[81]:	<pre>.eet's take a look at the data, which consists of two files:</pre>
	<pre>apps = apps_with_duplicates.drop_duplicates() # Print the total number of apps</pre>
	# Have a look at a random sample of 5 rows n = 5 print(apps.sample(n)) Total number of apps in the dataset = 9659 Unnamed: 0 App Category \ 3250 4103 Synd e-Passbook FINANCE
E n	CD View Lite TOOLS
a b v H	Price have a few special characters (+ , \$) due to the way the numbers have been represented. This prevents the columns from the price have a few special characters (+ , \$) due to the way the numbers have been represented. This prevents the columns from the price have a few special characters (+ , \$) due to the way the numbers have been represented. This prevents the columns from the price have a few special characters (+ , \$) due to the way the numbers have been represented. This prevents the columns from the prevents the column from the prevents of t
	<pre># Loop for each column in cols_to_clean for col in cols_to_clean: # Loop for each char in chars_to_remove for char in chars_to_remove: # Replace the character with an empty string apps[col] = apps[col].apply(lambda x: x.replace(char, '')) # Print a summary of the apps dataframe print(apps.info()) <class 'pandas.core.frame.dataframe'=""> Int64Index: 9659 entries, 0 to 9658 Data columns (total 14 columns): Unnamed: 0 9659 non-null int64 App 9659 non-null object Category 9659 non-null object Rating 8196 non-null float64 Reviews 9659 non-null float64 Reviews 9659 non-null float64 Installs 9659 non-null object Type 9659 non-null object Price 9659 non-null object</class></pre>
E v	Content Rating 9659 non-null object 9659 non-null object Last Updated 9659 non-null object Current Ver 9651 non-null object Android Ver 9657 non-null object dtypes: float64(2), int64(2), object(10) memory usage: 1.1+ MB None 3. Correcting data types From the previous task we noticed that Installs and Price were categorized as object data type (and not int or float) as a would like. This is because these two columns originally had mixed input types: digits and special characters. To know more about Pandata types, read this. The four features that we will be working with most frequently henceforth are Installs, Size, Rating and Price. While Size Rating are both float (i.e. purely numerical data types), we still need to work on Installs and Price to make them numeric. import numpy as np # Convert Installs to float data type apps['Installs'] = apps['Installs'].astype(float) # Convert Price to float data type
	<pre>apps['Price'] = apps['Price'].astype(float) # Checking dtypes of the apps dataframe print(apps.info()) <class 'pandas.core.frame.dataframe'=""> Int64Index: 9659 entries, 0 to 9658 Data columns (total 14 columns): Unnamed: 0</class></pre>
V bb cc	Last Updated 9659 non-null object Current Ver 9651 non-null object Android Ver 9657 non-null object Android Ver 9657 non-null object Atypes: float64(4), int64(2), object(8) Memory usage: 1.1+ MB None 4. Exploring app categories With more than 1 billion active users in 190 countries around the world, Google Play continues to be an important distribution platform build a global audience. For businesses to get their apps in front of users, it's important to make them more quickly and easily discoverage on Google Play. To improve the overall search experience, Google has introduced the concept of grouping apps into categories. This brings us to the following questions: Which category has the highest share of (active) apps in the market? Is any specific category dominating the market? Which categories have the fewest number of apps? We will see that there are 31 unique app categories present in our dataset. Family and Game apps have the highest market prevalence interestingly, Tools, Business and Medical apps are also at the top. Import plotly plotly graph_objs as go # Print the total number of unique categories num_categories = len (apps! 'Category').unique()) print ('Number of categories = ', num_categories)
	<pre># Count the number of apps in each 'Category'. num_apps_in_category = apps['Category'].value_counts() # Sort num_apps_in_category in descending order based on the count of apps in each category sorted_num_apps_in_category = num_apps_in_category.sort_values(ascending = False) data = [go.Bar(</pre>
A s ii	5. Distribution of app ratings After having witnessed the market share for each category of apps, let's see how all these apps perform on an average. App ratings (on scale of 1 to 5) impact the discoverability, conversion of apps as well as the company's overall brand image. Ratings are a key performar indicator of an app. From our research, we found that the average volume of ratings across all app categories is 4.17. The histogram plot is skewed to the indicating that the majority of the apps are highly rated with only a few exceptions in the low-rated apps.
