# Analytics using ClickStream Data

Path Optimization

Niranjan Kumar                    4/30/18                    BigData Expert

# CASE STUDY

**Input Data:**

Here's a summary of the data we're working with:

- Omniture logs – website log files containing information such as URL, timestamp, IP address, geocoded IP address, and user ID (SWID)

  - The Omniture log dataset contains about 4 million rows of data, which represents five days of clickstream data. Often, organizations will process weeks, months, or even years of data.

- Users– CRM user data (registered Users) listing SWIDs (Software User IDs) along with date of birth and gender.

- Products – CMS data that maps product categories to website URLs.

**Tools Used**:

- Hive – To perform the data analytics.
- Excel – To perform the data visualization.

## Creating Tables

Steps:

- Copy local data into Hadoop using command: Hadoop fs –put localPath HadoopPath
- Open Hive shell
- Create a Database alabs_db. CREATE DATABASE OF NOT EXISTS) alabs_db.

```
hive> Create database if not exists alabs_db
    > ;
OK
Time taken: 0.112 seconds
hive> show databases;
OK
alabs_db
computersalesdb
default
movieratings
zzniranjan
Time taken: 0.075 seconds
hive>
```

*Figure 1: Creating the Database*

- Now add tables into it.
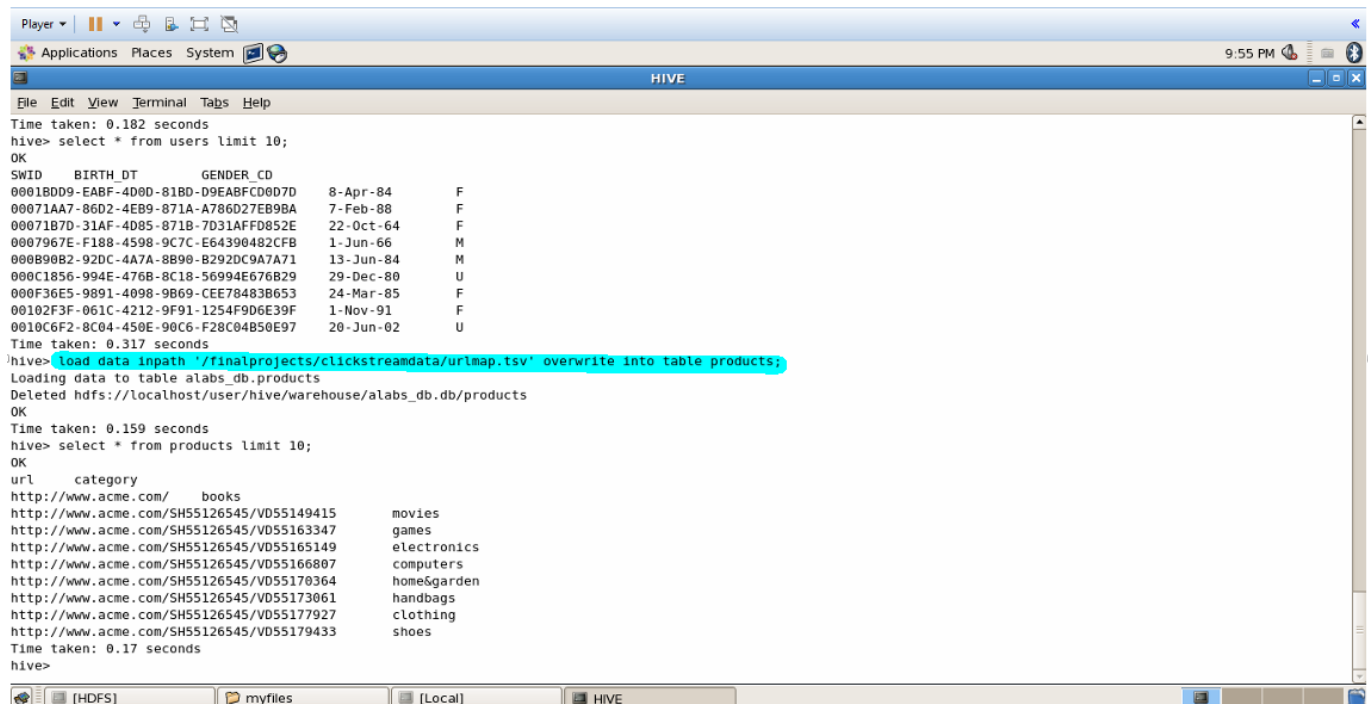- Now create the table schema for Users, Products and omniturelogs in hive



