# Certification CDSD Block 1 & 3

August 28th, 2025 – Louis Le Pogam

# Agenda

- **Block 1 - Build & Manage a Data Infrastructure – Kayak Project**

- **Block 3 – Unsupervised Machine Learning – Uber Pickups Project**

# Project Reminder

**Project**

**Kayak would like to create a holiday recommendation application based on :**
- **Weather**
- **Hotels in the area**
- **Based on real-time data**

**Goal**

**Get the needed data as following:**
- **Scrape data from destinations**
- **Get weather data and hotels' info from each destination**
- **Store all the information above in a data lake**
- **Extract, transform and load cleaned data from your datalake to a data warehouse**

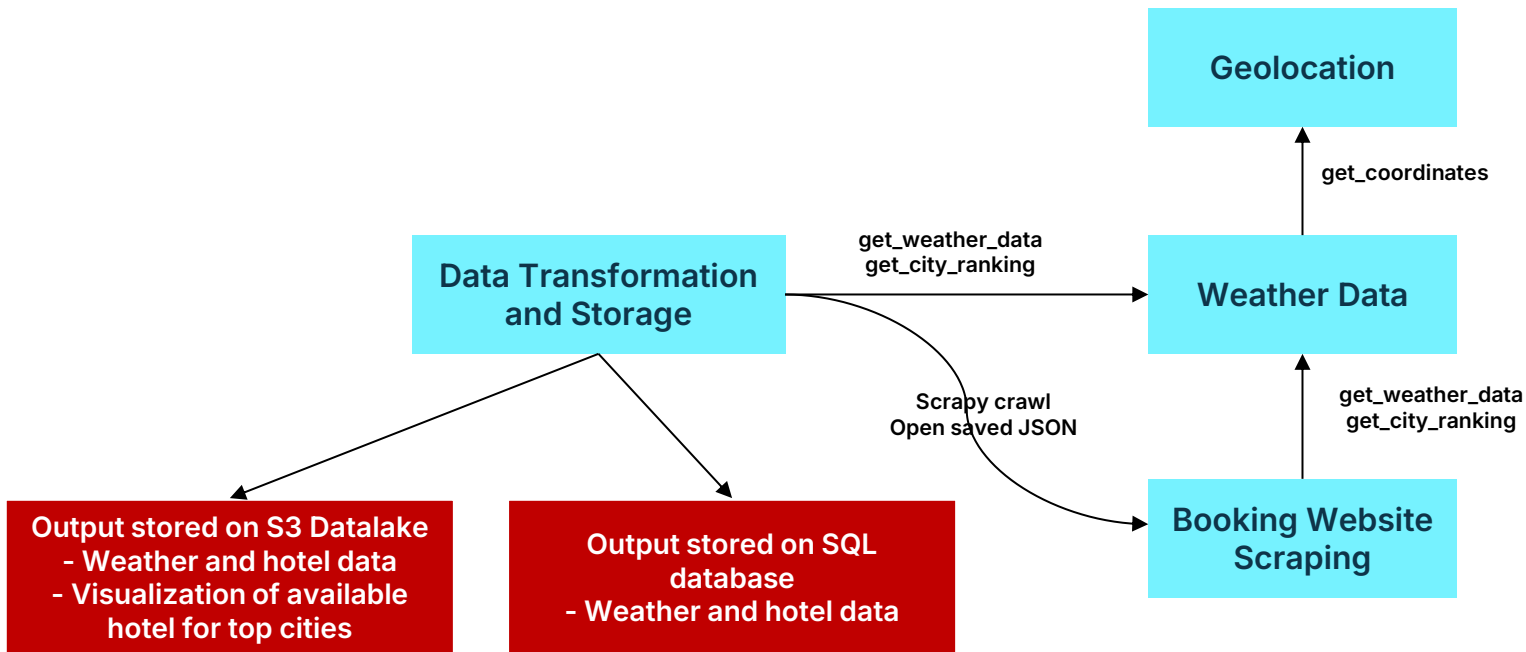# 4 building blocks for the data scraping model

## Description

**1** **Geolocation**
- Processing of list of cities by obtaining GPS coordinates and INSEE codes with API
- Saving all data to CSV for weather queries

