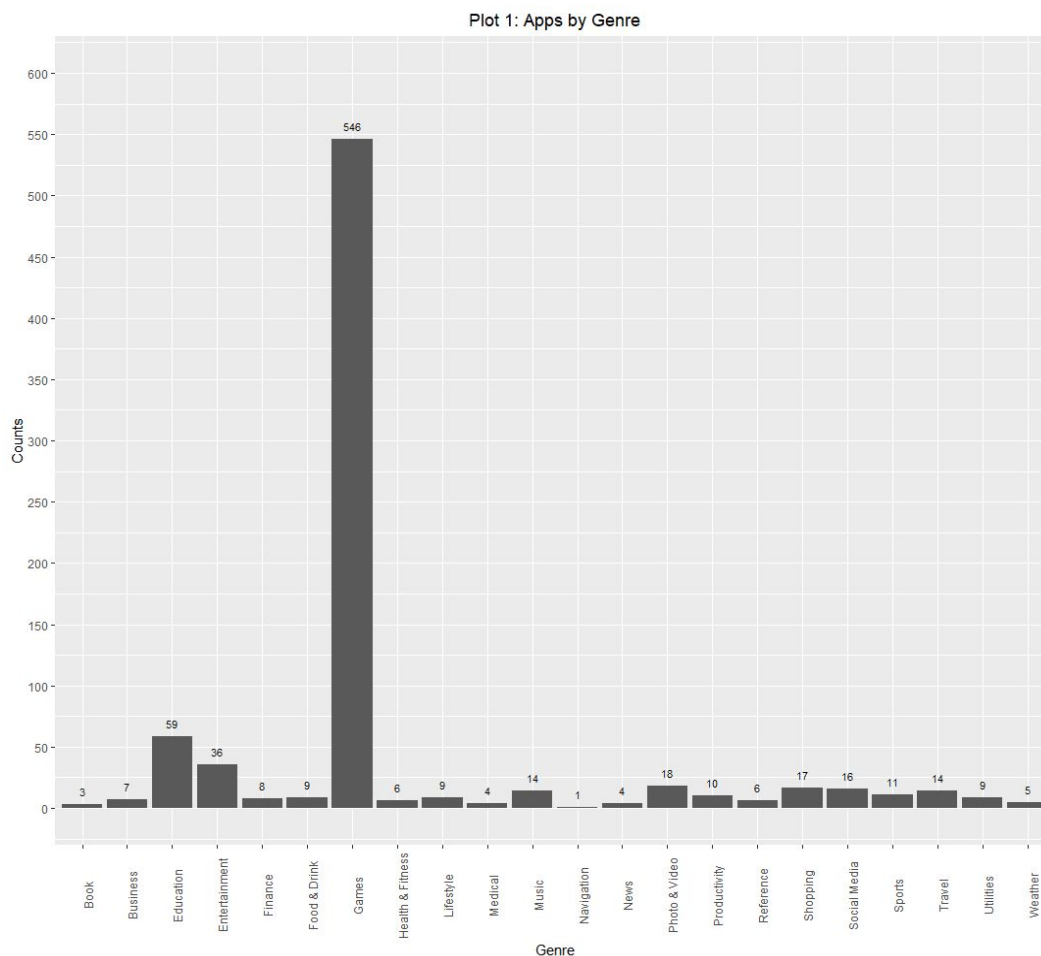


Module 6 Final Project  
Brad Viles  
ALY 6000  
October 25, 2020

## Abstract

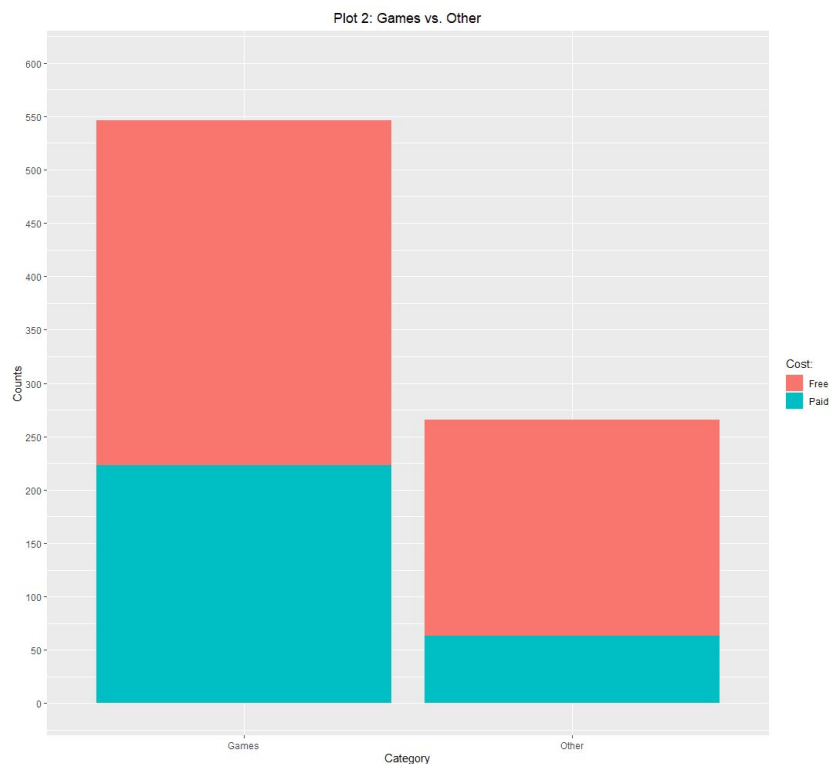
I chose a collection of three data sets to work with in this project that summarized characteristics of mobile applications available on the Apple store and the Google Play store. The main goal of this analysis is to understand the characteristics of apps that are present in both mobile application stores in order to increase potential audience for future app development. I will explore characteristics such as genre, the relationship between ratings and reviews, and content rating of the applications. The datasets that I worked with are as follows: Apple - 7197 objects of 17 variables, Google Play - 10841 objects of 13 variables, and Google Play 32K - 32000 objects of 12 variables (Prakash and Koshy, 2019) (Ramanathan 2017). In order to best understand the intersection of the two application stores, a cleaned data set combining the Apple data and the Google Play 32K data was created with 812 objects of 16 variables. It can be assumed that these do not encompass all apps available on each respective store, especially as the Apple data was collected in July of 2017 and the Google data was collected in April of 2019. Nonetheless, we will be able to make some starting generalizations to create a roadmap to create our own successful multi-platform app.

