# Uber Data Analytics Project

→ Date : 18 August, 2024

→ Reference : https://youtu.be/WpQECq5Hx9g?si=kxJC1Y8_FBDdIZLM
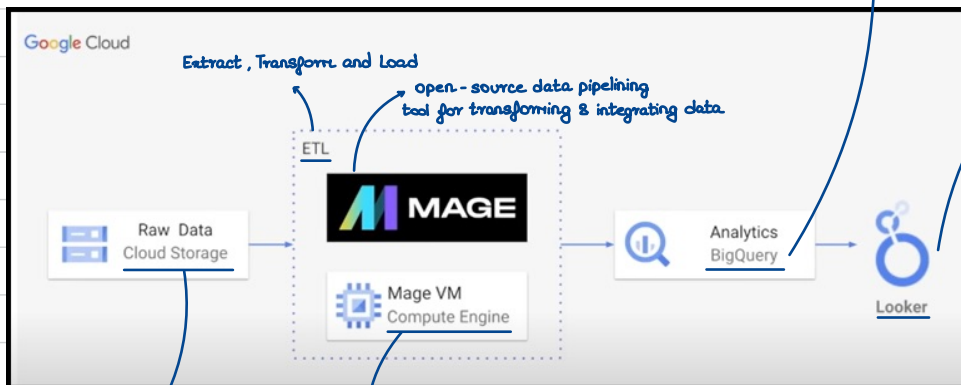
## → Architecture

**Differences from Operational Databases**

- **Operational Databases:** Used for day-to-day transaction processing (e.g., e-commerce transactions, inventory management). They are optimized for fast writes and real-time processing but may not perform as well with complex analytical queries.
- **Data Warehouses:** Designed for aggregating and analyzing large volumes of historical data. They support complex queries and reporting, and are typically read-heavy rather than write-heavy.

Google manages infrastructure, provisioning, backend, backups, patching, scaling...

* Fully - managed , Serverless Data Warehouse
    specialized database → querying & reporting
* SQL - based Querying   vs Traditional DBs optimized for transactional process
* Different from GCS (storing unstructured data) vs Analyzing & Querying Data

Google Cloud

Extract, Transform and Load

open-source data pipelining tool for transforming & integrating data

ETL — MAGE — Mage VM Compute Engine

Raw Data Cloud Storage → Analytics BigQuery → Looker

**Looker Studio** (formerly known as Google Data Studio) is a data visualization and business intelligence (BI) tool provided by Google. It allows users to create interactive dashboards and reports that can visualize data from various sources. Here's a detailed overview of Looker Studio:

**Key Features of Looker Studio**

1. **Data Integration:**
   - **Connectors:** Looker Studio supports a wide range of data connectors, including Google products like Google Analytics, Google Ads, BigQuery, and Google Sheets, as well as third-party sources such as SQL databases, cloud storage, and APIs.
   - **Data Blending:** Users can combine data from different sources to create comprehensive reports and dashboards.

2. **Customizable Dashboards:**
   - **Interactive Elements:** Create interactive elements like filters, date range selectors, and dynamic controls to allow users to explore data in various ways.
   - **Visualizations:** Offers a variety of visualization types including charts, graphs, tables, and geo maps.

3. **Collaboration:**
   - **Sharing and Permissions:** Share reports and dashboards with other users and control access levels (view, edit, comment).
   - **Real-time Collaboration:** Multiple users can work on the same report simultaneously, similar to Google Docs.

4. **Data Transformation:**
   - **Calculated Fields:** Create new metrics and dimensions by applying custom calculations to existing data.
   - **Data Blending and Aggregation:** Combine and aggregate data from different sources to generate comprehensive insights.

5. **User-Friendly Interface:**
   - **Drag-and-Drop:** Provides an intuitive drag-and-drop interface for designing reports and dashboards without needing extensive technical knowledge.
   - **Templates:** Offers pre-built templates and themes to streamline report creation.

6. **Customization and Branding:**
   - **Custom Themes:** Customize the appearance of reports to match your brand's look and feel.
   - **Embedded Reporting:** Embed interactive reports and dashboards in websites or applications.

