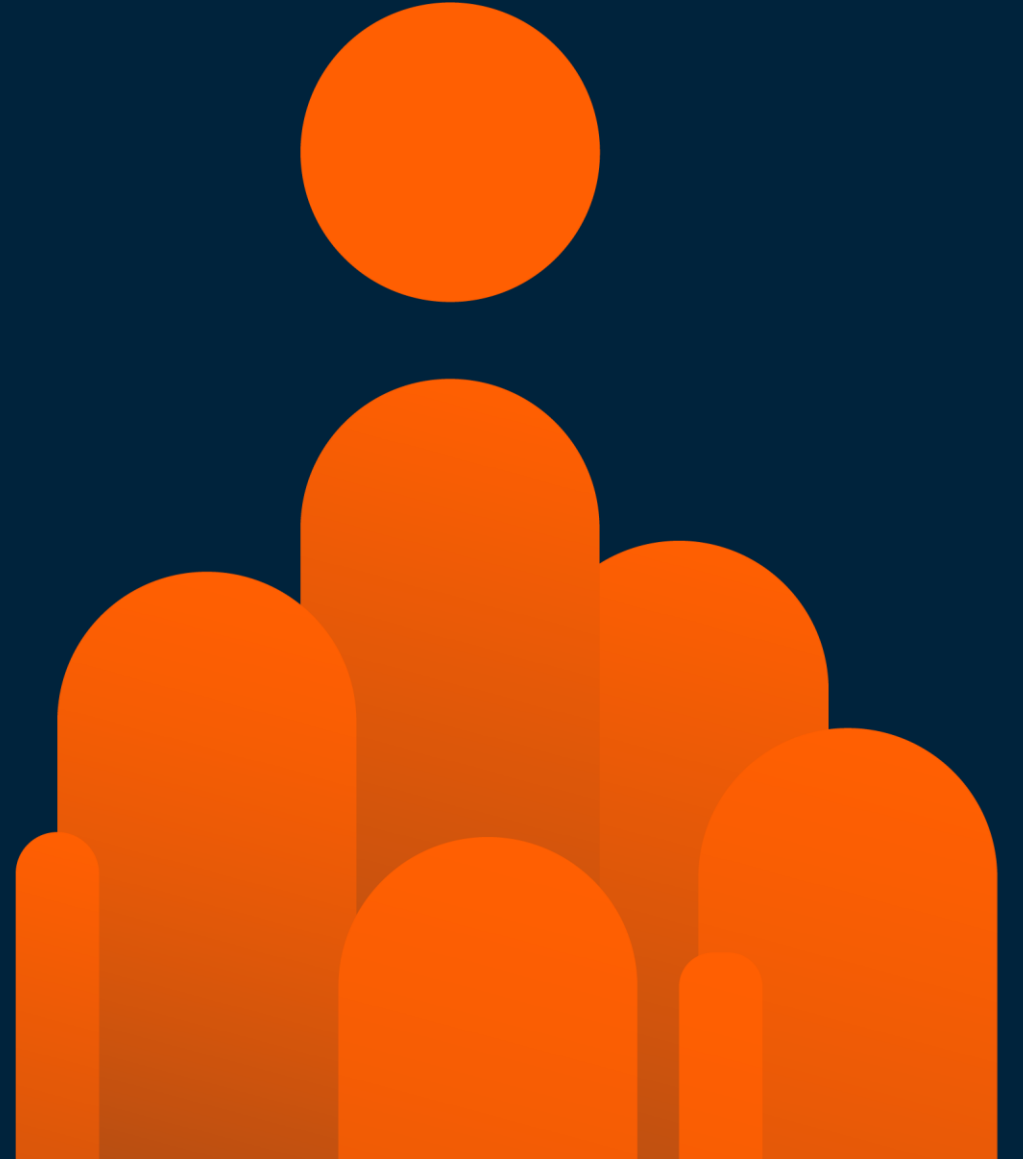27th Feb 2025

# Indian Railways Analysis Use case

**Yash Shah**
**Data Engineering Specialist**

# USECASE DESCRIPTION

## Overview:

The Indian Railways generates vast amounts of data daily, including train details, train schedules, delays, customer satisfaction, and operational performance. Analyzing this data can provide valuable insights to improve efficiency, enhance passenger experience, and optimize resource utilization.

