

Project 1 Report

Anthony Albanese
Binghamton University
Binghamton, New York, USA
aalbanese6@binghamton.edu

Ameen Kunnathu
Binghamton University
Binghamton, New York, USA
akunnat1@binghamton.edu

Kyle Enriquez
Binghamton University
Binghamton, New York, USA
kenriqu1@binghamton.edu

Andy Zheng
Binghamton University
Binghamton, New York, USA
azheng74@binghamton.edu

ACM Reference Format:

Anthony Albanese, Kyle Enriquez, Ameen Kunnathu, and Andy Zheng. 2018. Project 1 Report. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 3 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

For our project, we aimed to analyze data from different sources to determine general trends and correlations in fashion apparel. To do so, we closely analyzed Reddit and eBay as sources. For Reddit, we analyzed the subreddits r/findfashion and r/FashionReps to analyze discussion and trending topics in fashion. As for eBay, we focused on sales on apparel and clothing being auctioned recently to view their prices and categories. To gather the necessary information from these sources, we need to develop continuously functioning web crawlers which authorize access to these sites and then parse and formulate the data into a useful form. From there, the data needs to be inserted into a stable database to allow for its storage and viewing. During this process, various challenges and accommodations from our initial proposal will be taken to arrive at our ultimate goal of data analysis. In the end, through the data we gather, we can effectively visualize our findings and make preliminary remarks as we project the data that will be gathered over time.

2 IMPLEMENTATION

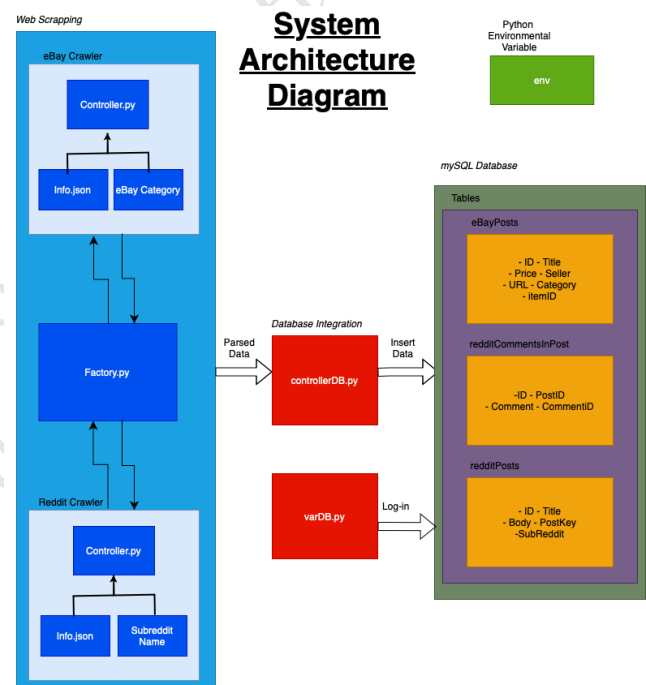
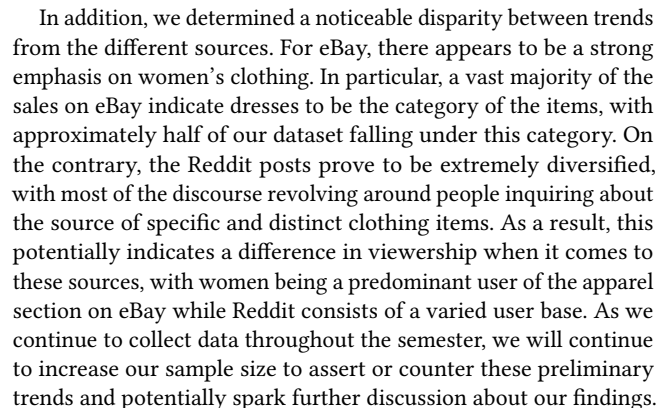


Figure 1: System Architecture Diagram

In order to effectively scrap data from different sites, several different resources were utilized. To begin, we predominantly utilized Python to build our web crawlers for Reddit and eBay. In the programs, we pass user credentials into the web scraper, allowing us to access posts from respective sites and abstract and parse the appropriate data. To do so, we imported the request module to allow us to submit login credentials, which are passed in through a JSON file. In addition, we utilized distinct tools between the Reddit and eBay crawlers to access the respective sites. For Reddit, we implemented a method to gain authentication into the Reddit API, and from there we created another method to gather comments and posts from Reddit based on specified subreddit name and post limit parameters. As for eBay, we imported the ebaySDK library and set it to gather data from a specific category based on a keyword parameter. From there, we utilized the HTML client BeautifulSoup