## Initial Analysis



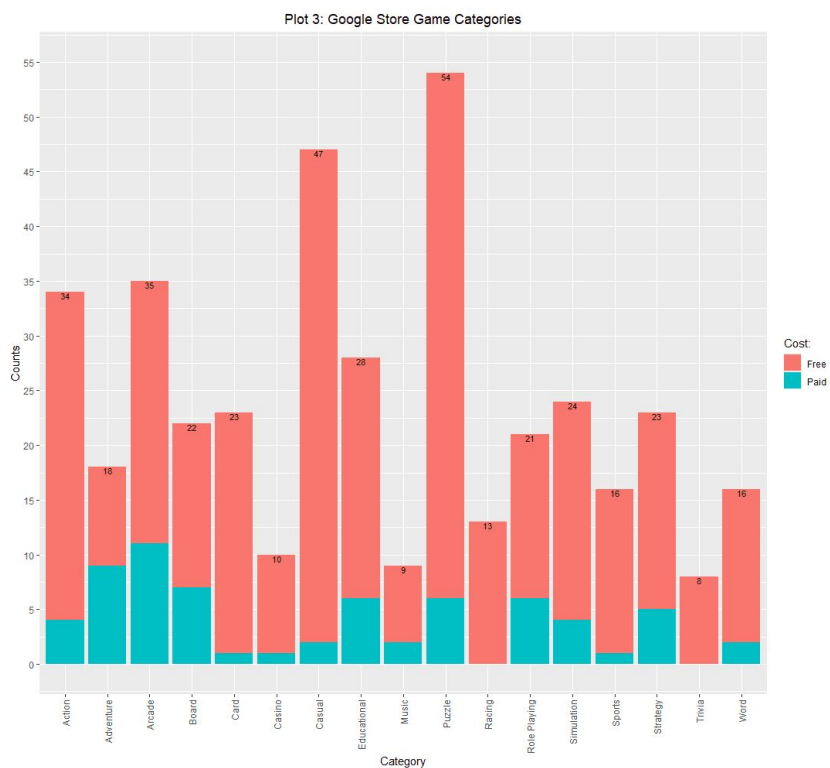
I began my analysis by breaking down the combined dataset into genres in order to see which of the 22 notated genres was most prevalent within the data. Unsurprisingly the dataset was dominated by games as shown in Plot 1. Upon seeing this I decided it would be more beneficial to separate games

into its own category compared against the rest of the population. Additionally, Plot 2 is segmented by free versus paid apps. This shows that even when compared against the

