



به نام خدا



دانشگاه تهران

دانشکده مهندسی برق و کامپیوتر

مقدمه ای بر علوم داده

تمرین کامپیوتری شماره ۲

نام و نام خانوادگی	محمد عسکری
شماره دانشجویی	۸۱۰۱۹۸۴۴۱
تاریخ ارسال گزارش	۱۴۰۴/۰۲/۱۱

Real-Time Payment Data Pipeline

Introduction

This report documents the design and implementation of a real-time data pipeline for processing financial transactions, developed as part of the "Introduction to Data Science" course at the University of Tehran. The pipeline simulates transaction flows from a fictional payment provider named "Darooqhe," aiming to handle large-scale streaming data, detect anomalies, and provide actionable insights.

The system integrates Apache Kafka for data ingestion, Apache Spark for both batch and stream processing, MongoDB and PostgreSQL for storage, and Python-based tools for validation and visualization. Each component was configured and tested in a Linux-based environment to emulate real-world conditions.

1. Environment Setup and Data Generation

1.1 Technology Stack

We utilized the following core technologies:

- **Java (OpenJDK 11):** Required for running Kafka and Spark.
- **Apache Kafka:** Acts as the real-time message broker.
- **PySpark:** Python API for Apache Spark, used for distributed processing.
- **Confluent Kafka (Python client):** Enables Kafka interaction from Python.
- **Virtual Environment:** Ensures isolated and reproducible Python dependencies.

1.2 Installation and Configuration

A virtual machine (Linux environment via WSL2) was set up. Java was installed using OpenJDK 11. Kafka was configured using standard scripts to start both the ZooKeeper and Kafka broker services. Spark was installed alongside PySpark.

Kafka topics were created manually and accessed using both CLI and Python scripts. A virtual Python environment was activated to manage the Python packages, including `confluent_kafka`, `pyspark`, `pymongo`, and `psycopg2`.

1.3 Transaction Generator

(`darooqhe_pulse.py`) generates both historical (backfilled) and live events using a Non-Homogeneous Poisson Process.

Key parameters were adjustable through environment variables, including the event rate, peak-hour multiplier, and fraud rate. This script produced 20,000 historical events and then continuously streamed transactions to the Kafka topic `daroooghe.transactions`.

Kafka CLI tools were used to validate that data was properly written to the topic.

