# amazon

**Team Amazon: Seller-Forums Information Filtering Tool (SIFT)**

**Project Plan**

**Spring 2015**

**Amazon Contacts**
John Marx
Darren Krape
Poornachandra Pesala

**MSU Team Members**
Max Goovaerts
Carl Johnson
Luke Pritchett
Benjamin Taylor
Johnny Zheng

# Table of Contents

# 1. Executive Summary

Amazon.com is an electronic commerce company based in Seattle, Washington. It started as an online bookstore, but quickly diversified and is currently the largest internet-based retailer in the United States.

Along with selling their own products, it features the ability for third party sellers to use Amazon as a web front to list and sell their merchandise. One of the features provided to third party sellers is Seller Forums, an online forum wherein users may post questions and replies related to selling on Amazon. Roughly sixty-five thousand questions and two million replies are posted annually. The topics of these posts and replies vary significantly.

There is currently no method to allow Amazon employees to monitor Seller Forums; short of manually scanning through each post. This is not feasible at the scale at which Amazon operates. If there is a widespread and trending issue, Amazon would like to know so that they may address it. Along with identifying topic spikes, a classifying system would be beneficial to help keep Seller Forums organized.

The purpose of our project, "Seller-Forums Information Filtering Tool," (SIFT) is to address these issues and unlock the value behind Seller Forums. This is done by identifying and classifying the topic and sentiment of conversations. SIFT's dashboard allows Amazon employees to drill down into specific date ranges to analyze topic trends and view intuitive graphs. The system also has a notification system to update employees of the current status of Seller Forums.

This dashboard, accessible by Amazon employees, streamlines analysis and searching of the Seller Forums. This allows Amazon to quickly address potential issues and provide a better quality of service to the many third party sellers on Amazon.

# 2. Functional Specifications

The functional specifications provide an overview of our project. They outline the problems at hand and our approach to solving them.

## 2.1 Overview

Since Amazon is a worldwide leader in e-commerce, it consequently provides a great service for third party sellers of all different experience levels. The sellers have a large community on Amazon Seller Forums, processing millions of posts annually. On these forums, the sellers discuss selling their goods through Amazon, whether it be novice questions, expert analysis, or somewhere in the middle.

Currently, Amazon has no way of analyzing these posts. Analyzing all of this data manually would be extremely tedious or perhaps impossible. However, the data is important enough to develop a solution to this problem. Through the forums, sellers are providing Amazon with great ideas for growth by identifying their pain points with Amazon's services, praising the great features, and even pinpointing areas of opportunity. All of this data is currently encrypted by its sheer size, and the decryption technique is to develop software to select only the useful information.

The focal point of the project is to unlock the value of this data. The application gives a visualization of the forums' topics and what sentiment the sellers are using. By identifying the users' sentiments, Amazon will be able to properly align their business plan to cater to the users.

## 2.2 Data Organization and Analysis

The forums data is stored in an Amazon Web Services relational database, and the application runs several algorithms against it. These algorithms will analyze and organize all of the data from the forums allowing Amazon to make sense of it.

Text clustering groups forum posts based on similarity. These posts are brought together on the basis of a subset of all words present in Seller Forums, called features. Words that are eligible to be used as features may be altered through the use of a stop list. The stop list is a collection of words that are distinct from the set of words that may appear in the features list. This allows words like "Amazon", "Selling", "Help", etc. that are present in most posts, and provide no uniqueness as compared to other posts, from causing noise within the clustering algorithm.

Through this process, new categories of discussion that the Seller Forums staff may not have been previously aware of may be brought to light, and addressed as seen fit. The title of each cluster is temporarily formed by the top three features found in that cluster by count, these words typically are a fairly descriptive label for the cluster e.g. "Shipping, Tax, Order." Clustering is more powerful than text classification because it allows new categories to be formed automatically based on popularity, as opposed to user defined categories present in classification. Twitter's trending topics is a great example of this in action.