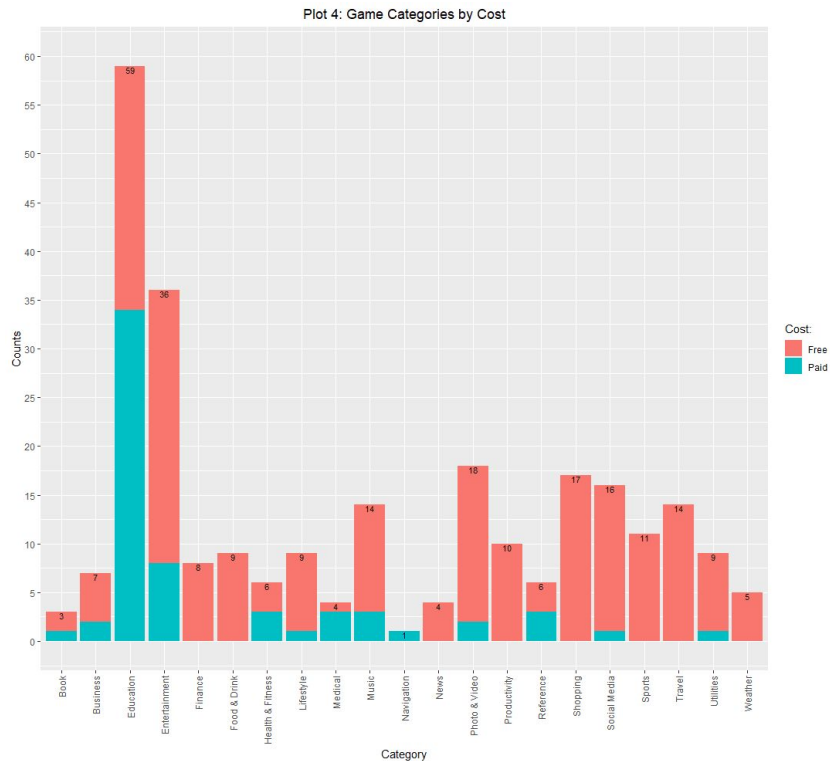


rest of the population games are far more prevalent, and there are almost as many paid games as there are other apps within the population. In total the combined data includes 266 non-games to 546 games. On initial inspection of the combined dataset it appears that games are the most successful application type on both stores, and notably the data does not notate “freemium” apps that include additional in-app purchases. Plot 3 breaks down the games data further by category.

Notably, this plot uses only data from the Google Play store, as the Apple data did not subset their games data by any specific subgenre or category. Puzzle and casual games top the charts with 54 and 47 apps respectively, but also have some of the largest counts of free apps. I will again reiterate that this data does not include apps with in-app purchases which can still generate large amounts of revenue. For strictly paid apps, arcade, and adventure games edge out the field, and notably have a high ratio of their sample as paid apps. On the other end of the spectrum are



racing and trivia games which logged no instances of paid apps on the Google store.

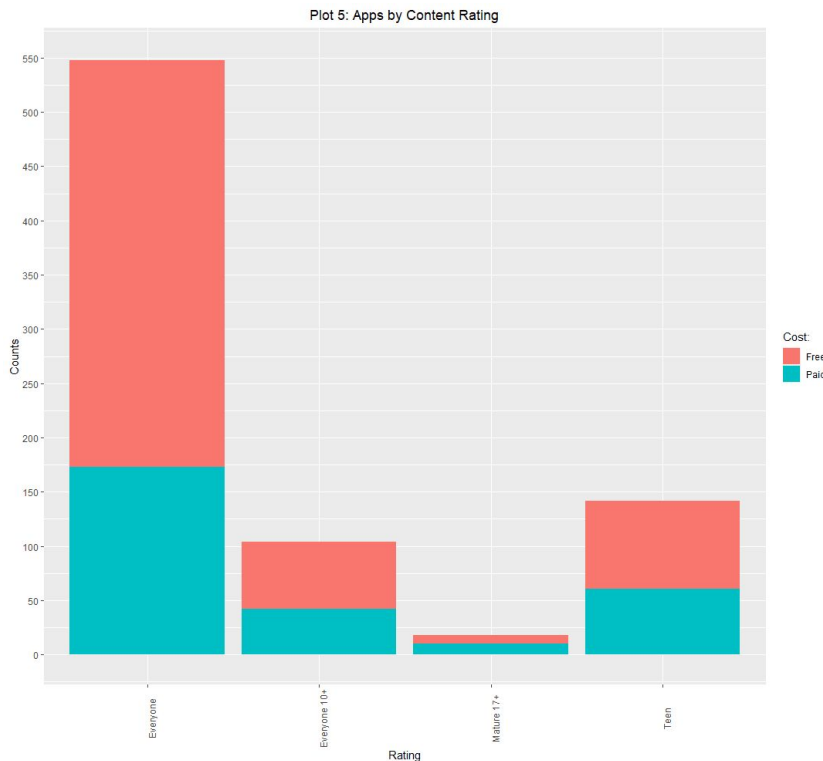


Plot 4 switches gears to explore the other side of Plot 2 within the combined dataset. Logging an enormous count of paid applications are educational apps. 34 of the 59 educational apps in the population are paid, which given its own category would still be second to entertainment apps by only 2 counts. There exists a large population of apps that are entirely or primarily free, with only paid entertainment apps getting close to double digits with 8 counts. Other notable features from this plot are

that 75% of medical apps, 50% of health & fitness apps, and 50% of reference apps are paid. The sole navigation app in the combined sample is paid as well, but due to the small sample size we may guess that each platform utilizes its own proprietary navigation application, and so there is no overlap within the respective application stores of navigation apps.

By breaking down the different categories of apps in this manner, we appear to be faced by the primary decision of whether we should aim to develop a game or a different type of app. From there the choice is how revenue within the app is generated. There is an option to simply charge a flat price for download, or create a free app with other revenue models like sponsored ads, data collection for resale, or freemium in-app purchasing. There are clearly certain genres both within and outside the games subset that lend themselves towards one revenue model over another.

