

## **Problem Setting**

Smartphones are undeniably an integral part of everyone's daily life, with additional capabilities being added to them frequently their importance in our lives only solidifies further. It is crucial to analyze their influence in all sectors of business so as understand how it can used to their advantage. The analysis for the purpose of this project will be centered around how smartphone usage affects the tourism industry, specifically how they affect the behavior of international tourists. Usage can be grouped into several forms such as social media, internet access, mobile operator used, etc. The goal of the analysis is the understand the extent of impact that smartphone usage has on tourists.

## **Data Source & Description**

The datasets that are used for this project are the following:

1. Questionnaire Data: Rusdi, Jack Febrian (2019), "Smartphone usage and International Tourist Behaviour", Mendeley Data, V1, doi: 10.17632/zwzb8hzc9j.1

The data in this data set was collected through the means of a questionnaires filled in by 302 participants from 52 countries who had travelled to the city of Bandung, Indonesia.

2. Trip Advisor Bandung, Indonesia "Things to Do" page: A web scraper tool was used to collect data about the 30 most popular things to do in the city of Bandung.

3. Trip Advisor Bandung Indonesia "Popular Hotels" page: A web scraper tool was used to collect data about the 30 most popular hotels in the city of Bandung.

Trip Advisor data was used as it was a source of information that the tourists looked at before making decisions about where they would stay and what they would do during their travel.

The questionnaire data is in .csv format where as the trip advisor data is in .json format.

The data can he found [here](#).

## **Problem Definition**

The raw data provides information about how tourists behave during the trip with respect to internet and smartphone feature usage as well as their behavior during the planning of the trip, specifically in relation to Trip Advisor reviews that may impact their decision. Dimensions derived from the data include:

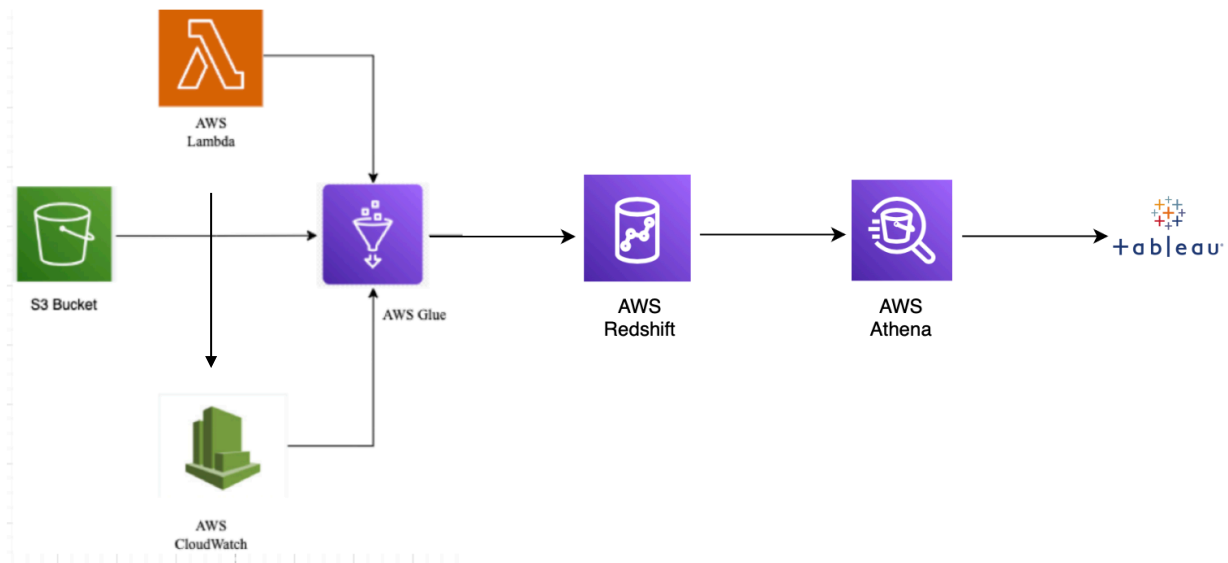
- stat (**stat\_id**, operator, number\_of\_smartphone\_functions, number\_of\_locations, purpose)
- social\_media (**scomed\_id**, fb\_usage, ig\_usage, twitter\_usage, youtube\_usage)
- person (**person\_id**, age, gender, education, country, *stat\_id*, *scomed\_id*)
- hotel (**hotel\_id**, name, no\_of\_ratings, rating)

- place (**place\_id**, tourist\_attraction, type, no\_of\_ratings, rating)
- behavior (**person\_id**, **hotel\_id**, **place\_id**, internet\_daily\_usage, ratings\_per\_review\_hotel, ratings\_per\_review\_place)

It is known that tourism can have a momentous effect on a country's economy and gaining insight into how it this be leveraged using smartphones as well as provide the tourist an enjoyable experience is ultimately a mutually beneficial endeavor.

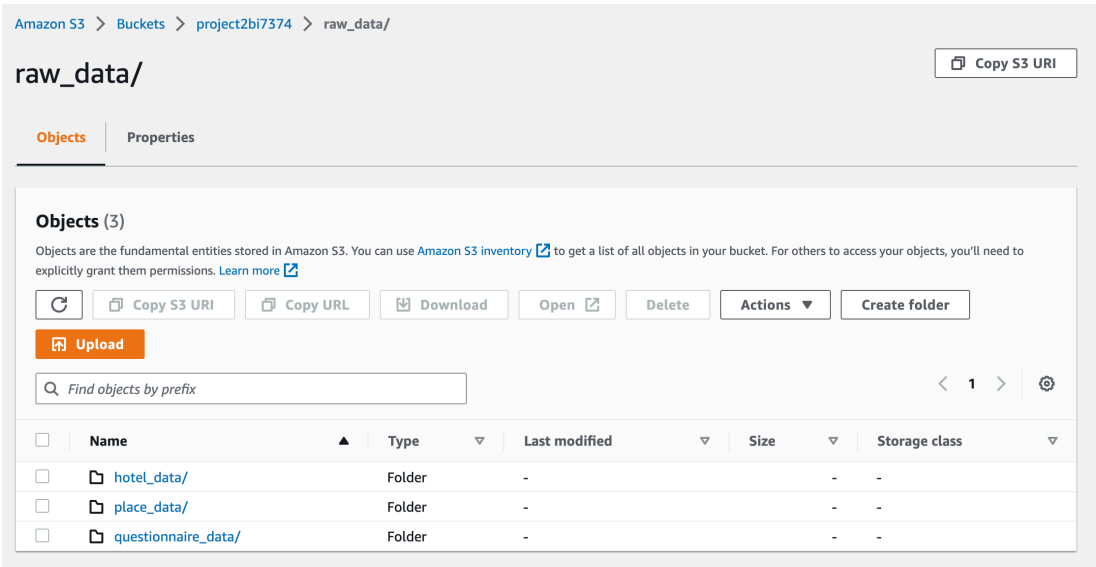
The goal is to analyze the relationship between smartphone usage and international tourist behavior in the city of Bandung, Indonesia, and to understand the extent of the impact that smartphone usage has on tourists. This will be achieved by analyzing the questionnaire data and Trip Advisor data to gain insights into how tourists use their smartphones during their trip, how they use social media, how they plan their trip using online resources, and how they rate hotels and tourist attractions.

