

TRƯỜNG ĐẠI HỌC CÔNG NGHỆ

VIỆN TRÍ TUỆ NHÂN TẠO

-----***-----



**BÁO CÁO MÔN HỌC KỸ THUẬT VÀ
CÔNG NGHỆ DỮ LIỆU LỚN**

ĐỀ TÀI

Finding Average Temperature of Each Year using Hadoop-HDFS

Giảng viên hướng dẫn: TS. Trần Hồng Việt

Nhóm sinh viên thực hiện

1. Thái Nguyễn Hoàng Bách - 22022672
2. Quách Đắc Chính - 22022518
3. Vũ Minh Đức - 22022587
4. Đào Duy Hưng - 22022589

HÀ NỘI, 12/2024

Mục lục

| | |
|---|----|
| MỞ ĐẦU | 1 |
| Chương 1: Tổng quan về dữ liệu lớn và Hadoop MapReduce | 2 |
| 1.1 Định nghĩa của Big Data..... | 2 |
| 1.2 Đặc trưng của Big Data | 2 |
| 1.3 Giới thiệu về Hadoop MapReduce | 2 |
| Chương 2: Cách triển khai phân tích dữ liệu khí tượng | 4 |
| 2.1 Cào và tạo file dữ liệu đầu vào | 4 |
| 2.2 Cấu trúc dữ liệu đầu vào | 5 |
| 2.3 Quy trình xử lý dữ liệu với MapReduce..... | 5 |
| 2.3.1 Mapper..... | 5 |
| 2.3.2 Reducer..... | 6 |
| 2.3.3 Driver | 7 |
| Chương 3: Đánh giá kết quả | 9 |
| 3.1 Kết quả chạy chương trình | 9 |
| 3.2 Đánh giá chương trình..... | 10 |
| Chương 4: Phát triển bài toán dự báo thời tiết | 11 |
| 4.1. Mục tiêu | 11 |
| 4.2. Model..... | 12 |
| 4.3. Kết quả..... | 12 |
| Chương 5: Kết luận và hướng phát triển | 14 |
| 5.1 Kết luận..... | 14 |
| 5.2 Hướng phát triển | 14 |
| BẢNG PHÂN CHIA CÔNG VIỆC CỦA TỪNG THÀNH VIÊN | 15 |

MỞ ĐẦU

Big Data ngày nay đã trở thành một phần không thể thiếu trong việc xử lý và phân tích dữ liệu lớn khi xu hướng sử dụng công nghệ ngày càng lớn đối với tất cả mọi người. Với sự phát triển của các công nghệ hỗ trợ như Hadoop, việc xử lý các tập dữ liệu lớn không chỉ trở nên khả thi mà còn hiệu quả. Điều đó sẽ được ứng dụng trong rất nhiều lĩnh vực như đời sống, giáo dục, ...

Chúng em chọn đề tài "**Finding Average Temperature using Hadoop-HDFS**" nhằm mục đích áp dụng thực tiễn công nghệ MapReduce để phân tích nhiệt độ từ dữ liệu khí tượng và cung cấp thông tin về biến động thời tiết qua các năm.

Báo cáo gồm 4 chương:

Chương 1: Tổng quan về dữ liệu lớn và Hadoop MapReduce

Chương 2: Cách triển khai phân tích dữ liệu khí tượng.

Chương 3: Đánh giá kết quả và hiệu quả xử lý.

Chương 4: Kết luận và hướng phát triển.

Chương 1: Tổng quan về dữ liệu lớn và Hadoop MapReduce

1.1 Định nghĩa của Big Data

- Theo Wikipedia: Dữ liệu lớn thường bao gồm tập hợp dữ liệu với kích thước vượt xa khả năng của các công cụ phần mềm thông thường để thu thập, hiển thị, quản lý và xử lý dữ liệu trong một thời gian có thể chấp nhận được.
- Dữ liệu lớn bao gồm các thách thức như phân tích, thu thập, giám sát dữ liệu, tìm kiếm, chia sẻ, lưu trữ, truyền nhận, trực quan, truy vấn và tính riêng tư.

