

TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI
TRƯỜNG CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG



BÁO CÁO BÀI TẬP LỚN

Đề tài: Lưu trữ và xử lý dữ liệu tiền ảo theo thời gian thực

Lớp: 151392

Học phần: Lưu trữ và xử lý dữ liệu lớn

Mã học phần: IT4931

Giảng viên hướng dẫn: TS. Trần Văn Đặng

Danh sách thành viên nhóm 5:

Họ và tên	Mã số sinh viên
Nguyễn Quốc Dũng	20215329

Hà Nội, tháng 6 năm 2024

CHƯƠNG 1. GIỚI THIỆU BÀI TOÁN

1.1 Giới thiệu

Trong thế giới tài chính hiện đại, tiền ảo hay crypto là một dạng tiền điện tử rất còn mới mẻ và đang đà phát triển. Việc lưu trữ và xử lý dữ liệu tiền ảo thời gian thực là rất quan trọng. Điều này không chỉ giúp các nhà đầu tư có được cái nhìn sâu sắc và kịp thời về thị trường, mà còn là cơ sở để phát triển các chiến lược giao dịch và quản lý rủi ro hiệu quả. Báo cáo này sẽ trình bày phương pháp và các công cụ lưu trữ dữ liệu lớn trong quá trình xây dựng hệ thống xử lý và phân tích dữ liệu tiền ảo trong thời gian thực.

1.2 Vấn đề

Giá tiền ảo thay đổi trong thời gian thực và thường khó dự đoán bằng các phương pháp truyền thống, tạo ra một thách thức lớn cho các nhà đầu tư và nhà phân tích. Sự biến động không thể dự đoán của giá tiền ảo là do nhiều yếu tố - từ tin tức kinh tế, báo cáo tài chính, đến hành động của các nhà đầu tư khác. Hơn nữa, các yếu tố dẫn đến thay đổi giá cổ phiếu phức tạp và đa dạng, bao gồm cả các yếu tố cơ bản của công ty và yếu tố kỹ thuật trên thị trường. Điều này đòi hỏi cần có phương pháp tiên tiến và đa dạng hơn trong việc lưu trữ và phân tích dữ liệu để có thể hiểu và dự báo chính xác các xu hướng giá cổ phiếu, đặc biệt là trong môi trường giao dịch nhanh và biến động cao của ngày nay.

CHƯƠNG 2. CÁC NGUỒN DỮ LIỆU & CẤU TRÚC

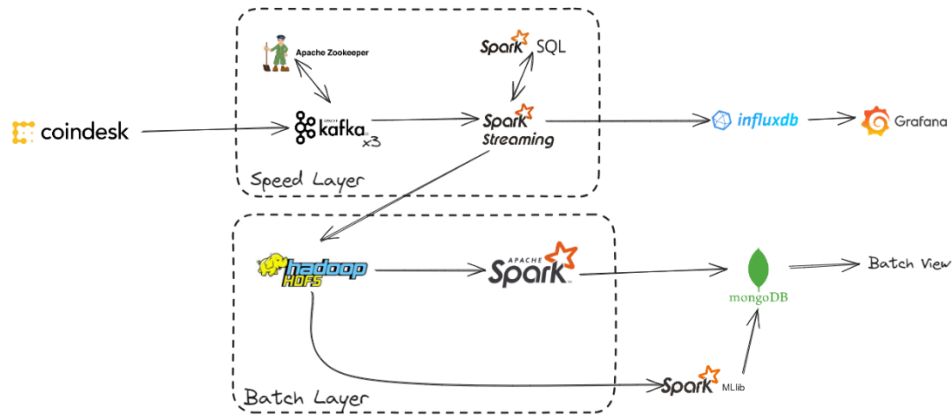
Nguồn dữ liệu cho việc phân tích và xử lý tiền thời gian thực trong dự án này chủ yếu đến từ nguồn là CoinDesk và dữ liệu được crawl trực tiếp từ API của trang web. Dữ liệu này bao gồm thông tin chi tiết về các mã chứng khoán, được cấu trúc như sau:

- ISO: Mã ISO của đồng tiền.
- Name: Tên đồng tiền.
- Date_time: Thời gian xuất hiện của bản ghi.
- Current_price: Giá hiện tại.
- Open: Giá mở cửa.
- High: Giá cao nhất.
- Low: Giá thấp nhất.
- Close: Giá đóng cửa

Thông tin từ nguồn dữ liệu này đóng vai trò quan trọng trong việc cung cấp dữ liệu chính xác và cập nhật cho hệ thống phân tích.