## Architecture Diagram



## Storage

The raw data sources (questionnaire.csv, hotel.json, place.json) were stored in an S3 bucket:



## Ingestion

The source crawlers source\_hotel, source\_place, source\_questionnaire are created to get the schema of the source data, target crawlers are created to get the schema of the warehouse tables.

### Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

[User preferences](#)

Add crawler Run crawler Action  Showing: 1 - 9 Refresh Help

<input type="checkbox"/>	Name	Schedule	Status	Logs	Last runtime	Median runtime	Tables updated	Tables added
<input type="checkbox"/>	source_hotel		Ready	<a href="#">Logs</a>	1 min	59 secs	0	1
<input type="checkbox"/>	source_place		Ready	<a href="#">Logs</a>	49 secs	47 secs	0	1
<input type="checkbox"/>	source_questionnaire		Ready	<a href="#">Logs</a>	42 secs	50 secs	0	2
<input type="checkbox"/>	tag_hotel		Ready	<a href="#">Logs</a>	1 min	1 min	0	1
<input type="checkbox"/>	tgt_behavior		Ready	<a href="#">Logs</a>	2 mins	2 mins	0	1
<input type="checkbox"/>	tgt_person		Ready	<a href="#">Logs</a>	2 mins	2 mins	0	1
<input type="checkbox"/>	tgt_place		Ready	<a href="#">Logs</a>	2 mins	1 min	0	1
<input type="checkbox"/>	tgt_scomed		Ready	<a href="#">Logs</a>	2 mins	3 mins	0	1
<input type="checkbox"/>	tgt_stat		Ready	<a href="#">Logs</a>	3 mins	2 mins	0	1

The source crawler names are used to create an event in CloudWatch to collect and log metric

data  
the

from

Project2ETL

Edit

Disable

Delete

CloudFormation Template ▼

Rule details

Rule name	Status	Event bus name	Type
Project2ETL	Enabled	default	Standard
Description	Rule ARN	Event bus ARN	
	arn:aws:events:us-east-1:28558:0274452:rule/Project2ETL	arn:aws:events:us-east-1:28558:0274452:event-bus/default	

Event pattern

Targets

Monitoring

Tags

Event pattern

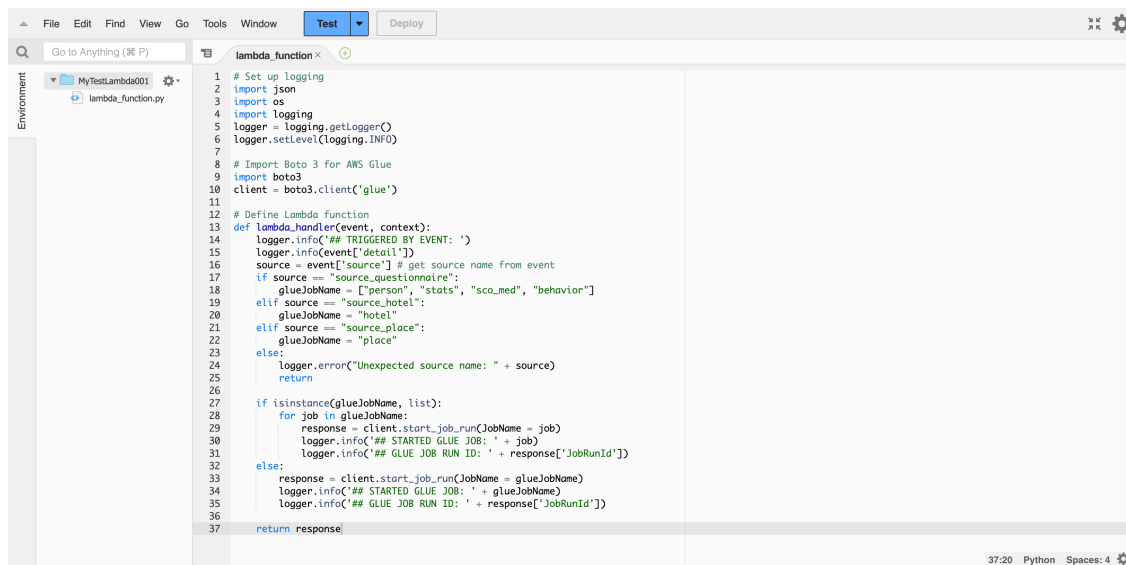
Edit

```
1 {
2   "detail-type": ["Glue Crawler State Change"],
3   "source": ["aws.glue"],
4   "detail": {
5     "crawlerName": ["source_questionnaire", "source_hotel", "source_place"],
6     "state": ["Succeeded"]
7   }
8 }
```

Copy

different sources.