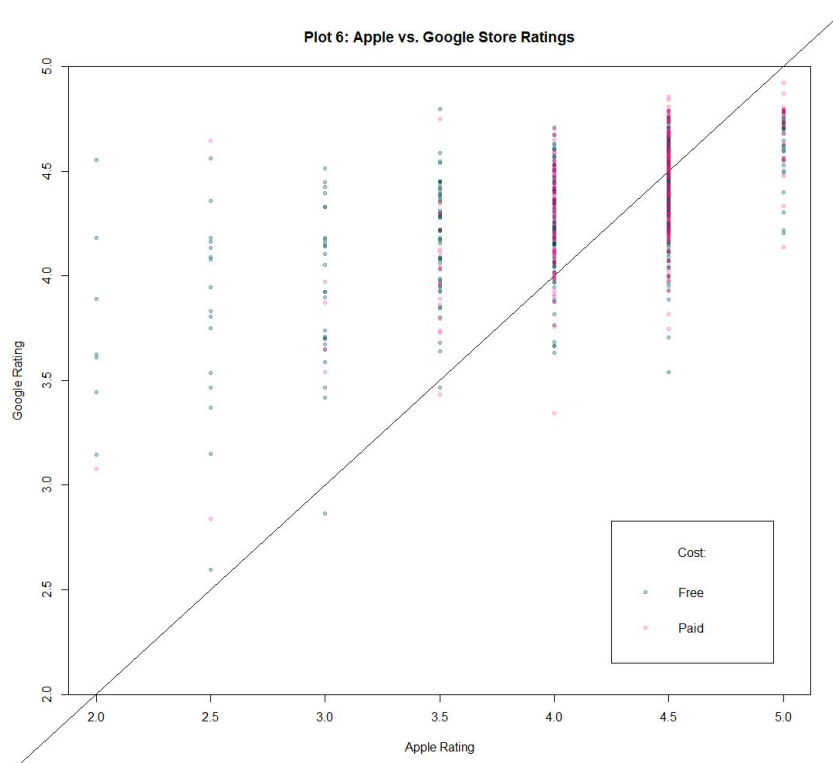
The final aspect of free versus paid apps that I was interested in exploring was a breakdown by content rating. Are you more likely to find a paid game in a rating bracket for teens or adults as they are a more likely target audience? Plot 5 shows that most apps in the combined dataset are rated for everyone, with about 30% of those being paid apps. The other age ranges are split fairly evenly down the middle of free versus paid, meaning that if we do decide to develop an app targeted at a higher age range or with a higher content rating, it is more viable to use a pay per download model than with an app

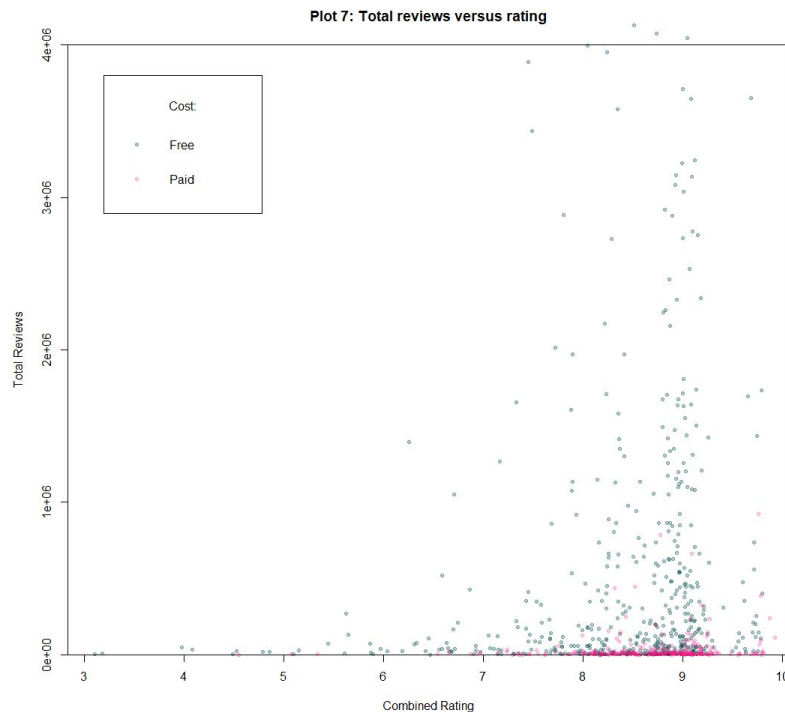


rated for all ages. The caveat to that is that there are far more apps designed for all ages, meaning there is a greater potential target audience with a lower rating. This could actually mean increased revenue if we were to implement a different revenue model into a free app rated for everyone.

The other major indicators of success that I wanted to explore in my combined dataset were ratings and reviews of apps. The Apple rating data was rounded to the nearest .5 as you can see in Plot 6, while the Google rating data was a raw decimal.