* Online file storage service
* Store and Retrieve Data from Cloud

* Service providing VMs for running applications & services
* Create, Configure and Manage VMs with various OS

## → Fact & Dimension Table

**Fact Table:**

- Contains quantitative measures or metrics that are used for analysis
- Typically contains foreign keys that link to dimension tables
- Contains columns that have high cardinality and change frequently
- Contains columns that are not useful for analysis by themselves, but are necessary for calculating metrics

**Dimension Table:**

- Contains columns that describe attributes of the data being analyzed
- Typically contains primary keys that link to fact tables
- Contains columns that have low cardinality and don't change frequently
- Contains columns that can be used for grouping or filtering data for analysis

https://builtin.com/articles/fact-table-vs-dimension-table

**Fact Table vs. Dimension Table Defined**

- **Fact table:** A fact table contains the primary keys of the referenced dimension tables along with some quantitative metrics. Examples of a fact table include customer orders or time-series financial data.
- **Dimension table:** A dimension table holds the descriptive information for the related fields that are in the fact table's records. It typically represents a physical entity like "customer" or "product."

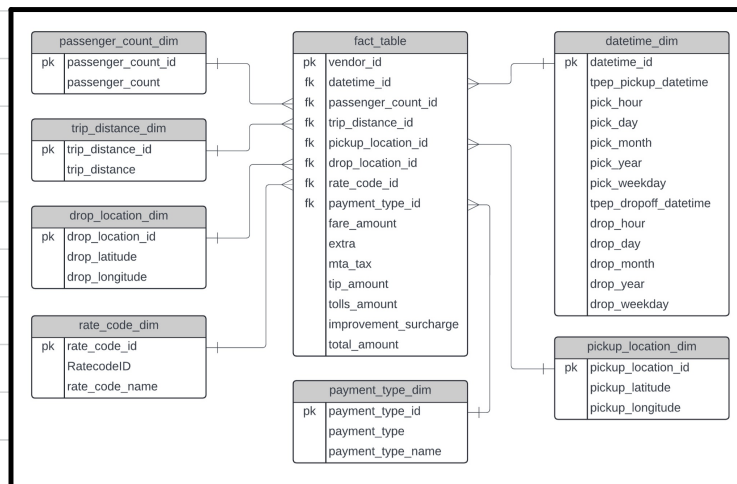## → Given Data : Taxi Trip Details

### Sample : uber_data.csv

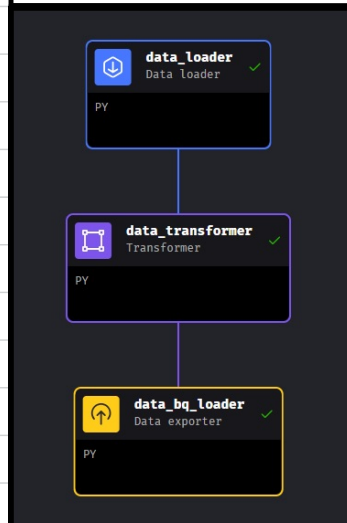| | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distance | pickup_longitude | pickup_latitude | RatecodeID | store_and_fwd_flag | dropoff_ |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2016-03-01 00:00:00 | 2016-03-01 00:07:55 | 1 | 2.50 | -73.976746 | 40.765152 | 1 | N | |
| 1 | 1 | 2016-03-01 00:00:00 | 2016-03-01 00:11:06 | 1 | 2.90 | -73.983482 | 40.767925 | 1 | N | |
| 2 | 2 | 2016-03-01 00:00:00 | 2016-03-01 00:31:06 | 2 | 19.98 | -73.782021 | 40.644810 | 1 | N | |
| 3 | 2 | 2016-03-01 00:00:00 | 2016-03-01 00:00:00 | 3 | 10.78 | -73.863419 | 40.769814 | 1 | N | |
| 4 | 2 | 2016-03-01 00:00:00 | 2016-03-01 00:00:00 | 5 | 30.43 | -73.971741 | 40.792183 | 3 | N | |