Because the content of Seller Forums is fairly volatile in terms of currently trending topics, users may adjust several aspects of the clustering algorithm such as number of clusters, date range, etc. and will receive results regarding the efficacy of the clustering. Due to the somewhat random nature of clustering and the capricious nature of new posts within the Seller Forums this allows the Seller Forums staff to dictate how clustering is run, ensuring the results are accurate, useful, and desirable. Once it is apparent that the accuracy of clustering has diminished, the user may experiment with new values until certain scores (detailed in technical specifications) are at a sufficient value. Upon this correction, clustering may be re-run over all posts by these user defined values.

## 2.2.2 Text Classification

Text classification organizes posts into a static set of predefined topics through keywords. The analysis will be sorting posts into the topics found by clustering.

Text Classification expands upon clustering by delegating each post to a previously created cluster. While clustering is useful in terms of determining previously unknown categories of posts, it will never generate the exact same output. With the combination of clustering and classification, the Seller Forums team can carefully select the number of clusters they wish to have, run clustering, ensure the results are to their liking, and from that point forward run classification on any new posts uploaded in the database.

If clustering were run each time new posts were added, the categories would constantly be changing, and require renaming and understanding of the underlying content in each category. With this approach the power of clustering's ability to find new topics may be paired with classification's ability to keep constant categories.

## 2.2.3 Sentiment Analysis

Sentiment Analysis refers to the use of natural language processing and computational linguistics to identify and extract disposition within text. SIFT provides sentiment as an additional metric for forum activity.

## 2.3 Report

To visualize the extracted data from text classification, text clustering and sentiment analysis, the application displays various charts and graphs. The home page displays information about the forums in a broad perspective: such as the volume of posts in topics over time, and the sentiment distribution for each individual topic. There are also individual cluster detail pages allowing the user to specify what time range they wish to view and displaying a sampling of posts within that topic over that period of time. Sentiment data for the selected topic will also be displayed.

## 2.4 Settings

Administrative settings are kept separate from SIFTs clean dashboard. For the users that wish to tinker with the inner workings of SIFTs clustering there are several pages that allow for this.

One of these pages allow changing of cluster names and also displays a list of the top 100 words per cluster with the ability to add each word to the list of "stop words." Another page allows users to run diagnostic clustering on a set date range with other specified settings. Once the diagnostic clustering is completed, the results will be populated in a table and an option to set these settings to the entirety of SIFTs data is available. On separate page, users can download a specified set of data into csv format for more freedom to analyze this set of data.

Our project also allows for notifications. Daily email notifications keep users up to date on the status of the topics, including sample posts, top key words, and the "health" of the clusters through silo score and normalized inertia. A notification page allows for easy adding and removal of users.

# 3. Design Specifications

The design specifications articulate the project's user interface features and elaborates on the design by providing and explaining screenshots.

## 3.1 Overview

SIFT is designed for Amazon employees who will use the data from the seller forums to make business decisions. It is accessible through any major browser. The dashboard will aid in visualizing the data analysis portion. Users can search for specific keywords and see what the sellers are saying about them via graphs and other reports.

## 3.2 User Interface

The index of the SIFT website is a general homepage, where the user can get an overview of how all forum topics are performing. The user can select an individual topic and it will link to that topic's page, where graphs specific to that topic and a table with sample posts (see Section 2.3) are shown. Additionally, there are several settings pages, where users can rename cluster titles, adjust settings for clustering, run diagnostic clustering, download samples of data and adjust the email list for SIFT's notifications each topic. The screenshots below (*Figure 1-7*) display all aspects of SIFT.

## 3.2.1 General Analytics

*Figure 1* displays SIFTs homepage. The left column is populated with available topics. The rest of the page provides information about the volume of posts in different topics and changes in topic volume over time, as well as per-topic sentiment information.
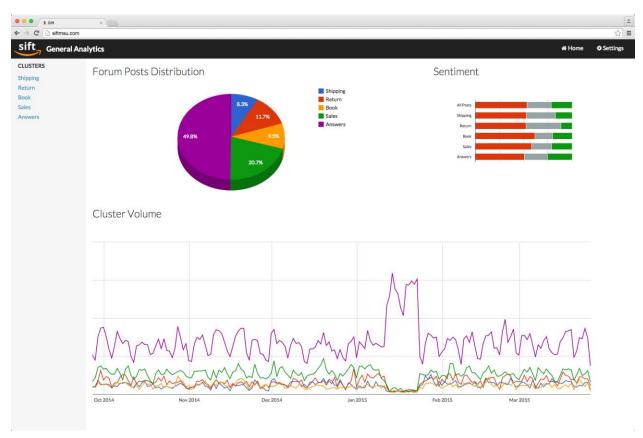