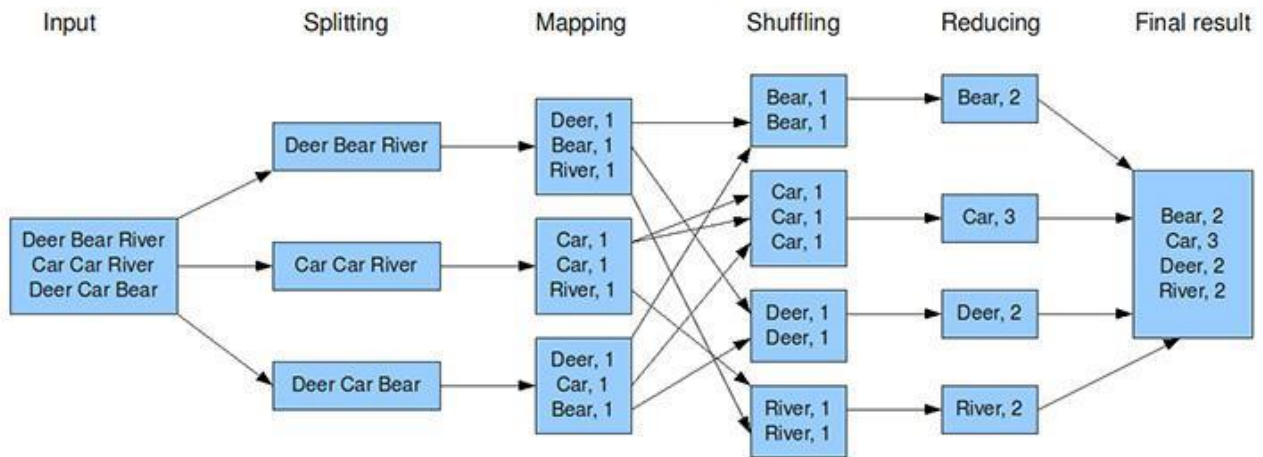
1.2 Đặc trưng của Big Data

- Khối lượng lớn (Volume): Khối lượng dữ liệu rất lớn và đang ngày càng tăng lên, tính đến 2014 thì có thể trong khoảng vài trăm terabyte.
- Tốc độ (Velocity): Khối lượng dữ liệu gia tăng rất nhanh.
- Đa dạng (Variety): Ngày nay hơn 80% dữ liệu được sinh ra là phi cấu trúc(tài liệu, blog, hình ảnh,...)
- Độ tin cậy/chính xác(Veracity): Bài toán phân tích và loại bỏ dữ liệu thiếu chính xác và nhiễu đang là tính chất quan trọng của bigdata.
- Giá trị(Value): Giá trị thông tin mang lại.

1.3 Giới thiệu về Hadoop MapReduce

- Hadoop là một framework mã nguồn mở dùng để lưu trữ và xử lý dữ liệu lớn trên các cụm máy tính thông thường.
- Hadoop bao gồm hai thành phần chính:
 - + HDFS (Hadoop Distributed File System): Hệ thống lưu trữ phân tán, giúp lưu trữ dữ liệu lớn bằng cách phân mảnh dữ liệu trên nhiều nút.
 - + MapReduce: Mô hình lập trình song song, xử lý dữ liệu qua các bước Map và Reduce.
- MapReduce là mô hình lập trình được thiết kế để xử lý dữ liệu song song trên các cụm máy tính phân tán. Mô hình hoạt động qua hai pha chính:
 - + Map: Chuyển đổi dữ liệu đầu vào thành các cặp key-value trung gian.
 - + Reduce: Tổng hợp và xử lý các cặp key-value để tạo ra kết quả cuối cùng.

The overall MapReduce word count process



Chương 2: Cách triển khai phân tích dữ liệu khí tượng

2.1 Cào và tạo file dữ liệu đầu vào

- Sử dụng web weather api để lấy api key và tiến hành cào dữ liệu

```
import requests
from datetime import datetime, timedelta

# WeatherAPI key (thay thế API_KEY bằng key thực tế của bạn)
api_key = 'API_KEY'

# Hàm lấy dữ liệu thời tiết từ API cho một ngày cụ thể
def fetch_weather_data(location, date):
    date_str = date.strftime("%Y-%m-%d")
    url = f'http://api.weatherapi.com/v1/history.json?key={api_key}&q={location}&dt={date_str}'
    response = requests.get(url)
    # (variable) status_code: int

    if response.status_code == 200:
        return response.json()
    else:
        print(f"Error fetching data for {location} on {date_str}: {response.status_code}")
        return None
```

- Sau khi cào thành công tiến hành ghi dữ liệu vào file

```
# Hàm ghi dữ liệu thời tiết vào file
def write_to_file(data, filename):
    with open(filename, 'w') as file:
        # Ghi tiêu đề cột
        file.write("city\tday/month/year\ttemperature\thumidity\ttrainfall\twind_speed\n")
        # Ghi từng dòng dữ liệu
        for entry in data:
            file.write(f"{entry['city']}\t{entry['date']}\t{entry['temp']:.2f} °C\t"
                      f"{entry['humidity']:.1f}%\t{entry['rainfall']:.1f} mm\t"
                      f"{entry['wind_speed']:.1f} kph\n")
    print(f"Weather data has been written to {filename}")
```

