

# Lab 2 Research Proposal

Oleg Ananyev, Oren Carmeli, Romain Hardy, Sam Rosenberg

Fall 2021

## 1 Introduction

Lab 2 Instructions: Your introduction should present a research question and explain the concept that you're attempting to measure and how it will be operationalized. This section should pave the way for the body of the report, preparing the reader to understand why the models are constructed the way that they are. It is not enough to simply say, "We are looking for product features that enhance product success." Your introduction must do work for you, focusing the reader on a specific measurement goal, making them care about it, and propelling the narrative forward. This is also a good time to put your work into context, discuss cross-cutting issues, and assess the overall appropriateness of the data.

In 2020, Statista [reported](#) a total of ~220 billion mobile application downloads globally. Another [study](#) found that ~70% of all US digital media is spent on mobile applications. As mobile applications become one of the the principal media through which organizations engage with their customers, it becomes increasingly valuable to understand the drivers that lead to greater usage and better customer experience.

The Google Play Store is one of the major hubs for Android mobile phone and tablet applications. Users can download any application on the Google Play Store for personal consumption across a wide range of categories. Making an application stand out in a sea of thousands of competing applications is no trivial task, however. To succeed, developers must carefully consider factors such as price, application size, and genre. Another variable that may be critical to an application's success is consumer rating. Today, smart algorithms play a key role in suggesting applications to consumers, creating a feedback loop that propels certain applications towards success and leaves many others behind. An understanding of the relationship between consumer rating and application success would be incredibly valuable to developers seeking to create the next viral application.

The following study a causal analysis of the relationship between average consumer rating and application downloads. Using a data set scraped directly from the Google Play Store, we will build a linear model that assesses the importance (or lack of importance) of average consumer rating on the number of downloads, with additional variables such as price, application size, and application category serving as controls. Confirming the existence of a causal pathway would signal application developers to invest heavily in improving their ratings, for instance by buying positive reviews from consumers. Our study may also be of interest to Google, who has a vested interest in preventing developers from unfairly inflating their ratings.

Given our prior beliefs on what factors motivate individuals to download applications, we believe there are omitted variables not included in our data set that may influence application downloads. These include brand awareness from marketing campaigns, differences in existing customer bases, seasonality, and individual life circumstances, among possible others. In spite of these limitations, our study should provide useful insight into the factors that cause an application's success.

The paper is structured as follows. Section 2 outlines our research question, the causal model we will use to contextualize our regression analysis, and the research design. Section 3 describes our data, the explanatory variables we will include in the modeling phase, and the transformations we will apply. Section 4 contains our statistical models, and Section 5 discusses their significance as well as their statistical validity. Finally, in

Section 6 we present our conclusions and discuss the implications of our results.

## 2 Data and Research Design

Lab 2 Instructions: After you have presented the introduction and the concepts that are under investigation, what data are you going to use to answer the questions? What type of research design are you using? What type of models are you going to estimate, and what goals do you have for these models?

The goal of this study is to assess the causal factors of application success. Specifically, we seek to determine if there is a statistically significant relationship between average consumer rating and application downloads. In the ensuing sections, we will answer the following question:

*Does a higher average consumer rating lead to greater downloads for Google Play Store applications?*

In order to answer this question, we will leverage publicly available data about applications available on the Google Play Store. This data was randomly scraped from the Google Play Store interface and uploaded to Kaggle.com in 2019. It contains key information about each sampled application, such as downloads, file size, rating, category, and price. In total, the data contains records of ~10,000 applications.

For the modeling step, we will use OLS regression. OLS regression is the plug-in estimator for the best linear predictor of a dependent random variable given the joint distribution of a set of independent random variables. In our case, the dependent variable is our conceptualization of application success (the number of downloads an application receives), and the independent variables are the various attributes of Google Play Store applications. Although OLS regression is often used with the goal of making predictions on new data, here it will be used to answer a causal question about the relationship between variables. By interpreting model coefficients within the context of a causal theory, we will develop a statistically valid argument that answers the research question.

Our analysis will make use the following features:

### Dependent Variable:

- **installs** — The accumulated number of installs since the application was uploaded to the Google Play Store.

### Independent Variables:

- **size** — The memory space occupied by the application.
- **reviews** — The accumulated number of reviews since the application was uploaded to the Google Play Store. This feature will be used to remove outliers (see the EDA section).
- **rating** — The average consumer rating for the application out of 5.
- **price** — The price of the application.
- **category** — The category tagged for the application (i.e Lifestyle, Game)
- **content\_rating** — The official content rating given to the application (i.e Teen, Everyone, Mature 17+)
- **type** — The price type of the application. Possible options are Free (free to download) or Paid (costs money to download).

