

Honours Computer Science Thesis

rally, a one stop-shop for all reddit data

by

Kevin J. Eger

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF

THE REQUIREMENTS FOR THE DEGREE OF

B.SC. COMPUTER SCIENCE HONOURS

in

Unit 5

(Computer Science)

THE UNIVERSITY OF BRITISH COLUMBIA

(Okanagan)

April 2016

© Kevin J. Eger, 2016

Abstract

Reddit is *the front page of the internet*, a slogan the company has coined and rightfully lived up to. It is a website which brings together members of all communities in a similar style to a typical forum but with much more structure and a lot more traffic. The open nature of Reddit generates a large amount of traffic, averaging over 200 million unique visitors a month. With such traffic screams the demand for data analysis through a human-interpretable medium which this thesis covers. Data analysis on reddit has been done before however this thesis focuses on bringing the data gathered in to a easily consumable format. Techniques for consuming less apparent analysis and alternative browsing techniques are covered. We will explore the implementation and results of querying the reddit API, generating aggregate statistics, querying large data dumps of historic reddit data with *Google BigQuery*, analyzing and labelling the content of Reddit using *Google Cloud Vision*'s image recognition and the use of unsupervised machine learning to draw powerful conclusions.

Table of Contents

Abstract	ii
Table of Contents	iii
List of Tables	v
List of Figures	vi
Acknowledgements	vii
Chapter 1: Introduction	1
1.1 Reddit	1
1.2 Motivation	1
Chapter 2: Background	3
2.1 Terms and Definitions	3
2.2 Reddit	3
2.2.1 History	3
2.2.2 Community	3
Chapter 3: Technical Stack	5
3.1 Laravel	5
3.1.1 MVC	5
3.2 Storage	8
3.2.1 MySQL	8
3.2.2 BigQuery	9
3.3 phpRaw	12
3.4 Clustering Script	13
3.4.1 Python	13
3.4.2 Scipy	13

TABLE OF CONTENTS

Chapter 4: Algorithms and Methods	14
4.1 Hierarchical Clustering	14
4.1.1 The Clustering Process	16
Bibliography	18

List of Tables

List of Figures

Figure 3.1	Example of Model	6
Figure 3.2	Example of View	6
Figure 3.3	Example of Controller	7
Figure 3.4	Example of Repository	7
Figure 3.5	MySQL Database schema	8
Figure 3.6	Query finding the best hours to post on Reddit	10
Figure 3.7	Query finding the best hours to post on Reddit	11
Figure 3.8	Registering the Google service provider	11
Figure 3.9	Requiring phpRaw as a dependency in composer. . . .	13
Figure 4.1	Dendogram of /r/movies	15
Figure 4.2	BigQuery for retrieving clustering data	17
Figure 4.3	Psuedocode for processing BigQuery data	17

Acknowledgements

Work on this thesis was widely facilitated with help from Dr. Ramon Lawrence through weekly meetings where ideas and progress were discussed extensively. It is also important to acknowledge Dr. Jeff Andrews for his support in advising on machine learning techniques which were implemented as described later.

Chapter 1

Introduction

High level overview and motivation for developing this thesis.

1.1 Reddit

Reddit is a a news and entertainment website whose content is sustained by members of the community. Users submit text posts or direct links similar to a typical forum setting. Registered users can vote on submissions bringing order to the posts which yields an ordered online bulletin board. Furthermore, what makes Reddit unique is that content is subsectioned into different areas of interest called “subreddits”. Some of the top subreddits include *movies*, *funny*, *AskReddit*, *food* and *news*. As of March 3rd, 2016 Reddit had 231,625,384 unique users a month viewing a total of 7,517,661,034 pages. The company was founded 10 years ago and has quickly become the most central place on the internet to partake in conversation or consume a wide array of content.

1.2 Motivation

For years data analytics has been used in many industries to give companies and organizations better business decisions and verification of their models and structures. Whether they are mining huge data sets, looking at specific use cases or aiming to prove or disprove a theory, companies and organizations alike aim to do one thing: identify and discover patterns, relationships and inferences that are not immediately apparent.

An early motivator for this thesis was some existing technology for Twit-

1.2. Motivation

ter insights. The community-content driven nature of Twitter parallels that of Reddit. There has already been a lot of academic research and production level software released for Twitter data management, pattern identification and tracking. The existing infrastructure in the Twitter space can be largely replicated and modified to suit Reddit, an effort which this thesis focuses on starting.