The CloudWatch event triggers an AWS Lambda function

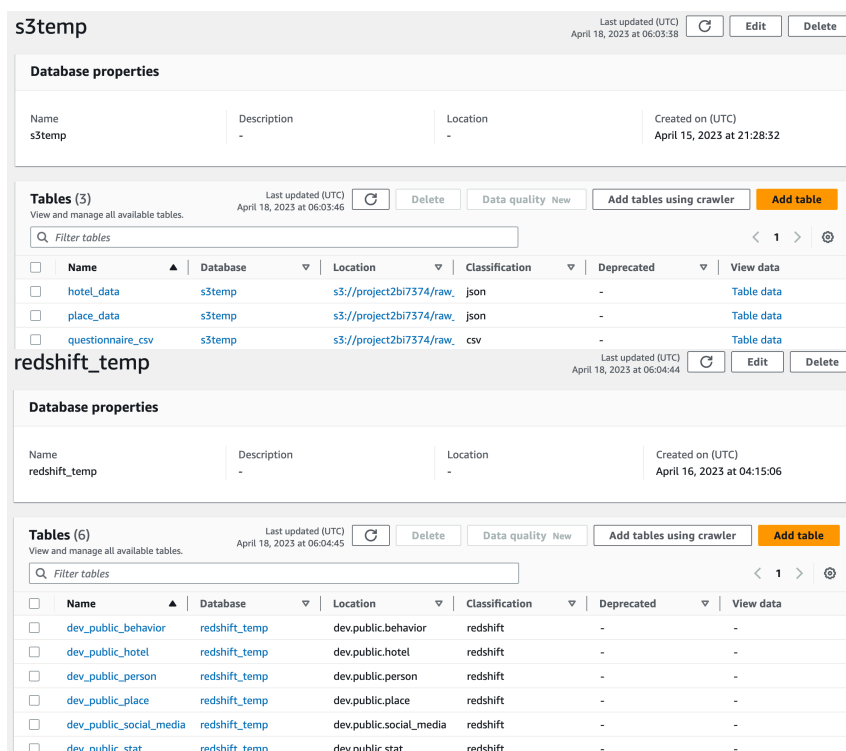


```
1 # Set up logging
2 import json
3 import os
4 import logging
5 logger = logging.getLogger()
6 logger.setLevel(logging.INFO)
7
8 # Import Boto 3 for AWS Glue
9 import boto3
10 client = boto3.client('glue')
11
12 # Define Lambda function
13 def lambda_handler(event, context):
14     logger.info("## TRIGGERED BY EVENT: ")
15     logger.info(event['detail'])
16     source = event['source'] # get source name from event
17     if source == "source_questionnaire":
18         glueJobName = ["person", "stats", "sco_med", "behavior"]
19     elif source == "source_hotel":
20         glueJobName = "hotel"
21     elif source == "source_place":
22         glueJobName = "place"
23     else:
24         logger.error("Unexpected source name: " + source)
25         return
26
27     if isinstance(glueJobName, list):
28         for job in glueJobName:
29             response = client.start_job_run(JobName = job)
30             logger.info("## STARTED GLUE JOB: " + job)
31             logger.info("## GLUE JOB RUN ID: " + response['JobRunId'])
32     else:
33         response = client.start_job_run(JobName = glueJobName)
34         logger.info("## STARTED GLUE JOB: " + glueJobName)
35         logger.info("## GLUE JOB RUN ID: " + response['JobRunId'])
36
37     return response
```

This lambda function will invoke the jobs “person”, “stats”, “sco\_med”, and “behavior” (jobs whose data source is the questionnaire data file) if source\_questionnaire crawler is run as it implies that changes have been made to the source questionnaire data. Similarly, changes to hotel data triggers the job “hotel” and changes to place data triggers to job “place”

## ETL Jobs

To load data data source crawlers are get the the source the required warehouse created in Redshift



s3temp			
Database properties			
Name	Description	Location	Created on (UTC)
s3temp	-	-	April 15, 2023 at 21:28:32

Tables (3)						
View and manage all available tables.						
Q Filter tables						
<input type="checkbox"/>	Name	Database	Location	Classification	Deprecated	View data
<input type="checkbox"/>	hotel_data	s3temp	s3://project2bi7374/raw	json	-	<a href="#">Table data</a>
<input type="checkbox"/>	place_data	s3temp	s3://project2bi7374/raw	json	-	<a href="#">Table data</a>
<input type="checkbox"/>	questionnaire_csv	s3temp	s3://project2bi7374/raw	csv	-	<a href="#">Table data</a>

redshift_temp			
Database properties			
Name	Description	Location	Created on (UTC)
redshift_temp	-	-	April 16, 2023 at 04:15:06


  

Tables (6)						
View and manage all available tables.						
Q Filter tables						
<input type="checkbox"/>	Name	Database	Location	Classification	Deprecated	View data
<input type="checkbox"/>	dev_public_behavior	redshift_temp	dev.public.behavior	redshift	-	-
<input type="checkbox"/>	dev_public_hotel	redshift_temp	dev.public.hotel	redshift	-	-
<input type="checkbox"/>	dev_public_person	redshift_temp	dev.public.person	redshift	-	-
<input type="checkbox"/>	dev_public_place	redshift_temp	dev.public.place	redshift	-	-
<input type="checkbox"/>	dev_public_social_media	redshift_temp	dev.public.social_media	redshift	-	-
<input type="checkbox"/>	dev_public_stat	redshift_temp	dev.public.stat	redshift	-	-

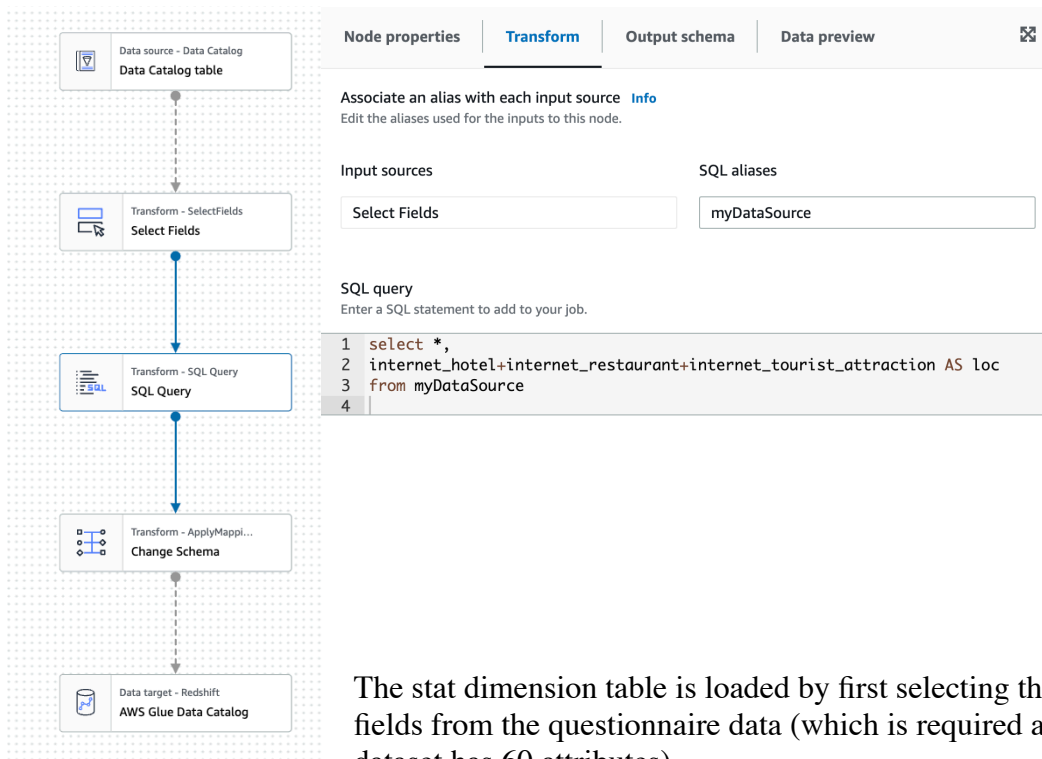
from the the source run so as the schema of data, then data tables are AWS which is