- Khi chạy chương trình tiến hành, lấy dữ liệu về thời tiết mong muốn từ người dùng

```
# Nhập dữ liệu từ người dùng
if __name__ == "__main__":
    # Lấy danh sách thành phố và ngày từ người dùng
    cities = input("Enter the locations separated by commas (e.g., Hanoi, Ho Chi Minh City): ").split(", ")
    start_date_str = input("Enter the start date (dd/mm/yyyy): ")
    end_date_str = input("Enter the end date (dd/mm/yyyy): ")

    # Chuyển đổi ngày sang định dạng datetime
    start_date = datetime.strptime(start_date_str, "%d/%m/%Y")
    end_date = datetime.strptime(end_date_str, "%d/%m/%Y")

    # Thu thập dữ liệu thời tiết
    weather_data = collect_weather_data(cities, start_date, end_date)

    # Ghi dữ liệu vào file
    write_to_file(weather_data, 'south_test.txt')
```

2.2 Cấu trúc dữ liệu đầu vào

- Dữ liệu đầu vào là 1 bảng bao gồm các cột city, day/month/year, temperature, humidity, rainfall, wind_speed
- Ý nghĩa của từng trường:
 - + City: Thành phố mà người dùng mong muốn thu thập dữ liệu
 - + Day/month/year: là các ngày thu thập dữ liệu từ ngày bắt đầu đến ngày kết thúc
 - + Temperature: Nhiệt độ của các ngày
 - + Humidity: Độ ẩm
 - + Rainfall: Lượng mưa
 - + Wind_speed: Tốc độ gió

2.3 Quy trình xử lý dữ liệu với MapReduce

2.3.1 Mapper

- Lớp Temperature_Mapper được thiết kế để xử lý các nhiệm vụ sau:
 1. **Đọc từng dòng dữ liệu đầu vào:** Mỗi dòng là một quan sát thời tiết từ trạm.
 2. **Trích xuất các trường dữ liệu cần thiết:** Gồm city, date, temperature, humidity, rainfall, windspeed
 3. **Phát các cặp key-value:** Với key là city và month sau khi đã kết hợp lại thành 1 chuỗi và value là các giá trị temperature, humidity, rainfall và windspeed kết hợp lại thành 1 chuỗi

```
1 package org.example;
2
3 import org.apache.hadoop.io.Text;
4 import org.apache.hadoop.mapreduce.Mapper;
5
6 import java.io.IOException;
7 import java.text.ParseException;
8 import java.text.SimpleDateFormat;
9 import java.util.Date;
10
11 public class TemperatureMapper extends Mapper<Object, Text, Text, Text> {
12
13     @Override
14     protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
15         String line = value.toString();
16         String[] fields = line.split("\\t");
17
18         if (fields.length == 6) { // Đảm bảo dòng có đúng 6 cột: D
19             String city = fields[0]; // Lấy tên thành phố
20             String date = fields[1]; // Lấy ngày
21             String temperature = fields[2].replace("°C", "").trim(); // Lấy nhiệt độ và loại bỏ ký tự '°C'
22             String humidity = fields[3].replace("%", "").trim(); // Lấy độ ẩm và loại bỏ ký tự '%'
23             String rainfall = fields[4].replace(" mm", "").trim(); // Lấy lượng mưa và loại bỏ ký tự 'mm'
24             String windspeed = fields[5].replace(" kph", "").trim(); // Lấy tốc độ gió và loại bỏ ký tự 'k'
25         }
26     }
27 }
```

```

// Chuyển đổi date thành định dạng tháng-năm
SimpleDateFormat inputFormat = new SimpleDateFormat("dd/MM/yyyy");
SimpleDateFormat outputFormat = new SimpleDateFormat("MM/yyyy");
try {
    Date parsedDate = inputFormat.parse(date);
    String month = outputFormat.format(parsedDate);

    // Gộp city và month thành key
    String cityMonthKey = city + "-" + month;

    // Gộp các giá trị lại thành chuỗi để truyền cho Reducer
    String combinedValues = temperature + "\t" + humidity + "\t" + rainfall + "\t" + windspeed;

    context.write(new Text(cityMonthKey), new Text(combinedValues));
} catch (ParseException e) {
    System.err.println("Error parsing date: " + date);
    e.printStackTrace();
}
}

```

2.3.2 Reducer

- Lớp Temperature_Reducer thực hiện các nhiệm vụ quan trọng nhằm tổng hợp và tính toán các thông tin từ dữ liệu nhiệt độ trung gian do Mapper phát ra:
1. **Nhận các cặp key-value từ Mapper:**
 - Các giá trị nhiệt độ được nhóm theo từng city-month (key) do Mapper phát ra.
 2. **Tính toán các thông số nhiệt độ:**
 - **Tính trung bình các thuộc tính:** Tính tổng các thuộc tính sau đó chia cho số lượng các bản ghi


```

for (Text val : values) {
    String[] fields = val.toString().split("\\t");
    double temperature = Double.parseDouble(fields[0]);
    double humidity = Double.parseDouble(fields[1]);
    double rainfall = Double.parseDouble(fields[2]);
    double windspeed = Double.parseDouble(fields[3]);

    sumTemperature += temperature;
    sumHumidity += humidity;
    sumRainfall += rainfall;
    sumWindspeed += windspeed;
    count++;
}

double avgTemperature = sumTemperature / count;
double avgHumidity = sumHumidity / count;
double avgRainfall = sumRainfall / count;
double avgWindspeed = sumWindspeed / count;