A terminal window with four tabs, all labeled 'm4skari@'. The terminal displays the output of a Kafka CLI command, showing a stream of JSON transaction records. The records include fields like 'transaction_id', 'timestamp', 'customer_id', 'merchant_id', 'merchant_category', 'payment_method', 'amount', 'location', 'device_info', 'status', 'commission_type', 'commission_amount', 'vat_amount', 'total_amount', 'customer_type', 'risk_level', and 'failure_reason'. The output shows several examples of transactions, some approved and some declined, with varying amounts and device information. At the bottom, it indicates that 20297 messages were processed.

```
{
  "transaction_id": "ec40c1e7-6990-4330-843b-951f63cb2ca2",
  "timestamp": "2025-04-15T08:16:49.038665Z",
  "customer_id": "cust_4463",
  "merchant_id": "merch_45",
  "merchant_category": "food_service",
  "payment_method": "mobile",
  "amount": 649311,
  "location": {
    "lat": 35.70630654600936,
    "lng": 51.3705068017449
  },
  "device_info": {
    "os": "iOS",
    "app_version": "3.1.0",
    "device_model": "iPhone 15"
  },
  "status": "approved",
  "commission_type": "tiered",
  "commission_amount": 12986,
  "vat_amount": 58437,
  "total_amount": 707748,
  "customer_type": "individual",
  "risk_level": 2,
  "failure_reason": null
}
{"transaction_id": "738dc608-f7b3-4874-afc7-496a659705be",
  "timestamp": "2025-04-15T08:16:51.910443Z",
  "customer_id": "cust_250",
  "merchant_id": "merch_66",
  "merchant_category": "retail",
  "payment_method": "online",
  "amount": 735754,
  "location": {
    "lat": 35.69310489684208,
    "lng": 51.339912555841885
  },
  "device_info": {
    "os": "Android",
    "app_version": "1.9.5",
    "device_model": "Google Pixel 6"
  },
  "status": "approved",
  "commission_type": "tiered",
  "commission_amount": 14715,
  "vat_amount": 66217,
  "total_amount": 801971,
  "customer_type": "individual",
  "risk_level": 1,
  "failure_reason": null
}
{"transaction_id": "26eaa58d-73ed-4fa2-aa05-a5a9cf2db775",
  "timestamp": "2025-04-15T08:16:52.137493Z",
  "customer_id": "cust_4126",
  "merchant_id": "merch_29",
  "merchant_category": "government",
  "payment_method": "pos",
  "amount": 1856686,
  "location": {
    "lat": 35.681901728793825,
    "lng": 51.31688019109848
  },
  "device_info": {},
  "status": "declined",
  "commission_type": "flat",
  "commission_amount": 37133,
  "vat_amount": 167101,
  "total_amount": 2023787,
  "customer_type": "CIP",
  "risk_level": 1,
  "failure_reason": "fraud_prevented"
}
{"transaction_id": "9548c4ca-7419-4f32-8b09-59285b38aa35",
  "timestamp": "2025-04-15T08:16:52.567015Z",
  "customer_id": "cust_2070",
  "merchant_id": "merch_81",
  "merchant_category": "entertainment",
  "payment_method": "online",
  "amount": 1741153,
  "location": {
    "lat": 35.7247898183667,
    "lng": 51.3686399820027
  },
  "device_info": {
    "os": "Android",
    "app_version": "2.4.1",
    "device_model": "Samsung Galaxy S25"
  },
  "status": "declined",
  "commission_type": "tiered",
  "commission_amount": 34823,
  "vat_amount": 156703,
  "total_amount": 1897856,
  "customer_type": "individual",
  "risk_level": 1,
  "failure_reason": "insufficient_funds"
}
{"transaction_id": "a4db3c46-710c-4c26-b5ba-8564a2069a2a",
  "timestamp": "2025-04-15T08:16:52.593358Z",
  "customer_id": "cust_1007",
  "merchant_id": "merch_94",
  "merchant_category": "food_service",
  "payment_method": "nfc",
  "amount": 1631786,
  "location": {
    "lat": 35.70384425004735,
    "lng": 51.32648660188084
  },
  "device_info": {},
  "status": "approved",
  "commission_type": "progressive",
  "commission_amount": 32635,
  "vat_amount": 146860,
  "total_amount": 1778646,
  "customer_type": "individual",
  "risk_level": 3,
  "failure_reason": null
}
^CProcessed a total of 20297 messages
m4skari@LAPTOP-06VNC6E0:~/kafka_2.13-3.5.1$ |
```

Sample output of real-time transaction generator published to `daroooghe.transactions` Kafka topic.

2. Data Ingestion and Validation

2.1 Kafka Consumer and Validator

The `transaction_validator.py` script acts as a Kafka consumer, reading from the `darrooghe.transactions` topic and applying a set of validation rules before forwarding invalid records to the `darrooghe.error_logs` topic.

2.2 Validation Logic

Each transaction was checked against three main rules:

- **Amount Consistency:** Ensures `total_amount = amount + vat_amount + commission_amount`. Violations trigger `ERR_AMOUNT`.
- **Time Warping:** Filters out transactions from the future or older than one day (`ERR_TIME`).
- **Device Mismatch:** Checks that mobile transactions originate from iOS or Android devices (`ERR_DEVICE`).