From the data we gathered, there are various noteworthy trends and facts we can determine. To begin, there is a clear correlation between big sales and fashion trends. For example, observations showed that the Bullseye deals section in eBay acted as a major source of sales in our data set. This is supplemented by several discussions in our targeted subreddits revolving around finding specific clothing items for a desired price. In addition, almost all of the listings on eBay fall at or below 100 USD, with a vast majority ranging between 5 USD and 30 USD. The data for this observation can be visualized from the graph below, which indicates the volume of eBay listings collected over time along with the prices associated with the listings:



In general, the premise of our project remained the same as our proposal. From the proposal, we analyzed data from the same sources using the same means of data extraction. Through this, we were able to work towards our goal from the proposal, which was to gather sufficient data from the sources to determine current trends

and correlations in fashion. To further solidify our project, we made a dedicated effort to incorporate the comments from the feedback to expand upon the original proposal. To begin, we expanded upon the details we included on the use of the Reddit API, ensuring we fully understood the endpoints we are using from the site. In addition, we made sure to have a detailed system architecture diagram that visualizes the various components of the system and shows how they communicate and interact with each other, allowing us to maintain a sense of organization and structure for our project. Through these enhancements, we were able to develop a firmer grasp of the implementation and sources we are handling and properly handle the data we gather.

5 CHALLENGES

Throughout the project, there were a myriad of problems that we faced while collecting data. To begin, during our plan phase, we attempted to experiment with various APIs to accomplish the different objectives we had in mind. However, we soon realized that a vast majority of the APIs we viewed required special user credentials or a paid membership to utilize. As a result, this severely limited the scope of our project. In general, much of the issues stem from this project being a predominantly new experience, as web crawling is seldom used in the computer science curriculum. As a result, the project required a significant amount of experimentation and research to complete. In particular, when we successfully arrived at a project idea with compatible APIs, there were several issues we encountered regarding our implementation. To begin, we had several issues utilizing the virtual machine for our project. In particular, we periodically had issues with the remote virtual machines failing to run and connect. This halted the development of our project in various parts, particularly when it came to making the Faktory file. In addition, we also had to figure out how to implement external libraries into the remote servers to accomplish the functions of our project. Eventually, we were able to resolve this issue through the use of a Python environment variable, which allows us to focus on a specified scope rather than the system-wide version of Python. Along with this, while we were able to affirm a quick understanding of the Reddit API, eBay caused issues in the earlier implementation stage. In particular, we initially struggled to find a reliable and affordable API that can web scrap eBay and its content. Thankfully, we were eventually able to find that enrolling in the eBay Developer Program allows us to have access to data from eBay for the extent of the project's scope, allowing us to import the eBaySDK library module to utilize its functions. Finally, after we developed a means for gathering data, it was difficult to integrate the data into a running database that can consistently support the flow of data it would receive. Eventually, we overcame this issue by utilizing mySQL Workbench, which provided a stable platform and interface to store and view our data. Through these solutions, we were able to resolve the overarching concern for our project, which was being able to collect a sufficient number of data to use for analysis. As such, we were ultimately able to overcome the various challenges we encountered to complete our objective for this project.

6 PROJECTION OF DATA

From the data we collected, we can project the volume of data that will be gathered over time. To begin, for Reddit we are projected to receive 100-200 subreddit posts per day, with about 5-10 comments per post on average. As for eBay, we are projected to have about 5000 API calls per day, which would pull around 88 posts each API call. After eliminating duplicate listings, this has proven to pull a comparable amount of listings to the number of subreddit posts per day. As a result, by the end of the semester, we can predict to have around 7000 to 14000 database items for redditPosts and corresponding comments to go along with them, as well as a comparable number of eBay listings as well. As for the storage, we initiated the AWS database to have 20 GB of total storage. Due to the predominance of text data over video or image media, we predict that this will be sufficient to collect data throughout the project length. This can be shown through our most recent storage update, which indicates that we have about 18.5 GB of storage left after all the data we've collected as seen here:

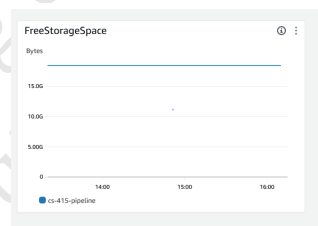


Figure 6: Database Storage Metric

Through these calculated projections, we gauge that we will gather a highly sufficient number of samples to be analyzed, as well as the storage capabilities to support the influx of data.