# Big Data Ingestion Using Hadoop : 1038-A

www.trendzzz4u.com

Faculty Advisor and Sponsor:
Dr. Rohit Aggarwal

# Team Members and Roles

Aditya Kannan

Software Engineer and Project Manager

Ekta Jaiswal

1st

Data Engineer

Vikal Gupta

Software Developer and Reporting Analyst

# Agenda

- Data Scenario today
- Why this project and Objective
- Business Questions to be answered
- Project Infrastructure
- Project Workflow
- Learnings and Conclusion

David Eccles
School of Business
THE UNIVERSITY OF UTAH

# Data Scenario today

- According to IBM, 80% of data generated today is unstructured

| Outsourcing Data | Video Streaming Data | Social Media Data | Logs Data |
|---|---|---|---|
|  |  |  |  |

- Need to process unstructured data to structured data
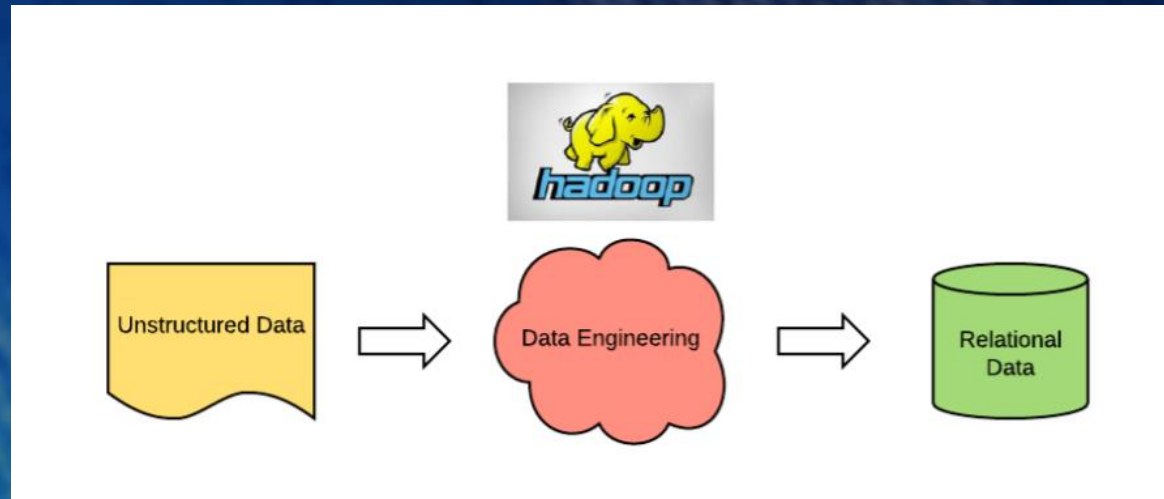
# Why this Project ?



- To learn and implement big data technologies which are used to process log data

- Clickstream Log Data as the data source

- What is Clickstream data?

  Clickstream Data is user navigation data on any website

```
66.249.66.127 - - [30/Apr/2017:04:24:41 -0700] "POST /swatches/ajax/media/ HTTP/1.1" 200 230 "http://www.trendzzz4u.com/women/
tops-women/jackets-women.html?material=38&size=167&style_general=126" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/
bot.html)"
```

- Big data world comprises of many technologies
- Focus of this project is to learn and implement Apache Hadoop ecosystem



- Hadoop is primarily used for Data Engineering by many companies notably Amazon, eBay, Walmart

# Business Questions to be answered

- Most popular browsing time
- Most popular product category
- Weekly Distribution of Clicks per page
- Customer Conversion Rate