**2** **Weather Data**
- Retrieving 7-day weather forecasts for cities based on INSEE codes
- Ranking cities based on customizable criteria and creates aggregated rankings
- Saving ranked list based on number of favorable days to CSV

**3** **Booking Website Scraping**
- Taking the top 5 cities based on previous analysis
- Searching for hotels in each city available on booking.com
- Saving in into a JSON file

**4** **Data Transformation & Storage**
- Removing low rated hotels and identifying consecutive night availability
- Creating interactive visualizations and storing results in AWS S3 (data lake) and an AWS RDS-hosted SQL relational database

# All 4 blocks are used to get the final output

**KAYAK**

**Geolocation**

↑ get_coordinates

**Data Transformation and Storage**

get_weather_data
get_city_ranking →

**Weather Data**

Scrapy crawl
Open saved JSON

get_weather_data
get_city_ranking

**Output stored on S3 Datalake**
- Weather and hotel data
- Visualization of available hotel for top cities

**Output stored on SQL database**
- Weather and hotel data

**Booking Website Scraping**

# The weather module scraps the weather by city and ranks them based on defined ideal conditions

## get_weather_data

- Reads a list of French cities from a file based on the INSEE code
- Connects to a weather service (Meteo Concept API) for each city
- Gathers 7-day forecasts including:
  - Daily rainfall predictions
  - Temperature highs and lows
  - Wind speed information
- Organizes everything in a DataFrame where each row represents one day's weather forecast for a specific city

**Output : DataFrame with weather prediction for eachs targeted cities**

## get_city_ranking

- Flags Less-Than-Ideal Weather Days when:
  - Too hot / Too cold / Too rainy / Too windy
  - Default value: 35°C / 20°C / 10mm / 50 km/h
  - Can be changed depending on the season
- Calculates Problem Days for each city over the forecast period
- Ranks Cities by:
  - Fewest problematic weather days
  - Lowest average rainfall as tiebreaker

**Output : DataFrame with a prioritized list of cities with optimal weather conditions for travelers**

KAYAK

# 700 pages scrapped to get the available hotels of the next 7 days

**Get Weather Data**

- Open file of ranked cities and select top 5

**Get URL for each city and date**

- Loop on each cities and for the next 7 days
- Get the URL of the search for each city x date

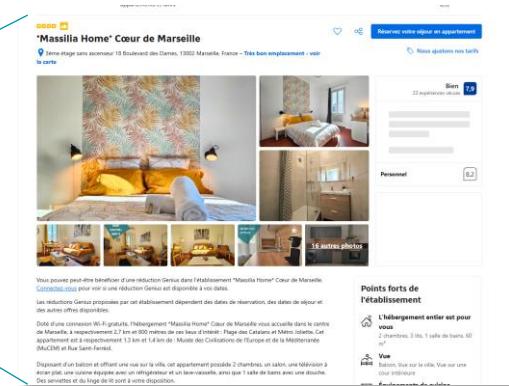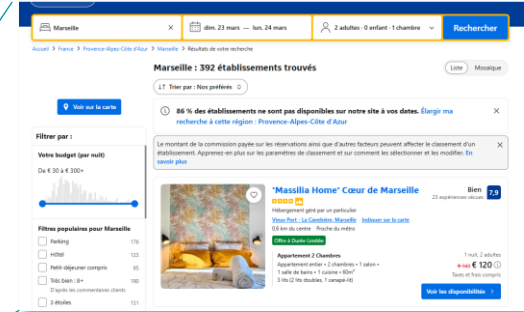**Get the Details from the Search Results**

- For each results on the search page, get the main info : Name, ranking, price, distance…
- Keep only 20 first results to limit output size

**Open each page and get detailed data**

- For each hotel, open the hotel page to get detailed info : Adress, latitude, longitude, description, URL…

**Store Data in a JSON file**

- Gather all data in one JSON
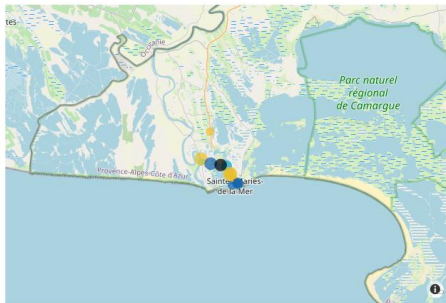- 700 pages scrapped (5 cities x 20 hotels x dates x 7 dates) in 6 minutes