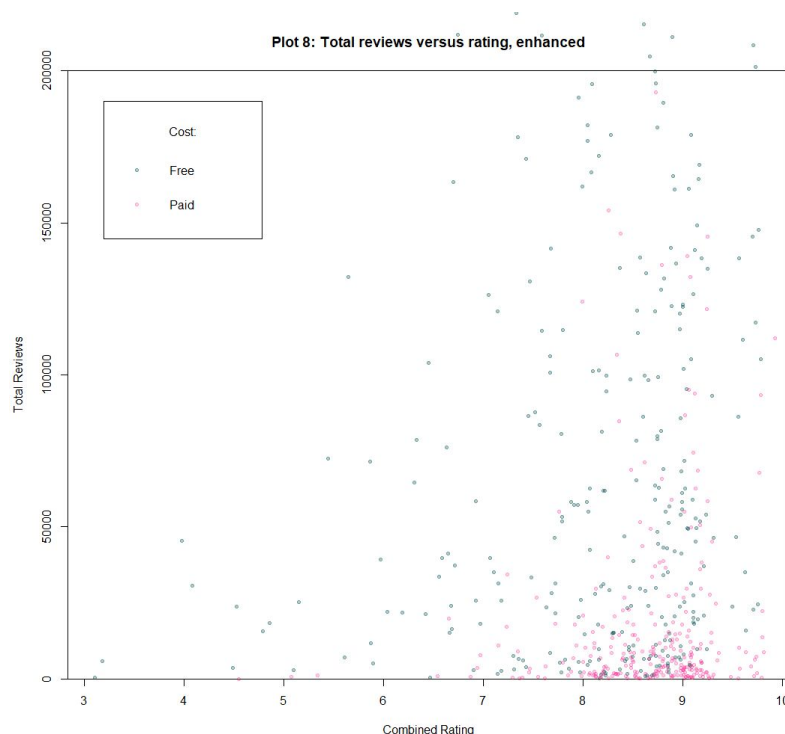
Included on the plot is a  $y=x$  reference line denoting what an equal rating between the two mobile stores would look like. Plot 6 shows that on average apps are higher rated in the Google Play store than in the Apple store. This seems to primarily be true for free apps. What we see from the paid apps, is a cluster in the top right corner of the graph, fairly equally bisected by the reference line. What this tells us is that while there is some variance between ratings between the two stores, paid apps are consistently highly rated on both platforms. Looking closer we can see that paid apps that are lower rated on the Apple store maintain a relatively higher performance in the Google



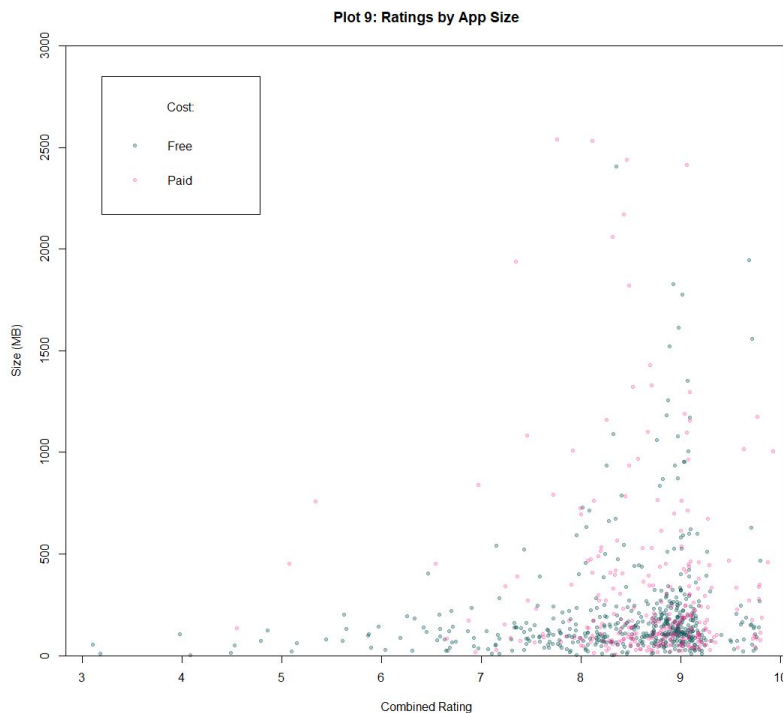


store, consistent with the trend shown by the free apps. This made me curious about the sample size of ratings and reviews for these apps, which is shown in Plot 7. This plot uses a combined rating of the sum of apple and google ratings, and a sum of the review count between the two platforms. This makes it vastly apparent that free apps in general have a much higher review count than paid apps. It also shows that apps with a high number of reviews are generally highly rated, with very few falling below 7.5/10 in a combined rating.