## Objective:

To analyze Indian Railways data to identify trends, improve decision-making, and optimize railway operations by leveraging big data processing frameworks like Azure Data Factory, Databricks, and Delta Lake.

# SOURCE DATASETS DETAILS

| railway_details.csv : | delay_details.json : | satisfaction_details.json : |
|---|---|---|
| Serves as the base dataset for mapping train operations. | Helps in delay pattern analysis and performance improvement. | Provides insights into passenger experience and service quality. |

**Schema:**

| | | |
|---|---|---|
| **Train_id** (String) | **Train_id** (String) | **Train_id** (String) |
| **Train_name** (String) | **Train_name** (String) | **Train_name** (String) |
| **Train_color** (String) | **Arrival_time** (String) | **Seats_available** (String) |
| **Distance** (String) | **Departure_time** (String) | **Cleanliness** (String) |
| **Src_Station_name** (String) | **Delay** (String) | **Status** (String) |
| **Dest_Station_name** (String) | | **Satisfaction** (String) |
| **Frequency** (String) | | |

# TOOLS USED

Azure Data Lake Storage Gen2

Azure Data Factory

databricks

Azure Logic Apps

---

# ACTIVITIES USED

**Lookup**

**Get Metadata**

**Copy**

**Set Variable**

**Fail**

**Web e-mail Notification**

**Databricks Notebook**

**Linked service**

**Dataset**

# USECASE DEMO
# &
# IMPLEMENTATION

# Azure Data Lake Storage Paths



# Databricks Workspace



# Azure Data Factory Pipeline

# Staging Layer :-

## (Source data in Parquet format)

**Location:** ys255066 / Input / delay_details

Search blobs by prefix (case-sensitive)

**Name**

- [ ] 📁 [..]
- [ ] 📄 delay_details.parquet

**Location:** ys255066 / Input / railway_details

Search blobs by prefix (case-sensitive)