Invalid records were serialized into JSON format and sent to the error topic. Console outputs and Kafka CLI confirmed correct logging behavior.

```
m4sl x m4sl x m4sl x m4sl x m4sl x + - □ ×
✓ Valid TXN: 8ccafc1f-d8d6-4b1d-b545-e34ce48cd663
✗ Invalid TXN: b20c4805-781a-4d60-ae2d-f6407cf41085 | Reason: ERR_AMOUNT
✓ Valid TXN: 5e9ba027-4ebd-47e6-ab5d-4f635e7a53e9
✗ Invalid TXN: c929e12b-245a-45fa-93a9-1350ffa9d731 | Reason: ERR_AMOUNT
✗ Invalid TXN: 8e8548ed-6ed8-44fb-8b59-1727bdb08296 | Reason: ERR_AMOUNT
✗ Invalid TXN: 939adde6-6428-432f-8f96-062877b2f4f4 | Reason: ERR_AMOUNT
✓ Valid TXN: bda54116-4da7-41a6-96fa-b574c5881b66
✗ Invalid TXN: b39138a3-c729-4fdf-8f95-0fd797c54f1c | Reason: ERR_AMOUNT
✓ Valid TXN: d8dda8c8-d722-477a-b710-8cdac364ceb6
✗ Invalid TXN: 8ccb158-b227-48b0-afee-75482bda52b4 | Reason: ERR_AMOUNT
✗ Invalid TXN: c1fd9bc8-1144-407b-9f83-ff4af2a16d16 | Reason: ERR_AMOUNT
✓ Valid TXN: ec3d3f82-a87b-41e5-a4df-fc423ef06cda
✓ Valid TXN: c19f9afb-0e75-471a-90dc-471d7221ca87
✓ Valid TXN: fc1df36f-eeab-4b23-8cc6-ac85151d38fb
✓ Valid TXN: 88f037c6-d8d0-4fd4-8404-599ea2492d99
✗ Invalid TXN: 7a3ec723-f6bc-4bb0-8032-e7ee06555f34 | Reason: ERR_AMOUNT
✓ Valid TXN: 77aa21cd-9a52-41d1-be34-155a83f93213
✓ Valid TXN: afe51c25-b445-4c63-a62e-d6039f1f86bd
✓ Valid TXN: 9d752178-c4d1-4075-8c2c-9359560f58d9
✗ Invalid TXN: 2b401139-be9e-4220-9850-52da1da351f6 | Reason: ERR_AMOUNT
✗ Invalid TXN: fdf11db0-b6e6-4df2-9f6e-7e5e8d478e98 | Reason: ERR_AMOUNT
✓ Valid TXN: 7edd44ec-5962-4642-a15d-5e90eea8fffd
✓ Valid TXN: daf5da2f-ec9a-4659-a9bd-5bdd71ebda9d
✓ Valid TXN: 39e16ec5-9819-423a-be5d-09e84cce0b72
✓ Valid TXN: 0eb9595f-5254-4d45-8436-196af213b990
✓ Valid TXN: aa58fde5-8613-46fa-aac8-2bdd1fdcb7fc
✓ Valid TXN: a46775da-1d1d-451e-b784-3dbf8389251e
✗ Invalid TXN: 95bd39a4-a72c-4d27-bc05-a58095fe9f1b | Reason: ERR_AMOUNT
✗ Invalid TXN: 198f04c7-b137-4118-8d9b-2c88d7d83068 | Reason: ERR_AMOUNT
✗ Invalid TXN: 83717ffa-9e68-4afb-9784-46227f13708a | Reason: ERR_AMOUNT
✗ Invalid TXN: b89083bc-d72d-43a1-b0ff-0bcda03c492a | Reason: ERR_AMOUNT
✓ Valid TXN: 85ca6b2f-9627-4c52-8a9b-de3929277666
✓ Valid TXN: 92fe6c9a-6942-42c3-92bd-4323e9ff65ea
✓ Valid TXN: 0369b30b-8d39-468c-81e2-73681076cbe4
✓ Valid TXN: 839869a0-4168-4323-95bf-147f0ba18414
✓ Valid TXN: 54374f4f-8b07-4698-98d9-ddd18d2a044e
✓ Valid TXN: f77c6ddd-e280-4aae-8db0-db6630c80104
✗ Invalid TXN: 73c23aa1-9004-4ff3-96d2-bd43981b0f5a | Reason: ERR_AMOUNT
✓ Valid TXN: 0ffa8838-0bcb-468a-ad77-77efb9a5bc34
✗ Invalid TXN: 0cfbc8a7-6fed-48e0-894a-2d2752676c94 | Reason: ERR_AMOUNT
```

