



UNIVERSITY OF MALAYA

Faculty of Computer Science & Information Technology

**Mastre Of Data Science
Semester II Session 2021/2022**

**WDQ7007 Big Data Management
Group Project**

**Title: The Big Data ETL Process Of Social Media MOMO's Behavior
Via Hadoop**

Lecturer: Dr. ALI FEIZOLLAH

Group Members

No.	NAME	MATRIX NO.
1	DI BEI	S2013283
2	ZITENG ZHANG	S2149768
3	ZUOGEB CHEN	S2125783
4	YINGGANG JIE	S2113868
5	JINXING LI	S2128176
6	HO WEI YAN	S2116489
7	JUNBO SHI	S2140997
8	HAN LI HUI	S2025432

Introduction

Social media applications have become essential to most individuals nowadays. Momo is a popular social media app in China. The distribution of the customer is necessary for the analysis to attract more people to use it. But for the thousands of data, how to perform it is a big problem for us. So in the research we will use hadoop technology to do the ETL and analyze the data for insight. For the tools, we will use MySQL, Sqoop, Hive and Hbase. Hadoop is an efficient platform to manage big data.

For the data, we search from a website. This dataset contains a total of 140,000 samples of user communication data with 20 columns that were collected from the social media app of Momoin in one day.

The dataset obtained from:

<https://blog.csdn.net/clark123432343/article/details/124409046>

Research Objectives

1. To set up the data warehouse and build a data pipeline to manage the social media data momo.
2. To do data cleaning and keep the data without noisy data and no value after EDA in hive.
3. To leverage a data access solution such as HBase so that we can better access and manage the data of social media data.
4. To do data visualization according to the behaviors of users in social media momo.

Problem statement

With the rapid growth of social media in recent years, Momo is one of the most popular apps among Chinese teenagers. Different from WeChat, Whatsapp, Facebook, qq, etc. Momo adds friends or joins interest groups through people's real time positioning. Customer losing is one of the important problems that Momo needs to consider for improving its ROI. Huge amount of user data generated every day. Integrating the user communication data and building a unified data warehouse can help it to classify, track, and monitor the business data to achieve the objectives of increasing ROI. Therefore, building the secure and reliable database or data warehouse and doing exploratory data analysis of the data records can help Momo to improve the performance of its products and avoid user loss.(Xu & Wu, 2019)(Chan, 2019)(Li, 2020)

Literature review

Stefan et al. proposed the use of social media app Momo to realize the process of user data collection. Momo will use centralized and distributed filtering methods to process and store data. Momo's data collection group mainly includes different user groups such as men and women. The system will process user data for visualization. Provide users with user behavior predictions and friend recommendations. In this paper, they also focus on the use of visual machine learning concepts on mobile and server. (Edlich & Vogler, 2013)

Xu & Wu examine a phenomenon that has become popular in mainland China in recent years, known as "stranger communication", focusing on an application known as Momo, a social discovery and dating platform widely used in China. The report begins with an overview of the post-reform and opening-up Chinese context,

which gave rise to urbanization, individualization and a sexual revolution. Thereafter, it turns to theories of strangers and cosmopolitanism, which are used to analyze stranger communication. Using cultural discourse analysis as an analytical procedure, the study presents the results of online and offline interviews conducted in both Beijing and Shanghai to examine how Momo users in urban metropolitan areas use the application in order to analyze the cultural radiation in their communication practices. Based on the data obtained, this paper analyzes MOMO Ren communication in the Chinese context. The paper analyzes MOMO Ren communication in the Chinese context and describes how Momo is designed to facilitate this communication. (Xu & Wu, 2019)

Lik et al. examined how heterosexual male users' endorsement of masculinity was related to the number of casual sex partners they encountered on Momo, based on a social constructivist view of gender. The study examined how heterosexual male users' endorsement of masculinity was related to the number of casual sex partners they encountered on Momo, based on a social constructivist perspective. The study also explored the mediating role of sexual motivation in using Momo. Analysis of survey data from 125 heterosexual male Momo users revealed an indirect positive relationship between endorsement of masculinity and the number of sexual partners, while at the same time, it was directly but negatively associated with the number of sexual partners. These contradictory associations are explained by different models of the various dimensions of masculine ideology. Since unsafe sex has been found to be associated with the use of geosocial networking applications, their study also calls for the integration of the concept of implementing safer sex with the cultural ideals of masculinity. (Chan, 2019)