**David Eccles
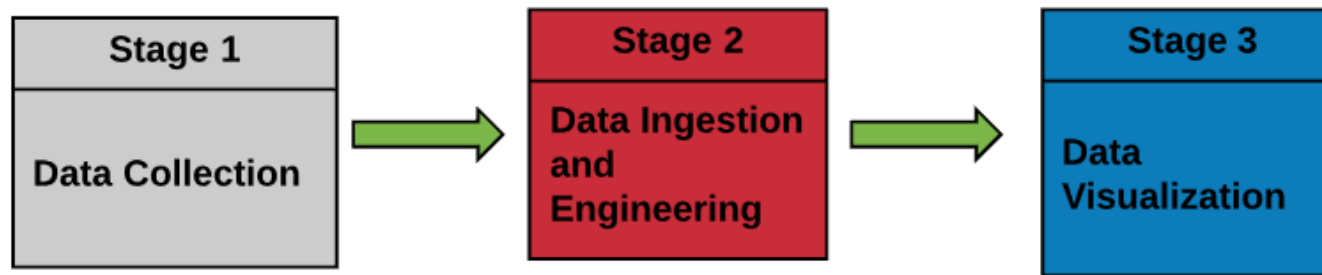School of Business**
THE UNIVERSITY OF UTAH

# Project Infrastructure

Major components of infrastructure

- eCommerce Website Setup: Magento eCommerce platform, goDaddy cPanel hosting

- Apache Hadoop Setup: Multi Node Cloudera Hadoop cluster on aws EC2

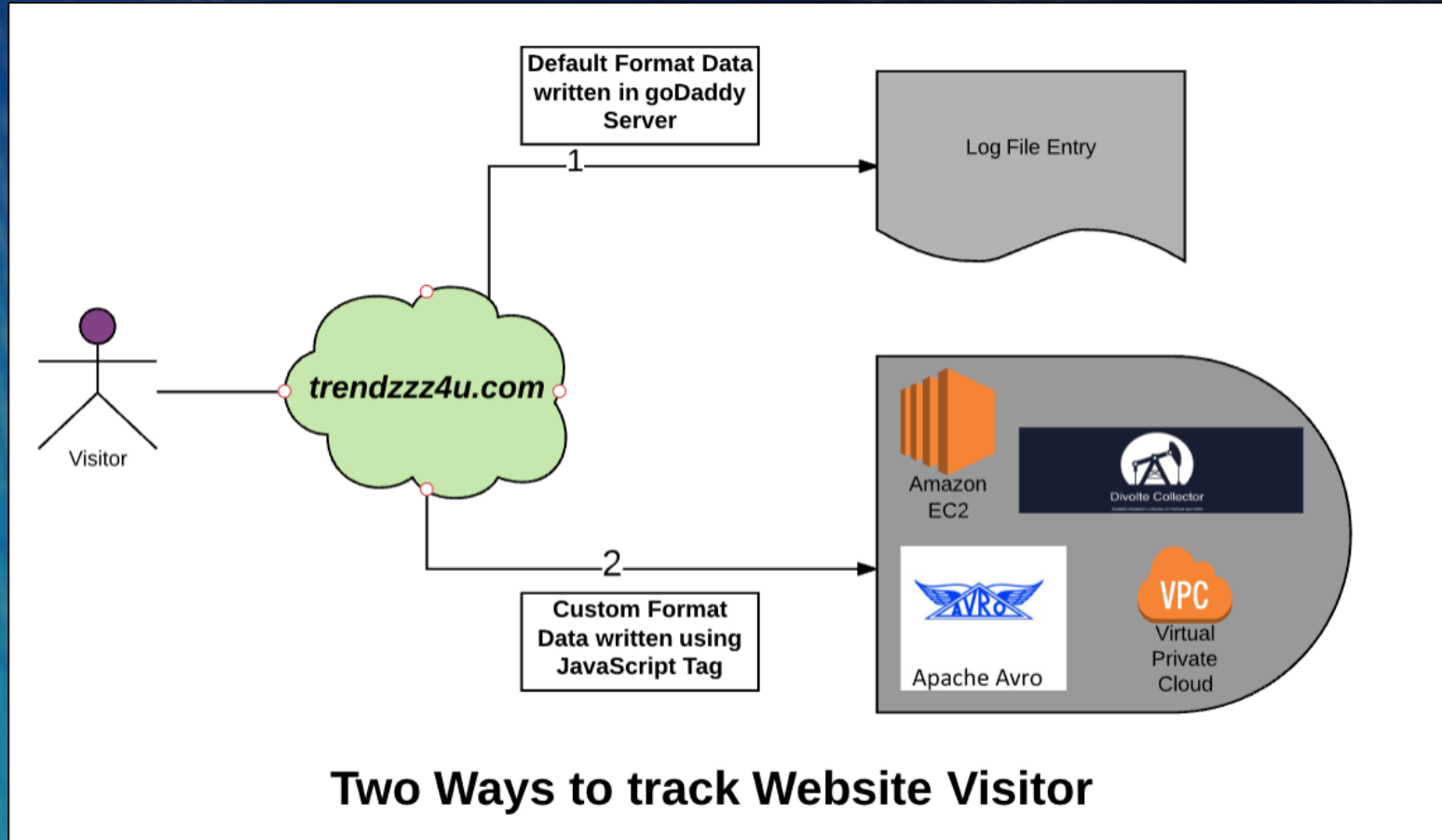- Data Collection Server: Divolte.js setup on aws EC2 to track custom events from website