*Figure 1 - General Analytics*

## 3.2.2 Topic Analytics

The Topic Analytics page (*Figure 2 and 3*) focuses on one of the topics selected from the left column. Users can query SIFT for a larger date range using the form and "Fetch Posts" button on the page. There are several graphs depicting how the topic is performing, such as the "Sentiment Analysis" bar at the top right and the "Top Ten Words" pie chart just below it.

The column graph, to the left, displays the sentiment breakdown of posts on a day-to-day basis. The table below, which contains a random sampling of five hundred posts from the selected date range, can be sorted by any column but defaults to show older posts at the top. Finally, the search bar allows users to search for posts containing a specific keyword.
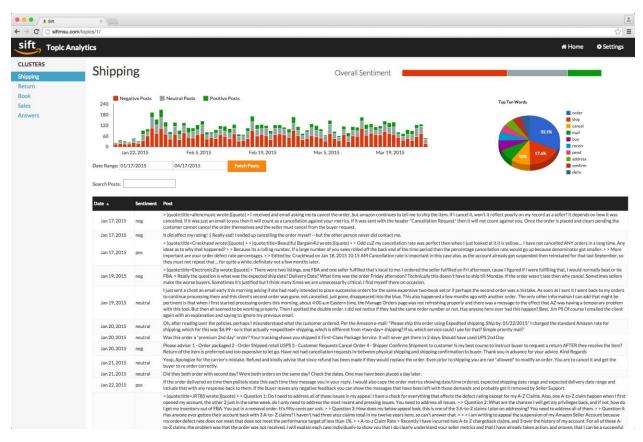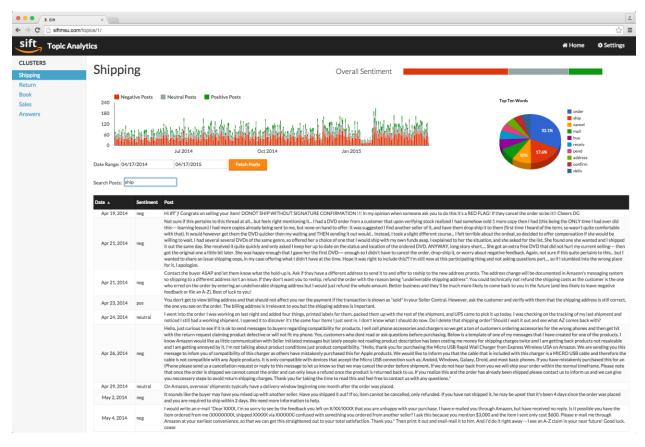


*Figure 2 - Topic Analytics*

*Figure 3 - Topic Analytics: Search + Date*

## 3.2.3 Settings: Clusters

The Clusters page (*Figure 4*) allows users to rename cluster topics. Since they are originally set as the top three features/words of the cluster this can be used in order to better represent them. This page also lists the top 100 features/words of each cluster, and allows users to easily select individual words to add to the stop words list.



*Figure 4 - Settings: Clusters*

## 3.2.4 Settings: Clustering

The Clustering page (*Figure 5*) allows users to run diagnostic clustering over a specified date range, number of clusters, clustering type, and number of features. All currently saved stop words are displayed and users can choose to add or remove words. Finally, a table showing the details of previous diagnostic cluster runs is available for easy comparison.



*Figure 5 - Settings: Clustering*

## 3.2.5 Settings: Export Data

The Export Data page (*Figure 6*) allows users to specify a set of data they would like to download. The page will compile these posts from the database and download them in csv format. This allows users to have more freedom in analyzing posts if they so choose.



*Figure 6 - Settings: Export Data*

## 3.2.6 Settings: Notifications

The Notifications page (*Figure 7*) allows users to add or remove emails from the mailing list of SIFTs notifications.



