# My Data Pipelines Project
# -Guoying Dai

This project is to help a e-scooter company Gans to improve the e-scooter-sharing system. As the populations and weather of a city, as well as the time when people arriving the airport are very important elements that influenced the sharing system. We should collect data from these sources that potentially help Gans e-scooter movement. Since data is needed every day, in real-time and accessible by everyone in the company, the challenge is to assemble and automate a data pipeline in the cloud. There are two major phases of the project: create local pipeline and then create the cloud pipeline.

## 1. Create Local Pipeline

In this phase, we should gather data from website by execute the scripts in Jupyter notebook, then store them to the local database-MySQL. There are several ways to dedicate data from website.

### 1) web scraping:

we use the python library :**BeautifulSoup** to access the information we needed by downloading and extracting the HTML code of the website (Fig1).

```
1  import requests
2  from bs4 import BeautifulSoup
3
4  def get_html_wiki_text_label(keyword, location):
5      response = requests.get(f'https://en.wikipedia.org/wiki/{location}')
6      soup = BeautifulSoup(response.content, 'html.parser')
7      results = []
8      i = 0
9      for item in soup.select("th.infobox-label"):
10         if(keyword in item.get_text()):
11             results.append(soup.select("td.infobox-data")[i].get_text())
12         i = i+1
13     return results
14
15 print('Country:', get_html_wiki_text_label('Country', 'Berlin')[0])
16 print('State:', get_html_wiki_text_label('State', 'Berlin')[0])
17 print('Population_City/State:', get_html_wiki_text_label('City/State', 'Berlin')[1])
18 print('Population_Urban:', get_html_wiki_text_label('Urban', 'Berlin')[1])
19 print('Population_Metro:', get_html_wiki_text_label('Metro', 'Berlin')[1])
20
21 Country: Germany
22 State: Berlin
23 Population_City/State: 3,769,495
24 Population_Urban: 4,473,101
25 Population_Metro: 6,144,600
```

Fig. 1 The code for getting the city's (Berlin) information: populations and state and country

**2) with APIs or particular API**

Application programming interface (API) is a way for two or more computer programs to communicate with each other. It is a type of software interface, offering a service to other pieces of software.(from Wikipedia) .

In this project, we use Python's **_request library_** to interact with APIs and collected weather and flights data by calling APIs.

- Weather:

Firstly sign up for an account in **OpenWeather**, then we got the _api keys._ Secondly, we make a API call. There are two ways to get the data by APIs ways: using a "for" loop (Fig 2), or  using '.json_normalize ()' (Fig 3).

```python
1  def get_weather_loop(cities):
2
3    API_key = 'api_key from OpenWeather'
4
5    tz = pytz.timezone('Europe/Berlin')
6    now = datetime.now().astimezone(tz)
7
8    weather_dict = {'city': [],
9                    'country': [],
10                    'forecast_time': [],
11                    'outlook': [],
12                    'detailed_outlook': [],
13                    'temperature': [],
14                    'temperature_feels_like': [],
15                    'clouds': [],
16                    'rain': [],
17                    'snow': [],
18                    'wind_speed': [],
19                    'wind_deg': [],
20                    'humidity': [],
21                    'pressure': [],
22                    'information_retrieved_at': []}
23
24    for city in cities:
25      url = (f"http://api.openweathermap.org/data/2.5/forecast?q={city}&appid={API_key}&units=metric")
26      response = requests.get(url)
27      json = response.json()
28
29      for i in json['list']:
30        weather_dict['city'].append(json['city']['name'])
31        weather_dict['country'].append(json['city']['country'])
32        weather_dict['forecast_time'].append(i['dt_txt'])
33        weather_dict['outlook'].append(i['weather'][0]['main'])
34        weather_dict['detailed_outlook'].append(i['weather'][0]['description'])
35        weather_dict['temperature'].append(i['main']['temp'])
36        weather_dict['temperature_feels_like'].append(i['main']['feels_like'])
37        weather_dict['clouds'].append(i['clouds']['all'])
38        try:
39            weather_dict['rain'].append(i['rain']['3h'])
40        except:
41            weather_dict['rain'].append('0')
42        try:
43            weather_dict['snow'].append(i['snow']['3h'])
44        except:
45            weather_dict['snow'].append('0')
46        weather_dict['wind_speed'].append(i['wind']['speed'])
47        weather_dict['wind_deg'].append(i['wind']['deg'])
48        weather_dict['humidity'].append(i['main']['humidity'])
49        weather_dict['pressure'].append(i['main']['pressure'])
50        weather_dict['information_retrieved_at'].append(now.strftime("%d/%m/%Y %H:%M:%S"))
51
```