# Project Workflow

- Primarily three stages
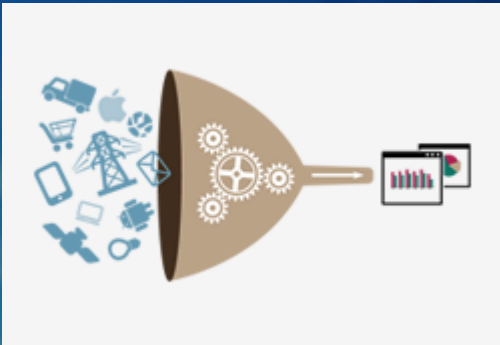- Performed several iterations of the below workflow

# Stage 1 : Data Collection



Two Ways to track Website Visitor

# Stage 2: Data Ingestion and Engineering

- The process of accessing and importing data for immediate usage or storage in a database is called as Data Ingestion



- Build Data Pipeline
- Transfer files from local Filesystem to HDFS (Hadoop File System)

E.g. Apache Sqoop is a popular tool used in big data ecosystem to transfer  bulk data

- Data engineering is a process of converting unstructured data to meaningful relational data using set of sophisticated tools or procedures
  - ✓ ELT Process (Extract Load Transform)
  - ✓ Focus on data transformation using Map Reduce

E.g. Components used in this project:
- Apache Hive
- Apache Pig
- Python

# Apache Hive

- It provides a SQL like interface to query data stored in various databases and file system

  Application in our project: Avro File (Data Source)

### Before



### Run Script



### After

# Apache Pig and Python

- It is a high level platform for creating programs on Hadoop. The language used is called as Pig Latin

  Application in our project: Web Server Log File (Data Source)

## Before

```
[30/Jun/2017:05:11:49 -0700] '
148&style_bottom=116 HTTP/1.1'
tp://www.google.com/bot.html)
 [30/Jun/2017:05:11:52 -0700]
53&material=39 HTTP/1.1" 200 1
tp://www.google.com/bot.html)'
[30/Jun/2017:05:11:59 -0700] '
99&style_bottom=116 HTTP/1.1"
tp://www.google.com/bot.html)'
[30/Jun/2017:05:12:16 -0700] '
```

## Run Script

```
A = LOAD '/home/training/Downloads/WebLog2' USIN
rarray,req_type:chararray,req_link:chararray,re
B = FOREACH A GENERATE ip_addr,temp1,timestamp,t
data = distinct B;
dump data;
STORE data INTO '/home/training/Downloads/Parsec
/* Calculate the number of web pages a user vis:
ip_data = GROUP data by ip_addr;
ip_count = FOREACH ip_data GENERATE group AS tir
dump ip_count;
/* Statistics where requests were successfull i.
time_data = GROUP data BY timestamp;
byte_count = FOREACH time_data GENERATE group AS
dump byte_count;
```