Li et al. identified three types of institutional privacy concerns, including information leakage, information tracking, and surveillance by apps and the Chinese government. Most participants took little action to address their privacy concerns, instead exhibiting an attitude of ignoring and accepting the privacy risks of GSNA. In contrast, few participants developed proactive strategies. In contrast, few participants developed proactive strategies to mitigate their concerns, such as abstaining or being cautious when using apps and features, and reporting privacy threats to app platforms to app platforms. Importantly, this study provides new insights into user privacy issues that are particularly relevant to the social and cultural context in China. This study provides new insights into users' privacy issues in the Chinese social and cultural context. However, since this study mainly focuses on users' privacy issues and strategies, they encourage future research to delve into solutions and strategies for mitigating users' privacy issues on social platforms in a global context.(Li, 2020)

Methodology

HDFS

I. Setting the root password in through my sql.

```
student@student-VirtualBox:~$ sudo cat /etc/mysql/debian.cnf
# Automatically generated for Debian scripts. DO NOT TOUCH!
[client]
host      = localhost
user      = debian-sys-maint
password  = C1F3y0vh6tLQXvX
socket    = /var/run/mysqld/mysqld.sock
[mysql_upgrade]
host      = localhost
user      = debian-sys-maint
password  = C1F3y0vh6tLQXvX
socket    = /var/run/mysqld/mysqld.sock
```

```
student@student-VirtualBox:~$ mysql -u debian-sys-maint -pCf3y6vh6tLQvX
mysql: [Warning] Using a password on the command line interface can be insecure
.
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 3
Server version: 5.7.38-0ubuntu0.18.04.1 (Ubuntu)

Copyright (c) 2000, 2022, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.
```

II. Create a database and a specific table of data with the correct column type from the dataset and check if the table framework has been created successfully in the correct database.

```
mysql> create table tb_msg (msg_time varchar(20),
-> sender_id varchar(15),
-> sender_sex varchar(8),
-> sender_ip varchar(20),
-> sender_os varchar(15),
-> sender_network varchar(5),
-> sender_gps varchar(30),
-> receiver_id varchar(15),
-> receiver_ip varchar(20),
-> receiver_os varchar(20),
-> receiver_network varchar(5),
-> receiver_gps varchar(30),
-> receiver_sex varchar(8),
-> msg_type varchar(5),
-> distance varchar(10));
Query OK, 0 rows affected (0.07 sec)

mysql> show tables;
+-----+-----+
| Tables_in_db2_msg |
+-----+-----+
| tb_msg           |
+-----+-----+
1 row in set (0.00 sec)
```

III. Load the data from local to the table created

```
mysql> LOAD DATA LOCAL INFILE '~/Downloads/data2.csv' INTO TABLE tb_msg FIELDS TERMINATED BY ',';  
Query OK, 68905 rows affected, 65535 warnings (1.36 sec)  
Records: 68905 Deleted: 0 Skipped: 0 Warnings: 274541
```

IV. Check if the data imports correctly and counts how much data has been loaded.

```
mysql> Select count(*) as cnt from tb_msg;
+-----+
| cnt   |
+-----+
| 68905 |
+-----+
1 row in set (0.00 sec)
```

V. Finished by mysql, in go back to hadoop, the system will use sqoop command to transfer data from mysql to HDFS. Lastly, ‘cat’ will be used to check if the data sent to HDFS correctly.

