# SKaMP. Tests

The goal of this report is to provide information on the performed testing of data acquisition, data pre-processing, batch processing and streaming processing.
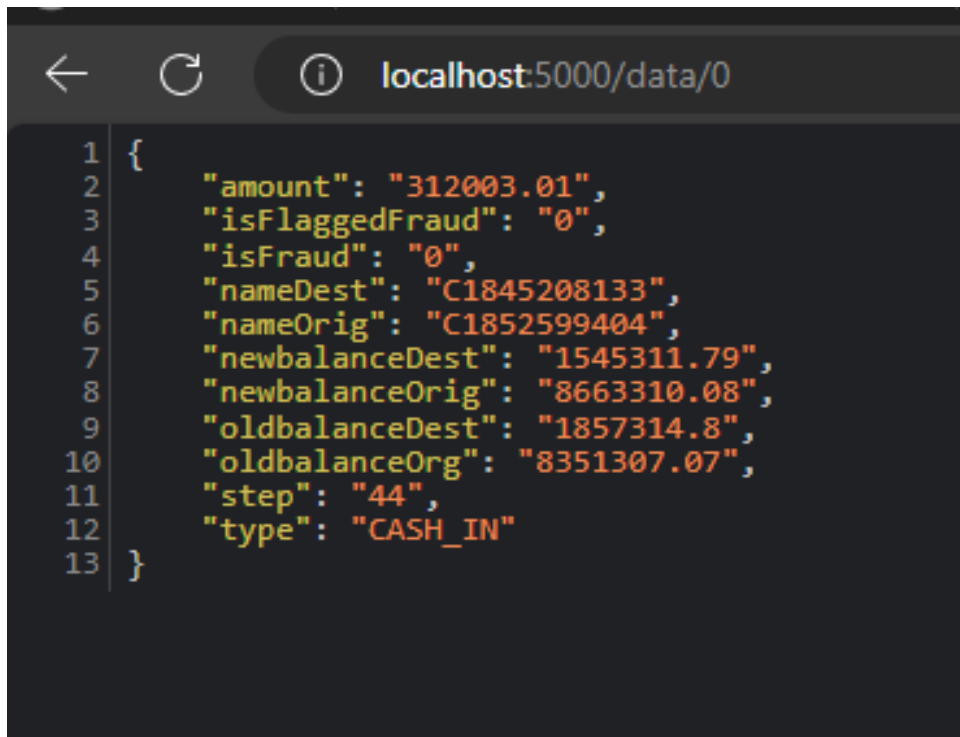
GitHub repository: https://github.com/salveendutt/Big-Data-Analytics.

# 1   Test Scenarios

| Test objective | Steps | Expected Result | Actual Result |
|---|---|---|---|
| Verify data incoming from stream API | 1. Start the server using start_containers.bat; 2. Navigate to http://localhost:5000 | Incoming data is available on /data/0 | Passed. The screenshot is provided in Figure 1 |
| Verify correct setup of the stream and data preprocessing functions | Run 'pytest' from the root folder | Data stream is configured as expected; Incoming data is not null; Returned status code - 200. Preprocessing utils return transformed data as expected | Passed. The screenshot is provided in Figure 10 |
| Verify the correct setup of Nifi - HDFS/Kafka flow | Run the containers - follow steps in README.md | Data flows from streamin API to Kafka topics and Hive tables | Passed. The screenshot is provided in Figure 2, 3, 4, 5 |
| Verify the correct setup of batch processing | Run the containers - follow steps in README.md | Views are available in the Cassandra tables | Passed. The screenshot is provided in Figure 6, 7, 8 |
| Verify the correct setup of streaming processing | Run the containers - follow steps in README.md | Views are available in the Cassandra tables | Passed. The screenshot is provided in Figure 9 |

Table 1: Test scenarios

| Test objective | Steps | Expected Result | Actual Result |
|---|---|---|---|
| Verify correct data pre-processing of dataset 1 | Run 'pytest' from the root folder | Feature 'type' is correctly transformed into numeric value (5 cases); Feature 'is-Merchant' is correctly prepared (2 cases) | PASSED. The screenshot is provided in Figure 10 |
| Verify correct data pre-processing of dataset 2 | Run 'pytest' from the root folder | Numeric boolean values are transformed to int from float (4 cases) | PASSED. The screenshot is provided in Figure 10 |
| Verify correct data pre-processing of dataset 3 | Run 'pytest' from the root folder | Feature 'entry_mode' is correctly transformed into numeric value (4 cases); Unnecessary features are omitted. | PASSED. The screenshot is provided in Figure 10 |
| Verify correct data pre-processing of dataset 4 | Run 'pytest' from the root folder | Features 'Amount', 'Class' are renamed to 'amount' and 'is-Fraud'; Extra features are removed | PASSED. The screenshot is provided in Figure 10 |

Table 2: Data pre-processing tests



Figure 1: Data incoming via the stream

Figure 2: Kafka Dataset1



Figure 3: Kafka Dataset2

Figure 4: Kafka Dataset3



Figure 5: Hive data



Figure 6: Cassandra batch processing views

Figure 7: Cassandra batch processing views



Figure 8: Cassandra batch processing views
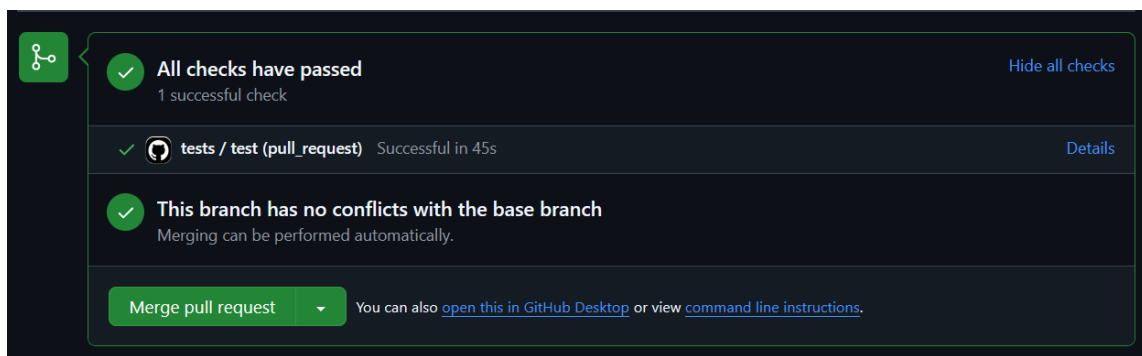


Figure 9: Cassandra stream processing views



Figure 10: Unit testing result

Figure 11: GitHub checks before merge