CHƯƠNG 3. KIẾN TRÚC HỆ THỐNG

Phân tích tiền ảo thời gian thực đòi hỏi một hệ thống có khả năng xử lý dữ liệu lớn một cách nhanh chóng và chính xác. Chương này sẽ mô tả chi tiết về kiến trúc hệ thống và các công nghệ được sử dụng.



3.1 Thu thập dữ liệu

Dữ liệu được thu thập real-time từ CoinDesk. Sau đó được chuyển đến Kafka, một nền tảng xử lý stream-data, để quản lý và phân phối dữ liệu đầu vào.

3.2 Tầng xử lý dữ liệu

Tầng xử lý dữ liệu được chia thành hai lớp chính: Batch Layer và Speed Layer.

3.2.1 Speed Layer

Lớp Speed xử lý dữ liệu thời gian thực sử dụng Spark Streaming và Spark SQL để cung cấp thông tin cập nhật liên tục và nhanh chóng. Dữ liệu sau khi transform được chuyển đến InfluxDB và Hadoop HDFS để lưu trữ.

3.2.2 Batch Layer

Lớp Batch xử lý dữ liệu lớn thông qua Hadoop HDFS và Apache Spark để thực hiện các tính toán theo lô (batch processing) với tần suất là 1 lần / ngày. Kết quả sau khi xử lý xong sẽ được lưu trữ vào MongoDB.

3.3 Phân tích và trực quan hoá

Để phân tích dữ liệu, MLlib được sử dụng để cài đặt mô hình học máy với dữ liệu được lưu trữ tại Hadoop HDFS.

Kết quả cuối cùng, đối với Speed Layer, dữ liệu tại InfluxDB sẽ được trực quan hóa dưới dạng stream thông qua Grafana. Đối với Batch Layer dữ liệu từ MongoDB sẽ được hiển thị dưới dạng Batch View vào cuối mỗi ngày

Như vậy, hệ thống sử dụng kết hợp giữa xử lý theo lô (batch processing) và xử lý theo dòng dữ liệu (stream processing). Với mục tiêu, cung cấp cái nhìn sâu sắc và cập nhật liên tục về thị trường tiền ảo.

CHƯƠNG 4. CHI TIẾT HƠN VỀ DỰ ÁN

4.1 Thu thập dữ liệu

```
URL = 'https://production.api.coindesk.com/v2/tb/price/ticker?assets='

2 usages  J Johnson +1

def default(args):
    if len(args) == 0:
        raise Exception("Missing argument...")

    final_url = URL + args
    result = get(final_url)
    json_data = result.json()

    individual_symbols = args.split(',')

    results = []

    for s in individual_symbols:
        coin_data = json_data['data'][s]
        results.append({
            "iso": coin_data['iso'],
            "name": coin_data['name'],
            "date_time": datetime.now().strftime("%Y-%m-%d %H:%M:%S%z"),
            "current_price": coin_data['ohlc']['c'],
            "open": coin_data['ohlc']['o'],
            "high": coin_data['ohlc']['h'],
            "low": coin_data['ohlc']['l'],
            "close": coin_data['ohlc']['c']
        })

    return results
```