**Name**

- [ ] 📁 [..]
- [ ] 📄 railway_details.parquet

**Location:** ys255066 / Input / satisfaction_details
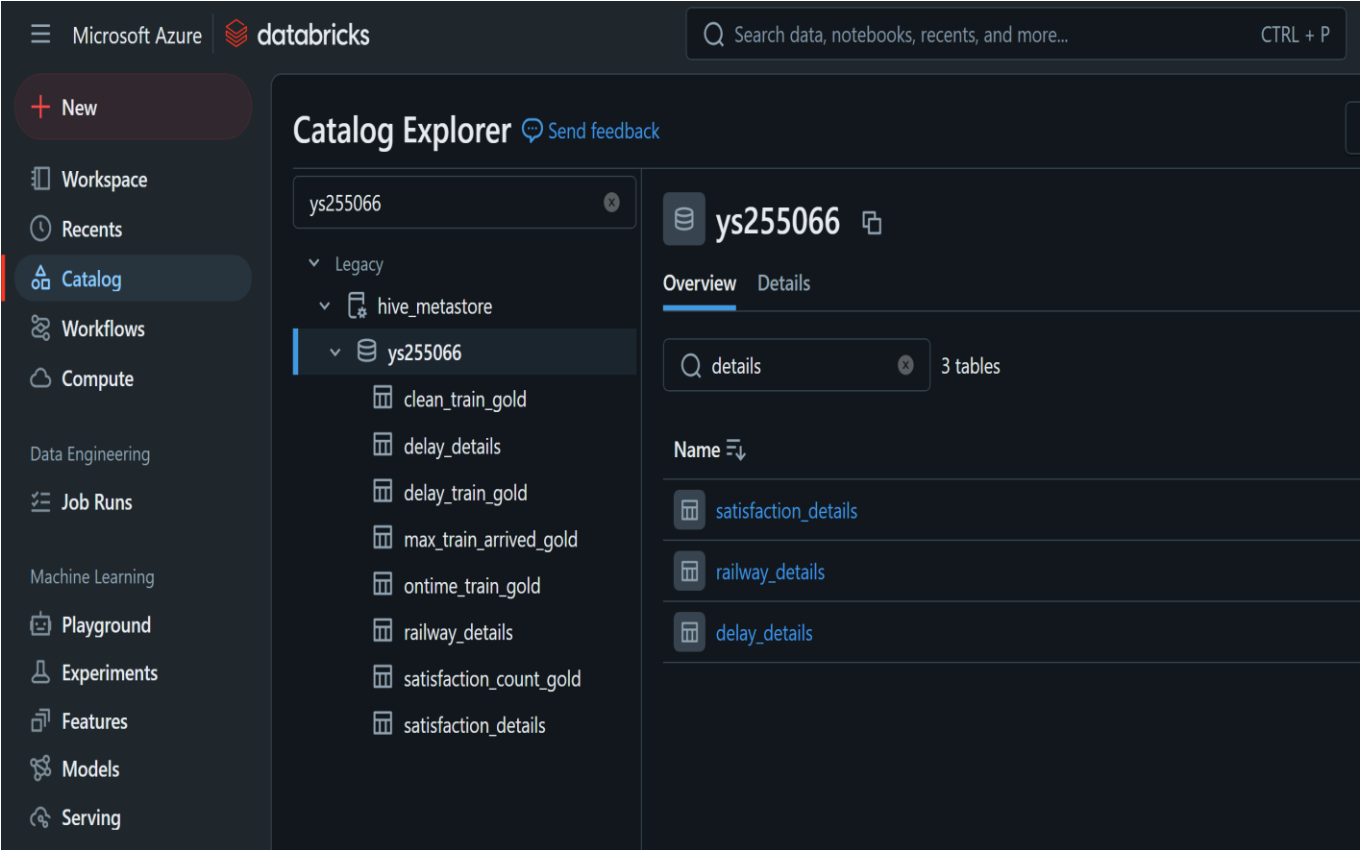
Search blobs by prefix (case-sensitive)

**Name**

- [ ] 📁 [..]
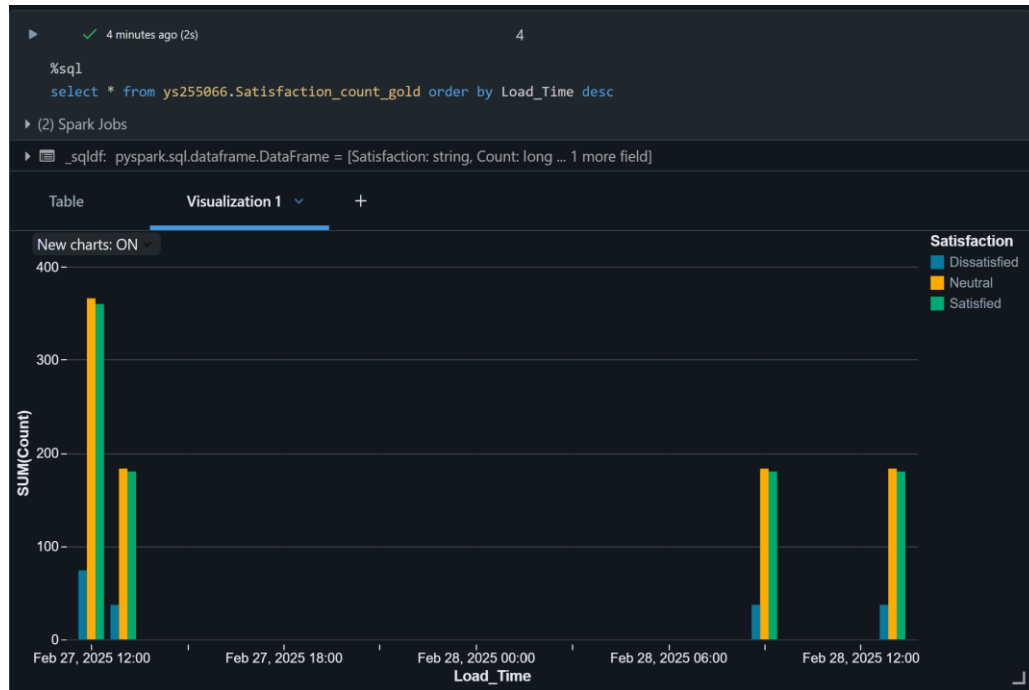- [ ] 📄 satisfaction_details.parquet