![KAYAK logo]

# All outputs are stored on an AWS S3 and SQL server and ready to be used

**Output 1 : All database**

**Output 2 : Visualization of available hotels for each top 5 cities**

**Output 3 : Overview of average temperature for the next 7 days of all targeted cities**

Hotels available in Saintes Maries de la mer



Average max temperature for the next 7 days per city



- **Dataframe with weather forecast**
- **Dataframe with City Ranking**
- **JSON with all booking.com data scrapped**

**All outputs are stored in a S3 and SQL database and ready to be extracted and used**

# Q&A

# Agenda

- **Block 1 - Build & Manage a Data Infrastructure – Kayak Project**

- **Block 3 – Unsupervised Machine Learning – Uber Pickups Project**

# Project Reminder

**Uber's Challenge**

- Drivers are not always located where and when riders need them

- This mismatch leads to longer wait times and reduced efficiency

**Project Goal**

- Detect real-time "hot zones" of high demand
- Provide actionable recommendations for driver reallocation
- Enable intuitive visualization for rapid decision-making

# 3 steps to define the hot zone

1. Exploratory Data Analysis

2. Model Selection

3. Real Time "Hot-Zone" Recommendation

# 3 steps to define the hot zone

1. **Exploratory Data Analysis**

2. **Model Selection**

3. **Real Time "Hot-Zone" Recommendation**

# Uber has provided a cleaned database ready to be used

## Dataset description

- Dataset of April 2014 used with 564k lines

- 4 columns:
  - Date
  - Latitude
  - Longitude
  - Base : Internal code, not used in the analysis

## Preprocessing

- Preprocessing limited to converting the date column into several sub-columns

- Focus on New York City inside this latitude and longitude line :
  - Latitude minimum = 40.4774
  - Latitude maximum = 40.9176
  - Longitude minimum = -74.2591
  - Longitude maximum = -73.7004