## After

| TIMESTAMP | TIMEZONE | REQUEST_TYPE |
|---|---|---|
| 30/Jun/2017:05:11:49 | 0700 | "GET /men/bottoms |
| 30/Jun/2017:05:11:52 | 0700 | "GET /men/bottoms |
| 30/Jun/2017:05:11:59 | 0700 | "GET /men/bottoms |
| 30/Jun/2017:05:12:16 | 0700 | "GET /men/bottoms |

David Eccles
School of Business
THE UNIVERSITY OF UTAH

# Sample of Python Script

```python
def sub_regular_exp1(str):
    match = []
    if re.search("GET",str):
        regex = re.compile(r"\"([\w]+)\s+\/([\w]+)\/(.+?html)(.*)\"");
        match = re.findall(regex, str);
    elif re.search("POST",str):
        regex = re.compile(r"\"([\w]+)\s+\/([\w]+)\/([\w]+)\/([\w]+)/(.*)\"");
        match = re.findall(regex,str);
    return match


def sub_regular_exp2(str):
    regex = re.compile(r"\"(http.*html)(.*)\"");
    match = re.findall(regex, str);
    if len(match) == 0:
        match = [' ',' ']
    return match
```

David Eccles
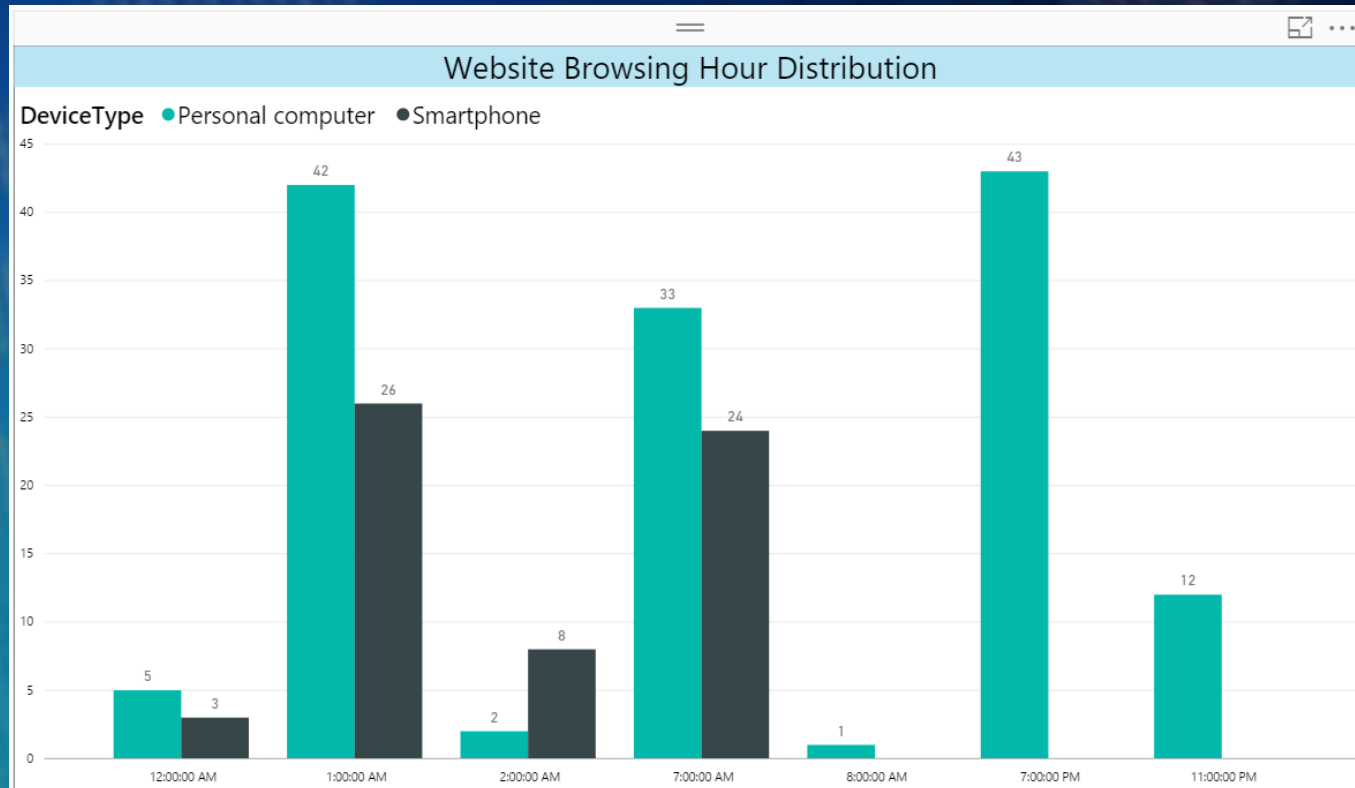School of Business
THE UNIVERSITY OF UTAH