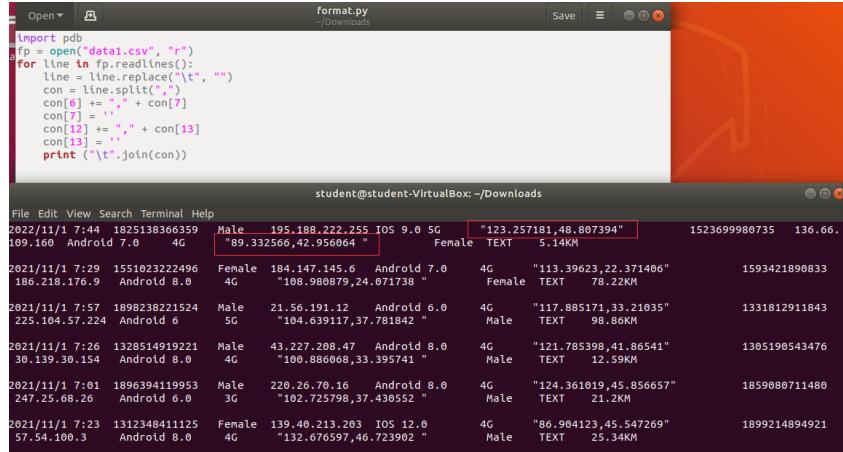
```
student@student-VirtualBox:~/home/WQD7007/sqoop/bin
File Edit View Search Terminal Help
student@student-VirtualBox:~/home/WQD7007/sqoop/bin$ sqoop import --connect jdbc:mysql://localhost/db2_msg -username root -password root --table msg --target-dir /home/WQD7007/sqoop/bin/sqoop_out --hive-import --hive-overwrite --create-hive-table
sqoop:ERROR [main]: java.lang.RuntimeException: HCatalog does not exist! HCatalog jobs will fail.
Please set SHACATE_HOME to the root of your HCatalog installation.
sqoop:ERROR [main]: java.lang.RuntimeException: Accumulo imports will fail.
Please set ACCUMULO_HOME to the root of your Accumulo installation.
sqoop:ERROR [main]: java.lang.RuntimeException: Accumulo imports will fail.
Please set ACCUMULO_HOME to the root of your Accumulo installation.
12/06/01 08:39:42 INFO sqoop.Sqoop: Running Sqoop version: 1.4.7
12/06/01 08:39:42 INFO sqoop.Sqoop: The parameter --password was passed on the command-line in an insecure. Consider using -P instead.
12/06/01 08:39:42 INFO Manager.MySQLManager: Preparing to use a MySQL streaming resultset.
12/06/01 08:39:42 INFO Manager.MySQLManager: Using the MySQL connector.
12/06/01 08:39:43 INFO Manager.MySQLManager: Establishing SSL connection without server's identity verification is not recommended. According to MySQL 5.5.4+, 5.6.2+ and 5.7.6+ requirements SSL connection must be established by default if explicit option is not set. For compliance with these requirements you must set the system variable 'mysql.secure-server-property' to set to 'True'. You need either to explicitly disable SSL by setting 'mysql.secure-server-property' to 'False' or provide truststore for server certificate verification.
12/06/01 08:39:43 SELECT t.* FROM tb_msg AS T LIMIT 1
12/06/01 08:39:43 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM tb_msg AS T LIMIT 1
12/06/01 08:39:43 INFO manager.SqlManager: Hadoop MapReduce home is [/home/WQD7007/hadoop]
12/06/01 08:39:43 INFO org.apache.hadoop.mapred.MapTask: user code uses or overides deprecated API. Please: Recompile with -Xlint:deprecation for details.
12/06/01 08:39:47 INFO org.apache.hadoop.mapred.MapTask: user code uses or overides deprecated API. Please: Recompile with -Xlint:deprecation for details.
12/06/01 08:39:47 INFO org.apache.hadoop.mapred.MapTask: Writing jar file: /tmp/sqoop-student/complier/b7c41d0ff415ac919af44
12/06/01 08:39:47 INFO org.apache.hadoop.mapred.MapTask: jar file length = 103200
12/06/01 08:39:47 MARN manager.MySQLManager: It looks like you are importing from mysql.
12/06/01 08:39:47 MARN manager.MySQLManager: This transfer can be faster by using the --direct
12/06/01 08:39:47 MARN manager.MySQLManager: option. It will skip reading the specific fast path.
12/06/01 08:39:47 MARN manager.MySQLManager: Setting zero DATETIME delimiter to convertnull (mysql)
12/06/01 08:39:47 INFO MapReduceJobBuilder: Beginning Import of tb_msg
student@student-VirtualBox:~/home/WQD7007/sqoop/bin
File Edit View Search Terminal Help
1..46,714-659,31..0,78_Female
1/1/2013 23:18,151798177574,Female,,151,37.182,209,Android 7.8,4G,,91.319474,29,0333363,,1569889396931,,240,156,169,Andro,4G,,95,,154,36,34,Female
11/1/2013 23:05,185163259240,Female,,167,240,213,72,Android 6.0,5G,,116,767744,,27,0086064,,178091715068,,58,250,218,Andro,4G,,115,36,34,Female
11/1/2013 23:54,1850805904,Female,,161,217,36,98,Android 6.0,5G,,113,39623,22,371406,,1556662187980,,5,.95,218,73,Andro,4G,,115,36,34,Female
11/1/2013 23:31,151,185082374951,Female,,182,113,283,52,IOS 9.0,5G,,195.54408,30,181359,,145,169,35,94,Andro,4G,,115,36,34,Female
11/1/2013 23:05,185082374952,Female,,161,233,150,137,Android 7.0,4G,,102,585031,40,478257,,154474784712,,31,224,120,Andro,4G,,126,56,49,14,Male
11/1/2013 23:05,151,185082374953,Female,,52,281,123,198,Android 8.0,4G,,120,067791,31,650084,,18813075984,,185,81,20,Andro,4G,,115,36,34,Female
11/1/2013 23:44,185082332369,Female,,216,145,15,96,Android 7.0,4G,,103,353279,37,3840227,,154474784743,,203,52,96,Andro,4G,,108,665,36,72,Female
11/1/2013 23:54,185082332370,Female,,116,250,145,190,IOS 9.0,4G,,110,89419,36,36727,,13492067847,,203,154,72,74,Andro,4G,,115,36,34,Female
11/1/2013 23:59,1735191259214,Male,,250,160,225,135,IOS 12,0,8G,,197,159034,24,745376,,18126815943,,89,154,255,135,Andro,4G,,104,271,30,30,Female
11/1/2013 23:59,1735191259215,Female,,83,229,79,81,Android 6,4G,,114,058533,39,39896,,183071342521,,46,141,201,236,Andro,4G,,104,271,30,30,Female
11/1/2013 23:59,1735191259216,Female,,83,229,79,81,Android 6,4G,,114,058533,39,39896,,183071342521,,46,141,201,236,Andro,4G,,104,271,30,30,Female
11/1/2013 23:06,18809910851,Male,,105,186,22,255,Android 7.0,5G,,111,041376,23,3580227,,177297338251,,148,126,22,Andro,4G,,100,886,,33,39,Female
9/1/2013 23:06,18809910852,Male,,105,186,22,255,Android 7.0,5G,,111,041376,23,3580227,,177297338251,,148,126,22,Andro,4G,,100,886,,33,39,Female
4..Andro,4G,,83,9065,31,01,Male
student@student-VirtualBox:~/home/WQD7007/sqoop/bin$ ls
```

