# Clickstream Data Processing

## Sreejith Sreekumar

December 10, 2021

#### Contents

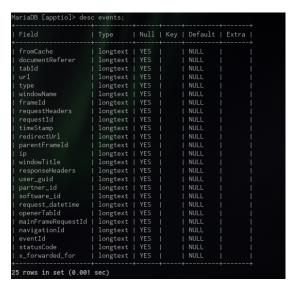
1	Notebooks & Code Overview	1						
<b>2</b>	2 Questions							
	2.1 How would you validate the data?	3						
	2.2 How would you normalize the data and make it representative?	4						

## 1 Notebooks & Code Overview

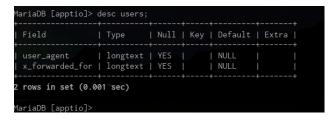
• 1.reading.ipynb: This notebook shows how the gzipped clickstream data is read and stored in a relational way for further processing and analysis.

Three tables(in a mysql/mariadb database) are created during the process of parsing data.

 events: Every row indicates the attributes associated with a click event. A description of the event table will look as follows.



 users: Information about the users to whom the ads were served. The same user can have same or a different ad served multiple times.



- status: Slightly less important table. Shows the description of various statuses (their codes and what they mean) during a click event.

```
[ariaDB [apptio]> select * from status limit 10;
 statusCode |
              statusLine
                                                             l error
              HTTP/1.1 200 OK
                                                               NULL
 200
 200
              HTTP/1.1 200
              HTTP/1.1 307 Internal Redirect
              HTTP/1.1 503 Service Temporarily Unavailable
                                                               NULL
              HTTP/1.1 302 Moved Temporarily
                                                               NULL
 302
 NULL
              NULL
 200
              HTTP/1.1 200 OK
                                                               NULL
              HTTP/1.1 200
              HTTP/1.1 302
                                                               NULL
 204
              HTTP/1.1 204
0 rows in set (0.000 sec)
```

## • 2.user-cleanup.ipynb

This notebook takes the user table as input and builds a profile for every user to whom the ads were served.



Additional attributes such as location(region, city, lat-long), ISP etc. could be extracted from the ip. However since I was unable to find a free API, this hasn't been done.

#### • 3.event-analysis.ipynb

This notebook takes the event table as input and analyses the ads served during each click event. For instance, where was the ad referred from? Where was the landing page etc:

	timeStamp	accept_language	user_guid	type	parentFrameld	windowName	statusCode	documentReferer	x_forwarded_for	user_agent	request_unixtime	country_code
1	.473207e+12	en-US,en;q=0.8	2f8b23ca273de94a51281b0697a126d7	sub_frame	0	None	200	None	69.207.104.248	Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.3	1473207176	US
1	.473207e+12	en-US,en;q=0.8	2f8b23ca273de94a51281b0697a126d7	main_frame	-1		200	http://www.imdb.com/title/tt0364845/episodes? r	69.207.104.248	Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.3	1473207176	US
1	.473207e+12	en-US,en;q=0.8	2f8b23ca273de94a51281b0697a126d7	sub_frame	0	None	200	None	69.207.104.248	Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.3	1473207159	US
1	.473207e+12	en-US,en;q=0.8	2f8b23ca273de94a51281b0697a126d7	sub_frame	0	None	200	None	69.207.104.248	Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.3	1473207159	US
1	.473207e+12	en-US,en;q=0.8	2f8b23ca273de94a51281b0697a126d7	main_frame	-1		200	http://www.imdb.com/title/tt0364845/? ref_=nv_sr_1	69.207.104.248	Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.3	1473207159	US

Snapshot of the data kept in events table

	from	count
0		1236120
26272	www.google.com	380009
25249	www.facebook.com	269685
21490	www.amazon.com	93385
24731	www.ebay.com	83527
31451	www.reddit.com	81046
35440	www.youtube.com	57912
35335	www.yahoo.com	35391
1580	apps.facebook.com	35306
665	accounts.google.com	27403
28673	www.mangahere.co	24309
24602	www.drudgereport.com	22520
10647	login.yahoo.com	19947
28328	www.linkedin.com	18962
10569	login.live.com	16100

Summary of click counts from a few of the top referrers

## 2 Questions

## 2.1 How would you validate the data?

- 1. Make sure that the zip files are not corruped, or check why this has happened if it did. The code takes in the variable  $log\_directory\_root$  where a report (failed.txt) will be written if any of the gzip files fail to read. It will also create the file (processed.txt) with a list of filenames that will be processed during the run.
- 2. <u>Validating Headers</u>: Every json record inside a gzipped text has the columns 'request\_keys' and 'server\_request\_key. However not all jsons have the same keys in the second level of hierarchy. During the first iteration of parsing, the code goes through all json instances and collects all the available keys (headers). From the files given, 36 columns were identified. These columns are written in the file schema.txt inside the configured log directory.

During the second iteration of parsing the data is decoded from json and is made to a dataframe, the columns that are missing in the dataframe are found and are attached to it with "None" as its values. This standardizes all the columns for all the data (from all the files).

#### Additional Validation

Ideally, the validation of data could be designed like a pipeline (chain) of conditions where each dataframe goes through, for example:

- a) Choose a subset of all the columns discovered in data
- b) Choose a date-range for the records from each dataframe
- c) Remove any instance where there is the error column is not null

These conditions have to be configurable or could be modified with very less code addition.

Note: Please see line 121 in the function standardize data function in the file utils.py

## 2.2 How would you normalize the data and make it representative?

The raw data can be split into 2 main parts. An **event** (a click) and a **user** (identified by an IP address -  $x_forwarded_for$ ). This redundancy can be reduced by splitting the whole data into two tables: events and users for separate analysis.

The attributes associated with the user comes from the  $user\_agent$  field in the request. The number of times a user has clicked an ad could also be added to this table through a groupby of  $x\_forwarded\_for$  and  $user\_agent$ . However, this needs to be done in a map reduce way since the data is too large to be held in RAM and processed.

The event table stores all the attributes associated with a clickevent. This table can be used to understand the details associated with an event such as the referrer, and the landing of a click. In additional, the table provides information on when the click happened (unix\_timestamp), when the data reached the data collection point (timeStamp), some details of the page such as frameId, tabId etc., and status of the click. The referrers (long urls) and the detination urls could be cleaned up to derive business insights on customer engagements.

The status table is a relatively small table that contains the description of every status that have appeared during a click event.

Note: These tables have been exported to a csv and has been shared as a link.