

```
In [ ]: Analysis of free User Centered Android and iOS Mobile Apps

This project is about building Android and iOS mobile apps that are free to download and install. Since the main source of revenue consists of in-app ads, it shows that our revenue for any app is highly dependent on the number of users who use our app.

Whenever more users see and engage with the ads, the more income would be generated. Therefore our goal for this project is to analyze data to help our developers understand what type of apps are likely to attract more users.

In [7]: %opening googleplaystore.csv

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
print(opened_file.readable())

True

In [12]: from csv import reader
opened_file = open('AppStore.csv', encoding='utf8')
read_file = reader(opened_file)
dataset = list(read_file)

print(opened_file.readable())

True

In [48]: #here we are exploring the two datasets using the explore_data() function

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]))

# below we are printing the first few rows of each data set

from csv import reader
opened_file = open('AppStore.csv', encoding='utf8')
read_file = reader(opened_file)
dataset = list(read_file)
print(dataset[0:5])

print("\n")

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
print(googledataset[0:3])

[['', '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'], ['1', '281656475', 'PAC-MAN Premium', '100788224', 'USD', '3.99', '21292', '26', '4', '4.3', '6.3.3', '4+', 'Games', '38', '5', '10', '1', '2', '281796108', 'Evernote - stay organised', '158578688', 'USD', '0', '161065', '26', '4', '3.5', '8.2.2', '4+', 'Productivity', '37', '5', '23', '1', '3', '281940292', 'WeatherBug - Local Weather, Radar, Maps, Alerts', '100524032', 'USD', '0', '188583', '2822', '3.5', '4.5', '5.0.0', '4+', 'Weather', '13', '1', '1', '1', '4', '282614216', 'eBay: Best App to Buy, Sell, Save! Online Shopping', '128512000', 'USD', '0', '262241', '649', '4', '4.3', '5.10.0', '12+', 'Shopping', '37', '5', '19', '1', '1']]

[['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', '17-Jan-18', '2.0.0', '4.0.3 and up'], ['Coloring book means', 'ART_AND_DESIGN', '3.9', '9.67', '14M', '500,000+', 'Free', '0', 'Everyone', 'Art & Design; Pretend Play', '15-Jan-18', '2.0.0', '4.0.3 and up']]

In [61]: # The number of rows and columns of each data set is as follows and the function assumes the argument for the dataset parameter doesn't have a header row

from csv import reader
opened_file = open('AppStore.csv', encoding='utf8')
read_file = reader(opened_file)
dataset = list(read_file)
#print(dataset)

print('Number of rows')
print(len(dataset))
print('\n')
print('Number of columns')
len(dataset[0])

number of rows
7198

number of columns

Out[61]: 17

In [60]: from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
print(len(dataset))
print('Number of rows')
print('\n')
print('Number of columns')
len(googledataset[0])

number of rows
10842

number of columns

Out[60]: 13

In [71]: # Below we are printing the column names and trying to identify the columns that could help us with our analysis

from csv import reader
opened_file = open('AppStore.csv', encoding='utf8')
read_file = reader(opened_file)
dataset = list(read_file)

column_names = dataset[0]

print(column_names)
print('\n')
print(print(column_names[1], ',', column_names[2], ',', column_names[4], ',', column_names[7], ',', column_names[8], ',', column_names[9]))

['', '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']

id , track_name , currency , rating_count_ver , user_rating , user_rating_ver
None

In [27]: #Based on the discussion section, the error is at row number 10472 and index number is 2 and it is indeed incorrect as the rating should not exceed 5. In this case the rating is 19 so we have to delete the row.

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

Wrong_data = googledataset[10472]

print(Wrong_data)

index = Wrong_data.index('19')
print('\n')
print('The index of 19:', index)

print(len(googledataset))
del googledataset[10472]

print(len(googledataset))

['Life Made Wi-Fi Touchscreen Photo Frame', '1.9', '19', '3.0M', '1,000+', 'Free', '0', 'Everyone', ' ', ' ', '11-Feb-18', '1.0.19', '4.0 and up', '']

The index of 19: 2
10841
10840

In [28]: #As we can see below the Google Play data set has duplicate entries which counts 1181 and names of some duplicate examples of app are also listed
#Examples of duplicate apps: ['Quick PDF Scanner + OCR FREE', 'Box', 'Google My Business', 'ZOOM Cloud Meetings',
# 'Join.me - Simple Meetings', 'Box']

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

duplicate_apps = []
unique_apps = []

for app in googledataset:
    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[:6])

Number of duplicate apps: 1181

Examples of duplicate apps: ['Quick PDF Scanner + OCR FREE', 'Box', 'Google My Business', 'ZOOM Cloud Meetings', 'Join.me - Simple Meetings', 'Box']
10841

In [36]: #here we created a dictionary key which is unique app name with a corresponding dictionary value having the highest count
#number of reviews of that app. looping through the googledataset, excluding the header row and we use
#the third index of the each row, by changing this particular review column (index 3) to float data type. After we check
#if name exists in the new dictionary or not. Finally we check if the length of the dictionary is 9659 rows after we subtract the
#duplicate entries as well as the actual length of the new dictionary we created is also 9659. Since the newly created dictionary is supposed to be the 9659 in length.