Hình 4.1 Code crawl dữ liệu

```
{'iso': 'BTC', 'name': 'Bitcoin', 'date_time': '2024-06-16 18:21:24', 'current_price': 66589.04965376819, 'open': 66269.23549705418, 'high': 66681.81572379959, 'low': 65846.1},
{'iso': 'ETH', 'name': 'Ethereum', 'date_time': '2024-06-16 18:21:24', 'current_price': 3560.8593261141546, 'open': 3532.42330951606, 'high': 3589.6116053697947, 'low': 3532},
{'iso': 'ETHFI', 'name': 'ether.fi', 'date_time': '2024-06-16 18:21:24', 'current_price': 3.9740091277, 'open': 3.8179212225, 'high': 4.0711812159, 'low': 3.7963240585, 'close': 3.9740091277},
{'iso': 'DOGE', 'name': 'Dogecoin', 'date_time': '2024-06-16 18:21:24', 'current_price': 0.1362600513, 'open': 0.1367315336, 'high': 0.13759965, 'low': 0.1348752542, 'close': 0.1362600513},
{'iso': 'ZETA', 'name': 'ZetaChain', 'date_time': '2024-06-16 18:21:24', 'current_price': 1.0400514633, 'open': 1.0366482172, 'high': 1.0501979546, 'low': 1.0030871901, 'close': 1.0400514633},
{'iso': 'BNB', 'name': 'BNB', 'date_time': '2024-06-16 18:21:24', 'current_price': 608.0293731664, 'open': 608.577255253, 'high': 610.6127641206, 'low': 601.958141428, 'close': 608.0293731664},
{'iso': 'SHIB', 'name': 'Shiba Inu', 'date_time': '2024-06-16 18:21:24', 'current_price': 2.07934e-05, 'open': 2.08183e-05, 'high': 2.09991e-05, 'low': 2.04228e-05, 'close': 2.07934e-05},
{'iso': 'SOL', 'name': 'Solana', 'date_time': '2024-06-16 18:21:24', 'current_price': 145.14669157718734, 'open': 143.8844168786579, 'high': 145.55328437014848, 'low': 143.05, 'close': 145.14669157718734},
{'iso': 'TON', 'name': 'Toncoin', 'date_time': '2024-06-16 18:21:24', 'current_price': 7.9735594273, 'open': 8.1217530927, 'high': 8.1910802899, 'low': 7.7492448139, 'close': 7.9735594273},
{'iso': 'WIF', 'name': 'dogwifhat', 'date_time': '2024-06-16 18:21:24', 'current_price': 2.540315364, 'open': 2.4087451698, 'high': 2.5660305859, 'low': 2.3102194783, 'close': 2.540315364},
{'iso': 'OKB', 'name': 'OKB', 'date_time': '2024-06-16 18:21:24', 'current_price': 46.1476339617, 'open': 45.5640742856, 'high': 46.4435857202, 'low': 45.5269185535, 'close': 46.1476339617},
{'iso': 'PEPE', 'name': 'Pepe', 'date_time': '2024-06-16 18:21:24', 'current_price': 1.19965e-05, 'open': 1.20502e-05, 'high': 1.23034e-05, 'low': 1.14667e-05, 'close': 1.19965e-05},
{'iso': 'RUNE', 'name': 'THORChain', 'date_time': '2024-06-16 18:21:24', 'current_price': 4.6837364756, 'open': 4.7150126065, 'high': 4.7802550631, 'low': 4.5839098332, 'close': 4.6837364756},
{'iso': 'AVA', 'name': 'Travala.com', 'date_time': '2024-06-16 18:21:24', 'current_price': 0.7394416146, 'open': 0.746540426, 'high': 0.7679110847, 'low': 0.7320058396, 'close': 0.7394416146},
{'iso': 'INJ', 'name': 'Injective', 'date_time': '2024-06-16 18:21:24', 'current_price': 25.1788217554, 'open': 26.0623597897, 'high': 26.6647242768, 'low': 24.4675209321, 'close': 25.1788217554},
{'iso': 'SUN', 'name': 'Sun', 'date_time': '2024-06-16 18:21:24', 'current_price': 0.0123356088, 'open': 0.0123959771, 'high': 0.0124604455, 'low': 0.0121975455, 'close': 0.0123356088}
```