**Limited preprocessing performed on the database**

# Database Overview: Each point corresponds to a single ride
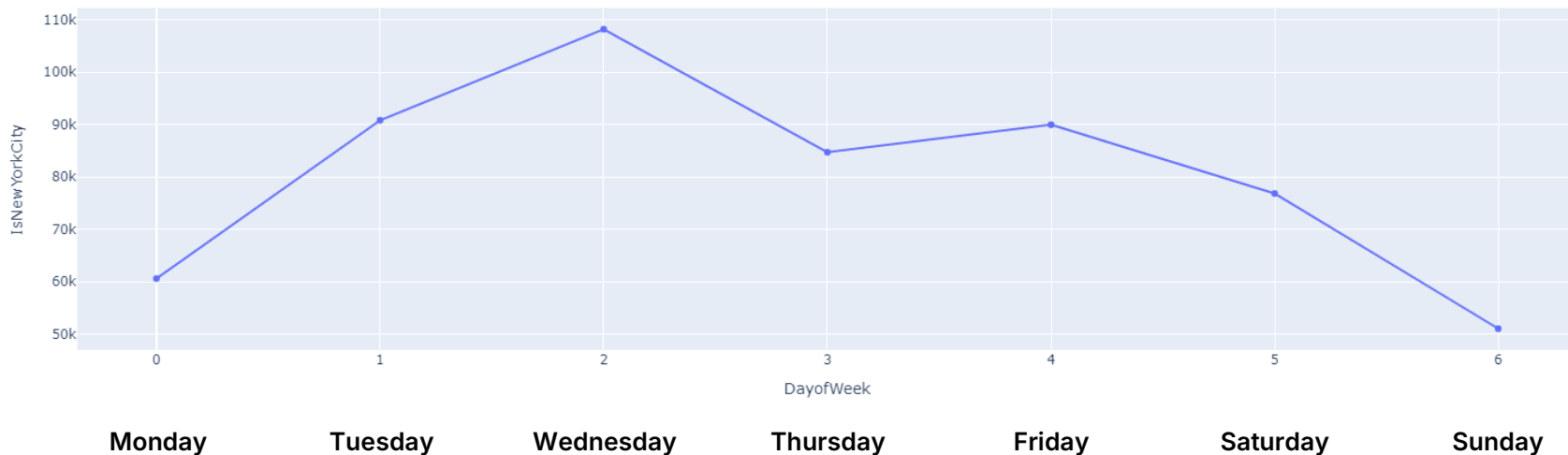


Pick up of the 30th at 5pm

Base
- B02512
- B02598
- B02617
- B02682
- B02764

# Monday and Sunday are the lowest day while the peak in on Wednesday

**Number of pickups per Day |** April 2014, NYC only, 0 = Monday



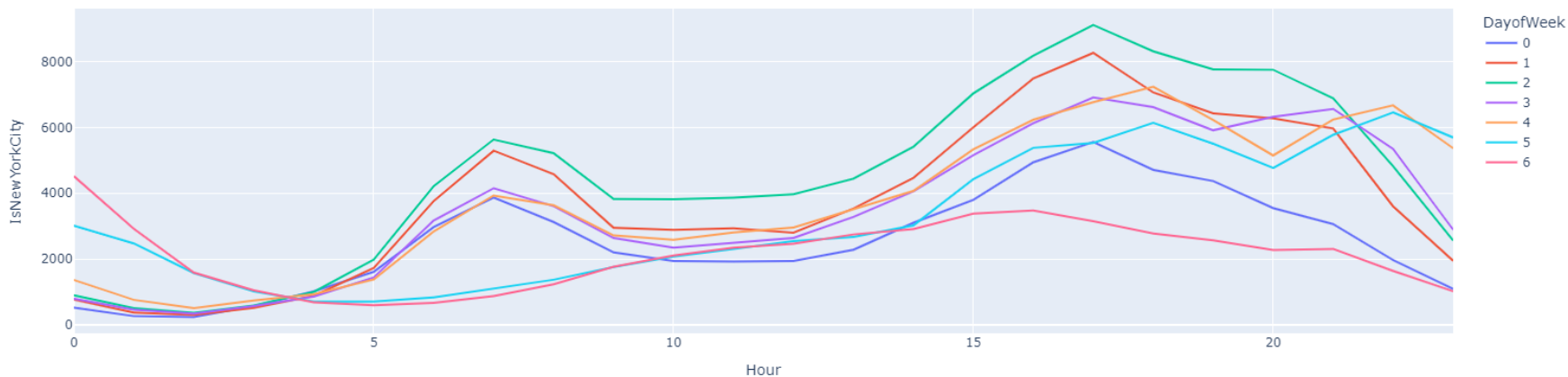| Monday | Tuesday | Wednesday | Thursday | Friday | Saturday | Sunday |

# During the week, there is a peak at 7am then between 5pm and 8pm while the night is busy during the weekend

**Number of pickups per Hour depending on the day |** April 2014, NYC only, 0 = Monday

# 3 steps to define the hot zone

1. Exploratory Data Analysis

2. Model Selection

3. Real Time "Hot-Zone" Recommendation

# 2 models can be chosen to define the hot zone

| | Kmeans Clustering | DBScan |
|---|---|---|
| **Principle** | • Centroid-based clustering algorithm | • Density-based clustering algorithm |
| **Shape** | • Spherical only | • No shape |
| **Strengths** | • Simple and fast<br>• Works well when clusters are spherical | • No need to know number of clusters<br>• Handle complex shapes |
| **Weaknesses** | • Number of clusters must be known<br>• Sensitive to initialization and outliers<br>• Poor performance on irregularly shape | • Commutationnally heavier |
| **Rationale for using it** | • Speed and simplicity when number of clusters is known | • Complex shapes and number cluster not known |