then crawled as well. The S3 bucket data is stored in a AWS Glue database “s3temp” and the Redshift table schema is store in another Glue database “redshift\_temp”

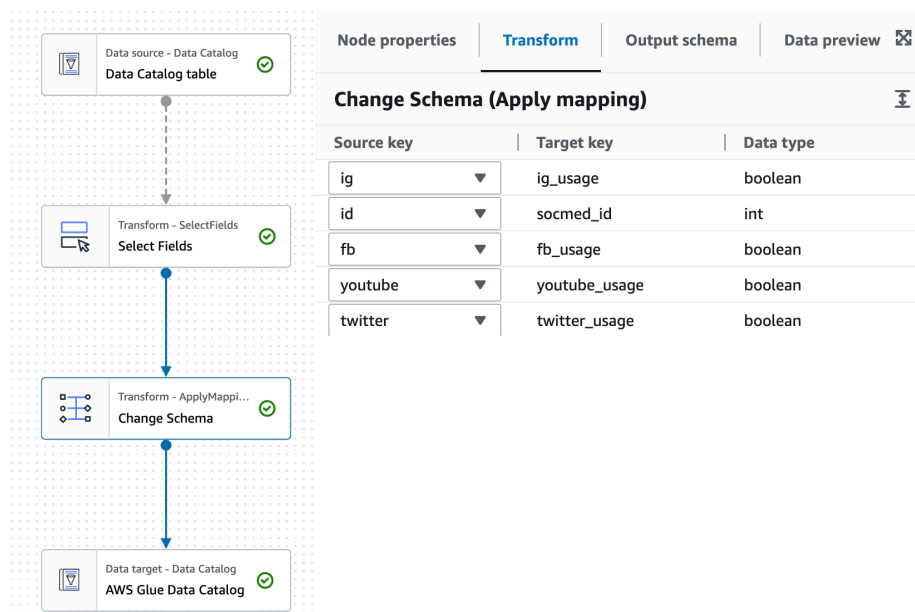
We use the visual job editor to configure nodes. We are reading data from source location, then applying a transform to the data and loading the data for AWS Redshift. The following ETL jobs were used to load data into the warehouse tables:

Your jobs (6) <a href="#">Info</a>						Actions ▾	Run job
<input type="text" value="Filter jobs"/>					< 1 > ⚙		
<input type="checkbox"/>	Job name ▾	Type	Last modified ▾	AWS Glue version ▾			
<input type="checkbox"/>	<a href="#">behavior</a>	Glue ETL	4/17/2023, 2:37:53 PM	3.0			
<input type="checkbox"/>	<a href="#">place</a>	Glue ETL	4/17/2023, 2:30:55 PM	3.0			
<input type="checkbox"/>	<a href="#">hotel</a>	Glue ETL	4/17/2023, 2:29:03 PM	3.0			
<input type="checkbox"/>	<a href="#">person</a>	Glue ETL	4/17/2023, 2:20:51 PM	3.0			
<input type="checkbox"/>	<a href="#">sco_med</a>	Glue ETL	4/17/2023, 2:08:26 PM	3.0			
<input type="checkbox"/>	<a href="#">stats</a>	Glue ETL	4/17/2023, 1:57:58 PM	3.0			


stat:



sco\_med:



The social\_media table is loaded the schema to the appropriate and data type in the warehouse.

Node properties	Transform	Output schema	Data preview 
<input type="checkbox"/>	sp_vc		double
<input type="checkbox"/>	sp_tl		string
<input type="checkbox"/>	sp_cc		double
<input type="checkbox"/>	sp_rn		double
<input type="checkbox"/>	sp_smp		double
<input type="checkbox"/>	sp_sp		double
<input type="checkbox"/>	sp_ob		double
<input type="checkbox"/>	sp_wa		double
<input type="checkbox"/>	phone_functions		double
<input type="checkbox"/>	sp_atg		string
<input type="checkbox"/>	place		string
<input type="checkbox"/>	hotel		string
<input checked="" type="checkbox"/>	youtube		double
<input checked="" type="checkbox"/>	twitter		double
<input checked="" type="checkbox"/>	fb		double
<input checked="" type="checkbox"/>	ig		double
<input checked="" type="checkbox"/>	id		long

by mapping columns


person:

For the dimension person, we are reading the date, then transform it by mapping the data and performing two join to combine multiple tables. We then drop the fields that are not required, mapping the changed schema and then loading the data. The person table contains a hierarchy to the stat and social\_media tables, hence it is joined with the respective tables from redshift so as to match each person to a stat\_id and scomed\_id

Node properties	Transform	Output schema	Data preview 
-----------------	-----------	---------------	--

Join type

Select the type of join to perform.

 Inner join

Select all rows from both datasets that meet the join condition.

Join conditions

Select a field from each parent node for the join condition.