Real-time output of Kafka consumer showing valid and invalid transactions with associated error codes.

3. Batch Processing and Analysis

3.1 Objective

Batch analysis focused on identifying patterns in historical data, such as temporal trends and customer segmentation.

3.2 Dataset and Tools

The dataset was exported from Kafka into a JSON file and loaded into PySpark. Tools used include PySpark DataFrame API, SQL functions, and timestamp transformation utilities.

3.3 Temporal Analysis

Using Spark SQL, transaction frequency was aggregated by hour of day and day of week. Results showed:

- Peak activity around noon and 19 PM.
- Highest transaction volume on Tuesday.

```
=== Most Active Hours ===
+-----+-----+
|hour_of_day|count|
+-----+-----+
|0          |830  |
|1          |816  |
|2          |849  |
|3          |824  |
|4          |825  |
|5          |879  |
|6          |828  |
|7          |898  |
|8          |819  |
|9          |823  |
|10         |878  |
|11         |834  |
|12         |7214 |
|13         |817  |
|14         |816  |
|15         |797  |
|16         |818  |
|17         |852  |
|18         |821  |
|19         |7061 |
|20         |4142 |
|21         |877  |
|22         |801  |
|23         |804  |
+-----+-----+
```

```

+-----+-----+
=== Most Active Days (Sunday=1) ===
+-----+-----+
|day_of_week|count|
+-----+-----+
|1          |2960 |
|2          |2863 |
|3          |18783|
|4          |2837 |
|5          |2858 |
|6          |2834 |
|7          |2788 |
+-----+-----+

```

3.4 Customer Segmentation

Customers were grouped based on their transaction frequency:

- **Frequent Users:** ≥ 10 transactions.
- **Normal Users:** 3–9 transactions.
- **Rare Users:** < 3 transactions.

The segmentation logic was implemented using the `when()` and `otherwise()` functions in PySpark. This information can support downstream marketing or anomaly detection.

```

=== Customer Segmentation by Activity ===
+-----+-----+-----+
|customer_id|txn_count|segment|
+-----+-----+-----+
|cust_1314|18|Frequent|
|cust_3234|17|Frequent|
|cust_1979|17|Frequent|
|cust_352|17|Frequent|
|cust_3890|17|Frequent|
|cust_1476|16|Frequent|
|cust_3624|16|Frequent|
|cust_3614|16|Frequent|
|cust_1618|16|Frequent|
|cust_905|16|Frequent|
|cust_2384|16|Frequent|
|cust_2634|16|Frequent|
|cust_2673|16|Frequent|
|cust_2779|15|Frequent|
|cust_1832|15|Frequent|
|cust_2616|15|Frequent|
|cust_2434|15|Frequent|
|cust_3882|15|Frequent|
|cust_1495|15|Frequent|
|cust_1161|15|Frequent|
+-----+-----+-----+
only showing top 20 rows

```

4. Real-Time Streaming and PostgreSQL Integration

4.1 Spark Structured Streaming

We built a real-time Spark streaming job (`realtime_processor.py`) to process Kafka data and apply simple fraud detection rules. The rules included:

- Late-night transactions (00:00–06:00).
- High-risk transactions with large amounts.
- Device anomalies in mobile transactions.