```

3. Xuất kết quả:

- Kết quả được định dạng thành chuỗi bao gồm thông tin city, month, avgTemperature, avgHumidity, avgRainfall, avgWindspeed

```

// Định dạng kết quả đầu ra đúng định dạng mong muốn
String result = String.format("%s\t%s\t%.2f°C\t%.2f%\t%.2fmm\t%.2fkph",
    city, month, avgTemperature, avgHumidity, avgRainfall, avgWindspeed);
context.write(new Text(city), new Text(result));
}

```

2.3.3 Driver

- Dùng để tạo file jars và chạy các lớp

```

public class TemperatureDriver {
    public static void main(String[] args) throws Exception {
        if (args.length != 2) {
            System.err.println("Usage: TemperatureDriver <input path> <output path>");
            System.exit(-1);
        }

        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "Temperature Average");

        job.setJarByClass(TemperatureDriver.class);
        job.setMapperClass(TemperatureMapper.class);
        job.setReducerClass(TemperatureReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(Text.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}

```

```

J AvgTemp_Mapper.java
1  import org.apache.hadoop.io.Text;
2  import org.apache.hadoop.mapreduce.Mapper;
3
4  import java.io.IOException;
5
6  public class AvgTemp_Mapper extends Mapper<Object, Text, Text, Text> {
7      private Text yearKey = new Text();
8      private Text tempValue = new Text();
9
10     @Override
11     protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
12         String line = value.toString();
13         // Tách các trường từ dữ liệu
14         String year = line.substring(15, 19); // Năm
15         String temp = line.substring(87, 92); // Nhiệt độ
16         String quality = line.substring(92, 93); // Chất lượng dữ liệu
17         if (!temp.trim().equals("+9999") && quality.matches("[01459]")) {
18             yearKey.set(year);
19             tempValue.set(temp);
20             context.write(yearKey, tempValue);
21         }
22     }
23 }

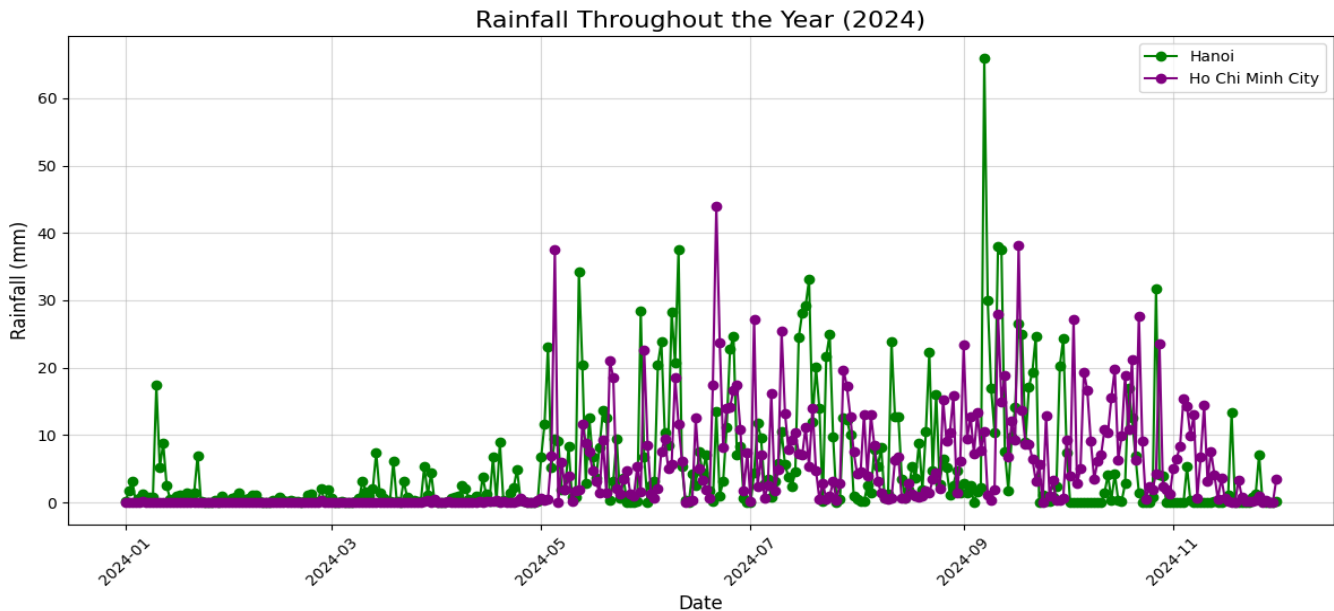
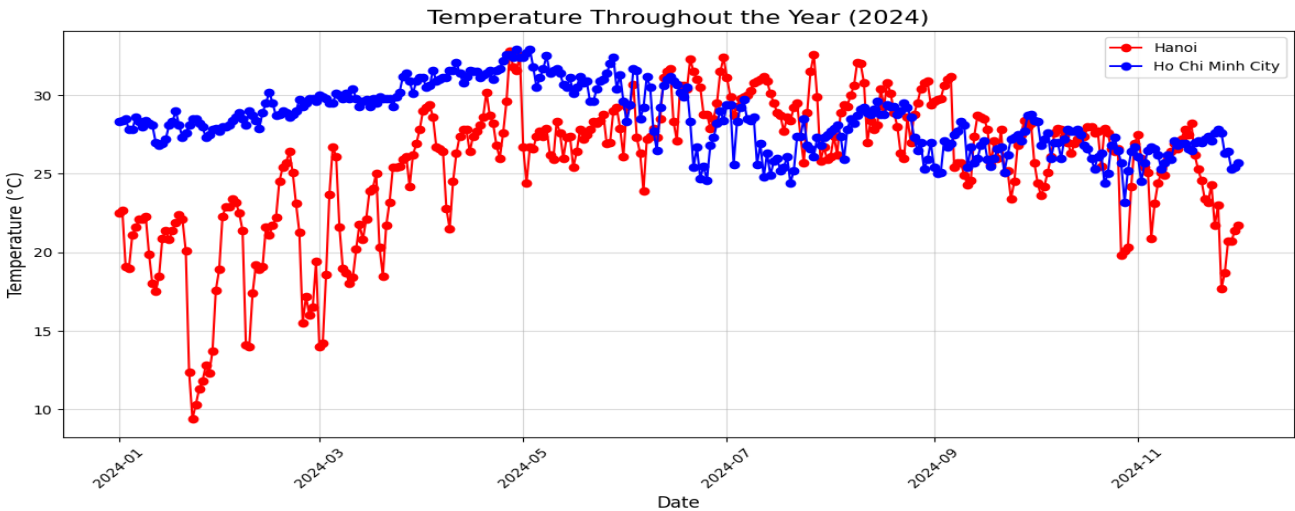
```

