

A German Dialogue Dataset for Tech Support

COMP 551

Applied Machine Learning

Thomas Page — thomas.page@mail.mcgill.ca — 260771672
David Tamrazov — david.tamrazov@mail.mcgill.ca — 26XXXXXXX
Alexander Wong — alexander.wong4@mail.mcgill.ca — 260602944

September 27, 2017

URL: www.oururl.com

1 Overview

Our dataset is gathered from Apple Support Community discussion boards, and would be useful for training machine learning systems specifically for the task of tech support in German. The applicability of this dataset is definitely most geared towards supporting Apple-related questions, but probably can be generalized to most technical customer service platforms.

2 Dataset Description

The dataset is made up of posts and responses from the German Apple Support Communities. The corpus includes a variety of formats such as block quotes, links, lists, and conversational dialogue between the questioner and the respondents. The data was collected over a large range of discussion posts, approximately 13000, to ensure that many of Apple's products and services were covered. Most of the individual dialogue pieces consist of a question with multiple responses from various users of the Apple Support Communities.

To add to the usefulness of the dataset, we extracted **score**, a feature beyond the project specifications. Score (or the number of likes for each post), reflects the helpfulness of a response, or in the case of the he original poster (OP), whether someone else had the same question). For the sake of the assignment we are submitting two datasets:

- *NotoriousApplePickers_GER.xml*
- *NotoriousApplePickers_GER_Scores.xml*

3 Discussion

Defining a Conversation One of the key design decisions made was how were we going to define a conversation and therefore what discussion threads would make the cut for the dataset? Our initial intuition was to only include discussions with at least one instance of a response from the OP, in addition to their original post. We then revised this to at least one instance of non-consecutive posts by any user, for the inclusion of the OP does not preclude a thread from being a conversation among others.

However, we eventually settled on the current definition of a conversation, which simply includes any thread with at least one response. While this may seem to disregard the objective of collecting data on conversations, the primary objective for this data set is to collect examples of good technical support, and one may argue that perhaps the best tech support takes the form of a singular, concise response (complexity and obscurity of the original question being a difference maker in this regard).

Data Characteristics While the responses are technical in nature, the benefit to scraping conversations from a website run on individual user input is the colloquialism in the responses. In contrast, a dataset developed from a formal setting like a FAQ would have a different type of language that would cause the automated tech support to sound robotic.

Design Choices To fit our the scraped data to the requisite .xml format, we made minor modifications to the extracted text. We changed all instances of '<' and '>' to '_open_angle_' and '_close_angle_', because they clearly would interfere with parsing xml. While parsing our corpus for basic analysis we found that the ampersand and at symbols would also provide unwanted interference, so changed them to '_and_' and '_at_' respectively.

4 Statement of Contributions

Thomas got the ball rolling by introducing us to the BeautifulSoup python library that played a large roll in simplifying the task from the coding side. He also wrote the initial scripts to test feasibility, and suggested sourcing our data from Apple tech support pages. Upon agreeing on our objective, we then worked cooperatively to create a program that extracted the data we wanted from the pages, without double counting or skipping posts. Dave and Alex worked on optimizing Thomas' initial program, and implemented the logic for testing each discussion thread to meet our definition of a conversation.

Thomas ran the program from home (approximately 7 hrs), and all three of us did some data analytics to include in the report. Alex wrote the first draft of the report, and from there everyone contributed to finalizing it for submission.

We hereby state that all the work presented in this report is that of the authors.