## Stage 3: Data Visualization

- Final stage of the process
- Structured Data exported from Hadoop to csv format
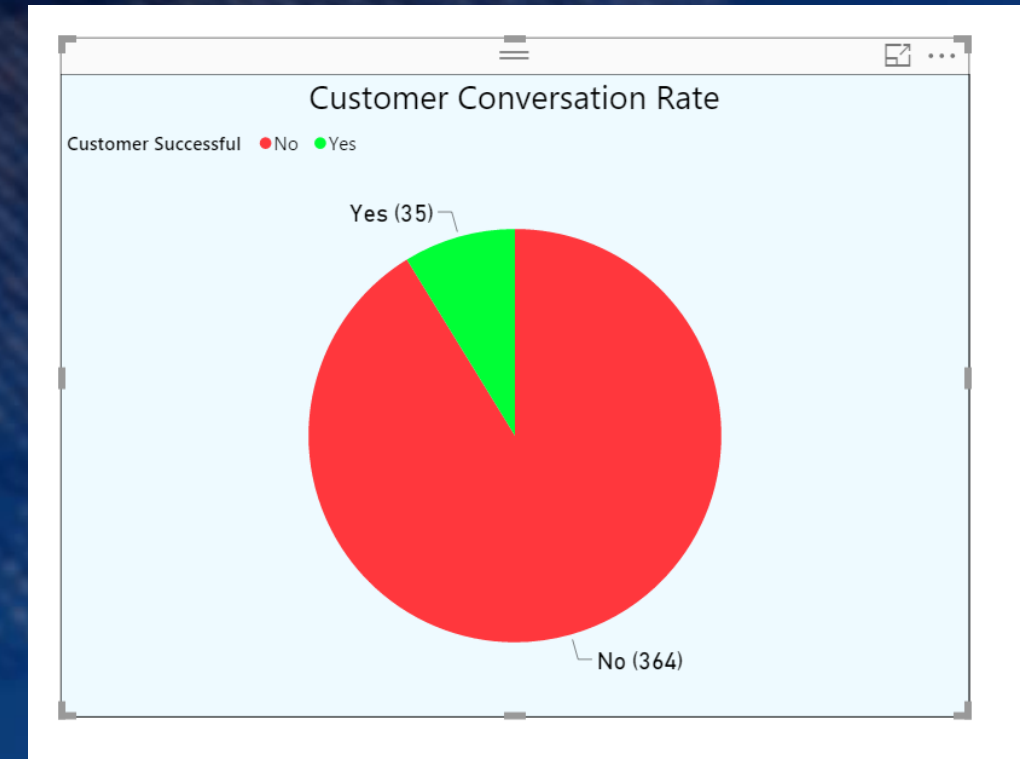- Use of Business Intelligence tools such as Tableau and Power BI