HIVE

For the second part, we will introduce the hive. hive help us to word with data.

I. Data format fix

After processed, we could the data format as below :



The screenshot shows a terminal window on an Ubuntu system. The title bar says "format.py". The code in the terminal is:

```
import pdB
fp = open("data1.csv", "r")
for line in fp.readlines():
    line = line.replace("\t", " ")
    con = line.split(",")
    con[6] += "," + con[7]
    con[7] = ","
    con[12] += "," + con[13]
    con[13] = ","
    print ("\"".join(con))
```

Below the code, the terminal displays a table of data with columns: Date, Time, ID, Gender, IP, Device Type, OS, RAM, CPU, Model, Latitude, Longitude, and a timestamp. The data is as follows:

Date	Time	ID	Gender	IP	Device Type	OS	RAM	CPU	Model	Latitude	Longitude	Timestamp	
2021/11/1	7:44	1825138366339	Male	195.188.222.255	iOS	9.0	5G		"123.257181,48.807394"			1523699980735	136.66
109.160	Android	7.0	4G	"89.332566,42.956664"					Female	TEXT	5.14KM		
2021/11/1	7:29	155102322496	Female	184.147.145.6	Android	7.0	4G		"113.39623,22.371406"			1593421890833	
186.218.176.9	Android	8.0	4G	"108.980879,24.071738"					Female	TEXT	78.22KM		
2021/11/1	7:57	1898238221524	Male	21.56.191.12	Android	6.0	4G		"117.885171,33.21035"			1331812911843	
225.104.57.224	Android	6	5G	"104.639117,37.781842"					Male	TEXT	98.86KM		
2021/11/1	7:26	1328514919221	Male	43.227.208.47	Android	8.0	4G		"121.785398,41.86541"			1305190543476	
30.139.30.154	Android	8.0	4G	"106.886068,33.395741"					Male	TEXT	12.59KM		
2021/11/1	7:01	1896394119953	Male	220.26.70.16	Android	8.0	4G		"124.361019,45.85657"			1859080711480	
247.25.68.26	Android	6.0	3G	"102.725798,37.430552"					Male	TEXT	21.2KM		
2021/11/1	7:23	1312348411125	Female	139.40.213.203	iOS	12.0	4G		"86.904123,45.547269"			1899214894921	
57.54.100.3	Android	8.0	4G	"132.676597,46.723902"					Male	TEXT	25.34KM		