ApplyMapping

stat\_table







id

=

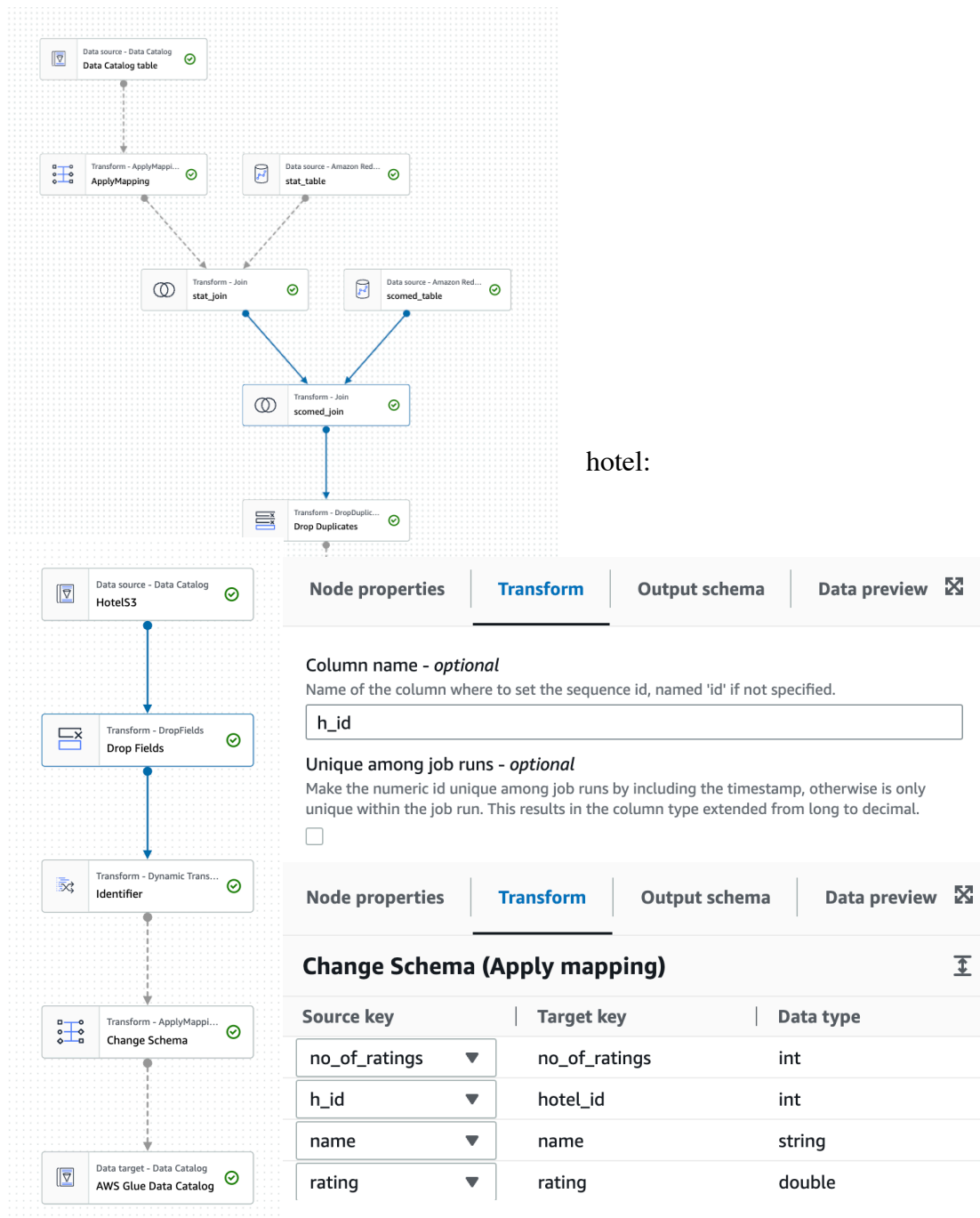
stat\_id



Add condition

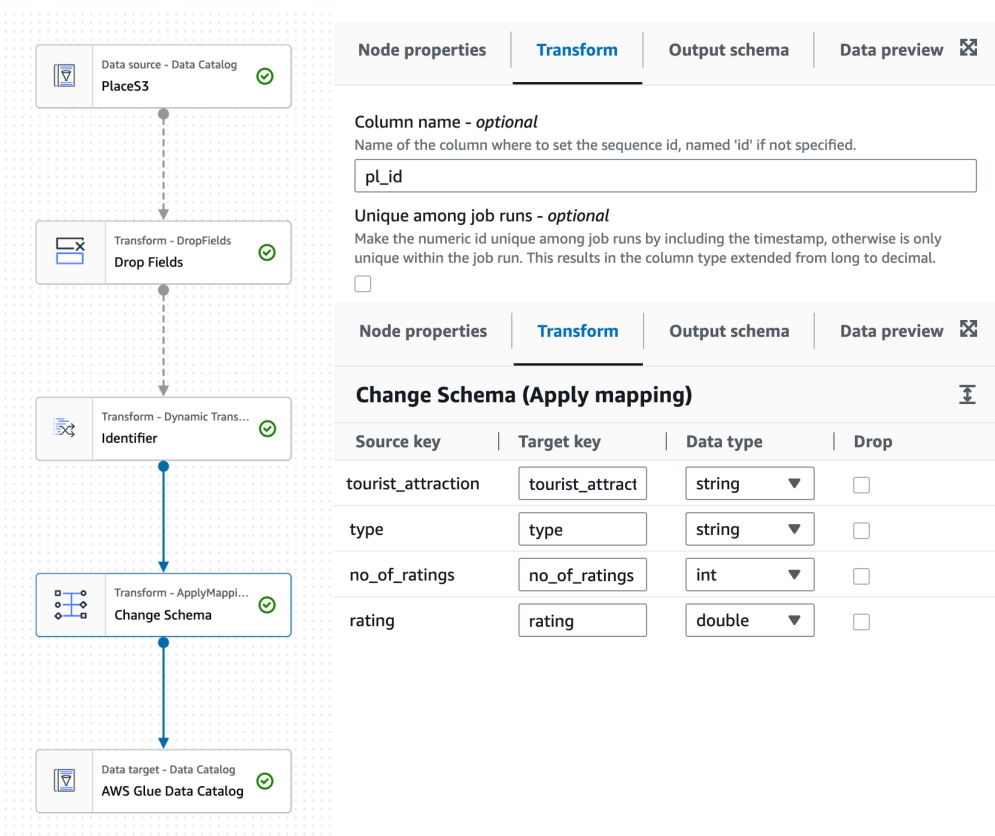
Node properties	Transform	Output schema	Data preview 
Change Schema (Apply mapping) 			
Source key	Target key	Data type	Drop
hotel_id	<input type="text" value="hotel_id"/>	int 	<input type="checkbox"/>
name	<input type="text" value="name"/>	string 	<input type="checkbox"/>
no_of_ratings	<input type="text" value="no_of_ratings_hotel"/>	int 	<input type="checkbox"/>
rating	<input type="text" value="rating_hotel"/>	double 	<input type="checkbox"/>





