

Digging into the FAQs: analyzing user queries to inform service and system development

Marcella E. Barnhart and Kevin A. Thomas

Lippincott Library of the Wharton School, University of Pennsylvania

Ask for a logo!

Ask for a logo!

Introduction

The Lippincott Library's Business FAQ, a knowledge base of 300+ questions that receives over 100,000 views a year, is an integral part of Lippincott's reference service offering. This poster reports on an analysis of user query data from the Business FAQ. User queries vary widely in their structure from long phrases like *'trust building strategies in inter-organizational negotiations'* to company names like *Nike*. Using machine learning techniques to analyze the data, we address a range of questions including:

- What types of queries are more likely to be successful as defined by the user selecting a question from the existing knowledge base?
- What user behaviors are common for successful versus unsuccessful queries?
- Does query complexity impact the success or failure of the match?
- Is there content we could add to the FAQ to improve the users' success?

Data

Business FAQ content is currently hosted on Springshare's LibAnswers platform. The Query Spy function captures the exact query entered, as well as whether the query led to the questioner selecting an answer from within the system.

- LibAnswers maintains 6 months of query data. Data used in this analysis were pulled from 8/31/2018 through 3/28/2019.
- To focus on patron behavior, Lippincott's business librarians identified IP addresses of machines they typically operate, and data were filtered to remove all searches from those addresses.

| Date | Query | Status | System | Business | URL |
|----------------------------|--|------------------------|--------|----------|--|
| 2019-05-17 14:30:30 | invest (add to FAQ?) (mark as reviewed) | Question not submitted | System | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 14:30:59 | How do I find analysis reports (investment bank research)? | Clicked on question | Widget | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 14:00:37 | Phy (add to FAQ?) (mark as reviewed) | Question not submitted | System | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 13:29:57 | Send (add to FAQ?) (mark as reviewed) | Clicked on question | System | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 13:27:07 | supply (add to FAQ?) (mark as reviewed) | Clicked on question | System | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 13:27:06 | supply (add to FAQ?) (mark as reviewed) | Question not submitted | Widget | Business | https://www.library.upenn.edu/businessfaq/faq45190 |
| 2019-05-17 108:130:219:219 | amazon competitors (add to FAQ?) (mark as reviewed) | Clicked on question | System | Business | https://www.library.upenn.edu/businessfaq/faq45190 |

The **Status** field reflects the user's next step as:

- **Clicked on question** when the user navigated to a FAQ on the Search Results page
- **Match via auto-suggest** when the user navigated to an FAQ in the suggested results before submitting the query
- **Not Submitted** when the patron did not select a result,
- **Question Submitted** when the user chose the *Submit Your Question* option.

This analysis considers a query successful if the patron navigated to an FAQ (either of the first two **Status** values) and unsuccessful if the patron submitted a query but did not identify a result (either of the last two **Status** values).

Methods

N-gram Conversion

Queries were converted, via a series of steps, into terms consisting of 1 to 3 consecutive words as follows:

1. Remove numbers
2. Remove punctuation (including apostrophes and dashes)
3. Convert all letters to lower case
4. Extract one-word, two-word, and three-word terms (spaces replaced by underscore)
5. Remove stop words identified by the Snowball stemmer project as implemented in the R package SnowballC (Bouchet-Valat 2019)
6. Perform stemming according to algorithms from the Snowball Stemming Library as implemented in the R package SnowballC (Bouchet-Valat 2019)
7. Remove one-letter terms

(Note that Query Spy exports the ampersand character as *amp;*, which the above process converts to *amp*. As a result, terms such as *M&A* and *S&P* appear as *mampa* and *sampp* in this analysis.)

Methods

k-Means Clustering

k-means analysis groups similar items, in this case terms from related queries. Five of the top query terms by group are shown below.

| 1 | 2 | 3 | 4 | 5 |
|--------------------------|----------------|--------------------------|-----------------|-------------------------|
| am_l | a_compan | financi | total | analyst_reports_invest |
| am_l_hav | adpt | financial | trading | trading |
| difficult | access_co | time | real | bank_research |
| difficulty_us | access_co_inve | real | bank_research | real_estate_j |
| difficulty_using_thomson | access_to_cas | do_l | find | find |
| having_difficult | account | find | analyst | report |
| having_difficult_in | acquist | activ | how | do_j |
| l_hav | advertis | one | l | find |
| l_having_difficult | atlas | invest | invest | invest |
| thomson | alum | investment_bank | investment_bank | research |
| thomson_on | america | investment_bank_research | reports_invest | reports_investment_bank |
| use | americans | reports_invest | research | research |
| using_thomson | amp | research | research | research |
| using_thomson_on | amp | research | research | research |
| why_am | analyst | research | research | research |
| why_am_j | and_acquist | research | research | research |
| why_am_l | and_report | research | research | research |
| | annual | research | research | research |
| | annual_report | research | research | research |
| | app | research | research | research |

Random Forest Analysis

In a random forest model, an algorithm randomly generates a large number of decision trees, each of which predicts an outcome. In this case, the trees used the presence or absence of groups of query terms to predict the success or failure of a query. For this analysis, 6376 total query records were divided into 5100 records used for training the model and 1276 records for testing the trained model's performance. Multiple variations on random forest and extra-random forest models were trained and tested for comparison. The ultimate model, chosen for its stability and test accuracy across both failed and successful searches, is a random forest consisting of 350 trees, validated through repeated cross-validation of four folds and four repetitions. Between 7 and 12 randomly-selected terms were available for splitting at each tree node in the final model, with the optimal model allowing 12 terms.