Chapter 2

Background

To best understand this thesis and the work done, it is necessary to first be introduced to the relevant technologies and key terms which will be heavily referenced and built upon.

2.1 Terms and Definitions

TODO

2.2 Reddit

2.2.1 History

The company was founded by two new graduates of the *University of Virginia*, Steve Huffman and Alexis Ohanian, in June 2005 [Gua05]. After a couple years of growth, Reddit's traffic exploded and the service went viral. The creators were quick to release Reddit Gold, which offered new features and usability improvements providing the company with a primary source of income.

2.2.2 Community

Reddit thrives on its open nature and diverse content fully generated by the community [Atl14]. The demographics Reddit serves allows for a wide range of subject areas thus having the ability for smaller communities to

2.2. *Reddit*

digest their niche content. Subreddits provide a very unique opportunity by raising attention and fostering discussion that may not be seen as mainstream and covered by other news or entertainment mediums.

Reddit as a company and as a community has been known for several philanthropic projects both short and long term. A few of notable efforts are as follows:

- Users donated \$185,356 to Direct Relief for Haiti after the earthquake that struck the country in January 2010
- Reddit donates 10% of it's yearly annual ad revenue to non-profits voted upon by its users [Red14]
- Members from Reddit donated over \$600,000 to DonorsChoose in support of Stephen Colbert's March to Keep Fear Alive [Don10]

Chapter 3

Technical Stack

Rally is a project that explores many different types of data access, processing techniques and display forms. Due to the nature of web applications, it is no surprise that *Rally* is implemented with modular programming in mind. Several key components outlined below are what will allow this project to be easily continued and built on. The technical stack is broken in to components as follows.

3.1 Laravel

Laravel is a *PHP* web application framework with expressive, elegant syntax [Lar14]. Laravel is designed primarily with the motive of removing the repetitive and often painful part of building trivial common tasks to a majority of web projects (ie: authentication, routing, sessions, etc.). Laravel aims to make the development process a pleasing one for the developer without sacrificing application functionality [Lar14]. The accessible and powerful framework was chosen for it's existing familiarity and power to implement a project spanning many domains.

3.1.1 MVC

Laravel follows the traditional Model-View-Controller design pattern. Models interact with the database through the *Eloquent* ORM providing an object oriented handle on information. Controllers handle the requests and retrieving data by leveraging the models. Views render the web pages and are returned to the user.

This intrinsic design pattern was followed tightly alongside the addition of

a repository layer. As discussed later, *Rally* interacts with several external resources such as the Reddit API and the *Google Cloud Platform*. These external resources house gigabytes of data thus storing them locally and accessing them through a model is counterproductive. To retain the structure of the MVC framework, a repository layer is built on top of the models. This allows for the convenience of a seemingly object oriented interaction with data outside of the application. Not only does it allow for convenient method calls but also abstracts logic away from the controllers, leaving them as slim as possible. This is a vital design philosophy to web development as it modularizes code to ensure a more rigid flow and testable code-base. Basic examples from *Rally* utilizing each level of the MVC framework as well as the repository layer are as follows:

```
$cluster_image = Cluster::where("name", $subreddit)->
    ↪ first();
```

Figure 3.1: A basic example of retrieving the first *Cluster* model where the name field matches.

```
<select name="labels" . . . multiple="">
  @foreach($labels as $label)
    <option value="{{_ $label_}}">{{ $label }}</option>
  @endforeach
</select>
```

Figure 3.2: A basic example demonstrating how objects passed to the view are utilized and iterated over to display the options for the index page of *Content Search*. *Laravel* leverages an HTML templating engine called *Blade* which allows for convenient variable dumping and interaction.