*Figure 7 - Settings: Notifications*

# 4. Technical Specifications

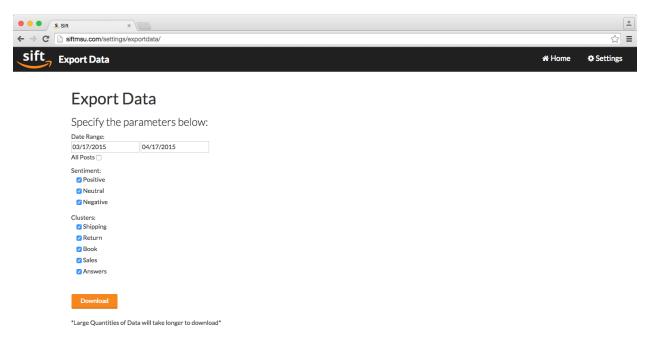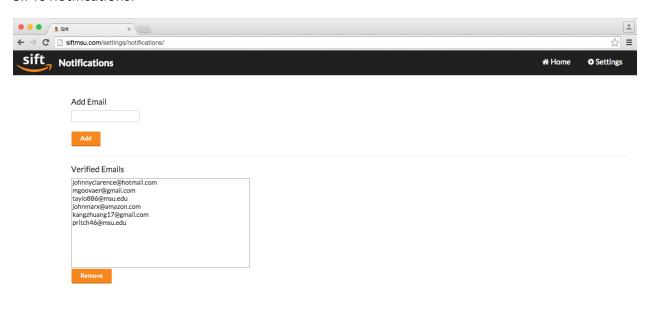The technical specifications articulate the project's system architecture and what technologies are used.

## 4.1 System Architecture



*Figure 8 – Overview of System Architecture*

*Figure 8* details the overall architecture of our system. There are three main components to SIFT: the database backend, the categorization module, and the Django front-end. These all interact through Django's Object Relational Models, and enable content to be easily transferred in-between.

## 4.1.1 Clustering and Classification Schedule

The clustering schedule details the process of new information being added to the database. Every night the cron job will kick off a series of processes. First, the parser script is used to take all data located in the S3 data dump, and convert it to rows in the MySQL Database. Following this, each post is stemmed, analyzed for sentiment and this new information is uploaded to the database. Following the addition of data to the database, these new posts are placed into categories using the classification script.

Upon completion of the classification script, diagnostic clustering is run. This is run over every post, and while it has no effect on the current clusters residing in the database, its output can be used to determine whether or not adjustments are necessary. This decision is based on the normalized inertia and silhouette score. The normalized inertia is the sum of squared distances from each post to its cluster center. While the silhouette score is another distance based metric that takes into account the distance from each post to its cluster and second nearest cluster. These are metrics regarding the accuracy of the clustering, and can determine if a larger or smaller number of clusters may be necessary.

Following these processes, a notification email is sent to all subscribed users detailing the results of diagnostic clustering.

## 4.1.2 MySQL Database

Each post featured in the MySQL database (*Figure 9*) contains a variety of information regarding the posts, and is either associated with a cluster, or will be via the nightly classification. The ClusterRun table displays metrics regarding the results of each clustering. This helps show the health of clustering each night, as well as give users the opportunity to run clustering for themselves with personalized input parameters. Finally the ClusterWord table contains the top features associated with each cluster – to assist in recognizing the topics being discussed in the cluster.



*Figure 9 - Database Schema*

## 4.2 Back-End Technologies

### 4.2.1 AWS S3 Text File Data
All data regarding posts in the third party forums is available via a key-value based text file. These files will be parsed on a scheduled basis (to be determined) and imported to the MySQL database. Since data regarding previous posts does not change, this only needs to occur for the most recent forum posts.

### 4.2.2 AWS Relational Database Service
RDS provides a scalable cloud solution on which our MySQL database resides. RDS can automatically scale the number or capability of servers the database is operating on if the number of accesses to it quickly rises or decreases.

### 4.2.3 Django ORM
Django Object Relational Mapping acts as an interface between the three non-database portions of the application and the database. Django ORM maps database tables to objects and provides a convenient format in which to access the database.