Fig.2 using a "for" loop to collect the weather information

```
1  city = 'Berlin'
2  API_key = 'api_key from OpenWeather'
3  url = (f"http://api.openweathermap.org/data/2.5/forecast?q={city}&appid={API_key}&units=metric")
4  response = requests.get(url)
5  json = response.json()
6
7  json_norm_df = pd.json_normalize(json['list'],
8                                   record_path=['weather'],
9                                   meta=['dt_txt', ['main', 'temp'], ['main', 'feels_like'], ['clouds', 'all'], ['rain', '3h'],
                                         ['snow', '3h'], ['wind', 'speed'], ['wind', 'deg'], ['main', 'humidity'], ['main', 'pressure']],
10                                  errors='ignore')
11
12 json_norm_df.head()
13
14
```

Fig.3 using '.json_normalize ()' to collect the weather information

- Flights:

  Firstly sign up for an account in *Rapid API*. Secondly go to the Aerodatabox api and subscribe. I choose the Basic and free plan. Thirdly go to the *Endpoints section* of the API and select the 'FIDS/Schedules: Airport departures and arrivals'. Then set the parameters and choose the 'Python-request', finally copy the code to my Jupter Notebook (Fig 4).



Fig. 4 setting the API from Aerodatabox

### 3) store the data in a database

Push the data from Jupyter Notebook to local MySQL.

We used two typical ways, here I choose the simple way, which could push the DataFrame data directly to MySQL database (Fig. 5).

```
 1  PUSH YOUR DATA FROM NOTEBOOK TO MySQL
 2
 3  !pip install SQLAlchemy
 4  !pip install PyMySQL
 5  from sqlalchemy import create_engine
 6  import pymysql
 7
 8  database = 'collected_data' #the database created in MySQL
 9  username = 'root'
10  password = 'password'
11  sqlEngine      = create_engine(f'mysql+pymysql://{username}:{password}@127.0.0.1/{database}', pool_recycle=3600)
12  dbConnection   = sqlEngine.connect()
13  tableName = 'flights_info'
14  |
15  try:
16
17      frame = tomorrows_flight_arrivals(icaos).to_sql(tableName, dbConnection, if_exists='append', index = False);
18
19  except ValueError as vx:
20
21      print(vx)
22
23  except Exception as ex:
24
25      print(ex)
26
27  else:
28
29      print("The data was pushed sucessfully to Table %s."%tableName);
30
31  finally:
32
33      dbConnection.close()
```

Fig. 5 Push the data in 'DataFrame' format  to the local MySQL

After pushing the data from the Jupyter notebook, I got the database in my local MySQL (Fig 6). Every time I execute the python code , the new collected data will be inserted to table. But in this way give it a disadvantage, that we need to refresh the data in the database manually if we need the real-time information of weather and flights,  which will be time consuming and maybe not so specific correct. Therefore, I started the second project, moving the data pipeline to the cloud, in which way, I could automatically gather data from website and insert to local database in MySQL.
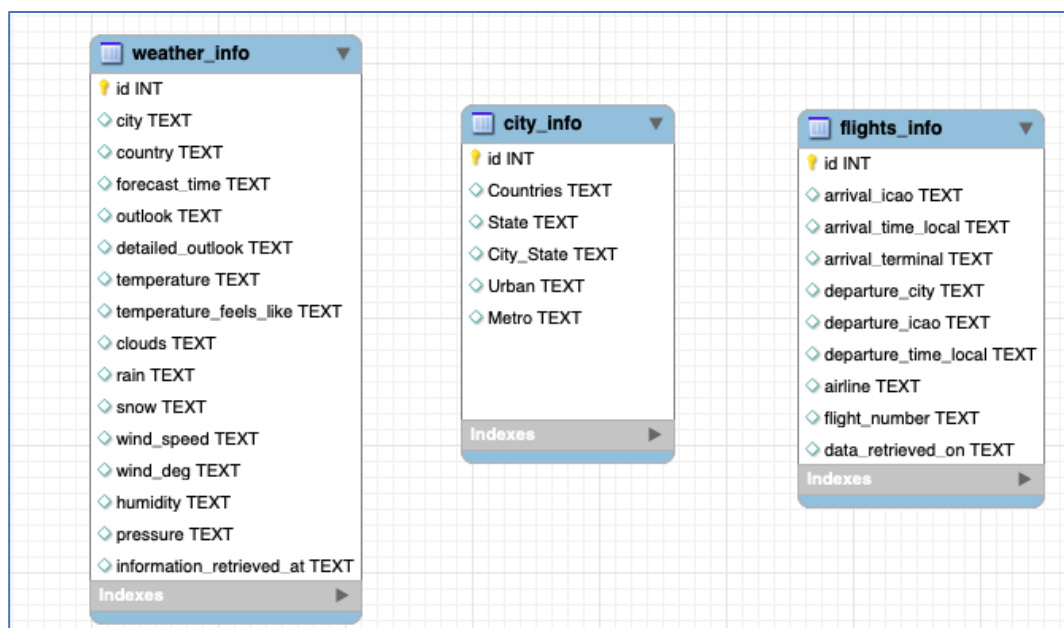
## 2. Cloud Pipeline Project