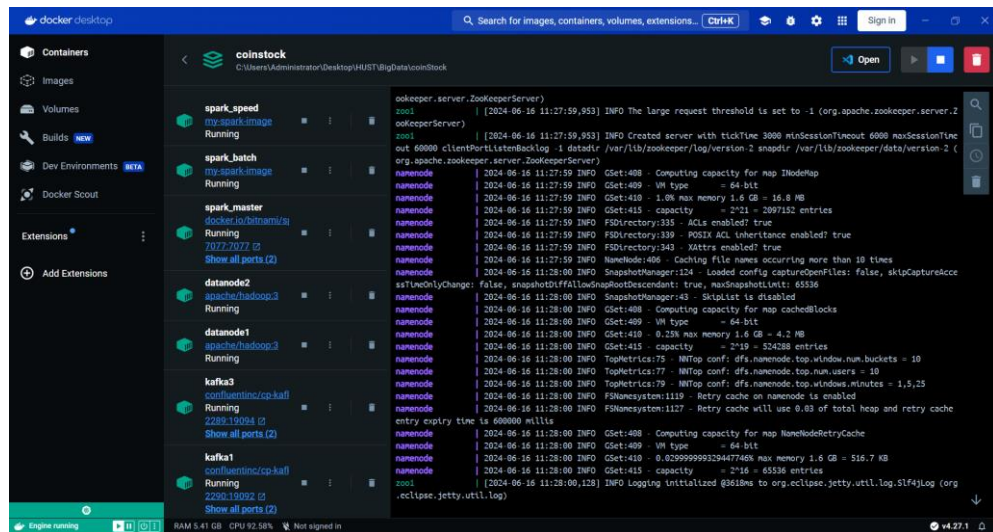
Hình 4.2 Kết quả dữ liệu sau khi crawl

Dữ liệu được crawl sẽ được sử dụng làm một nguồn producer cung cấp dữ liệu cho hệ thống.

4.2 Sử dụng Docker để thiết lập các thành phần hệ thống trên môi trường phân tán

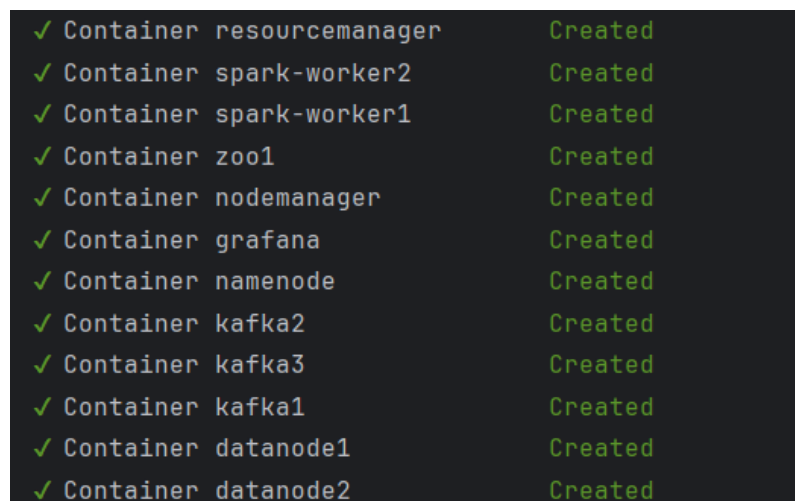
Để giả lập môi trường phân tán, các thành phần trong hệ thống đều được khởi chạy trên Docker thông qua docker-compose.yml. Trong đó có một image

phục vụ cho spark-submit được tạo lại có cài sẵn các thư viện cần thiết cho việc execute python file.

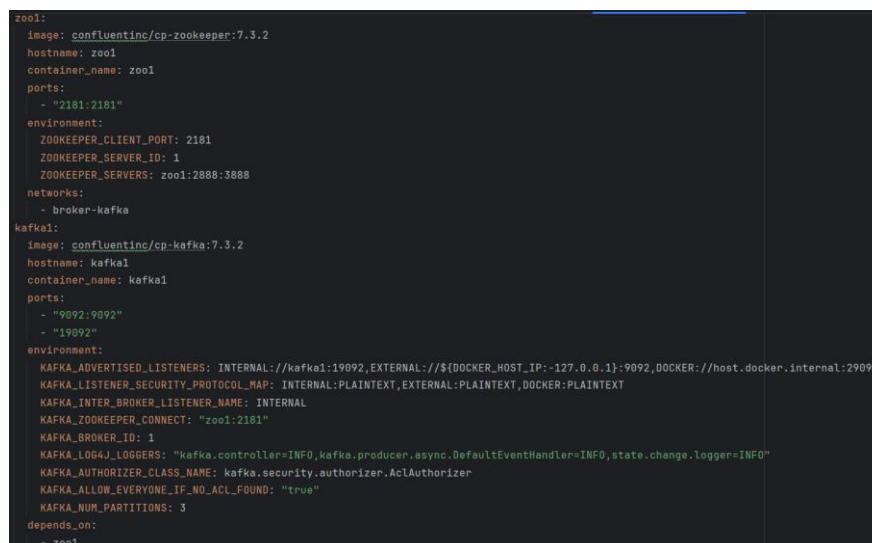


Hình 4.3 Các thành phần hệ thống được khởi chạy trên Docker

4.3 Thiết lập và sử dụng Multi Kafka brokers cùng ZooKeeper



Hình 4.4 Khởi động các container



Hình 4.5 Zookeeper và broker Kafka

3 broker kafka được cấu hình tại: kafka1:19092, kafka2:19093, kafka3:19094.