As mentioned above, the repository layer is utilized primarily to wrap auxiliary data sources. This gives them a similar feel and interaction as a traditional model. Seen in figure 3.4, a *RedditRepository* instance is injected in to the *RedditController* class which is then used in its internal functions to gather data using the *phpRaw* Reddit API wrapper in a chainable method technique identical to a traditional model.

```
public function show(Request $request)
{
    $subreddit = $request->get("subreddit");
    $about = $this->phpraw->aboutSubreddit($subreddit);

    return response()->view("subreddit.show", [
        "subreddit" => $subreddit,
        "about"      => $about->data,
        "tagline"    => "A_look_at_/r/" . $subreddit
    ]);
}
```

Figure 3.3: A basic example of the `show()` functions in the `subreddit` controller. This method retrieves the necessary data, then sends the data to a blade view (`subreddit/show.blade.php`) and returns a rendered instance of that view.

```
protected $redditor;

public function __construct(Repository $repository,
    ↪ $redditor)
{
    $this->redditor = $redditor;
}
...
public function show(Request $request)
{
    $user = $request->redditor;
    $subreddits = $this->redditor->getUserSubmitted($user
    ↪ )->getSubredditsList();
    ...
}
```

Figure 3.4: Code snippets from the `Redditor` Controller which leverages the power of a repository layer to make chain-able function calls to an auxiliary data source.

3.2 Storage

Databases used to house the necessary persistent information for the application.

3.2.1 MySQL

MySQL is an open-source relational database management system (DBMS). In Laravel, it is the default DBMS largely because of its *plug and play* nature. The MySQL database is what houses the caching layer as described in detail in the [insert section] section. A visual representation of the schema is depicted as follows:

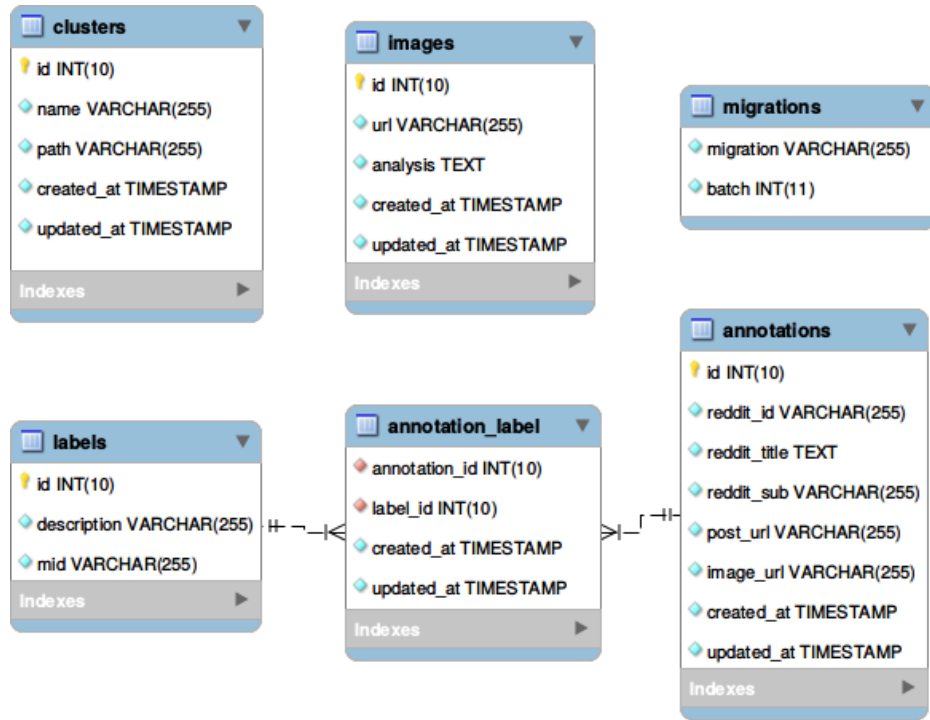


Figure 3.5: An ER diagram representing the MySQL database schema.

3.2.2 BigQuery

Querying massive datasets can not only be time consuming but expensive without the right hardware, infrastructure and software. *Google* alleviates this problem with *BigQuery*, an incredibly fast cloud-based storage platform. It is infrastructure as a service (IaaS) that handles all the hard work of both creating and accessing large data sets. Using the processing power of *Google*, a user can get up and running with *BigQuery* in a matter of minutes. The service can be used via their web UI, command-line tool or the REST API using one of the many client libraries.