4.2 Database Sink

A checkpointing mechanism was implemented to maintain fault tolerance. The final schema included transaction ID, timestamp, risk level, and alert reason.

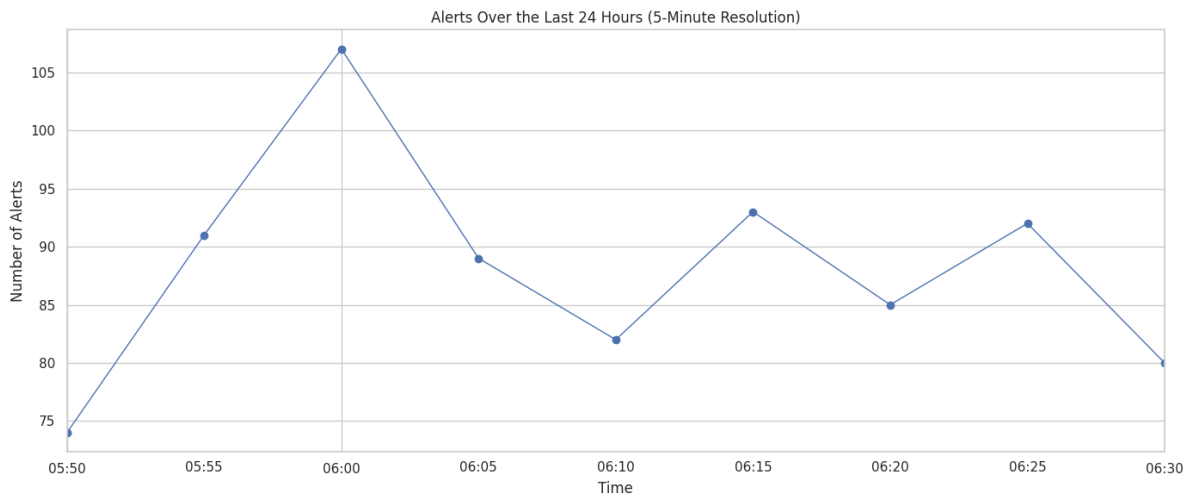
5. Data Visualization

5.1 Objective

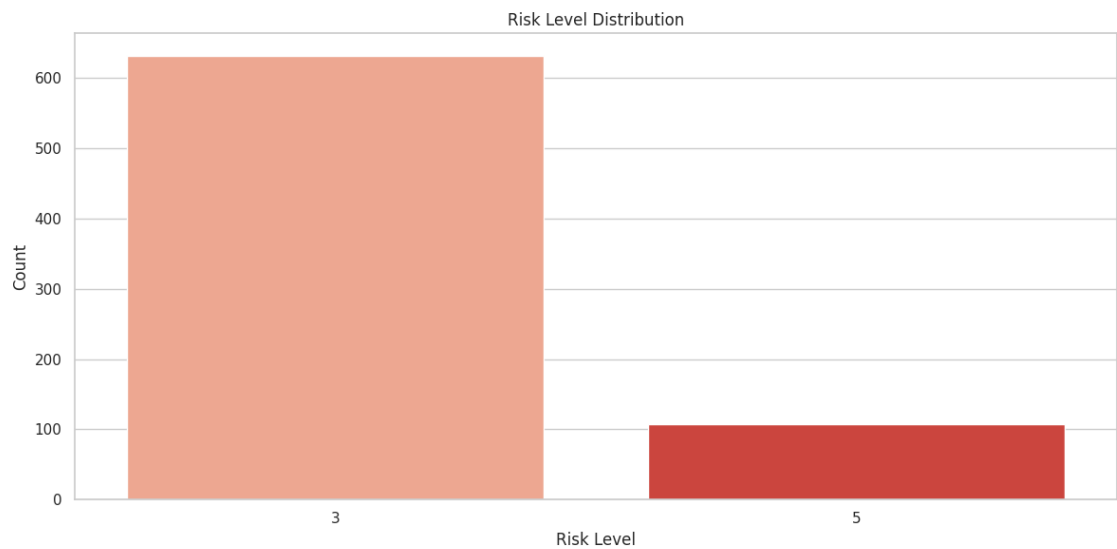
