# amazon

**Team Amazon**

**Michigan State University**

Amazon Review Confidence Tool

Project Plan

Fall 2022

## Amazon Contacts

Mitchell Morrell

Zach Arnold

Jacob Ackley-Smith

## Team Amazon

Collin Cole

Nikita Gupta

Cameron Hurley

Ashu Kher

Dylan Mccarroll

Ethan Strain

# Table of Contents

Team Amazon Project Plan

# Executive Summary

Amazon.com, founded by Jeff Bezos in July 1994, is renowned for its influence in the electronic commerce industry. Based in Seattle, Washington, the cooperation has grown to become part of the big five American Information Technology companies. Amazon was initially an online marketplace for books but eventually expanded into a multitude of product categories. The enterprise acts as a middleman between retailers and customers as well as aiming to provide convenience, affordability, and a wide selection. Amazon currently has over 300 million active customers and over 1.9 million seller partners worldwide. With so many people interacting with the platform, Amazon is trying to ensure that customers are having the most authentic, reliable, and safe experience while sellers are getting honest and accurate feedback.

One way that Amazon is working to reassure their clients about the standards of the products on their platform is by developing a review framework. Within this sector, shoppers can share their personal experience and analysis of the product. Through customer insight, future buyers can make more informed decisions and sellers can monitor if they have met their business goals. Additionally, by having a good review infrastructure in place, there will be an increase in buyer trust and loyalty as well as a boost in sales. Therefore, an abundance of illegitimate or low-quality reviews can hurt business for sellers and can cause brand damage. The presence of illegitimate reviews can mislead consumers and harm business' reputations by disseminating false information and opinions about products.

It is important that the reviews that are posted are honest and legitimate for both the customer and seller. To combat this issue, our software will help Amazon predict review authenticity and conduct review analysis. By utilizing machine learning, the software will be able to generate confidence scores for each review and calculate an adjustable total average rating after filtering out reviews with low confidence scores. Buyers and sellers will also have access to a more detailed analysis dashboard. This solution will not only help Amazon achieve its customer service goals, but also assist in avoiding customer confusion and preserving seller reputation.

# Functional Specifications

With the constant growth and dependency on the digital world, the e-commerce industry is rapidly expanding. Amazon is constantly adapting and innovating changes to its online marketplace to support and meet the expectations of its customers. The corporation insists on the highest standards and earning trust with their millions of customers and sellers. Amazon accomplishes this through a product review infrastructure. By allowing customer assessment and feedback, users can get insights on specific products and sellers can get a better understanding of the quality and reputation of their product. However, with the scale of Amazon, there are unavoidable bad actors who seek to gain from exploiting the review system.

From a consumer's perspective, online reviews offer security as they guide and reinforce a decision regarding the purchase of the product. Whereas for a seller, reviews provide feedback for improvement and help indicate if they are meeting their business goals. If this infrastructure becomes corrupted, it will be difficult for anyone to trust the reviews. Our goal for this project is to identify and eventually assist Amazon in preventing illegitimate product reviews by providing a supplementary tool that can help customers save time and money regarding their purchasing decisions. By doing this, we aim to improve the customer experience by bringing awareness to potential abuses of Amazon's review systems and promote a healthy skepticism when reading user experiences. By developing a visually intuitive summary of review-authenticity, Amazon will be able to increase trust and loyalty with both their customers and sellers.

The software will analyze Amazon's customer feedback under specific products to detect and predict fake reviews. While the user is on the tool's product page dashboard, they will be able to see confidence scores for each review as well as an average rating after filtering reviews with low confidence scores. The page provides visualizations to highlight distributional anomalies regarding review length and contradictory sentiment towards a product. This project aims to help bring more clarity and trust to the customer while maintaining the seller's reputation.