Chương 3: Đánh giá kết quả

3.1 Kết quả chạy chương trình

- Sau khi triển khai chương trình trên tập dữ liệu mẫu, kết quả đầu ra như sau:

```
kenkaw303@HB-QDT:~/Tempe$ hdfs dfs -cat /user/kenkaw303/output/part-r-00000
Hanoi Hanoi 01/2024 18.32°C 76.32% 1.92mm 14.16kph
Hanoi Hanoi 02/2024 20.83°C 75.62% 0.56mm 18.16kph
Hanoi Hanoi 03/2024 22.39°C 76.55% 1.34mm 18.41kph
Hanoi Hanoi 04/2024 27.93°C 73.27% 1.32mm 22.91kph
Hanoi Hanoi 05/2024 27.27°C 80.13% 7.94mm 16.69kph
Hanoi Hanoi 06/2024 29.31°C 77.70% 9.37mm 14.84kph
Hanoi Hanoi 07/2024 29.32°C 77.81% 10.40mm 15.13kph
Hanoi Hanoi 08/2024 29.01°C 78.45% 5.84mm 13.28kph
Hanoi Hanoi 09/2024 27.24°C 82.57% 13.78mm 15.17kph
Hanoi Hanoi 10/2024 25.98°C 65.90% 3.21mm 15.76kph
Hanoi Hanoi 11/2024 24.33°C 59.33% 0.99mm 15.67kph
Hanoi Hanoi 12/2024 21.70°C 65.00% 0.10mm 15.10kph
Ho Chi Minh City Ho Chi Minh City 01/2024 27.94°C 58.26% 0.00mm 18.86kph
Ho Chi Minh City Ho Chi Minh City 02/2024 28.90°C 59.07% 0.01mm 24.39kph
Ho Chi Minh City Ho Chi Minh City 03/2024 30.01°C 56.77% 0.03mm 26.63kph
Ho Chi Minh City Ho Chi Minh City 04/2024 31.53°C 56.83% 0.07mm 25.79kph
Ho Chi Minh City Ho Chi Minh City 05/2024 31.16°C 64.13% 6.19mm 18.68kph
Ho Chi Minh City Ho Chi Minh City 06/2024 28.79°C 73.40% 9.21mm 15.73kph
Ho Chi Minh City Ho Chi Minh City 07/2024 26.95°C 81.42% 7.96mm 18.55kph
Ho Chi Minh City Ho Chi Minh City 08/2024 28.21°C 76.61% 4.70mm 14.40kph
Ho Chi Minh City Ho Chi Minh City 09/2024 26.73°C 84.13% 9.36mm 16.54kph
Ho Chi Minh City Ho Chi Minh City 10/2024 26.59°C 84.39% 10.15mm 11.62kph
Ho Chi Minh City Ho Chi Minh City 11/2024 26.51°C 80.13% 4.02mm 11.12kph
Ho Chi Minh City Ho Chi Minh City 12/2024 25.70°C 82.00% 3.40mm 7.60kph
kenkaw303@HB-QDT:~/Tempe$
```