II. Create DB and table in hive and prepare the data table in this process. Load data into hive using sqoop, we could see we successfully load data into hive

```
student@student-VirtualBox:~$ sqoop import -connect jdbc:mysql://localhost/db2_msg -username root -password fsktm -table tb_msg -hive-import --hive-database tmp --hive-table tb_msg -m 1
Warning: /home/WQD7007/sqoop/bin/../.., hcatalog does not exist! HCatalog jobs will fail.
Please set SHCAT_HOME to the root of your HCatalog installation.
Warning: /home/WQD7007/sqoop/bin/../.., accumulo does not exist! Accumulo imports will fail.
Please set SACCUMULO_HOME to the root of your Accumulo installation.
22/06/04 11:09:05 INFO sqoop.Sqoop: Running Sqoop version: 1.4.7
22/06/04 11:09:05 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
22/06/04 11:09:05 INFO tool.BaseSqoopTool: Using Hive-specific delimiters for output. You can override
22/06/04 11:09:05 INFO tool.BaseSqoopTool: delimiters with --fields-terminated-by, etc.
22/06/04 11:09:06 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
22/06/04 11:09:06 INFO tool.CodeGenTool: Beginning code generation
Sat Jun 04 11:09:07 MYT 2022 WARN: Establishing SSL connection without server's identity verification is not recommended.
According to MySQL 5.5.45+, 5.6.26+ and 5.7.6+ requirements SSL connection must be established by default if explicit option isn't set. For compliance with existing applications not using SSL the verifyServerCertificate property is set to 'false'. You need either to explicitly disable SSL by setting useSSL=false, or set useSSL=true and provide truststore for s
```

III. Data problem

A. Empty data

Problem 1: in the current data, some data fields are empty and not legal data

B. Time format

Problem 2: in the demand, the daily and hourly message volume needs to be counted, but there are no day and hour fields in the data, only the overall time field, which is difficult to handle

C. GPS location problem

Problem 3: in the demand, it is necessary to build a visual map of the region for longitude and dimension, but the GPS longitude and latitude in the data is a field, which is difficult to handle

IV. Data cleaning

Requirement 1: filter illegal data with empty fields

• where filtering

2 requirement 2: build day and hour fields through time fields

• substr function

2 demand 3: extract longitude and dimension from longitude and latitude of GPS

• split function

2 requirement 4: save the results after ETL to a new hive table

Create another table to store the cleaning data

V. Show the result after data cleaning

```
Time taken: 00:00:03 seconds
hive> Select msg_time, dayinfo, hourinfo, sender_gps, sender_lng, sender_lat from tb_msg_etl limit 10;
OK
2021/11/1 7:44 2021/11/1      7      "123.257181,48.807394" 123.257181    48.807394
2021/11/1 7:29 2021/11/1      7      "113.39623,22.371406" 113.39623   22.371406
2021/11/1 7:57 2021/11/1      7      "117.885171,33.21035" 117.885171  33.21035
2021/11/1 7:26 2021/11/1      7      "121.785398,41.86541" 121.785398  41.86541
2021/11/1 7:01 2021/11/1      7      "124.361019,45.856657" 124.361019  45.856657
2021/11/1 7:23 2021/11/1      7      "86.904123,45.547269" 86.904123   45.547269
2021/11/1 7:57 2021/11/1      7      "119.651311,33.519113" 119.651311  33.519113
2021/11/1 7:42 2021/11/1      7      "114.058533,39.39896" 114.058533  39.39896
2021/11/1 7:17 2021/11/1      7      "117.07569,34.010827" 117.07569   34.010827
2021/11/1 7:15 2021/11/1      7      "112.439571,33.148464" 112.439571  33.148464
Time taken: 0.191 seconds, Fetched: 10 row(s)
```

VI. Using SQL to do the EDA

1. search the sender numbers about the sex, telephone system type and network.