4.4 Xử lý dữ liệu thời gian thực sử dụng Apache Spark Streaming

```
if __name__ == "__main__":
    spark = (
        SparkSession.builder.appName("KafkaInfluxDBStreaming")
            .master("spark://spark-master:7077")
            .config("spark.jars.packages", ",".join(packages))
            .getOrCreate()
    )

    spark.sparkContext.setLogLevel("INFO")

    stockDataFrame = spark \
        .readStream \
        .format("kafka") \
        .option(key: "kafka.bootstrap.servers", KAFKA_BOOTSTRAP_SERVERS) \
        .option(key: "subscribe", KAFKA_TOPIC_NAME) \
        .load()

    stockDataFrame = stockDataFrame.select(col("value").cast("string").alias("data"))
    inputStream = stockDataFrame.selectExpr("CAST(data as STRING)")

    stock_price_schema = StructType([
        StructField(name: "iso", StringType(), nullable: True),
        StructField(name: "name", StringType(), nullable: True),
        StructField(name: "current_price", DoubleType(), nullable: True),
        StructField(name: "open", DoubleType(), nullable: True),
        StructField(name: "high", DoubleType(), nullable: True),
        StructField(name: "low", DoubleType(), nullable: True),
        StructField(name: "close", DoubleType(), nullable: True),
        StructField(name: "date_time", StringType(), nullable: True)
    ])
})
```

Hình 4.6 Code Spark Streaming

```
24/06/16 12:02:07 INFO DAGScheduler: Job 52 finished: call at /opt/bitnami/spark/python/lib/py4j-0.10.9.5-src.zip/py4j/clientserver.py:617, took 0.537789 s
[["BTC", "Bitcoin", 66442.7846359773, 66280.83071709584, 66681.81572379959, 65846.13470482016, 66442.7846359773, "2024-06-16 19:02:06"]]
-----
[["ETH", "Ethereum", 3551.6567677362136, 3536.607104208684, 3589.6116053697947, 3535.836775962719, 3551.6567677362136, "2024-06-16 19:02:06"]]
-----
[["ETHFI", "ether.fl", 3.9510062371, 3.8335795795, 4.0711812159, 3.7963240585, 3.9510062371, "2024-06-16 19:02:06"]]
-----
[["DOGE", "Dogecoin", 0.1357487928, 0.1369720102, 0.13759965, 0.1348752542, 0.1357487928, "2024-06-16 19:02:06"]]
-----
[["ZETA", "ZetaChain", 1.0343858023, 1.0380019553, 1.0501979546, 1.0030871901, 1.0343858023, "2024-06-16 19:02:06"]]
-----
Batch processed 17 done!
```

Hình 4.7 Xử lý dữ liệu thời gian thực sử dụng Apache Spark Streaming

Từ hình 4.7, log của dòng dữ liệu (stream processing) đang được hiển thị theo các mini-batch.

4.5 Thực hiện các thao tác chuyển đổi và truy vấn dữ liệu bằng Spark SQL

```

# Remove duplicates
df = df.dropDuplicates()

# Basic Statistics for each stock
basic_stats = df.groupBy("iso").agg(
    F.mean("open").alias("avg_open"),
    F.mean("high").alias("avg_high"),
    F.mean("low").alias("avg_low"),
    F.mean("close").alias("avg_close"),
    F.stddev("close").alias("std_dev_close"),
    F.max("high").alias("historical_high"),
    F.min("low").alias("historical_low")
)

# Show basic statistics
basic_stats.show()

# Write basic_stats DataFrame to MongoDB
(basic_stats.write.format("mongo").mode("append")
 .option(key: "database", value: "bigdata").option(key: "collection", value: "basic_stats").save())

# Filter out Bitcoin data and prepare it as a feature
btc_df = df.filter(df.iso == 'BTC').select("date_time", df.current_price.alias("btc_price"))

# Join BTC data with the main dataframe on date_time
df = df.join(btc_df, on="date_time")

# Assemble features for ML model
assembler = VectorAssembler(inputCols=["btc_price"], outputCol="features")
df = assembler.transform(df)

```

Hình 4.8 Các thao tác trên dữ liệu sử dụng Spark SQL

Hình 4.8, là một số thao tác trong quá trình xử lý dữ liệu với Spark SQL như: dropDuplicate(), filter(), tính max, min, mean, std-dev, v.v

