

## Milestone 2: Airline Customer Holiday Booking (10%)

*Deadline: 2nd of December at 23:59*

In this milestone, we transition from predicting passenger satisfaction using tabular data to uncovering deep, verifiable insights into **airline booking and travel experience trends** through a structured **Neo4j Knowledge Graph (KG)**. The KG transforms raw operational and feedback data into an interconnected network that explicitly captures relationships among passengers, flights, routes, cabin classes, and satisfaction drivers. This structure enables the pre-computation of complex analytical features such as route popularity, traveler profiles, and sentiment patterns. The KG will serve as the verifiable knowledge source for a **Graph Retrieval-Augmented Generation (GraphRAG)** system in the next milestone, enabling users to ask complex, multi-criteria airline booking questions in natural language, such as “Which routes provide the highest satisfaction for frequent travelers in Economy Class?”, and receive accurate, transparent answers grounded directly in the graph’s structure.

The core challenge is transforming the tabular Airline data into a graph structure that explicitly models positional scoring rules, temporal flow, and complex relationships for efficient natural language reasoning.

You will be working with “**Airline\_survey.csv**”. This .csv file contains a sample from the original detailed survey .csv file, with each row representing a passenger’s experience on a particular flight. Columns include flight identifiers (flight\_number, origin\_station\_code, destination\_station\_code), passenger information (record\_locator, loyalty\_program\_level, generation), journey characteristics (number\_of\_legs, class, fleet\_type\_description, actual\_flown\_miles), and feedback metrics such as arrival\_delay\_minutes and food\_satisfaction\_score. You will use this dataset to build a knowledge Graph capturing passengers’ feedback, flights, and airports, and then run queries to extract insights on overall satisfaction and travel patterns.

## Milestone 2 Requirements:

This milestone requires you to construct a comprehensive Knowledge Graph database using Neo4j that represents the airline booking domain according to the specified schema. The knowledge graph will model the complex relationships between flights, their journey requirements, creating a rich semantic network that enables sophisticated querying and analysis.

Your primary tasks include:

### 1. Knowledge Graph Construction

Implement a script (using Python and the Neo4j driver) that reads the provided CSV files and constructs the knowledge graph according to the exact schema specification provided below. This involves creating nodes for Flight, Passenger, Airport, and Journey, and establishing the relationships between them as defined in the schema.

### 2. Schema Compliance

It is critical that you strictly adhere to the provided schema. You are **not permitted** to modify the schema in any way, including:

- Changing the type of any node (e.g., *you cannot rename "Flight" to "Plane"*)
- Changing the type of any relationship (e.g., *you cannot rename ": ARRIVES\_AT" to ": ARRIVED\_IN" or ": REACHED"*)
- Removing or modifying any properties specified for nodes
- Removing any node types or relationship types from the schema

### 3. Implementing The Given Scoring Rule

You need to apply the overall satisfaction scoring rules outlined below in order to identify the number of passengers whose satisfaction levels exceed the defined threshold.

These passengers are those with an **overall\_satisfaction\_score greater than 3**, calculated using weighted components from their journey data. Keep in mind that clamping is applied in order to ensure that the calculated values stay within specific ranges.

### 4. Cypher Query Development

After constructing the knowledge graph, you will develop Cypher queries to answer a comprehensive set of analytical questions. These queries will test your understanding of graph traversal, aggregation, filtering, and pattern matching in Cypher. The queries range from simple filtering operations to complex multi-hop traversals that require careful consideration of relationship directions and node properties.

## The Knowledge Graph Schema:

### 1. Nodes:

- **Passenger**: record\_locator (unique identifier), loyalty\_program\_level, generation
- **Journey**: feedback\_ID (unique identifier), food\_satisfaction\_score, arrival\_delay\_minutes , actual\_flown\_miles ,number\_of\_legs, passenger\_class
- **Flight**: flight\_number, fleet\_type\_description (uniquely identified by both properties)
- **Airport**: station\_code (e.g., ORD, LAX)

### 2. Relations:

- (Passenger)-[:TOOK ]->(Journey)
- (Journey)-[:ON ]->(Flight)
- (Flight)-[:DEPARTS\_FROM ]->(Airport)
- (Flight)-[:ARRIVES\_AT ]->(Airport)

## Knowledge Graph Verification Queries

**NOTE ON VERIFICATION:** These are example questions & answers for initial structure confirmation. Create queries that answer the following questions. If your queries return the same answers as below, then the structure of your knowledge graph is sound. The final assessment of the graph's design and robustness will involve a separate set of private queries to ensure the model is flexible and not hardcoded to these specific examples.

