

**Data gathering**

1. How did we collect the data?

We collected the data using a Python script running on a CSX server..

2. The Format we used to collect the data

(Number:: Dataset name::Link::Size:: Title)

1::google.com::http://google.com::206710::Google

- Delimiter – “ :: “ -non English characters

3. What type of data that is there :

Few notable titles:

- a. UNKNOWN – Any site that did not respond within 10 seconds or simply just unreachable, we set the size as -1, and the title as “UNKNOWN”. For all sites that returned the page, and allowed us to get the page size but did not have correct tag for the page title, they also got the title as “UNKNOWN”.
- b. Loading... – These websites had title as “Loading...” and some redirected to inappropriate sites. We looked at some sites that gave us “Loading” as the tile, and they had a JavaScript for a redirect.

**How was the data Collected?**

We used a python script (visist2.py). The script takes in two arguments: 1) Begin line# 2) End line#.

Example command to visit sites from line 1001 to 2000: “python visit2.py 1001 2000”

It reads the specified lines from top-1m.csv input file, and visits each site using the user agent of desktop Chrome browser on Windows 10. The use of user agent was done because some sites did not return a content if we did not provide a user agent.

It does not execute JavaScript, so if the site has moved and the page had JavaScript to redirect the visitor, it would not have been redirected. If the http request to the site resulted in http 301 redirect code, the script visited the redirected site. The script tries to visit site by appending “http://” in front of the hostname. If it results in error, it retries visiting the site with “https://”. Each visit had maximum timeout of 10 seconds. We divided the job of visiting all million sites into 100 jobs (one job is responsible of visiting 1000 sites). We ran 100 jobs on CSX0 at once. We were able to finish visiting all 100k sites within 2 hours. Each job generates one output file, so we have total of 100 output files.

Among all million sites, our script was able to get size of the page from 88397 sites.

**How to run on Hadoop ?**

- Put those file in the Hadoop cluster and unzip.

Commands used for transfer files form one server to another.

**//Transfer directory to another linux server**

**scp -r /home/username/unzip/ hadoop1:. :**

**// Copy directory containing files to hadoop cluster**

**hadoop fs -copyFromLocal unzip / username /**

- To run the code –Run the commands.

**// Run the jar file**

**hadoop jar Homework2.jar Question1 / username /unzip /username/output1/1**

**// Copy the output file from hadoop cluster to local machine**

**hadoop fs -copyToLocal / username /output1/1/part-r-00000**

**// Open the output file with use of nano editor**

**nano part-r-00000**

Question 1- List top 10 biggest websites

URL	Size in bytes
http://likealyzer.com	11298813
http://cartoonnetwork.com.co	9621305
http://cartoonnetworkla.com	9621297
http://cartoonnetwork.com	9621295
http://cartoonnetwork.com.br	9621295
http://cartoonnetwork.com.ar	9621288
http://cartoonnetwork.com.mx	9621288
http://cartoonnetwork.com.ve	9621288
http://ahu.go.id	7399705
http://diagnosisia.com	5649610

Question 2 – List top 10 most occurring unique titles (6<sup>th</sup> from the top is the sites with empty title.)

Title	# of occurrences
UNKNOWN	14842
Google	251
Loading...	110
Home	83
Not Found	77
	61
Yahoo	48
Welcome to nginx!	39
Index of /	38
You are being redirected...	36