In this section we use Amazon Web Service (AWS) to add a could pipeline. The process contains the following 3 steps:

**1) Sign up for an AWS account.**

Firstly we need sigh up an account on AWS, a tricky thing is that you need a valid credit card for it even you may not charge for anything. On AWS platform, there are lots of detail knowledge of AWS service, if you do not familiar with the service like Lambda, you could learn on AWS plateform. Meanwhile, it provides many specific training options, for example data science and machine learning etc.

Secondly we should set up a database on AWS RDS. The setting process contains many detailed steps, I will not list every steps here.

The last step will be to connect to the instance we created so that we can create a database with all of the tables to store the information that collected online.

**2) Automate the data pipeline**

After setting up a cloud MySQL instance, now we use AWS Lambda to run the Python scripts for data collection. We need three steps to create a Lambda function.

- Create a Role called 'LambdaAdminAccess' with the permission 'AdministratiorAccess'
- Create a Lambda function called 'weather_info' with the role 'LambdaAdminAccess'.
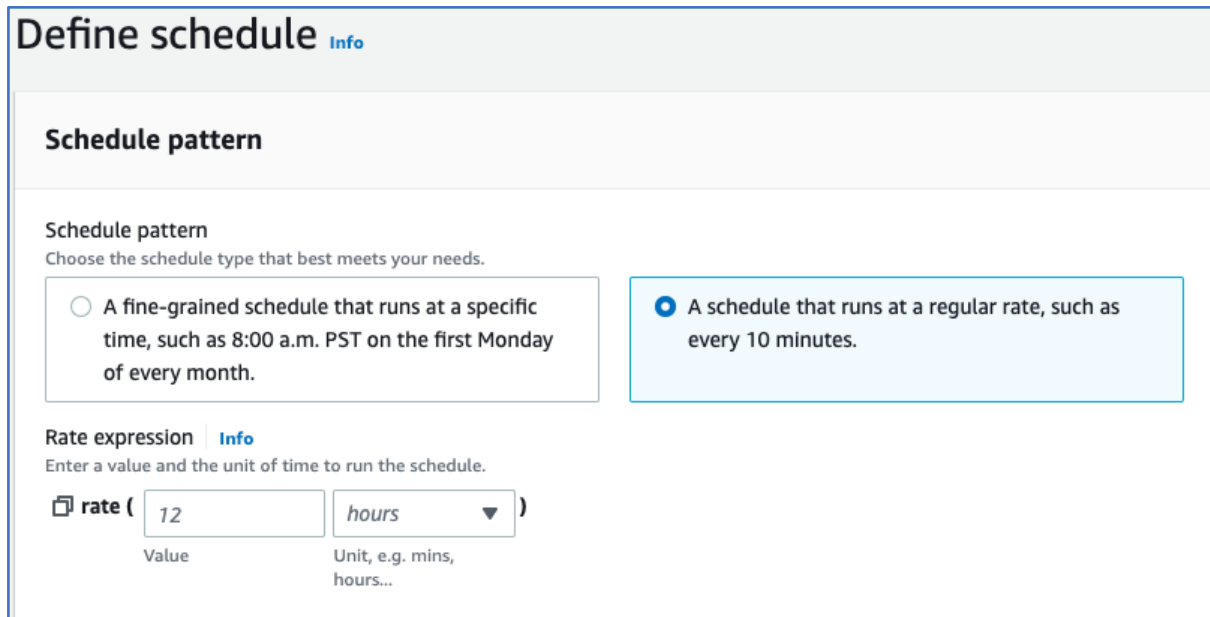- Connect the Lambda function to RDS instance.

The Python scripts we run in Lambda function is almost the same with the one we push to the local database manfully from Jupter notebook. There are two main differences: one is we should add 'def lambda_handler(event, context):" in it , another is adding 'return' at the end of the function.

By run the Lambda functions, we have got our data from the internet and insert it onto the RDS cloud instance.

### 3) Collection data from website to local database automatedly

We use AWS EventBridge to set up a trigger event.
The schedule of running the Lambda function could be stetted in the 'Rule-Define schedule' (Fig. 7).

Now the cloud pipeline was settled down. We could collect the data automatically and in real-time.



Fig. 7 define the scheldule for Lambda function to execute