We found most people use android.

2. search the numbers about the sex different between female and male

We found male would like to send messages rather than females, and the receivers are males too.

hbase

I. View the dataset.

The screenshot shows a Windows desktop environment. In the top-left, there's a file browser window titled '我的电脑' (My Computer) showing various local drives and network locations. In the center, a '打开文件' (Open File) dialog box is open, prompting for a CSV file to import. It lists several files on the desktop, including 'data1.csv'. In the bottom-right, a MySQL Workbench window is open, showing the 'test' database. A table named 'bookstore' is selected, and its data is displayed in a grid format. The data consists of 22 rows, each containing 13 columns of NULL values. The columns are labeled 'a' through 'o'.

II. Created a new table named hbase_test and a column family called info.

```
hbase(main):006:0> list
TABLE
users
1 row(s)
Took 0.0134 seconds
=> ["users"]
hbase(main):007:0> create 'hbase_test', 'info'
Created table hbase_test
Took 2.3067 seconds
```

III. Granted the permission to hbase_test.

```
=> Hbase::Table - hbase_test
hbase(main):008:0> grant 'root','RWCA','hbase_test';
hbase(main):009:0* user_permission 'hbase_test'
Took 0.1255 seconds
User           Namespace,Table,Family,Qualifier:Permission
atguigu      default,hbase_test,:, [Permission: actions=READ,WRITE,EXEC,CREATE,ADMIN]
root          default,hbase_test,:, [Permission: actions=READ,WRITE,CREATE,ADMIN]
2 row(s)
Took 0.0531 seconds
```

IV. Run the sqoop. Load the data into sqoop

```
[atguigu@hadoop102 sqoop]\$ bin/sqoop import --connect jdbc:mysql://hadoop102:3306/bookstore \
> --username 'root' --password '123321' \
> --split-by a \
> --hbase-table 'hbase_test' \
> --hbase-row-key 'a,b,c,d,e,f,g,h,i,j,k,l,m,n,o' \
> --create-table-if-not-exists \
> /opt/module/sqoop/bin/sqoop =====
> /opt/module/sqoop/bin
=====
> /opt/module/sqoop/bin
Warning: /opt/module/sqoop/bin/.../.HCatalog does not exist! HCatalog jobs will fail.
Please set $HCAT_HOME to the root of your HCatalog installation.
Warning: /opt/module/sqoop/bin/.../.accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
SLF4J: class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/module/hbase/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/module/hadoop-3.1.0/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory].
2022-06-14 12:22:52,161 INFO sqoop: Running Sqoop version: 1.4.6
2022-06-14 12:22:52,173 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
--table or --query is required for import. (Or use sqoop import-all-tables.)
```

V. Import dataset

```
[1] hadoop102          [2] hadoop102          [3] hadoop103          [4] hadoop104          [5]
2022-06-14 12:36:27,942 INFO [main] mapreduce.Job: The url to track the job: http://hadoop103:8080/proxy/application_1655169413131_0003/
2022-06-14 12:36:27,942 INFO [main] mapreduce.Job: Running: Job: job_1655169413131_0003
2022-06-14 12:36:36,099 INFO [main] mapreduce.Job: Job: job_1655169413131_0003 running in uber mode : false
2022-06-14 12:36:36,100 INFO [main] mapreduce.Job: map 0% reduce 0%
2022-06-14 12:36:36,100 INFO [main] mapreduce.Job: map 100% reduce 0%
2022-06-14 12:36:46,398 INFO [main] mapreduce.Job: Job job_1655169413131_0003 completed successfully
2022-06-14 12:36:46,506 INFO [main] mapreduce.Job: rations=1
HDFS: Number of large read operations=0
HDFS: Number of write operations=0
Job Counters
Launched map tasks=1
  Rack-local map tasks=1
Total time spent by all maps in occupied slots (ms)=13668
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=6834
Total map-milliseconds taken by all map tasks=6834
Total map-reduce-milliseconds taken by all map tasks=698016
Map-Reduce Framework
  Map Input records=71560
  Map Output records=704
  Input split bytes=939
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC Time (ms)=166
  CPU time spent (ms)=990
  Physical memory (bytes) snapshot=324119688
  Virtual memory (bytes) snapshot=2601510794
  Total committed heap usage (bytes)=201326592
  Peak Map Physical memory (bytes)=324119688
  Peak Map Virtual memory (bytes)=2668150784
Imported
  Bad Lines=70856
File Input Format Counters
  Bytes Read=12340265
File Output Format Counters
  Bytes Written=
```