Five months ago, user `/u/Stuck_In_the_Matrix` of reddit collected all Reddit submission data from 2006 to 2015. He had effectively bundled 200 million submission objects, each with score data, author, title, `self_text`, media tags and all the other attributes that are normally available via the Reddit API. The dataset complemented the Reddit comment corpus he released a couple months prior. When the data was initially made publicly available, he released it as a torrent where developers interested in using it could download their own local copies. Developers were all downloading the data for use either on their local machines or a cloud server. The problem with this is even with one of the most powerful desktop computers, loading the entire dataset into RAM was not feasible. Search times and joining (cross table) operations were expensive.

Conveniently soon after the release of this torrent, one of the lead engineers of *Google BigQuery*, Felipe Hoffa, uploaded the data to *BigQuery* and made the dataset publicly available. Each month, the dataset is updated with the latest information collected from the Reddit API.

With the convenience of BigQuery, it is now possible to query gigabytes of history Reddit data in a matter of seconds. Listed below are a few of the integral queries used in *Rally*, their sizes and the execution time.

Facades in Laravel with Google Services

In web programming, quite often developers will need access to static references to classes. Facades provide a static interface to such classes that are available in the application's service container. By default Laravel ships with several facades. These static proxies to underlying classes in the service


```
SELECT subreddit, total, sub_hour, num_gte_3000
FROM (
  SELECT
    HOUR(SEC_TO_TIMESTAMP(created - 60*60*5)) as
    ↪ sub_hour,
    SUM(score >= 3000) as num_gte_3000,
    SUM(num_gte_3000) OVER(PARTITION BY subreddit)
    ↪ total, subreddit,
  FROM [fh-bigquery:reddit_posts.full_corpus_201509]
  WHERE YEAR(SEC_TO_TIMESTAMP(created))=2015
  GROUP BY sub_hour, subreddit
  ORDER BY subreddit, sub_hour
)
WHERE total > 700
ORDER BY total DESC, sub_hour
```

Figure 3.6: The BigQuery SQL for finding the best hours to post on Reddit. This query processes 5.00GB across one table in roughly 8 seconds (1.5 seconds when cached)

container provide the benefit of a terse, expressive syntax while maintaining more testability and flexibility than traditional static methods.

The facade class itself only needs to implement a single method *getFacadeAccessor*. It is that method's job to define what to resolve from the container. Behind the scenes, the base facade class (which all facades must extend) makes use of a magic-method, *__callStatic()*, which defers calls from the facade to the resolved object.

The point of registering a facade may at times seem convoluted and unnecessary. It has always been a topic of discussion amongst the PHP world and a lot of the time boils down to personal preference and code readability. The facade approach was chosen particularly for *BigQuery* part of the project for a few main reasons:

- Expressive syntax without sacrificing testability of code
- Keep class responsibility narrow and well defined
- Clean constructor injection to automatically connect to *Google Services*

```
SELECT RIGHT( '0'+STRING(peak),2)+'-'+subreddit , hour , c
FROM (
  SELECT subreddit , hour , c , MIN(IF(rank=1,hour , null))
  OVER(PARTITION BY subreddit) peak
  FROM (
    SELECT subreddit , HOUR(SEC_TO_TIMESTAMP(created_utc
    ↪ )) hour , COUNT(*) c , ROW_NUMBER()
    OVER(PARTITION BY subreddit ORDER BY c ) rank
    FROM [fh-bigquery:reddit_comments.2015_08]
    WHERE subreddit IN (%subreddits)
    AND score>2
    GROUP BY 1, 2 )
  )
ORDER BY 1,2
```

Figure 3.7: Viewing activity (number of submissions) on subreddits over time. The wildcard *%subreddits* is replaced with a string comma-separated list of subreddits. This query processes 1.49GB across one table in roughly 2.5 seconds (1.1 seconds when cached)

```
public function register()
{
  $this->app->bind( 'google' , function () {
    $client = new Google_Client();
    $client->useApplicationDefaultCredentials();
    $client->addScope( Google_Service_Bigquery::BIGQUERY
    ↪ );

    return new GoogleAPI($client);
  });
}
```

Figure 3.8: Registering the Google service provider and binding the facade keyword *Google* to it.

and access the *BigQuery API*

- Explicit declaration defines what the class needs and what the class does

3.3 *phpRaw*

The Reddit API has several endpoints. It is through these endpoints where a client can retrieve posts specific to a subreddit, post a comment, moderate their account and all other actions that are normally available through the consumable web interface. For a single use or specific focus, calling the endpoints explicitly with cURL (or another client-side URL transfer) works fine but this strategy quickly fails as needs grow. Due to the wide array of endpoint calls utilized, it was necessary to develop an API wrapper that allows convenient calls to the API. Such a wrapper already existed for Python, Java, C and a few other languages but not PHP.

