# Predicting Collapsed Answers on Quora

Shannon Mitchell, Mason Montgomery, Akshat Verma, Mike Voltmer
SYS 6018 Final Project
10 Dec 2015

# Problem Description

**Quora**

# What are the most important Marketing Analytics tools?

For a book I'm writing on marketing analytics (scheduled to publish in May 2015), I'd like to know what the most important marketing analytics tools are. Which ones are absolutely essential to the marketing analytics process?

## 12 Answers

**Ania Skonieczny**
600 Views

My experience says that you should include the TEA Software in your list of Marketing Analytics tools. It certainly is an essential tool in the web analytics process.

TEA stands for Threat and Engagement Analytics Software that fights clickfraud, focuses on your best advertisers and improves the sales conversions of your website.

This new ecommerce solution ⧉ can analyze the behavior of the visitors arriving on your online store by keeping track on their mouse movement. Then, it assigns engagement scores to each visitor based on their performance with each keyword. This bifurcates each visitor into a valuable customer or a window shopper.

Written Sep 30 · View Upvotes

Upvote | 21    Downvote    Comment    Share    ...

6 Answers Collapsed (Why?)

**Answers are manually reviewed and collapsed if insufficient**

**Q: What are the most important Marketing Analytics tools?**    **Collapsed Answer: Marketing is a group of people**

# Objectives

# Objectives and Metrics

**Objective:**   Predict whether an answer will be collapsed

**Metrics:**   **Precision** = $\dfrac{\text{True Positives}}{\text{True Positives + False Positives}}$   (minimize false positives)

**Recall**   = $\dfrac{\text{True Positives}}{\text{True Positives + False Negatives}}$   (minimize false negatives)

**Q**: When does something count as a musical instrument?

**Collapsed Answer:** ... 1) You make music with it, and 2) Someone else asks you to show them how to do it ...

# Related Work

# Prior Work and Related Research



Flags <u>offensive</u> content for human review



Flags <u>low quality</u> content for human review

Q: Who was the worst woman in history

Collapsed Answer: Yoko Ono. She broke up the Beatles.

# Technical Approach

# Data

- Top 50 topics on Quora in 2015
- Dynamic content generation and user interaction ➡ Selenium webdriver
- Scraped 4635 answers from 183 questions
- Features collected:
    Question Text, Answer Text,
    Upvotes, Views,
    Profile of Answerer,
    Accepted/Collapsed

**Q: How can I run faster than Usain Bolt?**

**Collapsed Answer: Eat, drink sleep ..Keep doing this forever . But you cant run faster than Usain bolt**

# Feature Engineering

**Answer Text**
- Number of characters, words, sentences, spaces
- Percentage of uppercase letters
- Percentage of text abbreviations & stopwords
- Readability indices
- Entropy

**Question Text & Question/Answer Comparisons**
- Type of question
- Question/Answer length ratio
- Question/Answer semantic similarity index

Q: What is science?          Collapsed Answer: Solid is water

# Feature Selection

Feature Importance given by Random Forest



**rf.fit**

| Feature | MeanDecreaseGini |
|---|---|
| n.char | |
| entropy | |
| n.words | |
| Topic | |
| Followers.of.Answerer | |
| Answers.by.Answerer | |
| Upvotes | |
| n.spaces | |
| len.ratio | |
| semantic_similarity_index | |
| coleman.liau | |
| percen.stopwords | |
| fog | |
| percen.lower | |
| percen.upper | |
| percen.slang | |
| Questions.by.Answerer | |
| n.sentences | |
| Posts.by.Answerer | |
| what | |
| how | |
| why | |
| which | |
| who | |
| where | |
| when | |

Cut-off

MeanDecreaseGini

**Q: What is Physics?**                    **Collapsed Answer: xkcd nails it.**

# Evaluation

# Models

| | Precision | Recall |
|---|---|---|
| **Logistic Regression** | 0.75 | 0.57 |
| **Radial SVM** | 0.76 | 0.60 |
| **Random Forest** | 0.77 | 0.79 |

| | Ground Truth | |
|---|---|---|
| **Logistic** | Accepted | Collapsed |
| Accepted | 559 | 108 |
| Collapsed | 47 | 140 |

| **SVM** | Accepted | Collapsed |
|---|---|---|
| Accepted | 558 | 98 |
| Collapsed | 48 | 150 |

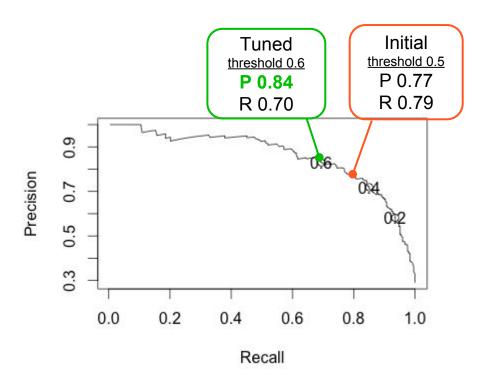| **Random Forest** | Accepted | Collapsed |
|---|---|---|
| Accepted | 549 | 52 |
| Collapsed | 57 | 196 |

**Q: What is mathematics?**

**Collapsed Answer: it is a subject where X is always missing and your duty is to find out Y…
basically it is a subject which only teachers understand. :)"**

# Performance



- Tuned threshold to maximize precision with acceptable recall

- **Beats precision target of 0.66**

**Q: What are the most important lessons to learn about personal finance?**

**Collapsed Answer: Earn more than you spend**

# Limitations

- Does the model hold for less popular topics with fewer answers?

- Readability indices for short answers < 100 words aren't reliable

  ex: "Alot" scores a 10th grade reading level

**Q: How many genres of music are there?**          Collapsed Answer: Alot

# Conclusions

# Recommendations

- Implement model to flag answers for manual review

- If manual review confirms model success, implement tool to:
    - automatically moderate answers
    - warn users if answers will be collapsed before they are posted

**Q: What are some of the best songs of all time?**   Collapsed Answer: Definitely Bohemian Rhapsody from Queen

# Future Work

- Increase sample size and question/answer diversity

- Add features: part of speech counts, syntactic analysis

- Predict number of upvotes

- Apply additional model types: gradient boosting, neural network classifiers

**Q: How did Rafael Nadal beat Novak Djokovic in the 2013 US Open Final?**

**Collapsed Answer: With his Tennis Racquet.**

# Questions?

Q: What are the best programming languages to learn today?
Collapsed Answer: Which is the best flower?

Q: Is there a liberal media bias?
Collapsed Answer: You betcha

Q: How do I understand physics easily?
Collapsed Answer: … 1. Drop an apple on your head 2. Try running on a slippery wet floor 3. Get down from a running bus …

Q: What is the difference between baking powder and baking soda?"
Collapsed Answer: The spelling.--- Edit ---Wow. Already flagged by someone for no reason. My answer is 100% correct.