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

Wrong_data = googledataset[10472]

print(Wrong_data)

index = Wrong_data.index('19')
print('\n')
print('The index of 19:', index)

print(len(googledataset))
del googledataset[10472]

reviews_max = {}
for app in googledataset:
    name = app[0]
    n_reviews = str(app[3])

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

print('Expected length:', len(googledataset) - 1181)
print('Actual length:', len(reviews_max))

['Life Made Wi-Fi Touchscreen Photo Frame', '1.9', '19', '3.0M', '1,000+', 'Free', '0', 'Everyone', ' ', ' ', '11-Feb-18', '1.0.19', '4.0 and up', '']

The index of 19: 2
10841
Expected length: 9659
Actual length: 9659

In [14]: #In this case the function takes in a string and returns False if the character of the string is greater than 127
#If not it returns True, which means it is a non-English language.

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

def is_english(string):
    for character in string:
        if ord(character) > 127:
            return False
        return True

print(is_english('Facebook'))
print(is_english('爱奇艺 (欢乐颂) 电视剧热播'))
print(is_english('Docs To Go® Free Office Suite'))
print(is_english('Instachat 0'))

True
False
False
False

In [15]: #What we did above does not work for all the English apps as some app names use images or symbols and the like which
#falls outside of the ASCII range (more than three non-ASCII characters). In addition to this we check if these app names are
#detected as English or non-English for each data set and explore the data sets.

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

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

print(is_english('Docs To Go® Free Office Suite'))
print(is_english('Instachat 0'))
print(is_english('爱奇艺 (欢乐颂) 电视剧热播'))

True
True
True
False

In [ ]: #The code below shows loop through each dataset to isolate the free apps in a separate list and check the length of each data set to check how many apps are remaining.

from csv import reader
opened_file = open('googleplaystore.csv', encoding='utf8')
read_file = reader(opened_file)
googledataset = list(read_file)
googledataset = googledataset[1:]

from csv import reader
opened_file = open('AppStore.csv', encoding='utf8')
read_file = reader(opened_file)
dataset = list(read_file)
dataset = dataset[1:]

android_final = []
ios_final = []

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

for app in dataset:
    price = app[4]
    if price == '0':
        ios_final.append(app)

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

In [ ]: #The reason why we want to find an app profile that fits both the App Store and Google Play is to determine the kinds of apps that are likely to attract more users as our revenue is highly influenced and depends on the number of people using our apps.
Based on the code above, which isolates both dataset in a sparse list we can inspect both data sets to identify the columns we can also generate frequency tables to find out what the most common genres in each market are. So according to the results shown below the App Store is dominated by apps designed for fun (game), while Google Play shows more balanced and scope of both practical and fun (game) apps. Therefore we would use the prime_genre column for the AppStore data set and the Genres and Category columns for the GooglePlaystore data set.