# Design Specifications

## Overview

The Amazon Review Confidence Tool is an application designed to display a specific product's review authenticity. It will be used as a way for customers browsing Amazon.com to identify the most accurate reviews after filtering out the ones that are not legitimate. Customers can install a browser extension that will allow them to have an increased awareness of the confidence for each review. It will present an adjusted review analysis of a specific product, display the star score without low-quality reviews, and include a link to the web application. This web application will display more details about the product's genuine reviews and allow Amazon sellers the opportunity of viewing an iterative graphical representation of the confidence scores, as well as a rating to show how accurate the adjusted score is for each individual review.
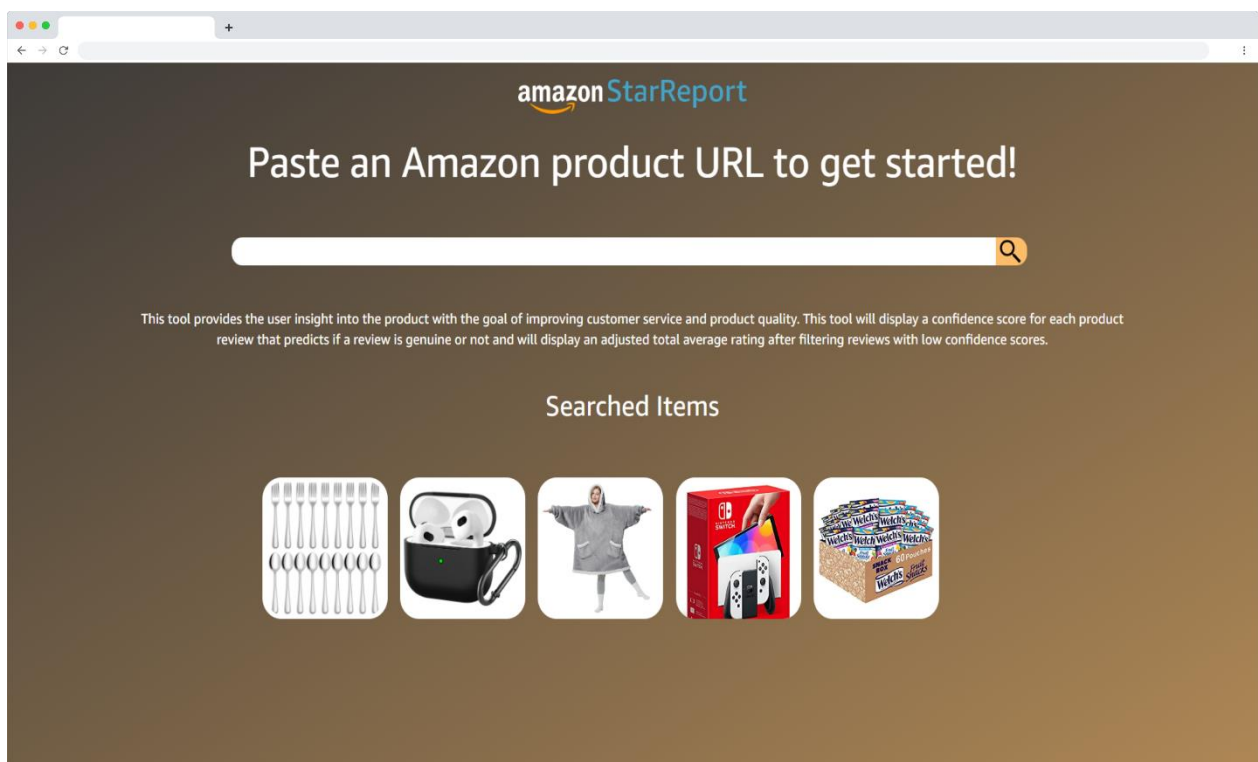
## Home Page



*Figure 1 - Web Application Home Page*

The web application starts on the home page, as shown in Figure 1, where the user can input their Amazon Product's URL at the top of the page. Once entered, this will redirect the user to the specified product's page within Amazon.com. Underneath the search bar

and description, the home page also has a history of previous products searched with the tool. The tool can show up to five previously searched Amazon products. When clicking on a previously access product, it will take the user to the product page similar to what is seen in Figure 4.

## Browser Extension

There will be two different browser extension displays, one for when a user is not viewing a product on Amazon.com and the other for when the user is viewing a product on Amazon.com.