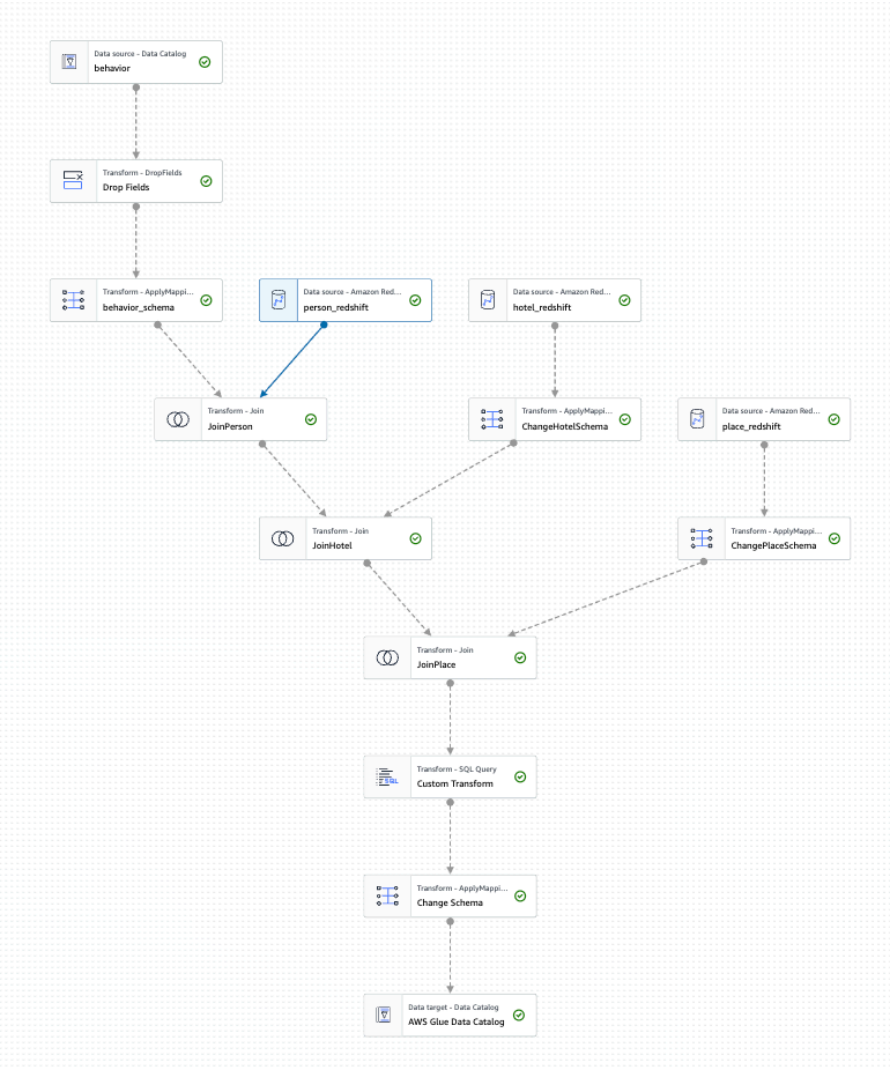
Hotel data is added to the redshift table first by dropping unnecessary field and then adding a unique identifier to each hotel using the transform 'Identity' then applying mapping.

place:



Place data is loaded in a similar way to hotel.

behavior:



Node properties

Transform

Output schema

Data preview

Change Schema (Apply mapping)

Source key	Target key	Data type
internet_daily_usage	internet_daily_usage	double
rpr_hotel	rating_per_review_hotel	double
hotel_id	hotel_id	int
rpr_place	rating_per_review_place	double
place_id	place_id	int
person_id	person_id	int

SQL query

Enter a SQL statement to add to your job.

```
1 select *,
2     round(rating_hotel / no_of_ratings_hotel, 4) AS rpr_hotel,
3     round(rating_place / no_of_ratings_place, 4) AS rpr_place
4 from myDataSource
```

Change Schema (Apply mapping)

Source key	Target key	Data type	Drop
hotel_id	hotel_id	int	<input type="checkbox"/>
name	name	string	<input type="checkbox"/>
no_of_ratings	no_of_ratings_hotel	int	<input type="checkbox"/>
rating	rating_hotel	double	<input type="checkbox"/>

Node properties

Transform

Output schema

Data preview

Join type

Select the type of join to perform.

Inner join

Select all rows from both datasets that meet the join condition.

Join conditions

Select a field from each parent node for the join condition.

JoinPerson

ChangeHotelSchema

hotel

=

name

Add condition

For the fact table behavior, data from Hotel, Place, and Person is read and an SQL statement is created to convert the data into columns “reviews\_per\_rating\_hotel”/ “reviews\_per\_rating\_place”

Data in their respective tables in Redshift:

Rows returned (30)

Export

Q Search rows

< 1 2 3 > ⓘ

hotel_id	name	no_of_ratings	rating
0	TheGaiaHotelBandung	632	5
1	ThePapandayan	3194	5
2	GrandTjokroPremiereBandung	4156	4.5
3	POPHotelFestivalCitylink	827	4
4	ibisStylesBandungGrandCentral	773	5
5	SwissBelresortDagoHeritage	1044	4.5
6	JanevallaBandung	1206	4.5
7	BelviuHotel	1929	5
8	TheTransLuxuryHotel	2173	4.5
9	TheNaripanHotel	1117	5

Rows returned (30)