An open source wrapper was discovered on GitHub but was no longer maintained, was not written to comply with the latest API security requirements (OAuth2) and was missing nearly half of the endpoints. Building on the work done on this API wrapper, a successful implementation was built and is what *Rally* utilizes and depends on for direct Reddit data access. The GitHub repository from the point at which it was forked and built on is linked in the appendix.

Listed below are a functions from *phpRaw* to give a feel for the wrapper.

Get the user submitted data.

```
$phpRaw->getUserSubmitted($user, $limit = 25, $after =  
    ↪ null);
```

Get the top 10 hottest listings for a specified subreddit.

```
$phpRaw->getHot('funny', 10);
```

phpRaw was then modified to serve as a standalone vendor service brought in through Laravel's default dependency manager *Composer*. By extracting the wrapper to a separate module, updating and maintaining the endpoints is simple as they are changed over time. Using the power of composer and package dependencies, by including the declaration as outlined in Figure 3.9, whenever composer is updated it automatically updates to the latest version of *phpRaw*.

```
...
"repositories": [
  {
    "name": "kevin/phpRAW",
    "type": "vcs",
    "url": "https://github.com/kevineger/phpRAW"
  }
],
...
```

Figure 3.9: Requiring *phpRaw* as a dependency for *rally* in composer.

3.4 Clustering Script

3.4.1 Python

[Insert chunk about python script]

3.4.2 Scipy

[Insert chunk about Scipy]

Chapter 4

Algorithms and Methods

The sets of operations integral to the key components of *Rally*.

4.1 Hierarchical Clustering

When observing an open environment, an equally interesting and powerful metric for how the community is distributed is with clustering. One of the biggest benefits of hierarchical clustering is that you don't need to already know the number of clusters in the data set. It is with hierarchical clustering that within a subreddit, we are able to detect sub-communities. Strategies for hierarchical clustering land within two groups: agglomerative and divisive. Agglomerative is a bottom up approach where each observation starts in it's own cluster and pairs are merged as you move up the hierarchy. Divisive is a top down approach where all observations begin in a single cluster and split recursively down throughout the hierarchy.

To best understand the hierarchical clustering process, we will begin by showing the end result in what is known as a dendogram. A clustering of users amongst the subreddit */r/movies* is shown as a dendogram in Figure 4.1. The dendogram is a visualization in the form of a tree that shows the order and distances of merges throughout the hierarchical clustering. It can be understood as snapshots throughout the linkage of observations. On the x axis are labels representing numbers of samples (if in brackets) or specific samples (without brackets). On the y axis are the relative distances (using the 'ward' method described later). Beginning at the bottom of the lines (near the labels), the height of the horizontal joining lines tells us about the distance at which that labelled group merged with another label or cluster.