### 4.2.4 AWS EC2 Virtual Server
SIFT is hosted on an Amazon Web Services Elastic Compute Cloud (AWS EC2) virtualized server running Ubuntu 14.04.  Static assets are served using the Nginx webserver, while the Django application is run by the Green Unicorn (Gunicorn) WSGI server.  Background tasks are enabled via Celery and the Django-Celery module, with task queueing handled by the Redis key-value cache store.  These technologies were chosen for their open source licenses and easy integration with Python and Django applications.

### 4.2.5 Scikit-learn Clustering
Clustering is performed using calls to Scikit-Learns's mini-batch k-means clustering. K-means clustering takes an input of n clusters, where n is the number of clusters to be created via the algorithm. Mini-batch k-means is an optimization on this algorithm wherein clustering is performed on separate groups, with the same cluster centroids. The result is a faster algorithm with almost identical results. As described previously, posts are fit into clusters via a list of features where features are words contained within posts in high frequency throughout the seller forums.

Once the feature list is created the module takes each post and converts it into an n-dimensional point in space where n is the number of features within the feature list. By doing this the clustering algorithm is able to perform its primary function, fitting points into clusters by optimizing for inertia. A cluster is defined as a group of points with an n-dimensional point defining its center. Inertia is the sum of squared distances from each point to its cluster. By optimizing for inertia it can be certain that at a large scale posts will be assigned to clusters that are "closest" to them in terms of n-dimensional feature space. Which, in this case means the post has words found often in the cluster.

### 4.2.7 Running Clustering

Clustering can be run on a scheduled or ad-hoc basis. Each night clustering is run over all posts via an on server application cron job, this does not reset all clusters, and is run simply as a diagnostic tool. However, upon user request clustering may be run to create entirely new clusters based on given input parameters. Upon being run, the module will consume all new forum posts present in the MySQL database, clustering them into relevant and popular topics. The now clustered posts are then stored in a table within the MySQL database to be accessed by the front end website.

### 4.2.8 Scikit-learn Classification

Classification uses Nearest-Neighbor classification provided by Scikit-Learn. By supplying the algorithm with a training set of posts already associated with a cluster, the Nearest-Neighbor algorithm determines the best category, of non-classified posts, based on likeness in words in the same way as the clustering algorithm.

### 4.2.9 Sentiment Analysis

The sentiment analysis is performed on a post by post basis. The post body is sent to mashape, a third party API service which uses NLTK's sentiment implementation. An overall sentiment tag is returned along with a positive, negative and neutral score.

### 4.2.10 AWS SES Email Sending Service

AWS Simple Email Service (SES) is an outbound-only email-sending service. SES is used for notifications within SIFT. When diagnostic clustering and routine classification complete it sends an email with details of the run, including the inertia score, and silhouette score.

## 4.3 Front-End Technologies

### 4.3.1 Django Web Application

The application utilizes the Model-View web framework Django. The web application queries the MySQL database for information regarding the top clustered posts, which are displayed on the web portal. The website is able to query the database in order to perform searches via user input. It uses this information to generate graphs and visualizations.

### 4.3.2 Front-End Libraries

Numerous libraries are used to make the user interface. Bootstrap, a front-end framework that contains UI templates, is used throughout the application. SIFT's graphs are displayed with Google Charts. Font-Awesome is used for icons on the navigation bar to give users a visualization of the link (i.e. a house icon for the Homepage).

## 4.4 Testing and Integration

### 4.4.1 Nose Unit Testing

Our software makes use of unit testing and the Nose testing framework for Python to ensure correct and reproducible behavior of each module and component of the software.

Nose allows us to write simple, loosely-coupled tests for our code, and automates the running of these tests. It presents the output in an easily readable format, and is capable of generating more detailed information about code coverage at the user's request.

# 5. Risks

Risks are events that could lead to delayed development or even significant changes in project plan. They can be either known or unknown and vary in severity. This section details a few risks that may affect our progress.

## 5.1 Unfamiliarity with Technologies

None of us have used SciKit or NLTK. We were not sure if they would work as we intend them to: for Clustering, Categorization and Sentiment Analysis. We mitigated this risk by setting up a test environment early on in order to test these unfamiliar technologies.