Q Search rows

place_id	tourist_attraction	type	no_of_ratings	rating
0	BragaStreet	Neighborhoods	1936	4
1	GeologyMuseum	ScienceMuseums	799	4
2	TransStudioBandung	AmusementampThemeParks	2189	4
3	ParisVanJava	ShoppingMalls	1283	4
4	GedungSate	ArchitecturalBuildings	527	4
5	NuArtSculpturePark	ArtGalleriesParks	227	4.5
6	MuseumofTheAsianAfricanConference	SpecialityMuseums	428	4
7	DusunBambuFamilyLeisurePark	PointsofInterestampLandmarks	792	4
8	BandungGrandMosque	PointsofInterestampLandmarksArchitecturalBuildings	616	4
9	RumahModeFactoryOutlet	FactoryOutlets	1271	4

Rows returned (302)

Q Search rows

< 1 2 3 4 5 6 7 ...

stat_id	operator	no_of_smartphone_func	no_of_locations	purpose
1	local	1	0	
2	local	9	3	work
3	local	10	3	
4	local	10	2	
5	local	9	2	
6	local	4	2	
7	local	10	1	
8	local	4	3	
9	local	8	2	
10	local	7	1	

Rows returned (302)

Q Search rows

socmed_id	fb_usage	ig_usage	twitter_usage	youtube_usage
1		true		true
2	true	true	true	
3	true	true	true	
4	true	true		
5	true	true		true
6	true	true		
7	true			
8	true			
9	true	true		
10	true			

Rows returned (302)

Q Search rows

< 1 2 3 4 5 6 7

person_id	age	gender	education	country	socmed_id	stat_id
18	34	F	D	Malaysia	18	18
38	60	M	PG	Netherlands	38	38
46	23	F	PG	Kazakhstan	46	46
73	34	F	ETC	Singapore	73	73
172	25	F	D	Australia	172	172
186	27	F	D	US	186	186
263	67	M	D	Mauritius	263	263
282	43	F	D	Malaysia	282	282
28	28	M	L	Australia	28	28
58	28	F	D	UK	58	58

Rows returned (262)

Q Search rows

< 1 2 3 4 5 6 7 ... 27 >

Export

person_id	place_id	hotel_id	rating_per_review_hotel	rating_per_review_place	internet_daily_usage
241	18	15	0.0079	0.0205	11
302	18	14	0.0101	0.0205	
74	18	6	0.0037	0.0205	4
192	18	4	0.0065	0.0205	5
153	18	4	0.0065	0.0205	3
254	18	20	0.125	0.0205	4
10	18	19	0.0367	0.0205	10
223	18	0	0.0079	0.0205	8
29	27	26	2.5	0.016	8
59	27	29	0.04	0.016	3

## **Data Visualization**

Analysis for the purpose of this project can be divided into two criteria:

### **I. Trip Information**

- Top hotels to stay and top place to go to in the city

### **II. Tourist Information while on the trip**

- Average number of smartphone functions used by a tourist
- Number of tourists who use social media during their travels
- Internet usage at places they travel to

#### **1. Average Number of Smartphone functions used by the tourists:**

```
SELECT AVG(no_of_smartphone_func) AS avg_phone_functions_used
FROM stat;
```

#### **2. Total number of international tourists who use Facebook, Twitter, YouTube, and Instagram during their trip:**

```
SELECT
    SUM(CASE WHEN fb_usage = true THEN 1 ELSE 0 END) AS fb_users,
    SUM(CASE WHEN twitter_usage = true THEN 1 ELSE 0 END) AS twitter_users,
    SUM(CASE WHEN youtube_usage = true THEN 1 ELSE 0 END) AS youtube_users,
    SUM(CASE WHEN ig_usage = true THEN 1 ELSE 0 END) AS ig_users
FROM social_media;
```

#### **3. Average internet usage of tourists ordered by the type of place they visited:**

```
SELECT place.type, ROUND(AVG(internet_daily_usage), 2) AS avg_usage
FROM behavior
INNER JOIN place ON behavior.place_id = place.place_id
GROUP BY place.type
ORDER BY avg_usage DESC;
```

4. Top 5 hotels to stay at based on number of ratings and of number of reviews

```
SELECT name, rating, no_of_ratings
FROM hotel
ORDER BY rating DESC, no_of_ratings DESC
LIMIT 5;
```

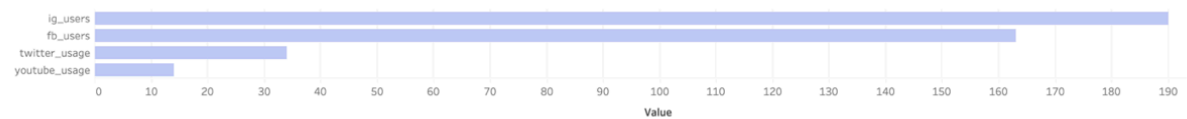
5. Top 5 places to visit based on number of ratings and of number of reviews

```
SELECT tourist_attraction, type, rating, no_of_ratings
FROM place
ORDER BY rating DESC, no_of_ratings DESC
LIMIT 5;
```

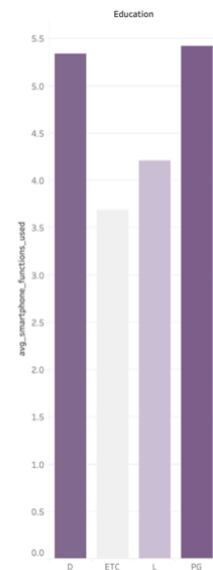
6. Number of tourists who visited a particular type of tourist attraction

```
SELECT place.type, COUNT(DISTINCT behavior.person_id) AS count
FROM place
INNER JOIN behavior ON place.place_id = behavior.place_id
GROUP BY place.type
ORDER BY count DESC;
```