89]:	<pre># Average rating of apps avg_app_rating = apps['Rating'].mean() print('Average app rating = ', avg_app_rating) # Distribution of apps according to their ratings data = [go.Histogram(</pre>
	<pre>plotly.offline.iplot({'data': data, 'layout': layout})</pre> Average app rating = 4.173243045387994
	7. Relation between app category and app price So now comes the hard part. How are companies and developers supposed to make ends meet? What monetization strategies can companies use to maximize profit? The costs of apps are largely based on features, complexity, and platform.
c c c k	of your customer to pay for your app. A wrong price could break the deal before the download even happens. Potential customers could be turned off by what they perceive to be a shocking cost, or they might delete an app they've downloaded after receiving too many add or simply not getting their money's worth. Different categories demand different price ranges. Some apps that are simple and used daily, like the calculator app, should probably be sept free. However, it would make sense to charge for a highly-specialized medical app that diagnoses diabetic patients. Below, we see Medical and Family apps are the most expensive. Some medical apps extend even up to \\$80! All game apps are reasonably priced belo (\$20. import matplotlib.pyplot as plt fig, ax = plt.subplots() fig.set_size_inches(15, 8) # Select a few popular app categories popular_app_cats = apps[apps.Category.isin(['GAME', 'FAMILY', 'PHOTOGRAPHY', 'MEDICAL', 'TOOLS', 'FINANCE', 'LIFESTYLE', 'BUSINESS'])] # Examine the price trend by plotting Price vs Category
	ax = sns.stripplot(x = popular_app_cats['Price'], y = popular_app_cats['Category'], jitter=True, linewidth ax.set_title('App pricing trend across categories') # Apps whose Price is greater than 200 apps_above_200 = popular_app_cats[opular_app_cats['Price'] > 200] apps_above_200[['Category', 'App', 'Price']] Category
	### App pricing trend across categories Hank Care Lam Rich Sab.99
	GAME FAMILY MEDICAL PHOTOGRAPHY TOOLS 0 50 100 150 200 250 300 350 40
95]:	**Examine price vs category with the authentic apps (apps_under_100) **Examine price vs category with the authentic apps (apps_under_100) **Examine price vs category with the authentic apps (apps_under_100) **ax = sns.stripplot(x = 'Price', y = 'Category', data = apps_under_100, jitter = True, linewidth = 1) ax.set_title('App pricing trend across categories after filtering for junk apps') **App pricing trend across categories after filtering for junk apps') **App pricing trend across categories after filtering for junk apps')
	FINANCE UFESTYLE GAME DATE FAMILY PHOTOGRAPHY TOOLS DEPORT OF THE Price O 20 40 Price O 80 80
F fi	9. Popularity of paid apps vs free apps For apps in the Play Store today, there are five types of pricing strategies: free, freemium, paid, paymium, and subscription. Let's focus of free and paid apps only. Some characteristics of free apps are: Free to download. Main source of income often comes from advertisements. Often created by companies that have other products and the app serves as an extension of those products. Can serve as a tool for customer retention, communication, and customer service. Some characteristics of paid apps are: Users are asked to pay once for the app to download and use it. The user can't really get a feel for the app before buying it. Are paid apps installed as much as free apps? It turns out that paid apps have a relatively lower number of installs than free apps, thoughed difference is not as stark as I would have expected! Utrace0 = go_Box(# Data for paid apps y = apps[apps['Type'] == 'Paid']('Installs'), name = 'Paid')
	<pre>tracel = go.Box(# Data for free apps y = apps[apps['Type'] == 'Free']['Installs'], name = 'Free') layout = go.Layout(title = "Number of downloads of paid apps vs. free apps", yaxis = dict(title = "Log number of downloads",</pre>
N S	10. Sentiment analysis of user reviews Mining user review data to determine how people feel about your product, brand, or service can be done using a technique called sentiment analysis. User reviews for apps can be analyzed to identify if the mood is positive, negative or neutral about that app. For example, positive words in an app review might include words such as 'amazing', 'friendly', 'good', 'great', and 'love'. Negative words might words like 'malware', 'hate', 'problem', 'refund', and 'incompetent'. By plotting sentiment polarity scores of user reviews for paid and free apps, we observe that free apps receive a lot of harsh comments, indicated by the outliers on the negative y-axis. Reviews for paid apps appear never to be extremely negative. This may indicate somethabout app quality, i.e., paid apps being of higher quality than free apps on average. The median polarity score for paid apps is a little
B ii a h	nigher than free apps, thereby syncing with our previous observation. n this notebook, we analyzed over ten thousand apps from the Google Play Store. We can use our findings to inform our decisions show
B ii a h	