4.6 Xử lý dữ liệu theo lô sử dụng Apache Spark

iso	avg_open	avg_high	avg_low	avg_close	std_dev_close	historical_high	historical_low
ZETA	1.037028455982353	1.0501979546	1.0030871901000002	1.0390431499705883	0.002630176427328519	1.0501979546	1.0030871901
ETH	3557.7287430509555	3606.8120380819387	3535.9136161289066	3584.5206784170923	0.771942363642621	3606.81203808194	3535.913616128907
TON	8.08076064222941	8.1105445807	7.749244813900002	7.985538694376472	0.007468504173688235	8.1260772589	7.7492448139
DOGE	0.13654611282941176	0.13759965000000002	0.1348752542	0.13652190748823526	0.817778846327358E-5	0.13759965	0.1348752542
INJ	26.09377705765882	26.66472427680001	24.4675209321	25.010940715082363	0.034772812080573284	26.6647242768	24.4675209321
OKB	45.82810460593529	46.4435857202	45.75748475134119	45.87240244494118	0.03124671123961511	46.4435857202	45.6662294656
RUNE	4.7262784346588225	4.7802550631	4.583909833199998	4.694624074082354	0.002616025338785165	4.7802550631	4.5839098332
PEPE	1.20637802353411...	1.2284399999999999	1.1468700000000000	1.212631764705882...	3.95573277598175E-9	1.22844E-5	1.14667E-5
SUN	0.0123540953524122	0.0124604455	0.01197545499999999	0.012382063888235293	6.007635291204116E-6	0.0124604455	0.01219795455
BTC	66186.50609890305	66728.5772907998	65846.13470482014	66550.38219661179	9.056651033070994	66728.5772907998	65846.13470482016
ETHFI	3.8886222592000004	4.071181215899999	3.8456058790099994	4.008263773535203	0.003055287322897...	4.0711812159	3.8456058791
SOL	144.61886221682622	146.62768550881395	143.05450288730887	146.18444374370134	0.00537512470299215	146.6276855088139	143.05450288730884
MTF	2.4010517801470583	2.5789173242	2.310219478299996	2.544759103911764	0.001638004098505...	2.5789173242	2.3102194783
BNB	606.1912657254236	611.511162389	601.958141428	609.8952116566942	0.13302537642279835	611.511162389	601.958141428
AVA	0.7566172723823529	0.7679110846999998	0.7320058395999999	0.7452492991176471	4.106511811496379...	0.7679110847	0.7320058396
SHIB	2.079458235294117...	2.1032400000000000	2.0422799999999999	2.091164705882352...	9.65564066087661E-9	2.10324E-5	2.04228E-5

Hình 4.9 Kết quả xử lý theo lô


```
import pyhdfs
from pyspark.sql import SparkSession
from pyspark.sql import functions as F
from pyspark.sql.types import StructType, StructField, StringType, DoubleType

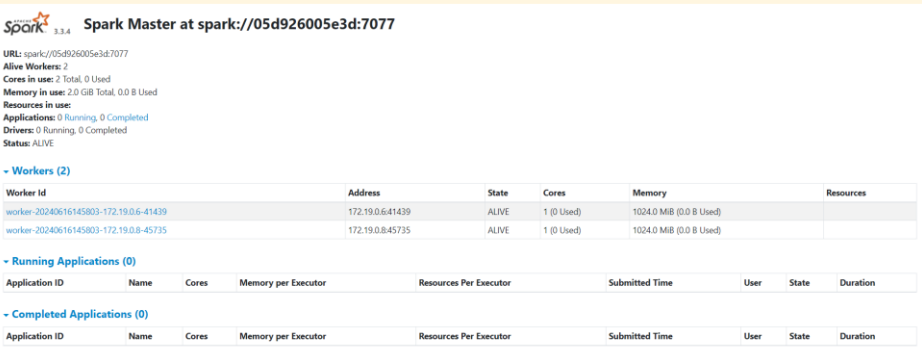
# Initialize Spark Session
spark = SparkSession.builder \
    .appName("Crypto Dependency Analysis") \
    .config("spark.mongodb.output.uri", "mongodb://root:admin@mongodb:27017/bigdata.stock2024") \
    .getOrCreate()

# Set up the HDFS client
hdfs = pyhdfs.HdfsClient(hosts="namenode:9870", user_name="hdfs")
directory = '/data'
if not hdfs.exists(directory):
    hdfs.mkdirs(directory)
files = hdfs.listdir(directory)
print("Files in '{}':".format(directory), files)

# Define the schema for the DataFrame
schema = StructType([
    StructField(name="iso", StringType(), nullable=True),
    StructField(name="name", StringType(), nullable=True),
    StructField(name="current_price", DoubleType(), nullable=True),
    StructField(name="open", DoubleType(), nullable=True),
    StructField(name="high", DoubleType(), nullable=True),
    StructField(name="low", DoubleType(), nullable=True),
    StructField(name="close", DoubleType(), nullable=True),
    StructField(name="date_time", StringType(), nullable=True)
])
```