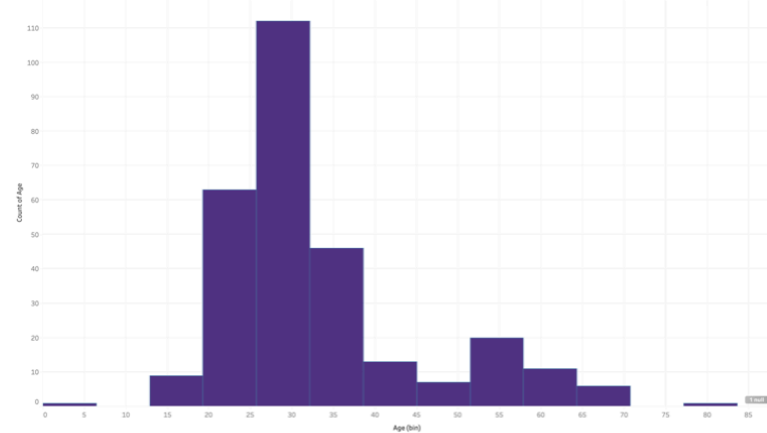
Social Media Users



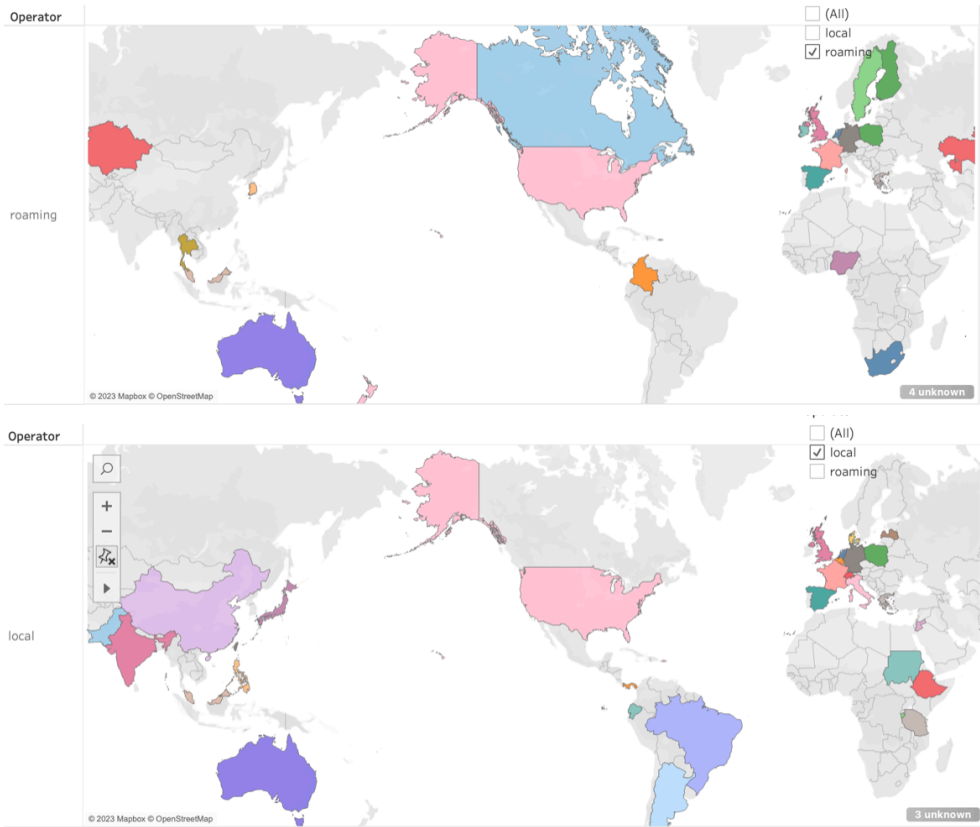
Avg. Smartphone Function by Education Level



Tourist Age Range



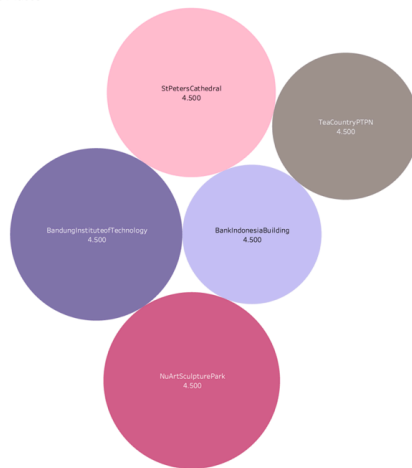
**Dashboards**



## Project

### 1. Tourist

Top Places



Internet Usage Hours by Place

Type (Avginternetus... F	
DepartmentStoresShoppi..	11.330
HotSpringsampGeysers	8.800
Waterfalls	7.800
ArchitecturalBuildings	7.750
Neighborhoods	7.480
PointsofInterestampLand..	7.250
FactoryOutlets	7.000
ScienceMuseums	6.530
SpecialityMuseums	6.470
ScenicWalkingAreas	6.430
Parks	6.000
PointsofInterestampLand..	5.830
PointsofInterestampLand..	5.820
NatureampWildlifeAreas	5.600
PointsofInterestampLand..	5.500
ShoppingMalls	5.100
AmusementampThemePa...	5.080
ArtGalleriesParks	3.500
ReligiousSites	2.500
NeighborhoodsHistoricW...	2.500

## Results

Top Hotels



Tourist Attraction (BestPlaces)

- BandungInstituteofTechnology
- BankIndonesiaBuilding
- NuArtSculpturePark
- StPetersCathedral
- TeaCountryPTPN

Name (Besthotels)

- BelviuHotel
- ibisStylesBandungGrandCentral
- TheGaiaHotelBandung
- TheNaripanHotel
- ThePapandayan

demographic: The majority of tourists are in their 30s and have a diploma, and they use social media platforms like Facebook and Twitter to share details of their trip. This information can be helpful for businesses that cater to tourists, such as hotels, restaurants, and tour operators, to target their marketing and advertising efforts towards this demographic.

2. Local operators vs. roaming: Tourists from India and the US are more likely to use local operators instead of roaming, while tourists from Australia and Canada are more likely to use roaming. This suggests that tourists from different regions may have different preferences and needs when it comes to mobile network services. Businesses in the telecommunications industry could use this information to tailor their services to better suit the needs of different tourist groups.

3. Top places and hotels: The top places and hotels have been calculated based on both the number of ratings and the highest reviews. This information can be useful for tourists who are looking for popular and highly rated places to visit or stay, as well as for businesses that want to improve their ratings and reviews by providing quality services and experiences to their customers.



4. Internet usage: The average hours of internet usage by places visited shows that internet usage is lowest in places like religious sites and highest in places like shopping malls. This information can be helpful for businesses that offer free Wi-Fi or other internet-related services, as they can target their efforts towards places where internet usage is higher. Additionally, it suggests that tourists may have different needs and preferences when it comes to internet access depending on the places they visit.