# Silver Layer :-

## (Clean & Transformed data in delta tables format)



Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P

+ New

- Workspace
- Recents
- Catalog
- Workflows
- Compute

Data Engineering

- Job Runs

Machine Learning

- Playground
- Experiments
- Features
- Models
- Serving

**Catalog Explorer** 💬 Send feedback

ys255066

⌄ Legacy
  ⌄ 📇 hive_metastore
    ⌄ 🗄 ys255066
      ▦ clean_train_gold
      ▦ delay_details
      ▦ delay_train_gold
      ▦ max_train_arrived_gold
      ▦ ontime_train_gold
      ▦ railway_details
      ▦ satisfaction_count_gold
      ▦ satisfaction_details

🗄 **ys255066**

Overview   Details

details   3 tables

**Name** ⇅

▦ satisfaction_details

▦ railway_details

▦ delay_details

# Gold Layer :- (Stores Aggregated & Business-ready Data)

## 1. Calculate satisfaction count on basis of satisfied and unsatisfied people



## 2. Find top 5 destinations with maximum train arrivals

# Gold Layer :- (Stores Aggregated & Business-ready Data)
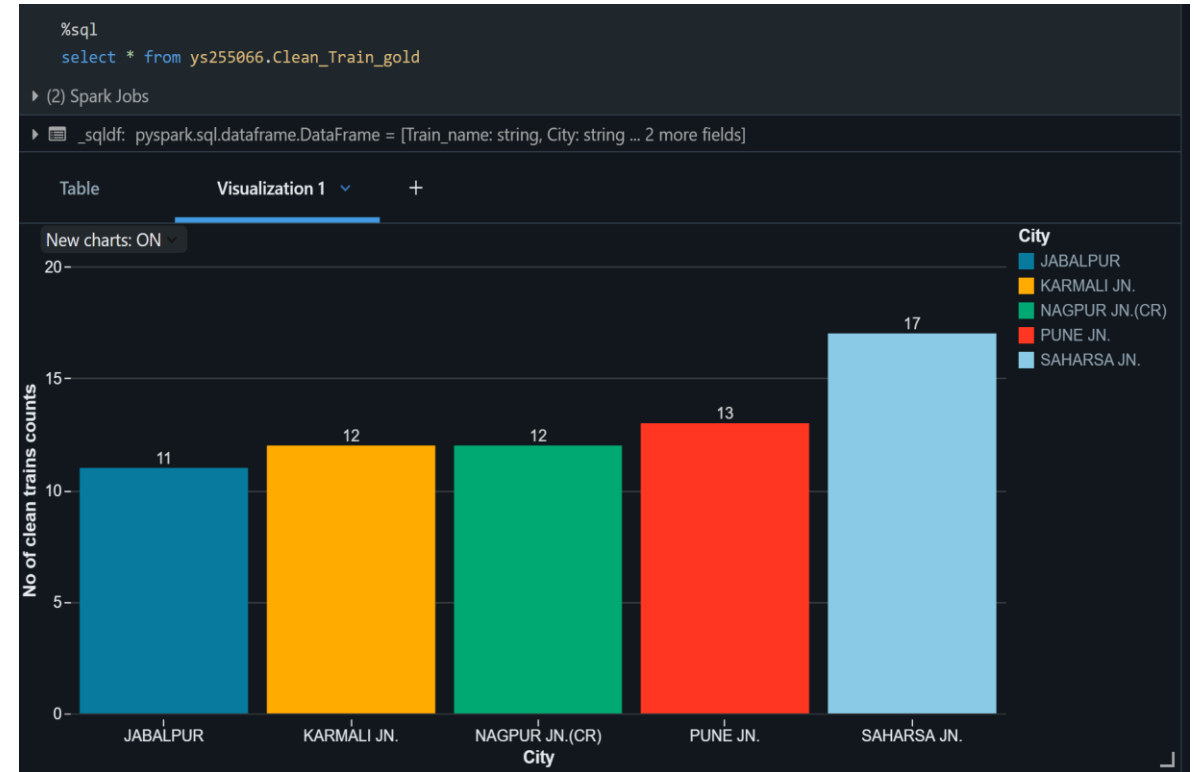
## 3. Analyze no. of trains delayed date wise

```sql
%sql
select * from ys255066.Delay_Train_gold order by load_time desc
```