3.2 Đánh giá chương trình

Ưu điểm:

- Có thể lấy dữ liệu từ thực tế với số lượng ngày tháng mong muốn
- Chương trình xử lý chính xác các dữ liệu nhiệt độ lớn và phức tạp.
- Khả năng mở rộng dễ dàng nhờ sử dụng Hadoop.
- Xử lý hiệu quả dữ liệu không hợp lệ thông qua bộ lọc Mapper.

Hạn chế:

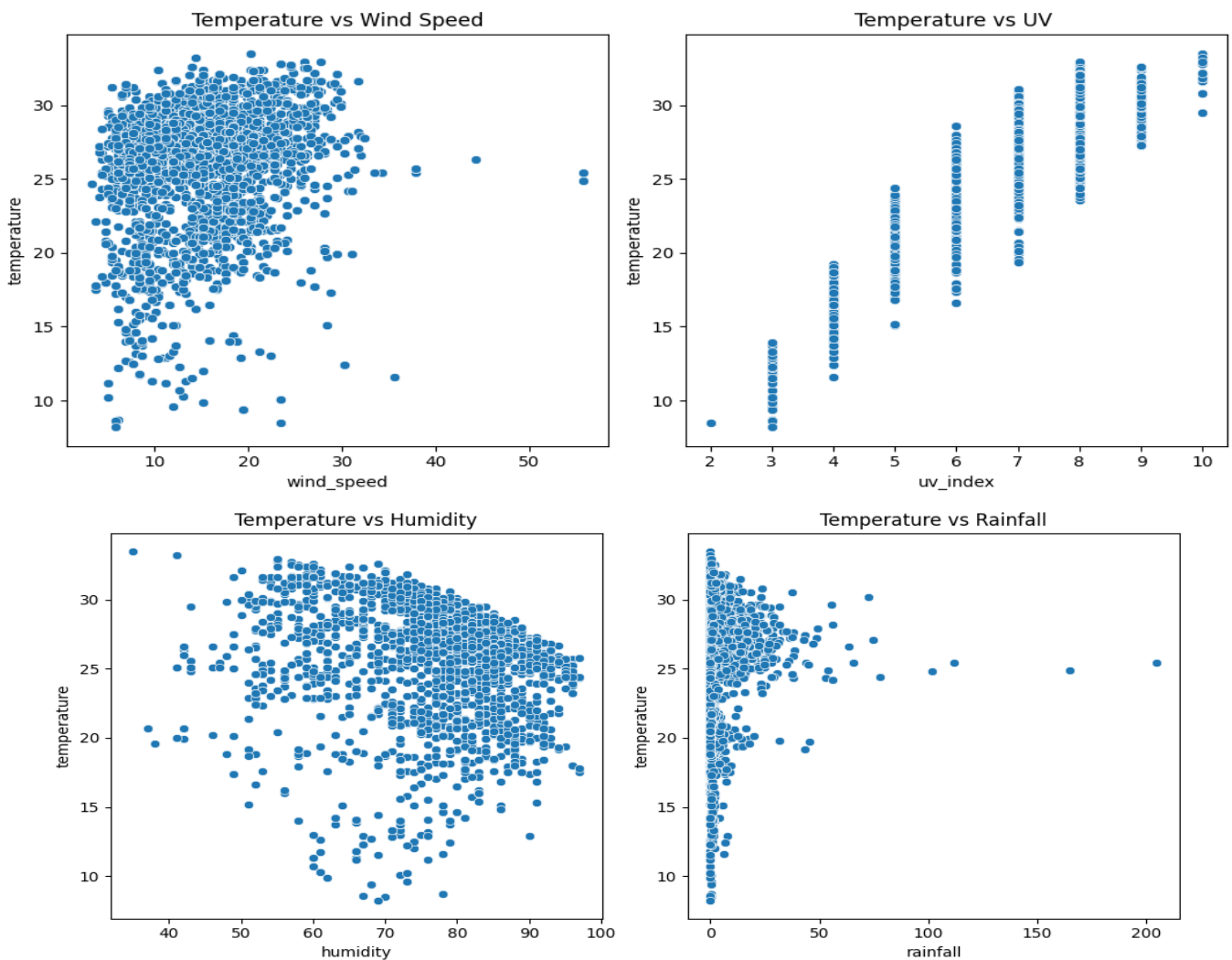
- Kết quả đầu ra chưa được tích hợp visualization để dễ hiểu hơn.