Hình 4.10 Một đoạn code trong xử lý theo lô

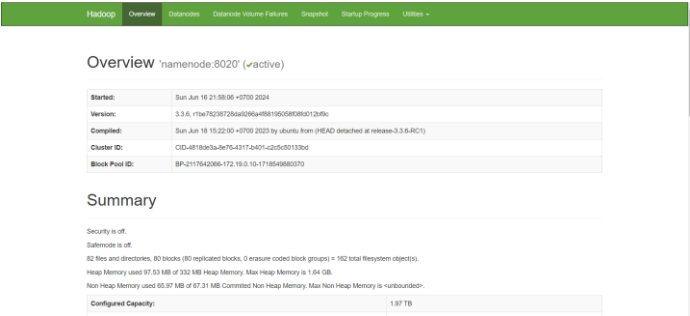
4.7 Thiết lập và sử dụng Multi-Node Spark Cluster



Hình 4.11 Giao diện hiển thị thông tin Spark Master

Từ Spark Master (hình 4.11) có thể thấy Apache Spark hoạt động với 2 worker node.

4.8 Lưu trữ dữ liệu trên Hadoop HDFS



Hình 4.12 Tổng quan của Hadoop HDFS

<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	107 B	Jun 16 21:59	<u>2</u>	128 MB	14875646-2bf1-11ef-93a5-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	135 B	Jun 16 21:59	<u>2</u>	128 MB	148d6234-2bf1-11ef-be8c-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	139 B	Jun 16 21:59	<u>2</u>	128 MB	1492bf96-2bf1-11ef-b3f3-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	114 B	Jun 16 21:59	<u>2</u>	128 MB	14983590-2bf1-11ef-b16b-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	111 B	Jun 16 21:59	<u>2</u>	128 MB	149d9e39-2bf1-11ef-bb33-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	114 B	Jun 16 21:59	<u>2</u>	128 MB	14a2f8d4-2bf1-11ef-925e-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	115 B	Jun 16 21:59	<u>2</u>	128 MB	14a8441c-2bf1-11ef-bb8e-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	109 B	Jun 16 21:59	<u>2</u>	128 MB	14ad6188-2bf1-11ef-8688-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	139 B	Jun 16 21:59	<u>2</u>	128 MB	14b29dce-2bf1-11ef-80be-9d734e00d044.json	
<input type="checkbox"/>	-rw-r--r--	hdfs	supergroup	111 B	Jun 16 21:59	<u>2</u>	128 MB	14b7c18d-2bf1-11ef-8161-9d734e00d044.json	

Showing 1 to 25 of 320 entries

Previous 1 2 3 4 5 ... 13 Next

Hình 4.13 Dữ liệu được lưu trữ trong HDFS

4.9 Cài đặt mô hình máy học thông qua Spark Mllib

```

# K-Means Model
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.clustering import KMeans

input_cols = ["avg_price_change", "std_dev_price_change"]
vec_assembler = VectorAssembler(inputCols=input_cols, outputCol="features")
df_kmeans = vec_assembler.transform(basic_stats)

kmeans = KMeans().setK(2).setSeed(1).setFeaturesCol("features")
model = kmeans.fit(df_kmeans)

# Predict clusters for each crypto
predictions = model.transform(df_kmeans)
predictions.select("iso", "prediction").show()