My assumption from this is not that high numbers of reviews create a good app, but rather good apps generate a high number of reviews. In order to dive deeper into the paid portion of this plot, Plot 8 scales the y limit down from 4 million to 200,000. We see



more of the same trend, but with far fewer total reviews for paid apps, with the bulk garnering fewer than 50,000 reviews. My hypothesis is that this is due to a greatly decreased number of downloads leading to a greatly decreased number of reviews. In comparison to the free apps shown in Plot 7, the bulk of the combined ratings for paid apps fall between 8/10 to 9.5/10, again indicating a very high quality in paid applications. Once again I find myself wishing there was more data on the



other revenue models or total revenue numbers for free download apps in order to better examine potential business models for our application development. A final characteristic between free and paid apps that I wanted to examine was application size. Plot 9 shows a breakdown of the size of each app compared against its combined rating from both platform stores. Most apps from the sample are under 500MB, but for those larger than 500MB there seems to be a range of ratings from 8/10 to 9.5/10. From this plot it is hard to differentiate much

between free versus paid apps, however the conclusion that can be drawn is that there is correlation between larger apps and a generally high rating.

## Conclusion

In conclusion, this data suggests that there are several directions to take with our potential app. I would caution that this data does not paint a full picture however, and there are attributes that were not appropriately present within the data sets that I would want to examine further to strongly recommend a direction. Specifically, more data regarding the revenue generation of free download apps would be very important to further analysis. There are many apps that are not classified as “paid” within this dataset that will still generate revenue. Additionally, a key metric regarding application success and audience that I would want to explore further is downloads. One of the datasets did include data on downloads, but in a largely segmented factor that was not easily analyzed or generalized across my sample. This alongside the ratings and reviews data that I explored here would give a better understanding of potential success of different categories of application.

If the goal is simply to create the best paid app, my recommendation based purely on this dataset is to proceed with the development of an education application. There is

strong precedent within that subset of the industry for a paid download, and educational apps hold a strong market share outside of games. If the goal is to maximize revenue with our application, I will wait on further data regarding revenue models in order to make a recommendation. I hypothesize that the best revenue model is a freemium model that makes use of in-app purchases based on the prevalence of games within the market, and the relatively small share of them that are paid to download.



## **Bibliography:**

Prakash, G. and Koshy, J. (April 2019). Google PlayStore App analytics (Version 1) [Data set].

<https://www.kaggle.com/gauthamp10/google-playstore-apps?select=Google-Playstore-Full.csv>

Ramanathan. (July 2017). Mobile App Store (Version 7) [Data set].

<https://www.kaggle.com/ramamet4/app-store-apple-data-set-10k-apps>

## Appendix A: R code for cleaning data

```
1 setwd("C:/NECPS/ALY6000/M6")
2 library(FSA)
3 library(FSAdat)
4 library(magrittr)
5 library(dplyr)
6 library(plotrix)
7 library(ggplot2)
8 library(moments)
9
10 ##### CLEANING DATA
11
12 ### Load 2 data sets from google play store and apple store
13 apple <- read.csv("Applestore.csv")
14 gplay <- read.csv("googleplaystore32k.csv")
15 gplay2 <- read.csv("googleplaystore.csv")
16
17 ### Convert google values to numeric
18 gplay$Rating <- as.character(gplay$Rating)
19 gplay$Rating <- as.numeric(gplay$Rating)
20
21 gplay$Price <- as.character(gplay$Price)
22 gplay$Price = as.numeric(gsub("\\$", "", gplay$Price))
23
24 ### Content rating cleaning
25 gplay$Content.Rating <- as.character(gplay$Content.Rating)
26 gplay$Content.Rating[gplay$Content.Rating %in% c("$0.99", "0", "100,000+")] <- NA
27 gplay$Content.Rating[gplay$Content.Rating %in% "Adults only 18+"] <- "Mature 17+"
28 gplay$Content.Rating <- as.factor(gplay$Content.Rating)
29
30 apple$cont_rating <- as.character(apple$cont_rating)
31 apple$cont_rating[apple$cont_rating %in% "4+"] <- "Everyone"
32 apple$cont_rating[apple$cont_rating %in% "9+"] <- "Everyone 10+"
33 apple$cont_rating[apple$cont_rating %in% "12+"] <- "Teen"
34 apple$cont_rating[apple$cont_rating %in% "17+"] <- "Mature 17+"
35 apple$cont_rating <- as.factor(apple$cont_rating)
36
37 ### Reviews
38 apple$rating_count_tot <- as.numeric(apple$rating_count_tot)
39
40 ### Slice columns to be merged
41 Capple <- apple[c(3,4,6,7,9,12,13)]
42 Cgplay <- gplay[c(1,3:5,7,8)]
43
44 ### Merge Capple and Cgplay
45 names(Capple)[1] <- "App Name"
46 names(Cgplay)[1] <- "App Name"
47 apps <- merge(Capple, Cgplay, by="App Name")
48 apps$`App Name` <- droplevels(apps$`App Name`)
49
50 ### Drop duplicate records
51 apps = apps[order(apps[, "App Name"], -apps[, "Reviews"]),]
52 apps = apps[!duplicated(apps$`App Name`),]
53
54 ### Additional cleaning
55 apps$prime_genre <- as.character(apps$prime_genre)
56 apps$prime_genre[apps$prime_genre %in% "Social Networking"] <- "Social Media"
57 apps$prime_genre <- as.factor(apps$prime_genre)
58
59 ### Create Free vs. Paid factor
60 apps$freepaid <- NA
61 apps$freepaid <- as.character(apps$Price)
62 apps$freepaid[apps$freepaid %in% "0"] <- "Free"
63 apps$freepaid[!apps$freepaid %in% "Free"] <- "Paid"
64 apps$freepaid <- as.factor(apps$freepaid)
65
66 gplay$freepaid <- NA
67 gplay$freepaid <- as.character(gplay$Price)
68 gplay$freepaid[gplay$freepaid %in% "0"] <- "Free"
69 gplay$freepaid[!gplay$freepaid %in% "Free"] <- "Paid"
70 gplay$freepaid <- as.factor(gplay$freepaid)
71
72 ### Clean name to fit on plots with LAS=2
73 apps$prime_genre <- as.character(apps$prime_genre)
74 apps$prime_genre[apps$prime_genre %in% "Social Networking"] <- "Social Media"
75 apps$prime_genre <- as.factor(apps$prime_genre)
```