*Figure 2 - Browser Extension Not Installed*

When the user selects the browser extension and they are not on Amazon.com, they will receive the display that is shown in Figure 2. It will display some text indicating where to navigate to in order to get the adjusted review data displaying in the extension window. Below that there will be two different links. The first link will direct the user to our web application's home page if the user clicks on "here." The second link, which is represented by the Amazon logo, will redirect the user to Amazon.com when clicked on.
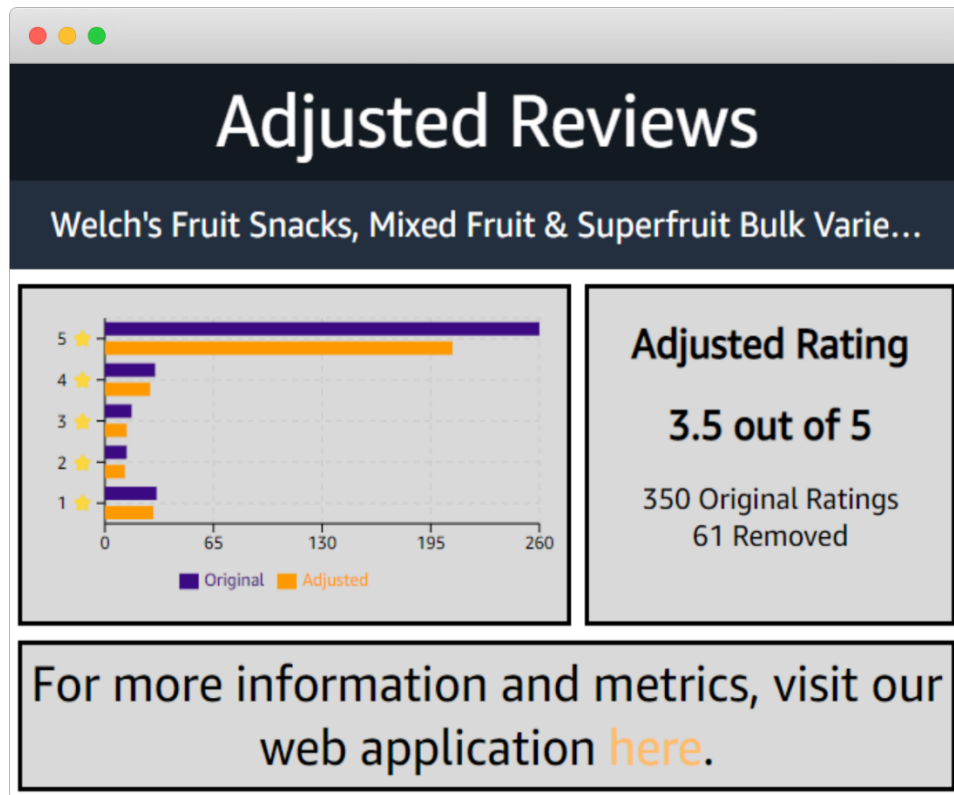
*Figure 3 - Browser Extension Installed*