**Seems more adapted to handle real-time data**
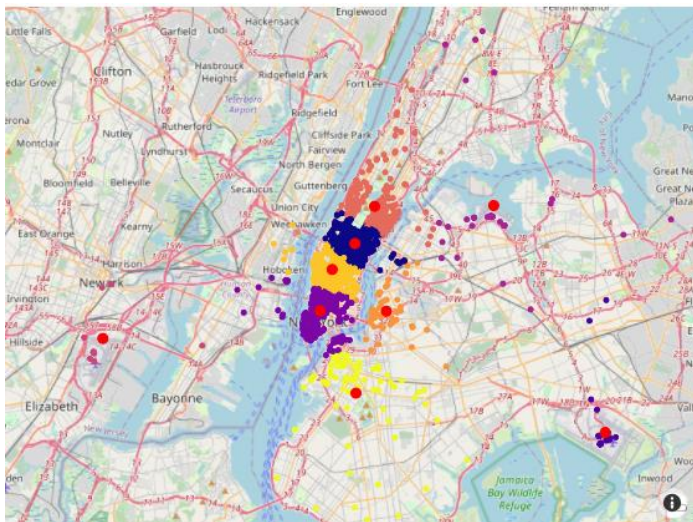
# DBScan was chosen for hot zone recommendation

## Kmeans Clustering

**Details**

- **Elbow and silhouette methods to get the optimal number of clusters**
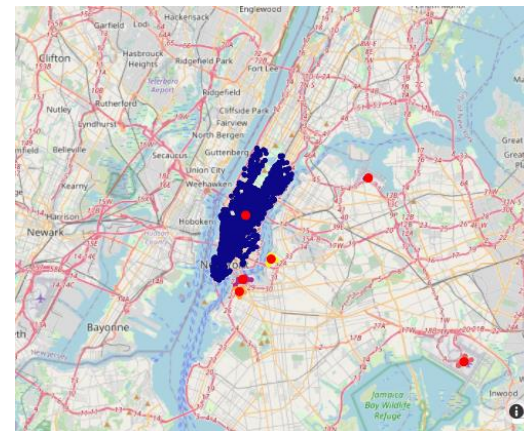- **9 clusters seems to be the best for April 30th**

**Cluster Overview**



## DBScan

**Details**

- **Parameters : Epsilon = 0.15 / Min Sample = 10**
- **6 clusters + outliers**

**Cluster Center Overview**



**Best model to handle variability in cluster location**

# 3 steps to define the hot zone

1. Exploratory Data Analysis

2. Model Selection

3. Real Time "Hot-Zone" Recommendation

# Hot zones are dynamically calculated with the DBScan algorithm

**Uber**

**Plot_hot_zone function**

| | |
|---|---|
| **Input** | • Dataset, day of the week, hour, dbscan parameters |
| **Calculation** | • For any given hours, would calculate the clusters center and plot them |
| **Output** | • A map with the hot zones of this hour |
| **Real-time visualization** | • For a given day, the evolution of hot spots can be shown by looping over different hours |

**Hot zones can be dynamically calculated and provided to the drivers at any given day and time of the week**
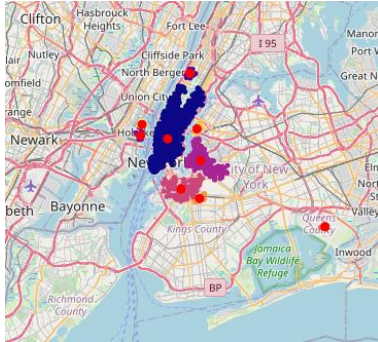
# The hot zone are changing every hour of the day

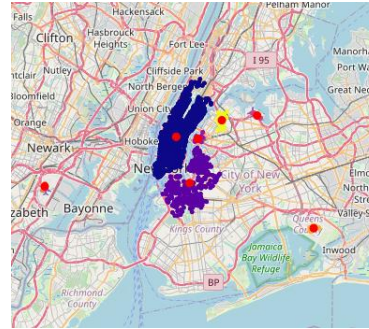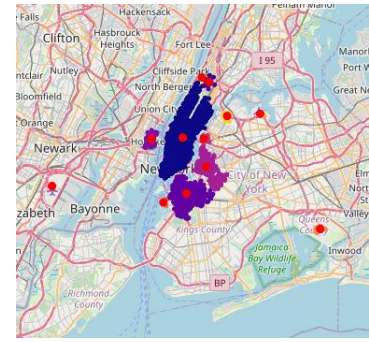Evolution of hot zones for a Saturday at 12am, 6am, 12pm, 6pm



| 12am | 6am | 12pm | 6pm |

# Q&A

# Thanks!