opened_file = open('googleplaystore.csv')
from csv import reader
read_file = reader(opened_file)
googledataset = list(read_file)

opened_file = open('AppStore.csv')
from csv import reader
read_file = reader(opened_file)
dataset = list(read_file)

googledataset_free_apps = []
dataset_free_apps = []

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

for app in dataset:
    price = app[4]
    if price == '0':
        dataset_free_apps.append(app)

display_table(dataset_free_apps, -5)

Games : 55.6459560749507
Entertainment : 8.234714003944774
Photo & Video : 4.117357001972387
Social Networking : 3.5256410256410255
Education : 3.2544378698224854
Shopping : 2.983234714003945
Utilities : 2.687376725838264
Lifestyle : 2.3175542406311638
Finance : 2.0710059171597637
Sports : 1.947731755424063
Health & Fitness : 1.8737672583826428
Music : 1.6518737672583828
Book : 1.6272189349112427
Productivity : 1.528599605226825
News : 1.4299802761341223
Travel : 1.3806706114398422
Food & Drink : 1.0601577909270217
Weather : 0.7642998027613412
Reference : 0.4930966469428008
Navigation : 0.4930966469428008
Business : 0.4930966469428008
Catalogs : 0.22189349112426035
Medical : 0.1972386587712032

display_table(dataset_free_apps, 1)

FAMILY : 17.739043824701195
GAME : 10.56772908366534
TOOLS : 7.6195219123505975
BUSINESS : 4.44231075697211
PRODUCTIVITY : 3.944223107569721
LIFESTYLE : 3.6155378486055776
SPORTS : 3.5856573705179287
COMMUNICATION : 3.5856573705179287
MEDICAL : 3.5258964143426295
FINANCE : 3.476096175298805
HEALTH_AND_FITNESS : 3.237051792828685
PHOTOGRAPHY : 3.11723880478098
PERSONALIZATION : 3.0776892430278884
SHOPPING : 2.98366533864542
NEWS_AND_MAGAZINES : 2.7598047808764938
TRAVEL_AND_LOCAL : 2.450199203187251
DATING : 2.2609561752988045
BOOKS_AND_REFERENCE : 2.0219123505976095
VIDEO_PLAYERS : 1.7031872509960162
EDUCATION : 1.5139442231075697
ENTERTAINMENT : 1.4641434262948207
MAPS_AND_NAVIGATION : 1.347410358565739
FOOD_AND_DRINK : 1.245019920318725
HOUSE_AND_HOME : 0.8764940239043826
LIBRARIES_AND_DEMO : 0.8366533864541833
AUTO_AND_VEHICLES : 0.8137672583826428
WEATHER : 0.7370517928286853
EVENTS : 0.6274900398406374
ART_AND_DESIGN : 0.6175298804780877
COMICS : 0.5976096143426295
PARENTING : 0.5776892430278884
BEAUTY : 0.5278884462151394

In [ ]: #The reason why we want to find an app profile that fits both the App Store and Google Play is to determine the kinds of apps that are likely to attract more users as our revenue is highly influenced and depends on the number of people using our apps.
Based on the code above, which isolates both dataset in a sparse list we can inspect both data sets to identify the columns we can also generate frequency tables to find out what the most common genres in each market are. So according to the results shown below the App Store is dominated by apps designed for fun (game), while Google Play shows more balanced and scope of both practical and fun (game) apps. Therefore we would use the prime_genre column for the AppStore data set and the Genres and Category columns for the GooglePlaystore data set.

opened_file = open('googleplaystore.csv')
from csv import reader
read_file = reader(opened_file)
googledataset = list(read_file)

opened_file = open('AppStore.csv')
from csv import reader
read_file = reader(opened_file)
dataset = list(read_file)

googledataset_free_apps = []
dataset_free_apps = []

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

for app in dataset:
    price = app[4]
    if price == '0':
        dataset_free_apps.append(app)

display_table(dataset_free_apps, -5)

Games : 55.6459560749507
Entertainment : 8.234714003944774
Photo & Video : 4.117357001972387
Social Networking : 3.5256410256410255
Education : 3.2544378698224854
Shopping : 2.983234714003945
Utilities : 2.687376725838264
Lifestyle : 2.3175542406311638
Finance : 2.0710059171597637
Sports : 1.947731755424063
Health & Fitness : 1.8737672583826428
Music : 1.6518737672583828
Book : 1.6272189349112427
Productivity : 1.528599605226825
News : 1.4299802761341223
Travel : 1.3806706114398422
Food & Drink : 1.0601577909270217
Weather : 0.7642998027613412
Reference : 0.4930966469428008
Navigation : 0.4930966469428008
Business : 0.4930966469428008
Catalogs : 0.22189349112426035
Medical : 0.1972386587712032

display_table(dataset_free_apps, 1)

FAMILY : 17.739043824701195
GAME : 10.56772908366534
TOOLS : 7.6195219123505975
BUSINESS : 4.44231075697211
PRODUCTIVITY : 3.944223107569721
LIFESTYLE : 3.6155378486055776
SPORTS : 3.5856573705179287
COMMUNICATION : 3.5856573705179287
MEDICAL : 3.5258964143426295
FINANCE : 3.476096175298805
HEALTH_AND_FITNESS : 3.237051792828685
```