▶ (1) Spark Jobs

▶ 🗒 _sqldf: pyspark.sql.dataframe.DataFrame = [Total_Trains: long, Delayed_Trains: long ... 3 more fields]

Table ⌄     Visualization 1     +

| | Total_Trains | Delayed_Trains | On_Time_Trains | Delay_Date | Load_Time |
|---|---|---|---|---|---|
| 1 | 400 | 364 | 36 | 2025-02-28 | 2025-02-28T09:05:25.964+00:... |
| 2 | 400 | 364 | 36 | 2025-02-27 | 2025-02-27T13:08:12.351+00:... |
| 3 | 400 | 364 | 36 | 2025-02-27 | 2025-02-27T12:43:41.845+00:... |
| 4 | 400 | 364 | 36 | 2025-02-27 | 2025-02-27T12:37:13.658+00:... |
| 5 | 400 | 364 | 36 | 2025-02-27 | 2025-02-27T12:35:31.273+00:... |
| 6 | 400 | 364 | 36 | 2025-02-27 | 2025-02-27T12:17:11.117+00:... |

## 4. Find top 5 cities with Clean trains

```sql
%sql
select * from ys255066.Clean_Train_gold
```

▶ (2) Spark Jobs

▶ 🗒 _sqldf: pyspark.sql.dataframe.DataFrame = [Train_name: string, City: string ... 2 more fields]

Table     Visualization 1 ⌄     +

New charts: ON ⌄

# Gold Layer :- (Stores Aggregated & Business-ready Data)

## 5. Identify trains that arrived on time

# Pipeline Monitoring With Success & Failure Notifications Over Email.

## Pipeline ran successfully : YS255066_Railway_Analysis_Demo

**Shah, Yash**
To   Shah, Yash

14:37

Retention Policy   Deletion Policy - All Mailboxes (3 Years) (3   Expires   28-02-2028

ⓘ This message was sent with Low importance.

Hi All,

Please find below details of successful pipeline execution for file: railway_details.csv, delay_details.json, satisfaction_details.json

The data factory name: td-aa-trng-adf

The pipeline: YS255066_Railway_Analysis_Demo

## Failed Pipeline YS255066_Railway_Analysis_Demo

**Shah, Yash**
To   Shah, Yash

Thu 19:12

Retention Policy   Deletion Policy - All Mailboxes (3 Years) (3   Expires   27-02-2028

ⓘ This message was sent with High importance.

Hi Team,

Please find below details for error pipeline:

Erro message: Missing Files: delay_details.json on ys255066/Inbound/ path

The Data Factory Name: td-aa-trng-adf

The Pipeline : YS255066_Railway_Analysis_Demo

Regards,