## Appendix B: Coding plots

```

80 ##### PLOTS #####
81
82 ### PLOT 1 ### CLEANING ###
83 genre <- apps$prime_genre
84 GENREcounts <- data.frame(genre)
85
86 ### Plot 1 ### genre counts - COMBINED DATA
87 p1 <- ggplot(data=GENREcounts, aes(x=genre)) +
88   geom_bar()+
89   geom_text(stat='count', aes(label=..count..), vjust=-1, cex=3) +
90   theme(plot.title = element_text(hjust = 0.5)) +
91   theme(axis.text.x = element_text(angle = 90)) +
92   scale_y_continuous(name="Counts",limits=c(0,600), breaks=seq(0,600,50)) +
93   labs(title="Plot 1: Apps by Genre",x="Genre",y="Counts")
94 print(p1)
95
96
97
98 ### PLOT 2 ### CLEANING ###
99 gamesdata <- as.character(genre)
100 gamesdata[!gamesdata %in% "Games"] <- "Other"
101 apps$gameother <- gamesdata
102
103 ### Plot 2 ### genre - games vs non games and free vs paid - COMBINED DATA
104 p2 <- ggplot(apps) +
105   geom_bar(aes(x=gameother, fill=freepaid)) +
106   labs(title="Plot 2: Games vs. Other",x="Category",y="Counts") +
107   theme(plot.title = element_text(hjust = 0.5)) +
108   scale_y_continuous(name="Counts",limits=c(0,600), breaks=seq(0,600,50)) +
109   scale_fill_discrete(name="Cost:")
110 print(p2)
111
112
113
114 ### PLOT 3 ### CLEANING ###
115 games <- gplay
116 gamecats <- c("GAME_ACTION", "GAME_ADVENTURE", "GAME_ARCADE", "GAME_BOARD", "GAME_CARD",
117   "GAME_CASINO", "GAME_CASUAL", "GAME_EDUCATIONAL", "GAME_MUSIC", "GAME_PUZZLE",
118   "GAME_RACING", "GAME_ROLE_PLAYING", "GAME_SIMULATION", "GAME_SPORTS",
119   "GAME_STRATEGY", "GAME_TRIVIA", "GAME_WORD")
120 games <- games %<>% filterD(Category == gamecats)
121 gamecatsclean <- c("Action", "Adventure", "Arcade", "Board", "Card",
122   "Casino", "Casual", "Educational", "Music", "Puzzle",
123   "Racing", "Role Playing", "Simulation", "Sports",
124   "Strategy", "Trivia", "Word")
125 levels(games$Category) <- gamecatsclean
126 GPfree <- games
127 GPfree <- GPfree %<>% filterD(freepaid == "Free")
128 GPpaid <- games
129 GPpaid <- GPpaid %<>% filterD(freepaid == "Paid")
130
131 ### Plot 3 ### ggplot with free/paid of game categories
132 p3 <- ggplot(data=games, aes(x=Category)) +
133   geom_bar(aes(x=Category, fill=freepaid)) +
134   labs(title="Plot 3: Google Store Game Categories",x="Category",y="Counts") +
135   theme(plot.title = element_text(hjust = 0.5)) +
136   scale_fill_discrete(name="Cost:") +
137   scale_y_continuous(name="Counts",limits=c(0,55), breaks=seq(0,60,5)) +
138   geom_text(stat='count', aes(label=..count..), vjust=1, cex=3) +
139   theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))
140 print(p3)
141
142
143
144 ### PLOT 4 ### CLEANING ###
145 nongames <- apps
146 nongames <- nongames %<>% filterD(prime_genre != "Games")
147
148 ### Plot 4 ### genre counts without games, free vs paid - COMBINED DATA
149 p4 <- ggplot(data=nongames, aes(x=prime_genre)) +
150   geom_bar(aes(x=prime_genre, fill=freepaid)) +
151   labs(title="Plot 4: Game Categories by Cost",x="Category",y="Counts") +
152   theme(plot.title = element_text(hjust = 0.5)) +
153   scale_fill_discrete(name="Cost:") +
154   scale_y_continuous(name="Counts",limits=c(0,60), breaks=seq(0,60,5)) +
155   geom_text(stat='count', aes(label=..count..), vjust=1, cex=3) +
156   theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))
157 print(p4)
158
159
160
161 ### PLOT 5 ### CLEANING ###
162 rated <- apps$Content.Rating
163 ratedcounts <- data.frame(rated)
164
165 ### Plot 5 ### genre counts - COMBINED DATA
166 p5 <- ggplot(data=apps, aes(x=Content.Rating, fill=freepaid)) +
167   geom_bar()+
168   theme(plot.title = element_text(hjust = 0.5)) +
169   theme(axis.text.x = element_text(angle = 90)) +
170   scale_fill_discrete(name="Cost:") +
171   scale_y_continuous(name="Counts",limits=c(0,550), breaks=seq(0,600,50)) +
172   labs(title="Plot 5: Apps by Content Rating",x="Rating",y="Counts")
173 print(p5)
174

