.0', '2974676', '212', '3.5', '3.5', '95.0', '4+', 'Social Networking', '37', '1', '29', '1']

['389801252', 'Instagram', '113954816', 'USD', '0.0', '2161558', '1289', '4.5', '4.0', '10.23', '12+', 'Photo & Video', '37', '0', '29', '1']

['529479190', 'Clash of Clans', '116476928', 'USD', '0.0', '2130805', '579', '4.5', '4.5', '9.24.12', '9+', 'Games', '38', '5', '18', '1']

Number of rows: 6183
Number of columns: 16

Isolating Free Apps
```

```
In [16]: android_final = []
ios_final = []

for app in android_english:
    price = app[7]
    if price == '0':
        android_final.append(app)

for app in ios_english:
    price = app[4]
    if price == '0.0':
        ios_final.append(app)

print(len(android_final))
print(len(ios_final))

8864
3222

Frequency Table For Genre of the App
```

```
In [17]: def freq_table(dataset, index):
    table = {}
    total = 0

    for row in dataset:
        total += 1
        value = row[index]
        if value in table:
            table[value] += 1
        else:
            table[value] = 1

    table_percentages = {}
    for key in table:
        percentage = (table[key] / total) * 100
        table_percentages[key] = percentage

    return table_percentages

def display_table(dataset, index):
    table = freq_table(dataset, index)
    table_display = []
    for key in table:
        key_val_as_tuple = (table[key], key)
        table_display.append(key_val_as_tuple)

    table_sorted = sorted(table_display, reverse = True)
    for entry in table_sorted:
        print(entry[1], ': ', entry[0])

In [18]: display_table(ios_final, -5)

Games : 58.16263190564867
Entertainment : 7.583352296718118
Photo & Video : 4.9658597144630665
Education : 3.662321539416512
Social Networking : 3.2898829608317814
Shopping : 2.60707635099311
Utilities : 2.5139664884469275
Sports : 2.1415270018621975
Music : 2.0484171322169147
Health & Fitness : 2.0173080909066205
Productivity : 1.738050900626732
Lifestyle : 1.5828677839851024
News : 1.3345747982619491
Travel : 1.2414649286157666
Finance : 1.1173184357541899
Weather : 0.8690254500310366
Food & Drink : 0.806952263602483
Reference : 0.558659217870949
Business : 0.5276225946617008
Book : 0.4345127250155183
Navigation : 0.186219739292365
Medical : 0.186219739292365
Catalogs : 0.12414649286157665
```

We can see that Games and Entertainment apps dominate. However that does not mean they have a large number of users.

```
In [19]: display_table(android_final, 1)

FAMILY : 18.997942238267147
GAME : 0.724729241877256
TOOLS : 0.461191335740072
BUSINESS : 4.591606498194946
LIFESTYLE : 3.963429628800966
PRODUCTIVITY : 3.892148014449433
FINANCE : 3.70836108303246
MEDICAL : 3.531137184115524
SPORTS : 3.905759122743682
PERSONALIZATION : 3.3167870936101084
COMMUNICATION : 3.2378158844765346
HEALTH_AND_FITNESS : 3.0798736462093865
PHOTOGRAPHY : 2.944494584837545
NEWS_AND_MAGAZINES : 2.7978339356180593
SOCIAL : 2.6624548736462095
TRAVEL_AND_LOCAL : 2.3352888986426
SHOPPING : 2.245936181083824
BOOKS_AND_REFERENCE : 2.1435018050541514
DATING : 1.861462093862816
VIDEO_PLAYERS : 1.7937725631708955
MAPS_AND_NAVIGATION : 1.3889169675090252
FOOD_AND_DRINK : 1.2409747292418771
EDUCATION : 1.1620036101083033
ENTERTAINMENT : 0.9589350816905415
LIBRARIES_AND_DEMO : 0.9363718411552346
AUTO_AND_VEHICLES : 0.9250962527075812
HOUSE_AND_HOME : 0.8235559566787804
WEATHER : 0.8069927797839394
EVENTS : 0.7107460722821661
PARENTING : 0.6543321299638989
ART_AND_DESIGN : 0.6430950451562455
COMICS : 0.6204873646209386
BEAUTY : 0.5979241877256317
```

In the android market, the "Family" genre tends to dominate.

Most Popular Apps on the App Store

One way to find out what genres are the most popular (most users) is to calculate the average number of installs for each app genre. For the Google Play data set, we can find this information in the Installs column, but for the App Store data set this information is missing. As a workaround, we can take the total number of user ratings as a proxy, which we can find in the rating\_count\_tot app.

```
In [21]: genres_ios = freq_table(ios_final, -5)

for genre in genres_ios:
    total = 0
    len_genre = 0
    for app in ios_final:
        genre_app = app[-5]
        if genre_app == genre:
            n_ratings = float(app[5])
            total += n_ratings
            len_genre += 1
    avg_n_ratings = total / len_genre
    print(genre, ': ', avg_n_ratings)

Social Networking : 71548.34905668378
ART_AND_DESIGN : 1986335.0877192982
AUTO_AND_VEHICLES : 647317.8170731707
BEAUTY : 513151.88679245283
BOOKS_AND_REFERENCE : 8767811.8947368441
BUSINESS : 1712290.1474201474
COMICS : 817657.2727272727
COMMUNICATION : 38456119.167247385
DATING : 854828.4303036303
EDUCATION : 1833495.145631068
ENTERTAINMENT : 11640705.88235294
EVENTS : 253542.22222222222
FINANCE : 13267892.475689756
FOOD_AND_DRINK : 4924897.7863636363
HEALTH_AND_FITNESS : 4188821.985479853
HOUSE_AND_HOME : 1331540.5616438356
LIBRARIES_AND_DEMO : 930563.734939759
LIFESTYLE : 1437816.2887861272
GAME : 15588015.603248259
FAMILY : 3695641.8198090694
MEDICAL : 120550.61980930671
SOCIAL : 23253652.1277118643
SHOPPING : 7036877.311557789
SPORTS : 17840110.40229885
TRAVEL_AND_LOCAL : 13984077.710144928
TOOLS : 10801391.298666677
PERSONALIZATION : 52061482.6122448975
PRODUCTIVITY : 16787331.344927534
PARENTING : 542603.6206896552
WEATHER : 5074486.197183099
VIDEO_PLAYERS : 24727872.452830188
NEWS_AND_MAGAZINES : 9540178.487741935
MAPS_AND_NAVIGATION : 4056941.7741935486
```

If we take the big companies such as Google, Facebook, Spotify, etc. out of the equation, they wouldn't have such high average of installs. So we can ignore them. If we look again, the Books category seems to be a good potential category.

Conclusion

For my final thoughts, I would suggest that building an app around a popular book can be profitable for both markets. It seems that taking a popular book (perhaps a more recent book) and turning it into an app could be profitable for both the Google Play and the App Store markets. Adding new a twist such as displaying daily quotes from the book, releasing an audio version of the book, providing quizzes on the book, or creating a forum where people can discuss the book could prove to be successful.