When the user selects the browser extension and they are an Amazon product page, they will receive the display that is shown in Figure 3. It will indicate which product they are currently viewing and below that there will be a couple of different sections of data regarding the product's reviews. The first section on the left will display the frequency of reviews for each star rating before and after adjusting The second section on the right will display the overall adjusted rating for the product and present the total number of reviews that were removed for not being authentic. Below that there will be a link that will direct the user to our web application's review page if the user clicks on "here." If the user clicks on the Amazon logo at the top of the page, it will redirect them to Amazon.com.
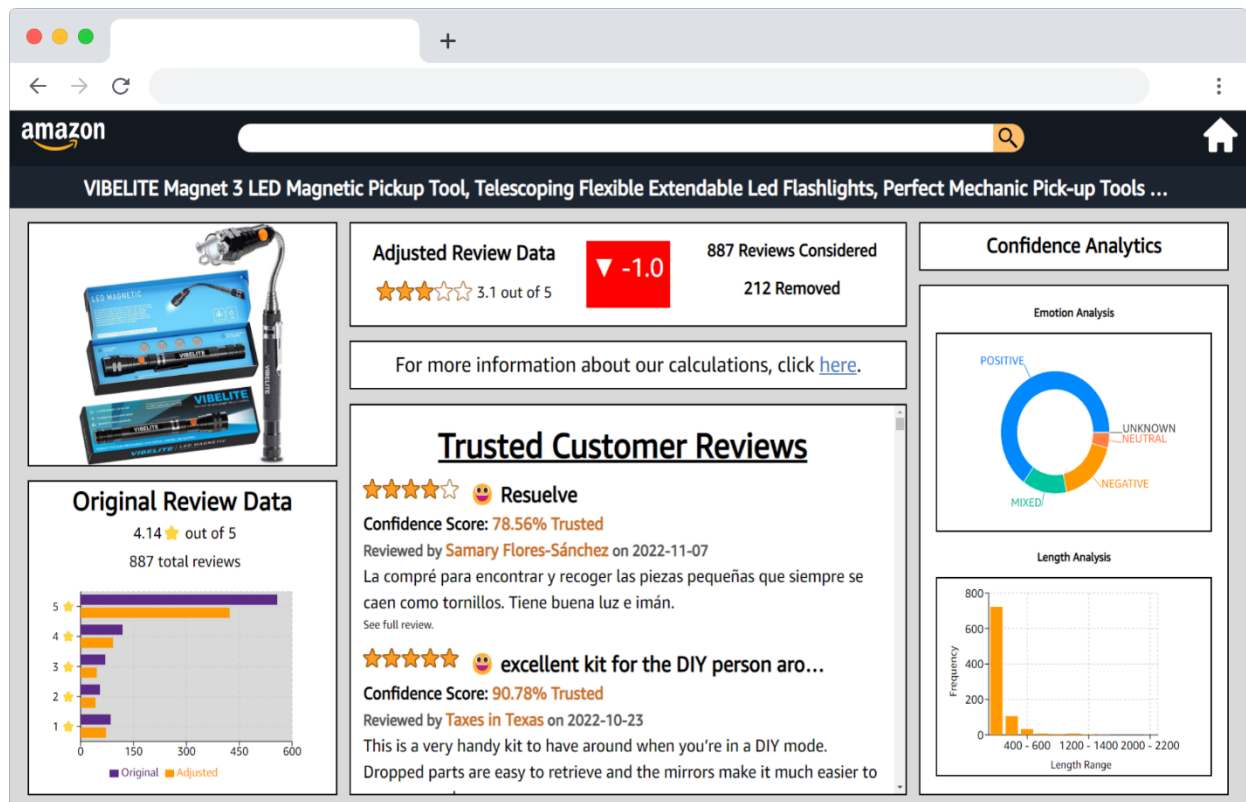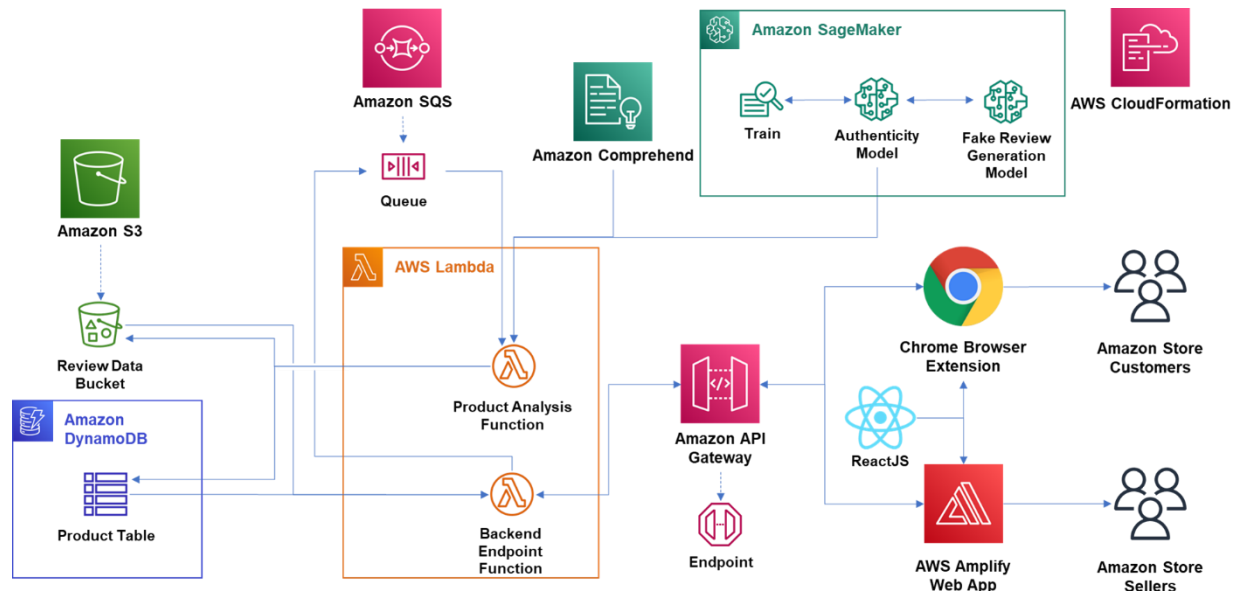
Team Amazon Project Plan

# Review Page



*Figure 4 – Web Application Review Page*

If the user selects the "here" link in the browser extension displayed in Figure 3, they will be redirected to our web application's review page, as shown in Figure 4. It will display a wide variety of different metrics, demonstrating the adjustments made by our tool. At the top of the page, we show the new rating after ignoring the reviews determined to be untrustworthy by our system. On the left-hand side, an image of the selected product will be shown above the original ratings taken directly from Amazon. On the right, two graphs will be present. One pie chart showing the breakdown of the sentiments for each review while the bar graph on the bottom has a histogram of review lengths. Finally, the main portion of the page pulls and sorts the top reviews that our tool has calculated to be trustworthy. This page is overall specific to the Amazon product that a user inputs into the search bar and will show details pertaining to that item.