```

```

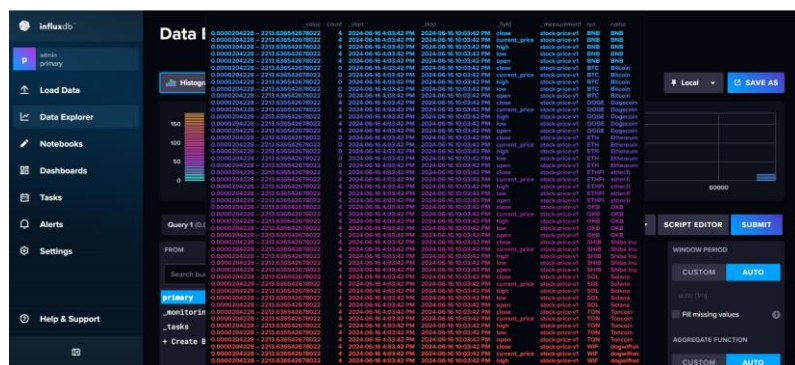
+-----+-----+
| iso|prediction|
+-----+-----+
| ZETA| 1|
| ETH| 1|
| TON| 0|
| DOGE| 1|
| INJ| 0|
| OKB| 1|
| RUNE| 0|
| PEPE| 1|
| SUN| 0|
| BTC| 1|
| ETHFI| 1|
| SOL| 1|
| BNB| 1|
| WIF| 1|
| AVA| 0|
| SHIB| 1|
+-----+-----+

```

Hình 4.14 Thuật toán K-means và kết quả phân cụm

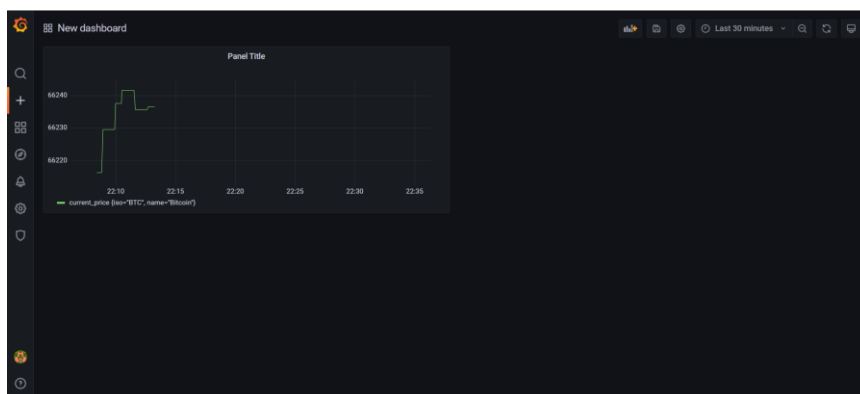
4.10 Triển khai và lưu trữ cơ sở dữ liệu NoSQL MongoDB

4.11 Triển khai và lưu trữ cơ sở dữ liệu time-series InfluxDB



Hình 4.15 Dữ liệu tại bucket primary trong InfluxDB

4.12 Trực quan hoá dữ dạng stream-data thông qua Grafana



Hình 4.16 Trực quan hoá dữ liệu với Grafana

Dữ liệu từ dạng time-series từ InfluxDB, sẽ được trực quan hoá theo dạng stream thông qua nền tảng Grafana.

4.13 Thiết kế kiến trúc hệ thống theo mô hình lambda

Việc thiết kế và khởi tạo kiến trúc lambda đã giúp cho hệ thống xử lý dữ liệu phân tán linh hoạt và hiệu quả hơn, có thể đáp ứng được các yêu cầu khác nhau. Điển hình của kiến trúc là 2 nhánh Batch Layer và Speed Layer như hình 3.1. Chi tiết đã được mô tả trong chương 3.

CHƯƠNG 5. KẾT LUẬN

Trong dự án này, em đã thành công trong việc xây dựng một hệ thống xử lý và phân tích dữ liệu tiền ảo thời gian thực. Kết hợp các công cụ hỗ trợ xử lý dữ liệu lớn một cách nhanh chóng và chính xác gồm: Kafka, Apache Spark, Hadoop cùng với các cơ sở dữ liệu NoSQL như MongoDB hay time-series như InfluxDB.

Em cũng đã gặp phải nhiều thách thức, và các vấn đề khác nhau. Tuy nhiên, từ đây nhiều bài học và trải nghiệm trong lưu trữ và xử lý dữ liệu lớn đã được rút ra. Dự án là bước đầu cho việc ứng dụng dữ liệu lớn trong các lĩnh vực tài chính, đặc biệt là trong việc phân tích và dự báo thị trường tiền ảo của em.

Em hy vọng rằng, dự án có thể hoàn thiện hơn trong tương lai. Cùng với các thông tin, phân tích đa dạng, thực tế và hữu ích hơn nữa cho người sử dụng.