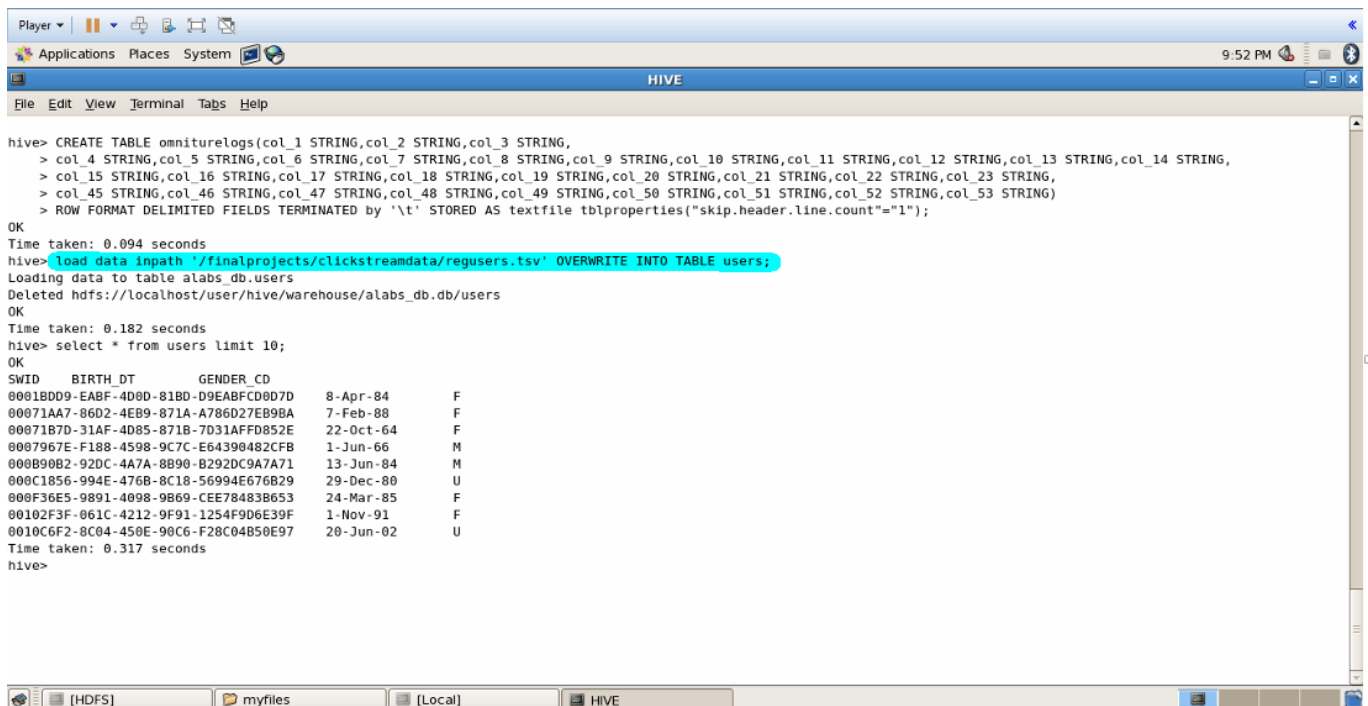*Figure 2: Loading Data into Products*



*Figure 3: Loading Data into Users*

Create View and final table for analysis



```
webanalytics.sql - Notepad
File   Edit   Format   View   Help
USE alabs_db;

--create view for analysis

CREATE VIEW omniture AS
SELECT COL_2 TIMESTAMP,
COL_8 IPADDRESS,
COL_13 URL,
COL_14 SW_ID,
COL_50 CITY,
COL_51 COUNTRY,
COL_53 STATE
FROM omniturelogs;

--create final table for analysis

CREATE TABLE webloganalytics AS
SELECT TO_DATE(o.timestamp) logdate,
o.url url,
o.ipaddress ipaddress,
o.city city,
UPPER(o.state) state,
o.country country,
p.category category,
CAST(DATEDIFF(FROM_UNIXTIME(UNIX_TIMESTAMP()),
FROM_UNIXTIME(UNIX_TIMESTAMP(u.birth_dt,'dd-MMM-yy')))/365 AS INT) age,
u.gender_cd gender_cd
FROM omniture o
INNER JOIN products p
ON o.url = p.url
LEFT OUTER JOIN users u
ON o.sw_id = CONCAT('{',u.sw_id,'}');
--end
```