## 2.1 Causal Model

In this section, we describe the causal model which will serve as the reference point for our analysis. We also discuss possible omitted variables and the impact they may have on our study.

### 2.1.1 Simple Model

We identify five factors that bear causal influence on the success of an application. These are (1) consumer rating, (2) production quality, (3) category, (4) age, and (5) size. Our proposed causal model is shown below.

## Consumer Rating

Our main variable of interest is consumer rating, which we are operationalizing using the `rating` field in the data set. We hypothesize that higher application ratings should cause greater downloads. This is because an application that has been highly rated is one that has been deemed worthwhile by other users. When new consumers decide whether or not to download an application, they will likely trust the opinions of their peers and download applications with positive reviews. Conversely, we expect applications with negative ratings to have less success, since other users have judged them poorly. There should not be causal pathways from rating to any of the other explanatory variables, since they are determined during the development phase of the application while consumer rating is decided once the application is published to the Google Play Store.

Although we have not included it in the diagram, there is the possibility that a reverse causal path exists from application success to consumer rating. Successful applications are those that are enjoyable to a large number of consumers. Therefore, it is possible that users will rate applications with high download counts higher than applications with low download counts because they are primed to believe that such applications are better—otherwise, the successful applications would not have received so many views. Generally, we expect the reverse path to be weaker than the forward path. If ratings and downloads both have positive effects on one another, however, our models may suffer from positive feedback. We will discuss the implications of this effect in the results section.

## Production Quality

We expect production quality to also have a causal effect on application success. Overall, we believe that applications with high production quality (such as nice graphics, fluid controls, and few bugs) should be more successful than low-quality applications since they have more appealing features. Because there is no exact measure of production quality in our data set, we will represent it using the `price` and `type` fields. This is valid assuming that price is a good proxy of production quality. Generally, paid applications should be associated with higher quality, but this may not always be the case. For instance, some free applications could have high production quality while certain paid applications could have low production quality. We recognize that this assumption is imperfect, and we will discuss the possible ramifications later in our analysis.

We also anticipate a causal pathway from production quality to rating, for similar reasons as the pathway from quality to success. Specifically, we believe that consumers are likely to rate high-quality applications positively since those applications have more desirable features. Conversely, low-quality applications will be reviewed negatively. Again, this causal pathway may be difficult to assess if the relationship between price and quality is weaker than we assume.

## 2.2 Model Building

Lab 2 Instructions: You will next build a set of models to investigate your research question, documenting your decisions. Here are some things to keep in mind during your model building process: 1. *What do you want to measure?* Make sure you identify one, or a few, variables that will allow you to derive conclusions relevant to your research question, and include those variables in all model specifications. How are the variables that you will be modeling distributed? Provide enough context and information about your data for your audience to understand whatever model results you will eventually present. 2. What [covariates](#) help you achieve your modeling goals? Are there problematic covariates? either due to *collinearity*, or because they will absorb some of a causal effect you want to measure? 3. What *transformations*, if any, should you apply to each variable? These transformations might reveal linearities in the data, make our results relevant, or help us meet model assumptions. 4. Are your choices supported by exploratory data analysis (*EDA*)? You will likely start with some general EDA to *detect anomalies* (missing values, top-coded variables, etc.). From then on, your EDA should be interspersed with your model building. Use visual tools to *guide* your decisions. You can also leverage statistical *tests* to help assess whether variables, or groups of variables, are improving model fit. At the same time, it is important to remember that you are not trying to create one perfect model. You will create several specifications, giving the reader a sense of how robust (or sensitive) your results are to modeling choices, and to show that you're not just cherry-picking the specification that leads

to the largest effects. At a minimum, you need to estimate at least three model specifications: The first model you include should include *only the key variables* you want to measure. These variables might be transformed, as determined by your EDA, but the model should include the absolute minimum number of covariates (usually zero or one covariate that is so crucial it would be unreasonable to omit it). Additional models should each be defensible, and should continue to tell the story of how product features contribute to product success. This might mean including additional right-hand side features to remove omitted variable bias identified by your casual theory; or, instead, it might mean estimating a model that examines a related concept of success, or a model that investigates a heterogeneous effect. These models, and your modeling process should be defensible, incremental, and clearly explained at all points. Your goal is to choose models that encircle the space of reasonable modeling choices, and to give an overall understanding of how these choices impact results.

Based on our exploratory data analysis we decided on various interaction terms and transformations to use in our linear model. Our focus was to maximize its prediction accuracy ( $R^2$ ) as well as maintaining a model that is explainable.

## 2.3 Model Limitations

### 2.3.0.1 5a. Statistical limitations of your model

Lab 2 Instructions: As a team, evaluate all of the large sample model assumptions. However, you do not necessarily want to discuss every assumption in your report. Instead, highlight any assumption that might pose significant problems for your analysis. For any violations that you identify, describe the statistical consequences. If you are able to identify any strategies to mitigate the consequences, explain these strategies. Note that you may need to change your model specifications in response to violations of the large sample model.