```



```

177 ### PLOT 6 ### CLEANING ###
178 free <- apps
179 free <- free %<>% filterD(freepaid == "Free")
180 paid <- apps
181 paid <- paid %<>% filterD(freepaid == "Paid")
182
183 ### Plot 6 ### Ratings for free apps versus paid apps - COMBINED DATA
184 plot(Rating~user_rating, data=free, pch=20, col=rgb(0,0.3,0.3,0.3),
185       main="Plot 6: Apple vs. Google Store Ratings",
186       ylab="Google Rating", xlab="Apple Rating",
187       xlim=c(2,5), ylim=c(2,5))
188 abline(coef=c(0,1))
189 points(Rating~user_rating, data=paid, pch=20, col=rgb(1,0,0.5,0.2))
190 legend("bottomright", title="Cost:", inset=.05, c("Free", "Paid"),
191       pch=20, col=c(rgb(0,0.3,0.3,0.3),rgb(1,0,0.5,0.2)))
192
193
194
195 ### PLOT 7 ### CLEANING ###
196 free$totalreviews <- free$rating_count_tot + free$Reviews
197 paid$totalreviews <- paid$rating_count_tot + paid$Reviews
198 free$tenrating <- (free$user_rating + free$Rating)
199 paid$tenrating <- (paid$user_rating + paid$Rating)
200
201 ### Plot 7 ### Number of reviews for free apps vs paid apps
202 plot(totalreviews~tenrating,
203       data=free,
204       pch=20,
205       col=rgb(0,0.3,0.3,0.3),
206       main="Plot 7: Total reviews versus rating",
207       xlab="Combined Rating",
208       ylab="Total Reviews",
209       ylim=c(0,4000000),
210       las=0)
211 points(totalreviews~tenrating, data=paid, pch=20, col=rgb(1,0,0.5,0.2))
212 legend("topleft", title="Cost:", inset=.05, c("Free", "Paid"),
213       pch=20, col=c(rgb(0,0.3,0.3,0.3),rgb(1,0,0.5,0.2)))
214
215
216
217 ### Plot 8 ### Zoom in on plot 7
218 plot(totalreviews~tenrating,
219       data=free,
220       pch=20,
221       col=rgb(0,0.3,0.3,0.3),
222       main="Plot 8: Total reviews versus rating, enhanced",
223       xlab="Combined Rating",
224       ylab="Total Reviews",
225       ylim=c(0,200000),
226       las=0)
227 points(totalreviews~tenrating, data=paid, pch=20, col=rgb(1,0,0.5,0.2))
228 legend("topleft", title="Cost:", inset=.05, c("Free", "Paid"),
229       pch=20, col=c(rgb(0,0.3,0.3,0.3),rgb(1,0,0.5,0.2)))
230
231
232
233 ### PLOT 9 ### CLEANING ###
234 apps$tenrating <- (apps$user_rating + apps$Rating)
235 apps$sizeMB <- (apps$size_bytes/1000000)
236
237 ### Plot 9 ### size x rating
238 plot(sizeMB~tenrating,
239       data=free,
240       xlab="Combined Rating",
241       ylab="Size (MB)",
242       ylim=c(0,3000),
243       main="Plot 9: Ratings by App Size",
244       pch=20,
245       col=rgb(0,0.3,0.3,0.3)
246       )
247 points(sizeMB~tenrating, data=paid, pch=20, col=rgb(1,0,0.5,0.2))
248 legend("topleft", title="Cost:", inset=.05, c("Free", "Paid"),
249       pch=20, col=c(rgb(0,0.3,0.3,0.3),rgb(1,0,0.5,0.2)))
250

```