For the example shown in Figure 4.1 there are 4265 samples (users) being

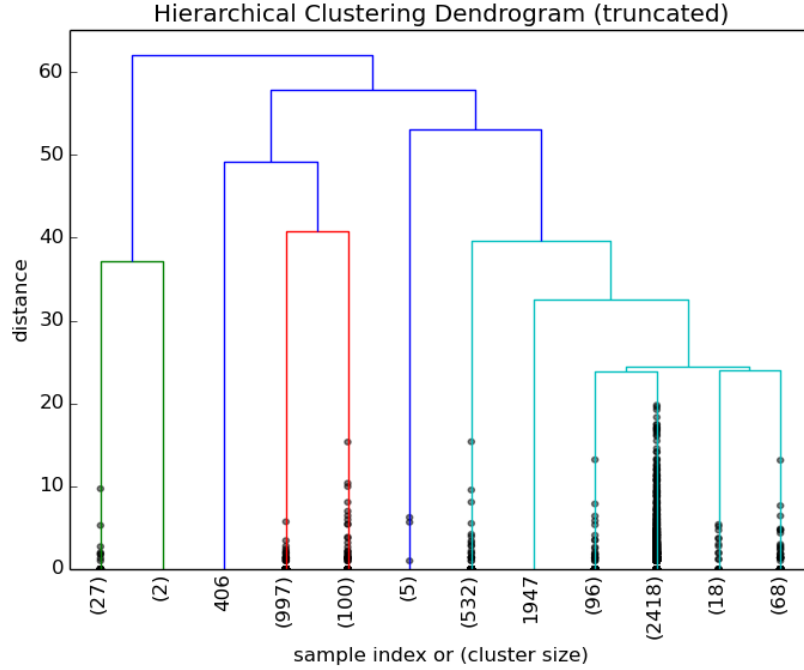


Figure 4.1: A dendrogram representing the hierarchical clustering amongst the subreddit */r/movies*.

processed. What is actually being shown is a truncated dendrogram, showing only the last 12 merges. The small black dots along the vertical lines represent joins that happened prior to the final 12. Truncation is an incredibly useful tool when plotting dendrograms. More often than not, we are only interested in the last few merges amongst the samples. The merge that carries the largest vertical distance will be the merge that attaches the most segregated groups. Again with the example in Figure 4.1 we can see three distinct groups being formed, identified by their green, red and teal colours.

Before we start summarizing the process, it is sometimes easy to forget what variables map to. A list is provided here:

- X: samples ($n \times m$ array), or data points or "singleton clusters"
- n: number of samples

- m: number of features
- Z: cluster linkage array
 - Contains the hierarchical clustering information
- k: number of clusters

4.1.1 The Clustering Process

To begin the clustering, we first gather the necessary data from *Google BigQuery*. The query retrieves the most recent 300 posts for the specified subreddit. A join is then made with the `link_id` from the inner query and a `UNION ALL` with the comment shard tables over the past 3 months. *BigQuery* does not directly support the `UNION ALL` syntax familiar to most sql languages, but instead supports comma separated tables wrapped in a `SELECT *`. After joining up the relations, user accounts that were deleted or made by an auto moderator are filtered out. The remaining authors are grouped by the `link_id` and selected out by the number of times they commented on each link. The query as it is executed in the application can be seen in figure 4.2. The query processes 9.95GB of data across a total of 4 tables and is completed between 5 and 10 seconds (depending on the subreddit under consideration).

Upon retrieving the data, the X matrix needs to be generated. As a reminder, in the matrix there are n samples and m features. Our samples are authors of comments on listings and our features are each of the listings. Outlined in figure 4.3 is the algorithm for processing the raw *BigQuery* response in to a usable matrix. Because the matrices can become very large in size, we are currently limiting the data gathered by using only the most recent 300 posts however as mentioned later, future work could focus on coming up with a preprocessing check to determine the anticipated size of response data from *BigQuery*.

The "ward" method is one of several that can be used to calculate the distance between newly formed clusters.

```

SELECT author , link_id , COUNT(link_id) as cnt
FROM (
  SELECT *
  FROM
    [fh-bigquery:reddit_comments.2016_01] ,
    [fh-bigquery:reddit_comments.2015_12] ,
    [fh-bigquery:reddit_comments.2015_11]
)
WHERE link_id IN (
  SELECT posts.name
  FROM [fh-bigquery:reddit_posts.full_corpus_201512] AS
    ↪ posts
  WHERE posts.subreddit = (%subreddits)
  AND posts.num_comments > 0
  ORDER BY posts.created_utc DESC LIMIT 300
)
AND author != '[deleted]'
AND author != 'AutoModerator'
GROUP BY author , link_id
ORDER BY author

```

Figure 4.2: The query executed on *BigQuery* to retrieve all cluster data.

```

for Each row in BigQuery response do
  end
  Do things

```

Figure 4.3: The psuedocode for processing the *BigQuery* response data.

Bibliography

- [Atl14] Ama: How a weird internet thing became a mainstream delight, 2014 [cited March 3, 2016]. → pages 3
- [Don10] Welcome redditors!, 2010 [cited March 3, 2016]. → pages 4
- [Gua05] A new website makes it easier to sift the mountains of news content online - and learns what you like, 2005 [cited March 3, 2016]. → pages 3
- [Lar14] Laravel documentation, 2014 [cited March 13, 2016]. → pages 5
- [Red14] Decimating our ads revenue, 2014 [cited March 3, 2016]. → pages 4