VI. Validate the data

```

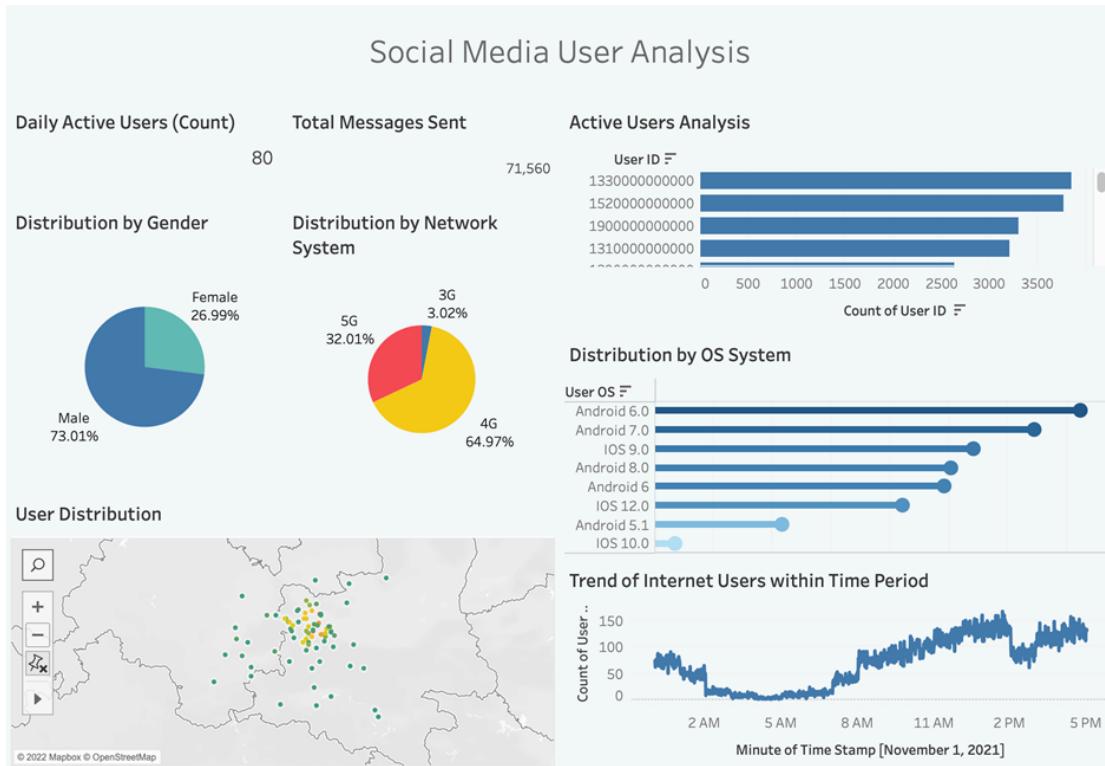
2021/11/1 8:48
2021/11/1 8:49
2021/11/1 8:52
2021/11/1 8:53
2021/11/1 8:55
2021/11/1 8:59
2021/11/1 9:01
2021/11/1 9:02
2021/11/1 9:03
2021/11/1 9:04
2021/11/1 9:05
2021/11/1 9:06
2021/11/1 9:08
2021/11/1 9:11
2021/11/1 9:12
2021/11/1 9:13
2021/11/1 9:18
2021/11/1 9:21
2021/11/1 9:22
2021/11/1 9:24
2021/11/1 9:27
2021/11/1 9:28
2021/11/1 9:30
2021/11/1 9:31
2021/11/1 9:34
2021/11/1 9:37
2021/11/1 9:38
2021/11/1 9:39
2021/11/1 9:40
2021/11/1 9:41
2021/11/1 9:43
2021/11/1 9:44
2021/11/1 9:45
2021/11/1 9:47
2021/11/1 9:49
2021/11/1 9:51
2021/11/1 9:52
2021/11/1 9:54
2021/11/1 9:57
2021/11/1 9:58
2021/11/1 9:59

452 rows(s)

Took 0.3801 seconds
hbase(main:013:0>
```