Team Amazon Project Plan

# Technical Specifications

## System Architecture



## Overview

The Review Confidence Tool will rely on Amazon Web Services (AWS) to provide a fully cloud-hosted infrastructure to the end-user.

To analyze review authenticity, the tool will rely on machine learning models for classification and sentiment analysis of text. New labeled review data can either be entered into a DyanmoDB table, or directly into an S3 bucket. Amazon SageMaker will use the given data to train a machine learning model that can provide real-time inference on reviews to determine if the review is legitimate. The baseline for generation and detection models are GPT-2 by OpenAI and C-Support Vector Classification (SVC) by Scikit-learn respectively, which are both pretrained on other text speech relevant datasets prior to review training. SageMaker will also generate fake data in combination with open-source Amazon review datasets provided by Tensorflow for a supervised model training approach. Detection classifiers for the 10 most popular product categories on Amazon will be hosted on a multi-model endpoint in order to reduce cloud-computing costs and improve classification accuracy within a particular product category. By training separate classifiers for each product category, the appropriate language patterns can be adopted and applied to future inference instances. For example, the high rate of technical language found on electronic products differs greatly and may not apply well to the high rate of emotional language found on toy & games products. If a new product does not match any of the accounted for categories, a default classifier trained on a multitude of different products will be used.

On the backend, a series of REST APIs hosted by API Gateway point to Lambda functions. The product analysis function will analyze a review with the SageMaker model and returns a rating of authenticity for the review, as well as caching the review analysis into a DynamoDB table to avoid redundant calculations. This additionally improves loading times for previously analyzed products other users have run.

This analysis is displayed to the user in two ways. For Amazon.com customers, a browser extension will automatically analyze reviews and display them in a quickly digestible format without having to navigate to another web page. For sellers, a more detailed analysis will be provided through a web app hosted on AWS Amplify, using browser cookies for user state management. Both the web app and the browser extension use React as a web framework.

The final version of the tool will be made into a CloudFormation template to be given to the client so that they can automatically provision and configure the tool on their own AWS account.

## System Components

### Software
### Amazon Web Services

**Lambda**

AWS Lambda will be the backbone of the App Layer of the project, providing serverless event-driven functions that can interact with other AWS services natively.

**Simple-Queue System (SQS)**

Amazon SQS is the backbone of the messaging system within the architecture, helping compile and process API calls to scrape for information. This will allow web scraping reviews to be separated from the API call, allowing the API to be within an acceptable response time.

**DynamoDB**

Amazon DynamoDB provides fast and scalable performance in a fully managed NoSQL database. Several tables will be used to contain new training data and cache product reviews, aggregated product scores, and seller overviews.

**S3 Buckets**

Amazon S3 acts as the large-scale storage for training data. S3 buckets will be used to contain the large amount of review data used to train the ML model and can be updated with new data as it is acquired.

**SageMaker**

Amazon SageMaker is a machine learning platform to quickly build, train, and deploy machine learning models. Two key types of models that will be built are the illegitimate review generation model which will provide labeled training data to the model that will be used to detect illegitimate reviews and calculate confidence scores for each review.

### Comprehend

Amazon Comprehend is a natural language processing (NLP) service that uses machine learning to gather insights from text. The primary NLP tool that will be used from Comprehend is sentiment analysis through its API.

### API Gateway

Amazon API Gateway manages the creation of APIs for Lambda functions to interact with the front-end on both the browser extension and web application. Since data will be passed into these APIs, REST APIs will be the most commonly used.

### CloudFormation

AWS CloudFormation is an Infrastructure as Code (IaC) service that allows the provisioning of AWS resources through JSON or XML templates. This service will primarily be used to generate a final deliverable to the client.

### Amplify

AWS Amplify is a simple web application deployment solution, allowing for easy deployment of the web app using React. The web app can be updated as needed using Amplify's auto-deploy feature.

## Development Environments/Languages

### Python

Python will act as our main scripting and model development language due to its large data science capabilities and wide team familiarity. Experiments and testing can be done on our own integrated development environments as well as Amazon's web environment.

### React.js Framework

React.js is an open-source JavaScript framework developed by Meta released in 2013. It will serve as the front-end for both the web application and the browser extension. React will be used to display adjusted average reviews and confidence scores for each review after filtering out illegitimate reviews.