To validate that OLS is the most appropriate method to analyze our data, we will validate all 5 assumptions; I.I.D, no perfect collinearity, linear conditional expectation, homoskedastic errors, normally distributed errors.

- **I.I.D:** According to the Kaggle authors, this data set was collected by randomly scraping the Google Play Store. Since no clusters of applications were specifically targeted, we can reasonably use the entirety of the store as our reference population. We recognize that applications likely have some degree of interdependence, especially within genres. For example, the success of one application probably has a negative impact on other applications of the same type. Due to the large size of this data set, however, we expect any dependencies to be negligible. We also have reason to believe that the data are identically distributed, as they are drawn from the same population of applications. One could argue that since the Google Play Store changes over time, the distribution also shifts in response. Because the authors do not mention the time frame across which the data was collected, we will assume that they originated from a single snapshot of the Play Store and that no shifts in the underlying distribution occurred.
- **No Perfect Collinearity:** We can immediately conclude that `log_installs`, `log_reviews`, `rating`, and `log_size` are not perfectly colinear as otherwise the regression above would have failed. We can also assess near perfect collinearity for these variables by observing the robust standard errors returned by the regression model. In general, highly colinear features will have large standard errors. Since the standard error of the coefficients are small relative to their magnitude, we can reasonably conclude that they are not nearly colinear.
- **Linear Conditional Expectation:** To verify the assumption of linear conditional expectations, we seek to show that there is no relationship between the model residuals and any of the predictor variables. That is, the model does not systematically underpredict or overpredict in certain regions of the input space. Plots 1 through 3 show the relationships between the model residuals and individual predictors. The residuals are generally well-centered around zero, although the model seems to underpredict when `log_reviews` is high and `rating` is low. The fourth plot shows the model residuals as a function of the model predictions. Here, the model seems to underpredict in the left-most and right-most regions, and slightly overpredict in the middle. Overall, there are no strong non-linear relationships between the

model residuals and the input features, and we do not find enough evidence to reject the assumption of linear conditional expectation.

- **Homoskedastic Errors:** When assessing homoskedastic errors, we seek to determine if there is a relationship between the variance of the model residuals and the predictors. If the homoskedastic assumption is satisfied, then we should observe a lack of relationship; conversely, if the data are heteroskedastic then the conditional variance will depend on the predictors. The first plot is an eyeball test of homoskedasticity, showing the model residuals as a function of the model predictions. We notice that the spread of the residuals is mostly consistent throughout the data, although the right-hand side is somewhat narrower. As a more concrete assessment, we also perform a Breush-Pagan test with the null hypothesis that there are no heteroskedastic errors in the model. Since the  $p$ -value falls below our significance threshold of 0.001, we find enough evidence to reject the null hypothesis. In response to this failed assumption, we report robust standard errors (adjusted for heteroskedasticity) instead of non-adjusted errors.
- **Normally Distributed Errors:** When assessing the normality of the error distribution, we seek to determine if the model residuals are approximately Gaussian. If so, then the sample quantiles of the residuals should closely match the theoretical quantiles of a normal distribution in a Q-Q plot. Below, we plot the Q-Q plot associated with our model. In general, the residuals seem to follow a normal distribution, as the middle quantiles match the corresponding theoretical quantiles. However, the tails of the residual distribution are fatter than expected; the first quantiles occur at smaller than expected values, and the last quantiles occur at larger than expected values. Overall, the assumption of normally distributed errors seems imperfect but reasonably justified.

### 2.3.0.2 5b. Structural limitations of your model

Lab 2 Instructions: What are the most important *omitted variables* that you were not able to measure and include in your analysis? For each variable you name, you should *reason about the direction of bias* caused by omitting this variable and whether the omission of this variable calls into question the core results you are reporting. What data could you collect that would resolve any omitted variables bias?

## 3 Results

Lab 2 Instructions: You should display all of your model specifications in a regression table, using a package like `stargazer` to format your output. It should be easy for the reader to find the coefficients that represent key effects near the top of the regression table, and scan horizontally to see how they change from specification to specification. Make sure that you display the most appropriate standard errors in your table. In your text, comment on both *statistical significance* and *practical significance*. You may want to include statistical tests besides the standard t-tests for regression coefficients. Here, it is important that you make clear to your audience the practical significance of any model results. How should the product change as a result of what you have discovered? Are there limits to how much change you are proposing? What are the most important results that you have discovered, and what are the least important?

## 4 Conclusion

Lab 2 Instructions: Make sure that you end your report with a discussion that distills key insights from your estimates and addresses your research question

## 5 Appendix