5.2 Tools

5.3 Charts and Insights

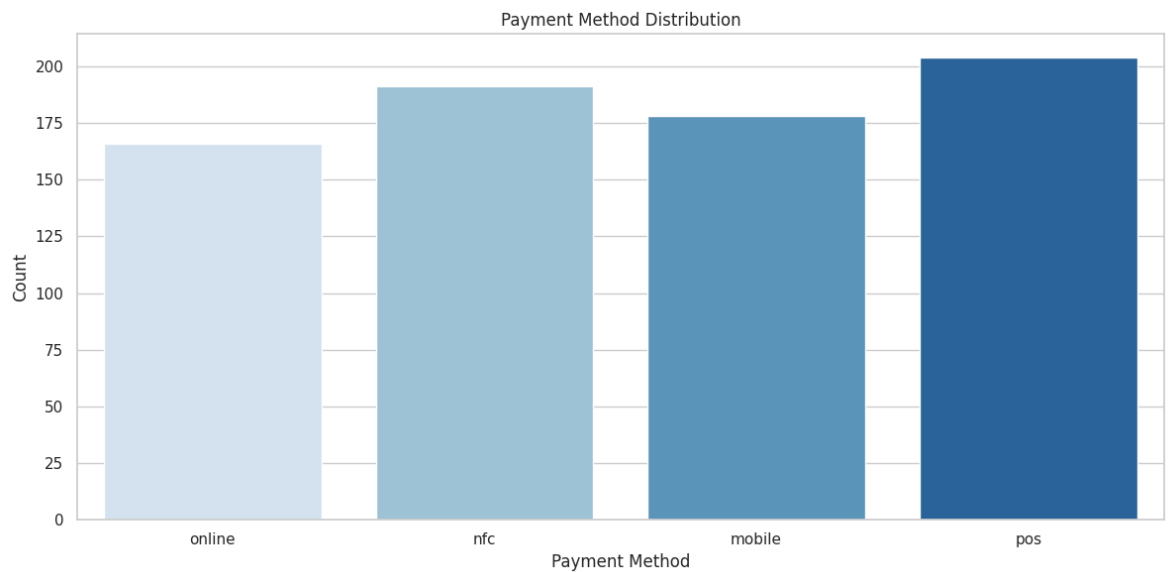
- **Alerts Over Time:** Time-series plot showed real-time fraud activity.



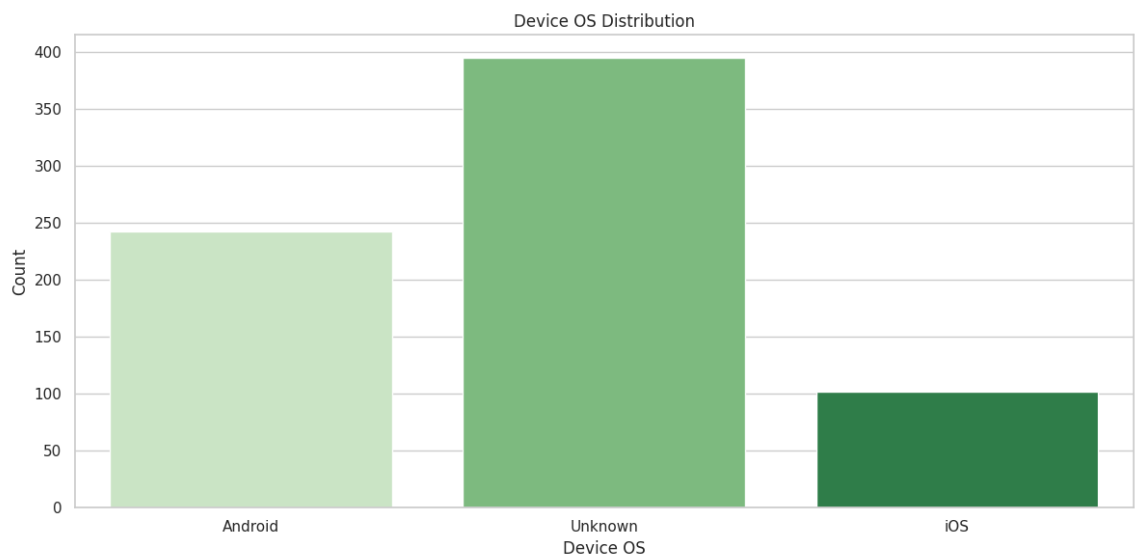
- **Risk Level Distribution:** High concentration in risk level 3.