# Sample Reports

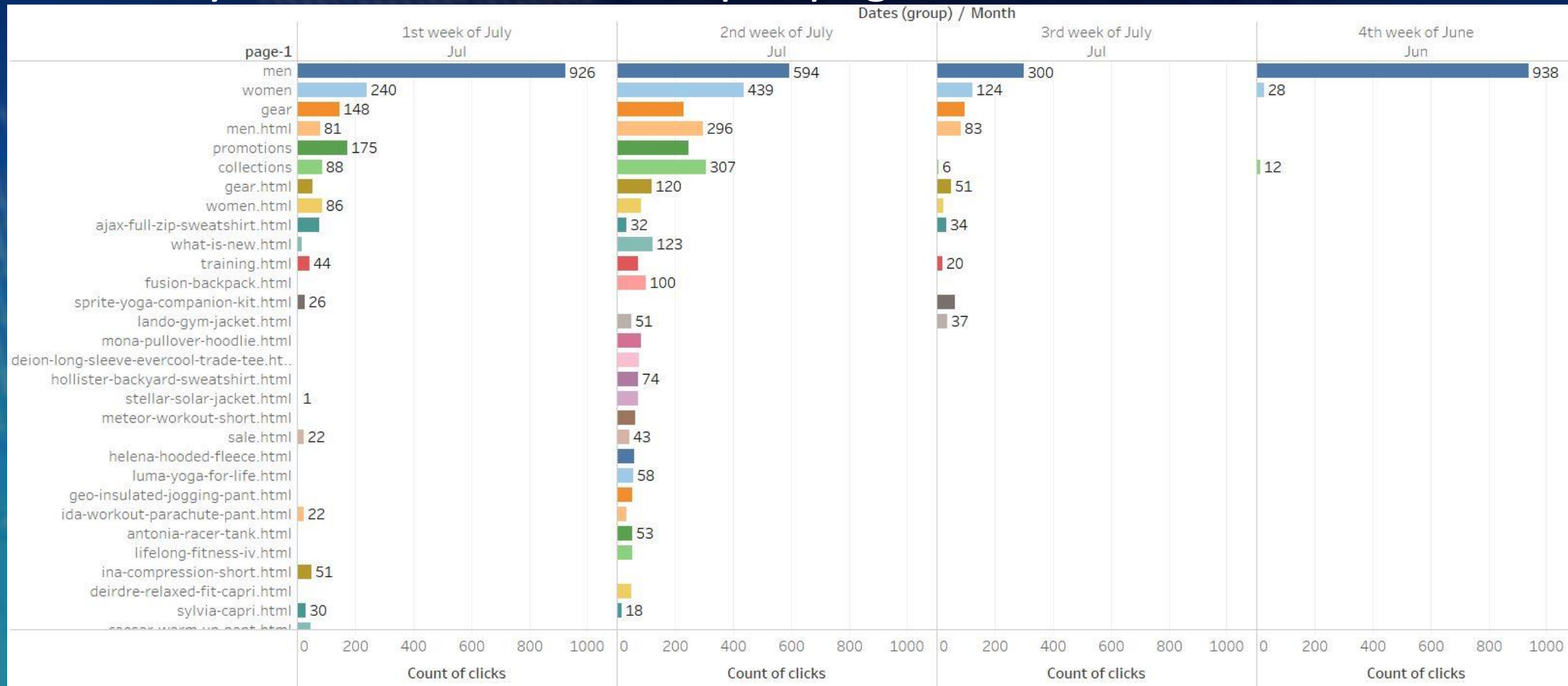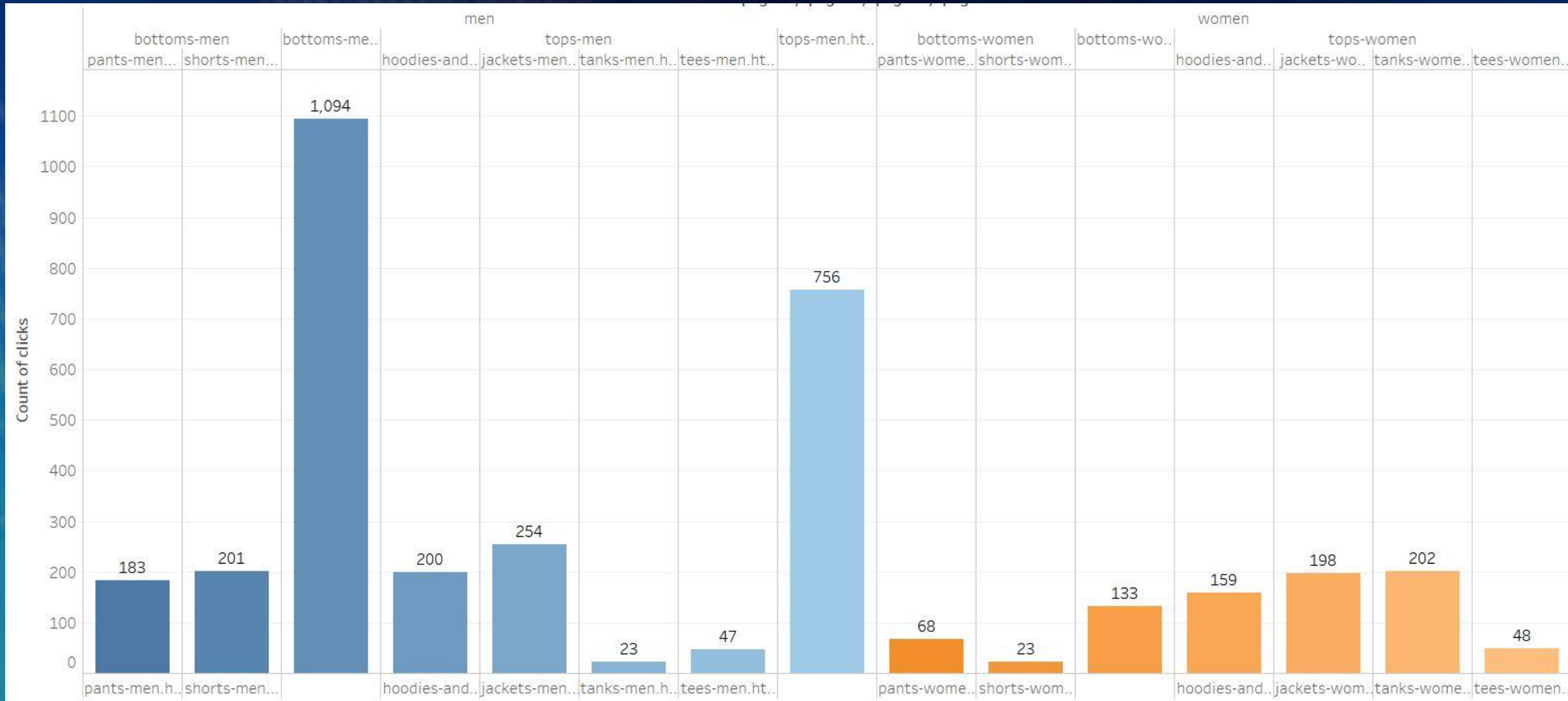## Website Browsing Hour Distribution

## Customer Conversion Rate

# Weekly Distribution of Clicks per page

# Most Popular Product Category

# Challenges Faced



- eCommerce Website Setup
- Multi-node cluster creation on Amazon Ec2
- Cloudera Hadoop Installation on cluster
- Parsing links to get detailed information about Product Category and sub-category

David Eccles
School of Business
THE UNIVERSITY OF UTAH

# Learnings and Conclusion

- Implementation of an end to end project
- Technology Stack worked on:

## Special Thanks

- Thanks to our entire Capstone Faculty team and Sponsor for timely guidance
- Bi-weekly progress reports helped us to get a reality check of the project
- Great learning experience

**David Eccles School of Business**
THE UNIVERSITY OF UTAH

# Any Questions?

# Thank You!!!