## 5.2 Machine Learning

Machine Learning is a tricky topic. It can get difficult and we did not exactly know how it tied into our project. We mitigated this risk by using existing libraries, NLTK, rather than developing our own implementations of Machine Learning.

## 5.3 Feature Creep

Feature creep is the ongoing expansion or addition of new features in our project. We mitigated this risk by developing wireframes early and using them to solidify future project features and capabilities. Along with sticking to our detailed timeline (see Section 6).

## 5.4 Scalability

SIFT has to be scalable in order to properly load and interact with the large quantity of data provided by Amazon. We mitigated this risk by making sure that the technologies we used are scalable, and that SIFT does not attempt to load more data than is necessary. For instance, on the Topic Analytics page, SIFT loads a sampling of data (500 posts) over the specified date range rather than all posts in order to keep page load times reasonable.

# 6. Timeline

Below is our plan for development with tasks organized in a week by week manner.

**Week 1: (1/12 - 1/16)**
- Receive project description
- Individual research on project technologies
- Meet team and assign roles
- First client meeting: January 14th

**Week 2: (1/19 - 1/23)**
- Set up code repository
- Explore given data
- Install and configure computers with specific technologies
- Complete status report presentation and rough draft of Design Specification

**Week 3: (1/26 - 1/30)**
- Status Report Presentation: January 26th
- Meet with client, solidify technical specifications/technologies to use. Show first wireframes and get initial thoughts.
- Complete Project Plan Rough Draft, email to Amazon on January 28th
- In person visit with Amazon on January 30th go over Project Plan, fix any issues

**Week 4: (2/2 - 2/6)**
- Project Plan due: February 2nd
- Develop skeleton of front end
- Setup backend: database and simple python script for data analysis

**Week 5: (2/9 - 2/13)**
- Make website live
- Hookup front end and backend
- Rework front end based on input from Amazon
- Work on *classification* – both front end and backend

**Week 6: (2/16 - 2/20)**
- Refine Project
- Work on *clustering* – both front end and backend
- Test all features prior to presentation
- Prepare for Alpha Presentation

**Week 7: (2/23 - 2/27)**
- Alpha Presentation due: February 23rd

- Add *sentiment* functionality – both front end and backend

**Week 8: (3/2 - 3/6)**
- Add *notification* functionality
- Include "admin" functionality in dashboard to adjust notification features
    - Would include logging in – develop simple authentication system

**Week: (3/9 - 3/13)**
- Spring Break
- Continue work from previous week

**Week 9: (3/16-3/20)**
- Identify current standing in project
    - If on schedule begin work on "nice to have"
    - Otherwise complete essential functionality

**Week 10: (3/30-4/3)**
- Begin rigorous second round of testing, bring in unfamiliar users
- Prepare for Beta Presentation

**Week 11: (4/6-4/10)**
- Beta Presentation due: April 6th
- Time cushion – plan for this time to be taken for emergencies and unforeseen problems

**Week 12: (4/13-4/17)**
- Wrap up any final glitches and features
- Do final round of testing and acceptance

**Week 13: (4/20-4/24)**
- Only critical defects to be worked on this week, no new functionality
- Develop script for project video and record
- Prepare for Design Day presentation

**Week 14: (4/27-5/1)**
- All Deliverables due: April 27th and 29th
- Design Day Setup: April 30th
- Design Day: May 1st

# 7. Additional Information

## 7.1 Point of Contact

For further information regarding this document and project please contact Dr. Wayne Dkysen at Michigan State University. All materials in this document are free of proprietary data. The students and instructor gratefully acknowledge the participation of our industrial collaborators.

## 7.2 Acknowledgements

We would like to acknowledge the Amazon team for their contribution and effort spent working with us on our Senior Design Capstone Project at Michigan State University. We would like to specifically thank John Marx, Darren Krape and Poornachandra Pesala.

Additionally, we would like to thank the Department of Computer Science at Michigan State University, Malcom Doering and Dr. Wayne Dyksen for preparing us for applying our skills in a business environment.