## → Google Cloud Steps

→ Google Cloud Storage → Create Bucket → Upload file → Fine-grained access
  → Edit Access → Public → URL

→ Compute Engine → Create Instance → Hardware + OS Configuration
  → Allow HTTP / HTTPS traffic
  → SSH Connect , Install Packages and Dependencies

→ Virtual Private Cloud → Create Firewall Rule
  → Target, Source IP : 0.0.0.0/0 (open for any incoming traffic)
  → Open TCP port 6789 → Mage Access

→ Service Account → Create → Grant BigQuery Admin Role
  → Create Key → Download JSON file

→ BigQuery → Create Dataset → Note Dataset ID

## → Mage Pipeline  (Access: external_ip:6789)

Planned
  Data Structure
    Transformation



Pipeline

→ Code Explanation

```python
import io
import pandas as pd
import requests
if 'data_loader' not in globals():
    from mage_ai.data_preparation.decorators import data_loader
if 'test' not in globals():
    from mage_ai.data_preparation.decorators import test


@data_loader
def load_data_from_api(*args, **kwargs):
    """
    Template for loading data from API
    """
    url = 'https://storage.googleapis.com/my_uber_data_analytics_project/uber_data.csv'
    response = requests.get(url)

    return pd.read_csv(io.StringIO(response.text), sep=',')


@test
def test_output(output, *args) -> None:
    """
    Template code for testing the output of the block.
    """
    assert output is not None, 'The output is undefined'
```