Data Visualization

Data visualization was performed using Tableau. The tool was chosen because Tableau allows users to connect Apache Hadoop as a datasource through Cloudera's driver directly. This allowed us to complete analytics on top of massive data rapidly. We generated 6 graphs and 2 summary cards to visualize the overall performance of Momo using our dataset. We tried to conclude some key performance indicators that could assist the app management team to better keep track of their users and further instruct operation decisions.



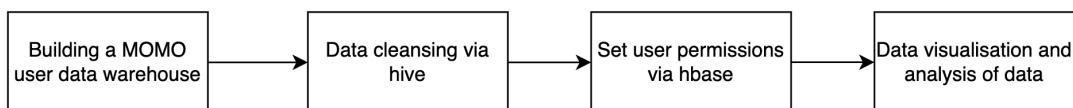
Results & discussion

No.	Data description	Description of the conclusions
1	Male: 73% Female: 27%	MOMO users are predominantly male
2	3G: 3.02% 4G: 64.97% 5G: 32%	MOMO users are predominantly on 4G networks
3	Number of Android users: 73%	MOMO users are mainly on Android OS System
4	Most active time: 13pm	Users' internet access is concentrated in the 13pm time slot

With the final visual analysis, it is easy to see that we have a total of 80 users who sent 71,560 messages on November 1, 2021. Over 73% of users are male and approximately 27% are female.

4G is also by far the most widely distributed, with nearly 73.01% of males online, which is a relatively large percentage. In addition, users of the Stranger platform are predominantly Android OS users.

The three most active users sent over 3,000 messages on 1 November 2021. In addition, the data shows that platform users are mainly located in the western province of Shaanxi, with internet access time slots concentrated around 1pm, with an upward trend from morning to 2pm, then gradually declining and rising.



In this assignment, we first built the momo user data warehouse through hadoop to manage momo's social media data. Then we cleaned it with hive to remove the dirty data. Then we set up data user permissions via hbase for better access and management of social media data. Finally, the momo user data was analysed through visualization of the data.

In addition, in using Hadoop, we found that the Hadoop ecosystem provides various open source tools for collecting, storing and processing data, as well as cluster deployment, monitoring and data security. Also, this research on momo user behavior data has benefited from the three components of Hadoop which are HDFS, MapReduce and HBase. and, we found that it can be applied to offline computing and mass data storage in addition to being well suited for analysing mass data. We will continue to explore the advantages and application scenarios of hadoop in future projects to really bring out its value.

Conclusions

Reference.

1. Edlich, S., & Vogler, M. (2013). Smart apps for applied machine learning on mobile devices: the MOMO project. *Https://Doi.Org/10.1117/12.2005173*, 8667, 224–233. <https://doi.org/10.1117/12.2005173>
2. Xu, D., & Wu, F. (2019). Exploring the cosmopolitanism in China: examining mosheng ren (“the stranger”) communication through Momo. *Https://Doi-Org.Ezproxy.Um.Edu.My/10.1080/15295036.2019.1566629*, 36(2), 122–139. <https://doi.org/10.1080/15295036.2019.1566629>
3. Chan, L. S. (2019). Paradoxical Associations of Masculine Ideology and Casual Sex Among Heterosexual Male Geosocial Networking App Users in China. *Sex Roles 2019 81:7*, 81(7), 456–466. <https://doi.org/10.1007/S11199-019-1002-4>
4. Li, H. (2020). Negotiating Privacy and Mobile Socializing: Chinese University Students’ Concerns and Strategies for Using Geosocial Networking Applications: *Https://Doi-Org.Ezproxy.Um.Edu.My/10.1177/2056305120913887*, 6(1). <https://doi.org/10.1177/2056305120913887>