# Risk Analysis

## Responsive Browser Extension

Difficulty: Low

Description: For customers to use the browser extension, they will expect it to be responsive and not load for long periods of time. Reducing load time will be a priority for the extension.

Mitigation: There are several solutions that could help mitigate loading times. Firstly, if a new pass of the algorithm is occurring, the extension can alert the user and redirect them to the web app version of the service. In addition, prior analysis data will be cached into a database to avoid repetitive and unnecessary runs of the authenticity model.

## Acquiring Labeled Review Data for Machine Learning

Difficulty: Medium/High

Description: The most difficult part of creating a machine learning model is acquiring well-labeled review data to train the model with. Product reviews, which both have natural language features as well as metadata, will require multifaceted metrics including sentiment, review length, reviewer account information, and many more to determine authenticity.

Mitigation: To obtain relevant training datasets, artificial reviews will be generated by a separate pretrained model in several product categories. Additionally, widely available open-source datasets will be collected to assist with correcting weak predictive areas in the detection model. Amazon has been able to provide a small sample to reference as a template for dataset structure thus far.

## Complex CloudFormation Stack

Difficulty: Medium

Description: As much of the project is on AWS, a simple git repository won't capture everything required to spin up the project. AWS offers CloudFormation, their Infrastructure as Code (IaC) service, allowing the provisioning of AWS resources using JSON or XML formatting. However, we have not done research on the limits of the service, including if code such as our proposed browser extension or our machine learning algorithm can be put onto the template, and if the AWS SDK would be needed for this.

<u>Mitigation</u>: Mitigating this risk involves assessing the limits of CloudFormation early in the project. The first step is to study what can be put into a stack, and produce a stack of the alpha to prove the viability of fully packaging the tool.

# Schedule

## Week 1 – (8/29 – 9/4)

- Initial meeting with team members
- Designate roles into key areas
- Research vital technologies (AWS)

## Week 2 - (9/5 – 9/11)

- Initial client meeting (schedule, office hours)
- Initial triage meeting with TM Griffin Klevering
- Status Report Presentation

## Week 3 – (9/12 – 9/18)

- Design screen mockups and system architecture
- Work on Project Plan Document and Presentation
- Prototype the machine learning classifiers

## Week 4 – (9/19 – 9/25)

- Project Plan Presentation
- Complete the layout for the web application home page
- Load training datasets into S3/Dynamo
- Start model preprocessing scripts

## Week 5 - (9/26 – 10/2)

- Complete the layout for the web application review page
- Training models, feature engineering
- DynamoDB -> S3 pipeline
- Lambda function call to SageMaker and return results, while caching to DynamoDB

## Week 6 – (10/3 – 10/9)

- Develop a browser extension

- Generate reasonable confusion matrix from the model
- Create Lambda Functions for Gathering and Analyzing Review Data (integrating confusion matrix)
- Present alpha deliverable to client

## Week 7 – (10/10 – 10/16)

- Alpha Presentation
- Link browser extension to both web pages
- Generate API from Lambda Function
- Limit-test dataset sizes

## Week 8 - (10/17 – 10/23)

- Caching recently viewed products to DynamoDB from Lambda
- Develop user state for web app
- Integrate API Calls for Web App
- Create Product and Seller APIs
- Finalize list of classifiers used for detection (top 10 product categories)

## Week 9 – (10/24 – 10/30)

- Integrate Cookies for recently analyzed products
- Integrate API Calls into front-end
- Individual review confidence score optimizations
- Review generation model optimizations

## Week 10 – (10/31 – 11/6)

- Refine frontend to Amazon-compliant formatting
- Incorporate more real data into the Frontend environment
- Develop SageMaker endpoints for front-end/Lambda
- Work on CloudFormation/CDK stack

## Week 11 - (11/7 – 11/13)

- Testing/Bug Fixing
- Model Optimization
- Test out UX for all web pages with client and other sources for feedback

## Week 12 – (11/14 – 11/20)

- Beta Presentation
- Clean up front-end and implement final client design changes
- Incorporate feedback
- Continue training new product category classifiers

### Week 13 – (11/21 – 11/27)

- Bug Fixing
- Final CloudFormation Stack
- Finalize all front-end pages
- Finalize Browser Extension

### Week 14 - (11/28 – 12/4)

- Final CloudFormation Stack
- Create Video
- Present final deliverables to client

### Week 15 – (12/5 – 12/11)

- Design Day
- Final Deliverables Due