Term importance is measured by the difference in gini coefficient, which measures how dissimilar data are, with versus without the given term. Terms with higher gini importance scores have more impact on results.

Search Behavior & Follow-up Searches

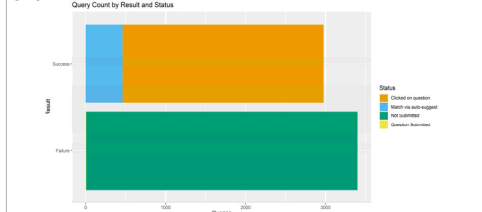
Queries do not always stand alone, but sometimes form part of a larger, iterative process. This research defines a search as a series of one or more queries from the same IP address that occurred within 30 minutes of another query. Although SpringShare does not provide information that follows a patron's entire search process, we can use a query's referring URL to identify when a search originated from a given FAQ. Comparing these instances to the quantity of visits to the same FAQ during the same timeframe (*Export Views by Month* report), we can arrive at a follow-up search rate.

Regression

Regression analysis can help indicate which factors contribute toward or detract from query success. This analysis uses a logistic regression model to estimate the relative odds of query success vs. failure under a given set of conditions. For search queries, which tend to be brief and need not be formatted as a sentence, typical text complexity measures may not be highly valid. See Liu, Croft, et al for research assessing the validity of several approaches, which found a substantially valid hybrid approach.

Results

Overall Query Success

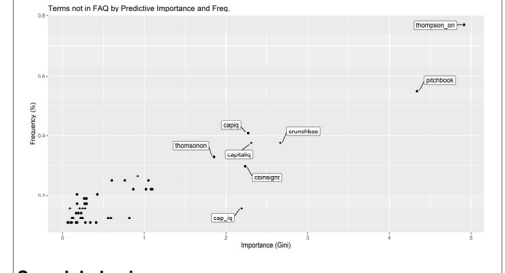


Results & Discussion

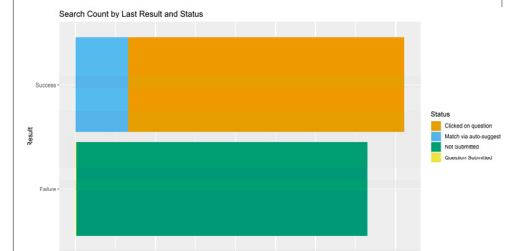
Predictive Importance & Term Frequency

| Term | Importance (Gini) |
|----------------|-------------------|
| analyst | 25.98 |
| report | 23.27 |
| investor | 20.28 |
| invest | 20.09 |
| economist | 20.09 |
| analyst_report | 20.09 |
| propm | 20.09 |
| thomson | 20.09 |
| one | 19.22 |
| l | 18.47 |
| l | 18.47 |
| gamer | 18.47 |
| is | 18.47 |
| felicia | 18.47 |
| proquest | 18.47 |
| mg | 18.47 |
| capital | 18.47 |
| cap_j | 18.47 |
| statist | 18.47 |
| find | 18.47 |
| where | 18.47 |
| research | 18.47 |
| where | 18.47 |
| research | 18.47 |
| booming | 18.47 |
| can_l | 18.47 |
| find | 18.47 |
| report | 18.47 |
| where | 18.47 |
| where | 18.47 |
| can | 18.47 |

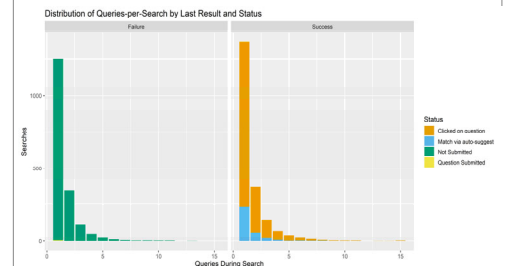
Terms that are important to and frequent in queries, yet not present in any FAQs present a clear first target for expanding FAQ content.



Search behavior



53% of searches ultimately leading the patron to view an FAQ, an improvement over queries considered individually. Searches that ended in success tended to include more queries. While the difference accounts, on average, for only a fraction of an additional query per search, there is a strong indication that the difference is more than coincidence. Could patrons be struggling for the specific words or spellings needed to return a relevant result? We might consider follow-up searches an indicator of the referring FAQ's failure to address the patron's query and explore potential improvements for FAQs with both high traffic and high follow-up search rates.



Regression

The generalized linear model used attempts to predict a query's success based on its word count. Applied to a test sample, this model demonstrates 52.5% accuracy; it is better than a coin toss, but only slightly.