Chương 4: Phát triển bài toán dự báo thời tiết

Sau khi sử dụng Hadoop để xử lý và tổng hợp dữ liệu thời tiết lớn từ nhiều nguồn khác nhau, chúng ta có thể phát triển bài toán dự báo thời tiết bằng cách tích hợp các kỹ thuật học máy hoặc học sâu. Hadoop đã giúp giải quyết bài toán tiền xử lý dữ liệu, như làm sạch, chuyển đổi định dạng, và tính toán các chỉ số trung bình hay tổng hợp theo từng khu vực và thời gian.

4.1. Mục tiêu

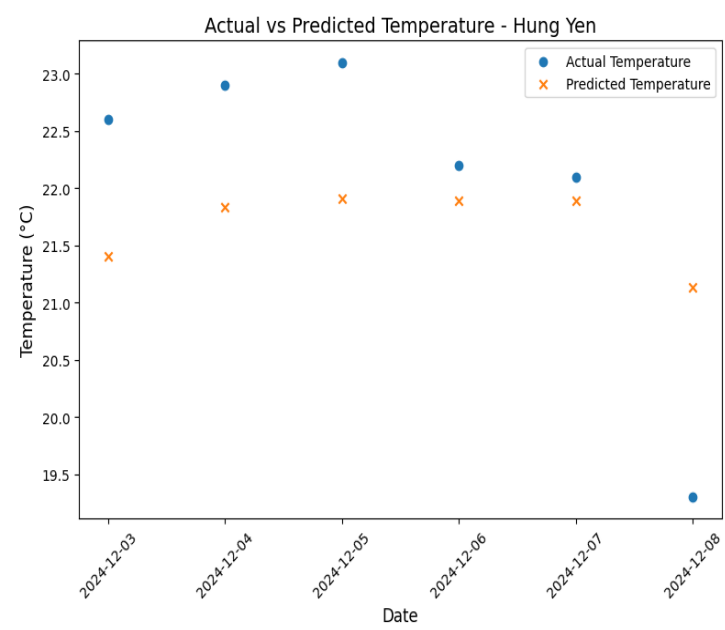
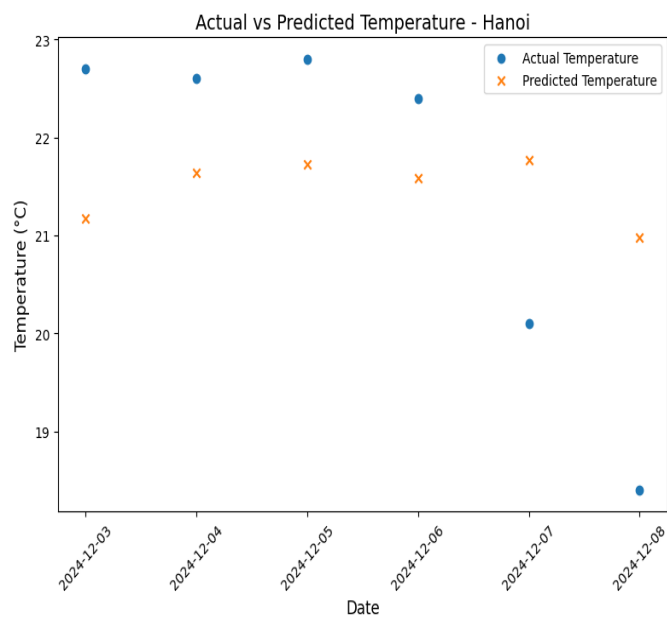
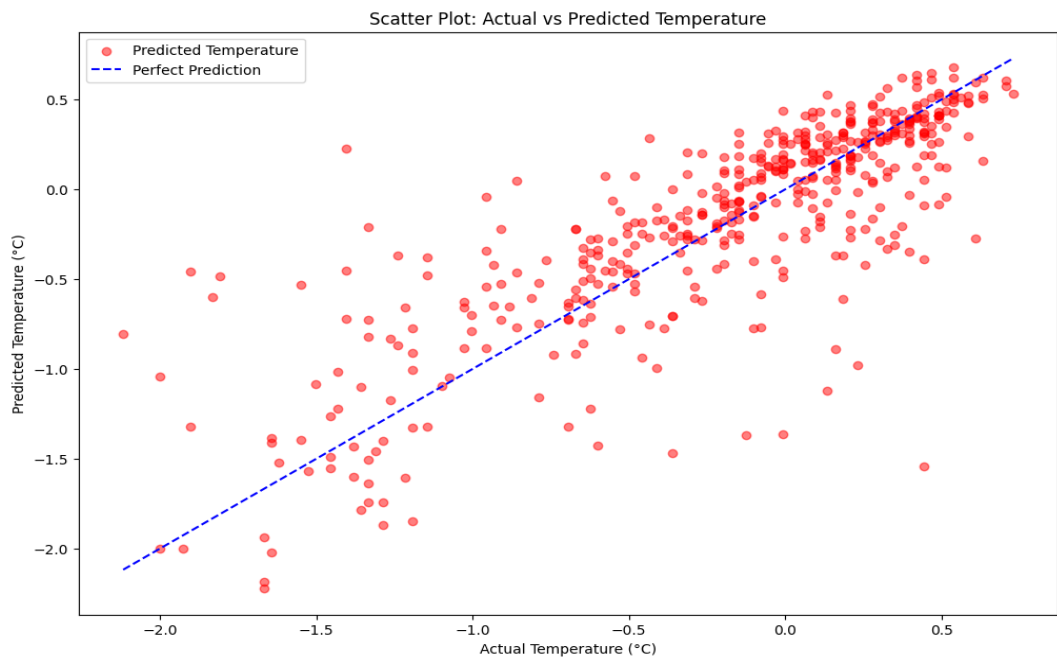
Với các đặc trưng đã thu thập được bao gồm: Rainfall, Humidity, Wind_speed, UV_index. Xây dựng một mô hình dự đoán thời tiết, cụ thể hơn trong trường hợp này sẽ là thời tiết của các tỉnh thành Việt Nam.