1. Identify the 5 airport-to-airport routes with the highest number of flights. This simply shows which routes have the most scheduled flights. Your query should return:
  - **origin** – the departure airport code
  - **destination** – the arrival airport code
  - **flight\_count** – the total number of flights on that route

Results must be **sorted in descending order of flight\_count**.

2. Identify the top 10 Flights with the Most Passenger Feedback. It is important to highlight that “most feedback” does not imply “positive feedback”, these flights simply have the largest volume of responses, regardless of whether the ratings were good or bad. Your query should return the Flight IDs as well as the passenger counts per flight, **sorted in descending order**.
3. Calculate the average food satisfaction score for all multi-leg journeys (journeys with more than one flight leg), grouped by passenger generation. Your query should return:
  - **generation**: the passenger’s generation group
  - **multi\_leg\_count**: how many multi-leg journeys that generation took
  - **avg\_score**: the average food satisfaction score for those journeys

Results must be **sorted in descending order of multi\_leg\_count**.

4. Calculate the average arrival delay for all flights and identify the 10 flights with the shortest delays. Your query should return:
  - **flight\_id**: the unique identifier of the flight
  - **avg\_arrival\_delay**: the average delay in minutes across all journeys for that flight

Results must be **sorted in ascending order of avg\_arrival\_delay**

5. Calculate the average flown miles for passengers, grouped by loyalty program level. Your query should return:
  - **loyalty\_level**: the passenger’s loyalty program category
  - **avg\_actual\_flown\_miles**: the average flown miles for that level

Results must be **sorted in descending order of average flown miles**

## Overall Satisfaction Score Rule Calculation:

Identify the number of passengers with an overall satisfaction score higher than 3, where **overall\_satisfaction\_score** represents the passenger's overall score based on how satisfied they are.  $\text{overall\_satisfaction\_score} = 0.5 * \text{food} + 0.35 * \text{delay\_score} + 0.1 * \text{legs\_score} + 0.05 * \text{miles\_score}$ .

Component	Source in KG	Rule
<b>Food Score</b>	Journey.food_satisfaction_score	Direct value from 1–5
<b>Delay Score</b>	Journey.arrival_delay_minutes	$\text{delay\_score} = 5 - \text{clamp}(\text{abs}(\text{arrival\_delay\_minutes})/2, 0, 5)$ where $(\text{abs}(\text{arrival\_delay\_minutes})/20)$ is expressed as a decimal value rounded to one decimal place.
<b>Legs Score</b>	Journey.number_of_legs	$\text{legs\_score} = 5 - \text{clamp}(\text{number\_of\_legs} * 1.5, 0, 5)$ where $(\text{number\_of\_legs} * 1.5)$ is expressed as a decimal value rounded to one decimal place.
<b>Miles Score</b>	Journey.actual_flown_miles	$\text{miles\_score} = 5 - \text{clamp}(\text{actual\_flown\_miles} / 3000, 0, 5)$ where $(\text{actual\_flown\_miles} / 3000)$ is expressed as a decimal value rounded to one decimal place.

Answer: **1856**

## Deliverables:

You are required to submit the following by filling out the [form](#):

### 1. Create\_kg.py

- Python script that creates the Knowledge Graph in Neo4j.
- Should read CSV files located in the same directory
- Should follow the exact schema specification.
- The script should use the config.txt file with Neo4j credentials. Not doing so might result in 0.

### 2. config.txt

URI=neo4j://localhost:7687

USERNAME=neo4j

PASSWORD=your\_password

### 3. rule.txt

A text file containing the implemented rule for the scoring task.

### 4. queries.txt

A text file containing the queries for the above questions with their numbers.

//1. Query 1

Match(...)

Return (....)

//2. Query 2

Match(..)

Return (....)

Etc..