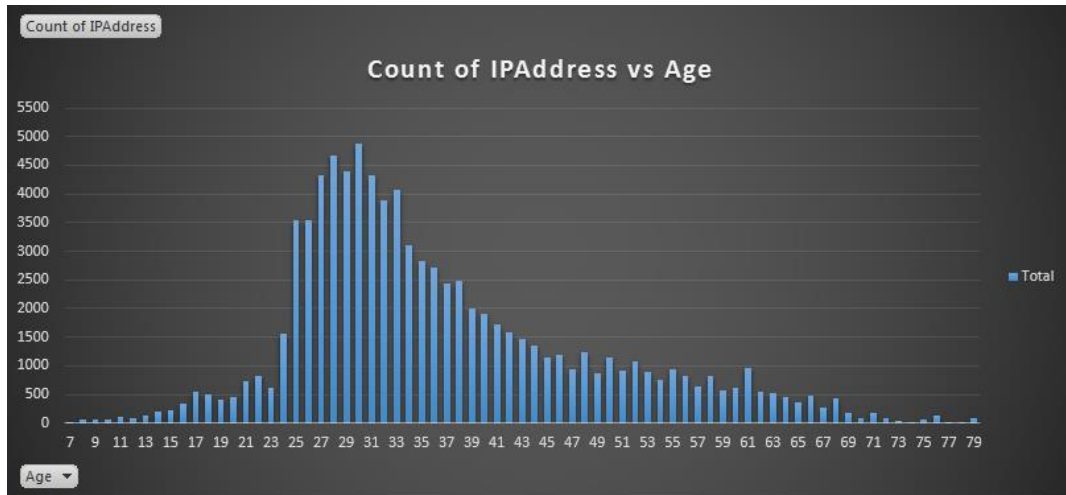
#Omniture VIEW and Webloganalytics TABLE created



```
Player ▼ | ❚❚ ▼ 🖧 🖫 ⛶ ▨
🔲 Applications  Places  System  🖼 🌐
🔲                                              training@localhost:~
File  Edit  View  Terminal  Tabs  Help
[training@localhost ~]$ hive
Hive history file=/tmp/training/hive_job_log_training_201804300448_174889098
3.txt
hive> use alabs_db;
OK
Time taken: 3.877 seconds
hive> show tables;
OK
omniture
omniturelogs
products
users
webloganalytics
Time taken: 0.859 seconds
hive>
```
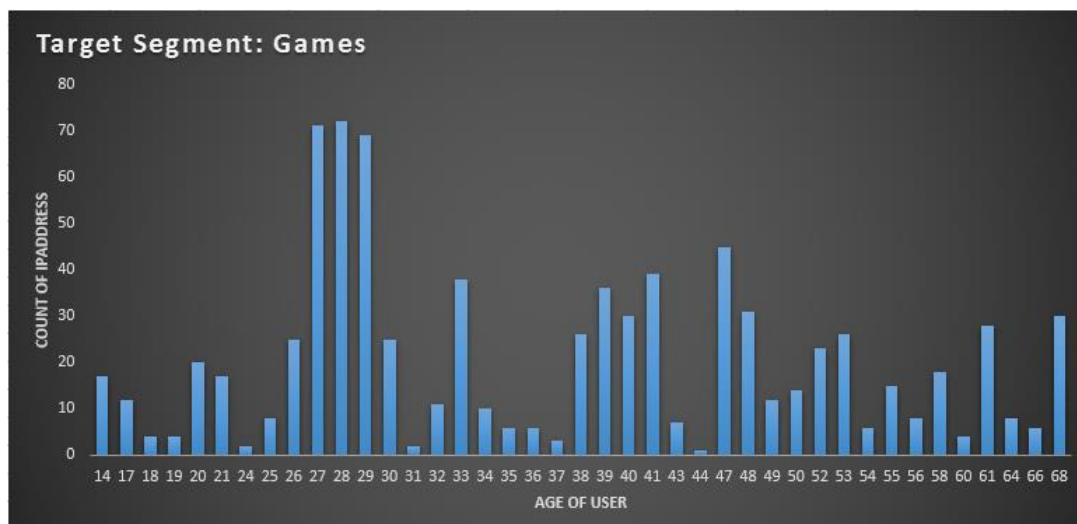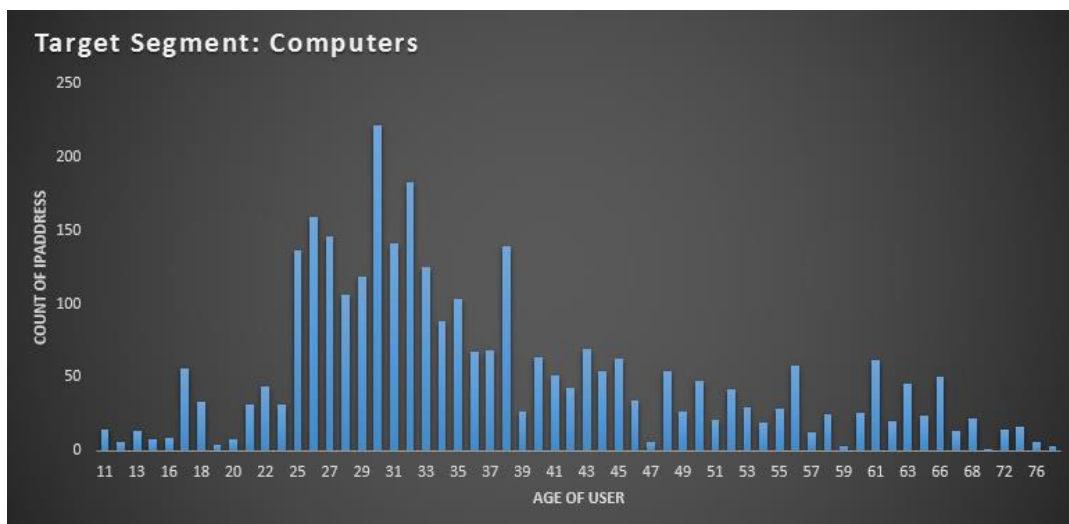
Export the data from 'Webloganalytics' to csv file for Visual Analytics.

# Analysis

From the Output file we can analyse the website usage based on IpAddress against the age of the user.



Identifying the Customer behaviour per target segment – Computers and Games

# Identifying the webpages with highest bounce rates