*URL for Cloud stored data.*

```python
import pandas as pd
if 'transformer' not in globals():
    from mage_ai.data_preparation.decorators import transformer
if 'test' not in globals():
    from mage_ai.data_preparation.decorators import test


@transformer
def transform(df, *args, **kwargs):
    """
    Template code for a transformer block.

    Add more parameters to this function if this block has multiple parent blocks.
    There should be one parameter for each output variable from each parent block.

    Args:
        data: The output from the upstream parent block
        args: The output from any additional upstream blocks (if applicable)

    Returns:
        Anything (e.g. data frame, dictionary, array, int, str, etc.)
    """
    # Specify your transformation logic here

    # Converting objects to datetime for required fields

    df['tpep_pickup_datetime']  = pd.to_datetime(df['tpep_pickup_datetime'])
    df['tpep_dropoff_datetime'] = pd.to_datetime(df['tpep_dropoff_datetime'])

    df = df.drop_duplicates().reset_index(drop=True)
    df['trip_id'] = df.index
```

```python
32
33    # Creating datetime_dim table
34
35    datetime_dim = df[['tpep_pickup_datetime','tpep_dropoff_datetime']].reset_index(drop=True)
36    datetime_dim['tpep_pickup_datetime'] = datetime_dim['tpep_pickup_datetime']
37    datetime_dim['pick_hour']            = datetime_dim['tpep_pickup_datetime'].dt.hour
38    datetime_dim['pick_day']             = datetime_dim['tpep_pickup_datetime'].dt.day
39    datetime_dim['pick_month']           = datetime_dim['tpep_pickup_datetime'].dt.month
40    datetime_dim['pick_year']            = datetime_dim['tpep_pickup_datetime'].dt.year
41    datetime_dim['pick_weekday']         = datetime_dim['tpep_pickup_datetime'].dt.weekday
42
43    datetime_dim['tpep_dropoff_datetime'] = datetime_dim['tpep_dropoff_datetime']
44    datetime_dim['drop_hour']            = datetime_dim['tpep_dropoff_datetime'].dt.hour
45    datetime_dim['drop_day']             = datetime_dim['tpep_dropoff_datetime'].dt.day
46    datetime_dim['drop_month']           = datetime_dim['tpep_dropoff_datetime'].dt.month
47    datetime_dim['drop_year']            = datetime_dim['tpep_dropoff_datetime'].dt.year
48    datetime_dim['drop_weekday']         = datetime_dim['tpep_dropoff_datetime'].dt.weekday
49
50    datetime_dim['datetime_id'] = datetime_dim.index
51    # Ordering the Columns
52    datetime_dim = datetime_dim[['datetime_id', 'tpep_pickup_datetime', 'pick_hour', 'pick_day', 'pick_month', 'pick_year', 'pick_weekday',
53                                 'tpep_dropoff_datetime', 'drop_hour', 'drop_day', 'drop_month', 'drop_year', 'drop_weekday']]
54
55    # Creating passenger_count_dim table
56
57    passenger_count_dim = df[['passenger_count']].reset_index(drop=True)
58    passenger_count_dim['passenger_count_id'] = passenger_count_dim.index
59    passenger_count_dim = passenger_count_dim[['passenger_count_id','passenger_count']]
60
61    # Creating trip_distance_dim table
62
63    trip_distance_dim = df[['trip_distance']].reset_index(drop=True)
64    trip_distance_dim['trip_distance_id'] = trip_distance_dim.index
65    trip_distance_dim = trip_distance_dim[['trip_distance_id','trip_distance']]
```

```python
67    # Rate Code Type Dictionary
68
69    rate_code_type = {
70        1:"Standard rate",
71        2:"JFK",
72        3:"Newark",
73        4:"Nassau or Westchester",
74        5:"Negotiated fare",
75        6:"Group ride"
76    }
77
78    rate_code_dim = df[['RatecodeID']].reset_index(drop=True)
79    rate_code_dim['rate_code_id'] = rate_code_dim.index
80    rate_code_dim['rate_code_name'] = rate_code_dim['RatecodeID'].map(rate_code_type)
81    rate_code_dim = rate_code_dim[['rate_code_id','RatecodeID','rate_code_name']]
82
83    # Creating pickup_location_dim table
84
85    pickup_location_dim = df[['pickup_longitude','pickup_latitude']].reset_index(drop=True)
86    pickup_location_dim['pickup_location_id'] = pickup_location_dim.index
87    pickup_location_dim = pickup_location_dim[['pickup_location_id','pickup_latitude','pickup_longitude']]
88
89    # Creating dropoff_location_dim table
90
91    dropoff_location_dim = df[['dropoff_longitude','dropoff_latitude']].reset_index(drop=True)
92    dropoff_location_dim['dropoff_location_id'] = dropoff_location_dim.index
93    dropoff_location_dim = dropoff_location_dim[['dropoff_location_id','dropoff_latitude','dropoff_longitude']]
94
```

```python
# Payment Type Dictionary

payment_type_name = {
    1:"Credit card",
    2:"Cash",
    3:"No charge",
    4:"Dispute",
    5:"Unknown",
    6:"Voided trip"
}

# Creating payment_type_dim table

payment_type_dim = df[['payment_type']].reset_index(drop=True)
payment_type_dim['payment_type_id']   = payment_type_dim.index
payment_type_dim['payment_type_name'] = payment_type_dim['payment_type_id'].map(payment_type_name)
payment_type_dim = payment_type_dim[['payment_type_id','payment_type','payment_type_name']]

fact_table = df.merge(passenger_count_dim, left_on='trip_id', right_on='passenger_count_id') \
                .merge(trip_distance_dim      , left_on='trip_id', right_on='trip_distance_id') \
                .merge(rate_code_dim          , left_on='trip_id', right_on='rate_code_id') \
                .merge(pickup_location_dim    , left_on='trip_id', right_on='pickup_location_id') \
                .merge(dropoff_location_dim   , left_on='trip_id', right_on='dropoff_location_id')\
                .merge(datetime_dim           , left_on='trip_id', right_on='datetime_id') \
                .merge(payment_type_dim       , left_on='trip_id', right_on='payment_type_id') \
                [['trip_id','VendorID', 'datetime_id', 'passenger_count_id',
                'trip_distance_id', 'rate_code_id', 'store_and_fwd_flag', 'pickup_location_id', 'dropoff_location_id',
                'payment_type_id', 'fare_amount', 'extra', 'mta_tax', 'tip_amount', 'tolls_amount',
                'improvement_surcharge', 'total_amount']]

return {"datetime_dim":datetime_dim.to_dict(orient="dict"),
        "passenger_count_dim":passenger_count_dim.to_dict(orient="dict"),
        "trip_distance_dim":trip_distance_dim.to_dict(orient="dict"),
        "rate_code_dim":rate_code_dim.to_dict(orient="dict"),
        "pickup_location_dim":pickup_location_dim.to_dict(orient="dict"),
        "dropoff_location_dim":dropoff_location_dim.to_dict(orient="dict"),
        "payment_type_dim":payment_type_dim.to_dict(orient="dict"),
        "fact_table":fact_table.to_dict(orient="dict")}
```

```python
from mage_ai.data_preparation.repo_manager import get_repo_path
from mage_ai.io.bigquery import BigQuery
from mage_ai.io.config import ConfigFileLoader
from pandas import DataFrame
from os import path

if 'data_exporter' not in globals():
    from mage_ai.data_preparation.decorators import data_exporter


@data_exporter
def export_data_to_big_query(data, **kwargs) -> None:
    """
    Template for exporting data to a BigQuery warehouse.
    Specify your configuration settings in 'io_config.yaml'.

    Docs: https://docs.mage.ai/design/data-loading#bigquery

    """
    config_path = path.join(get_repo_path(), 'io_config.yaml')
    config_profile = 'default'

    for key, value in data.items():
        table_id = 'project-uberdataanalytics.Uber_Transformed_Dataset.{}'.format(key)
        BigQuery.with_config(ConfigFileLoader(config_path, config_profile)).export(
            DataFrame(value),
            table_id,
            if_exists='replace',  # Specify resolution policy if table name already exists
        )
```