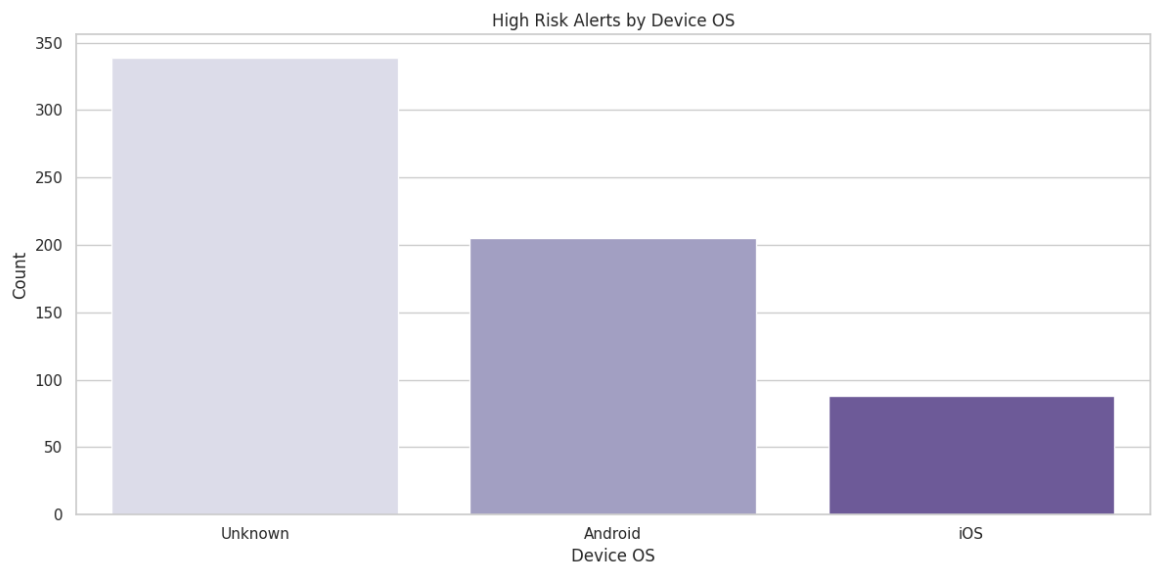
- **Payment Method Distribution:** Fraudulent activity was spread across POS, NFC, and Mobile.



- **Device OS Breakdown:** Android was the most common among flagged mobile transactions.



- **High-Risk Devices:** Certain OS types were repeatedly linked to high-risk alerts.



Conclusion

This project successfully demonstrated the architecture and implementation of a real-time data pipeline for transaction processing. By leveraging Kafka, Spark, and PostgreSQL, the system supports continuous ingestion, validation, fraud detection, and visualization.

Each component was modular and extensible, making the pipeline suitable for scaling in production environments. Potential improvements include advanced fraud detection rules, dynamic pricing strategies, and real-time dashboards.

The pipeline met its functional goals and proved to be a reliable platform for streaming financial data analytics.