# Tests Description

| Test Name | Objective | Steps | Expected Results | Actual results |
|---|---|---|---|---|
| Sentiment analysis | Check if sentiment analysis works correctly. | 1. Ensure that subscribed topic in kafka contains messages with articles<br>2. Run spark code with sentiment analysis code | New columns with with averaged sentiment in table is created | ```
----------------------------------------
Batch: 3
----------------------------------------
+--------------------------------------------+--------------+-------------------+
|window                                      |avg_sentiment |win_start          |
+--------------------------------------------+--------------+-------------------+
|{2023-12-18 20:10:00, 2023-12-18 20:15:00}|5.676335010691815E11|2023-12-18 20:10:00|
+--------------------------------------------+--------------+-------------------+
``` |
| Streaming rates | Check if streaming data from kafka to spark works correctly | 1. Ensure that NiFI processors from crypto_rates _api_to_kafk a are running<br>2. Run script test_streami ng_rate.py, which is a copy of our solution that prints stream to terminal instead of joining data and dropping | Within 30 seconds after running the program spark printed table with columns: id, symbol, changePercent24H r, and priceUSD. This table also can't be empty. | ```
----------------------------------
Batch: 1
----------------------------------
+-------+------+-------------------+----------------+--------------------+-------------------------------------------+
|id     |symbol|changePercent24Hr  |priceUsd        |current_timestamp   |window                                     |
+-------+------+-------------------+----------------+--------------------+-------------------------------------------+
|bitcoin|BTC   |-0.8768706096683539|41806.72134461311|2023-12-18 20:38:23.027|{2023-12-18 20:35:00, 2023-12-18 20:40:00}|
+-------+------+-------------------+----------------+--------------------+-------------------------------------------+
``` |

| | | | | |
|---|---|---|---|---|
| | | it to cassandra<br>3. Ensure that within 30 seconds after running the program spark printed table with columns: id, symbol, changePercent24Hr, and priceUSD. Also ensure that this table isn't empty. | | |
| Streaming rates | Check if streaming data from kafka to spark works correctly | 1. Ensure that NiFI processors from NewsAPI_news_api_to_kafka are running<br>2. Run script simulate_news.py, which takes newest article from | Within 15 minutes after running the program spark printed non empty table with expected columns |  |



```
-------------------------------------------
Batch: 3
-------------------------------------------
+--------------------------------------------+------------------+---------------------+
|window                                      |avg_sentiment     |win_start            |
+--------------------------------------------+------------------+---------------------+
|{2023-12-18 20:10:00, 2023-12-18 20:15:00}|5.676335010691815E11|2023-12-18 20:10:00|
+--------------------------------------------+------------------+---------------------+
```

| | | | | |
|---|---|---|---|---|
| | | kafka topic waits 30 seconds and puts it  back simulating data ingestion<br><br>3.  Run script test_streaming_news.py, which is a copy of our solution that prints stream to terminal instead of joining data and dropping it data to cassandra<br><br>4.  Ensure that within 15 minutes after running the program spark printed table with columns: window, sentiment. Also ensure | | |

| | | that this table isn't empty. | | |
|---|---|---|---|---|
| Machine Learning predictions | Checking if models generate predictions and if they are saved to Cassandra | 1. Ensure that data is retrieved from Kafka.<br><br>2. Ensure that models generate correct predictions (not a constant value).<br><br>3. Check if data is saved to Cassandra. | After a few minutes after running a script, the predictions should be available in Cassandra. |  |
| Rates and sentiment join | Checking if joining streams of crytocurrency rates and aggragated sentiment from the articles on window works correctly. | 1. Run the script for calculating sentiment from incoming news articles<br><br>2. Run the script for reading the crypto rates stream | Results of successful join are saved to kafka topic. (new column avg_sentiment and window in rates data) |  |

| | | 3. Run the script for joining the stream ans saving it to kafka topic | | |
|---|---|---|---|---|
| Raw ingestion | Checking if new currency rates are ingested into HDFS using NiFi and if API works correctly | 1. Ensure that NiFi process group called crypto_rates _api_to_mas ter is working<br><br>2. Edit script test_raw_in_ hdfs.py changing the date to current one<br><br>3. Run script test_raw_in_ hdfs.py it will return the name of the newest file in HDFS. Remember this name<br><br>4. Run script again no sooner than 1 minute | After running program again after at least 1 minute name of newest file will be different | Results after running program once<br><br>```\n>>> print(paths[-1:])\n['hdfs://namenode:9000/master_dataset/crypto_rates/2024-01-14/crypto_rates_2024-01-14-13-09-00.avro']\n```<br><br>Results after running program one minute later<br><br>```\n>>> print(paths[-1:])\n['hdfs://namenode:9000/master_dataset/crypto_rates/2024-01-14/crypto_rates_2024-01-14-13-10-00.avro']\n``` |

| | | later. Check if the name of newest file has changed. | | |
| --- | --- | --- | --- | --- |