→ SQL Query

```sql
1   CREATE OR REPLACE TABLE `project-uberdataanalytics.Uber_Transformed_Dataset.AnalysisReport` AS (
2     SELECT
3       f.VendorID,
4       dt.tpep_pickup_datetime,
5       dt.tpep_dropoff_datetime,
6       p.passenger_count,
7       td.trip_distance,
8       rc.RatecodeID,
9       rc.rate_code_name,
10      f.store_and_fwd_flag,
11      pl.pickup_latitude,
12      pl.pickup_longitude,
13      dl.dropoff_latitude,
14      dl.dropoff_longitude,
15      pt.payment_type,
16      f.fare_amount,
17      f.extra,
18      f.mta_tax,
19      f.tip_amount,
20      f.tolls_amount,
21      f.improvement_surcharge,
22      f.total_amount
23    FROM
24      `project-uberdataanalytics.Uber_Transformed_Dataset.fact_table` f
25      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.passenger_count_dim` p
26        ON f.passenger_count_id = p.passenger_count_id
27      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.trip_distance_dim` td
28        ON f.trip_distance_id = td.trip_distance_id
29      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.rate_code_dim` rc
30        ON f.rate_code_id = rc.rate_code_id
31      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.pickup_location_dim` pl
32        ON f.pickup_location_id = pl.pickup_location_id
33      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.dropoff_location_dim` dl
34        ON f.dropoff_location_id = dl.dropoff_location_id
35      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.payment_type_dim` pt
36        ON f.payment_type_id = pt.payment_type_id
37      JOIN `project-uberdataanalytics.Uber_Transformed_Dataset.datetime_dim` dt
38        ON f.datetime_id = dt.datetime_id
39  );
40
```

→ Data Visualization Tool — Looker Studio

# Uber

## Filters

| Vendor ID ▾ | Rate Code Name ▾ | Network Error ▾ |

**Trip Distance**

0 ———————————————————————————————— 184.4

## Quick Statistics

| Total Revenue | Record Count | Avg. Trip Distance | Avg. Fare Amount | Avg. Tip Amount |
|---|---|---|---|---|
| 1.6M | 100.0K | 3.0 | 13.3 | 1.9 |

## Pick-Up Location



**Rate Code Name**  ● Negotiated fare  ● Newark  ● Standard rate  ● Nassau or Westchester  ● JFK  ● Group ride

## Drop-Off Location



**Rate Code Name**  ● Negotiated fare  ● Newark  ● Standard rate  ● Nassau or Westchester  ● JFK  ● Group ride