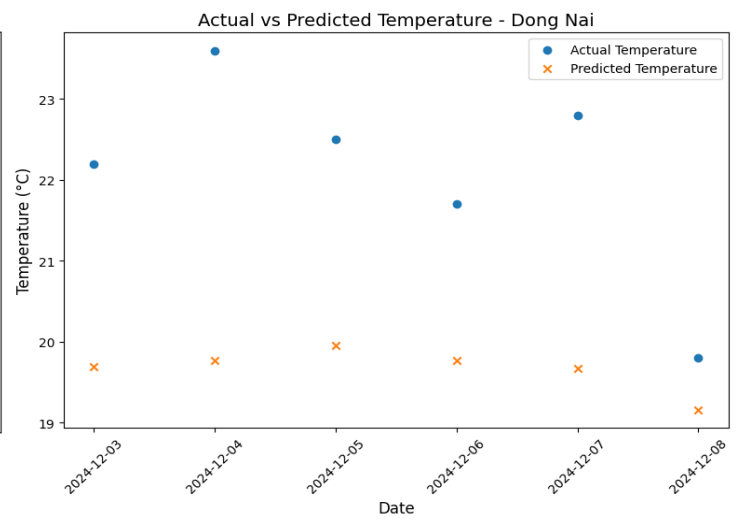
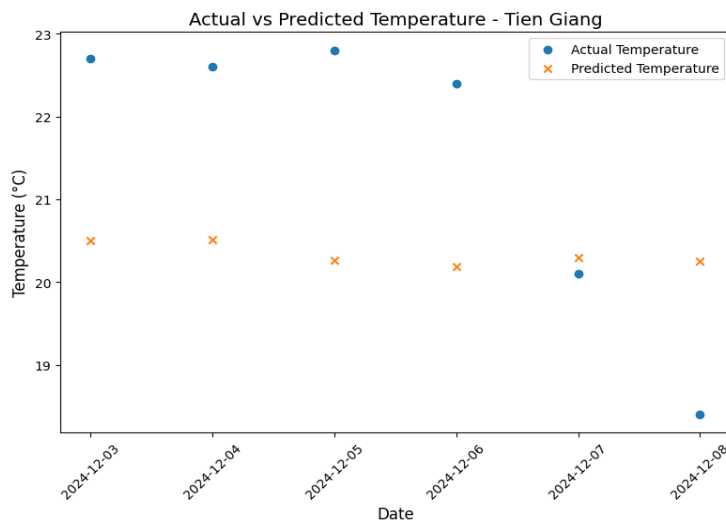


4.2. Model

Linear Regression trong PySpark là một thuật toán hồi quy tuyến tính, được sử dụng để dự đoán giá trị của một biến phụ thuộc (label) dựa trên một hoặc nhiều biến độc lập (features).

4.3. Kết quả





Chương 5: Kết luận và hướng phát triển

5.1 Kết luận

- Đề tài "Phân tích dữ liệu khí tượng bằng Hadoop MapReduce" đã hoàn thành các mục tiêu đề ra, bao gồm:
 - Triển khai thành công hệ thống phân tích nhiệt độ từ dữ liệu lớn.
 - Xử lý hiệu quả dữ liệu không hợp lệ và cung cấp các thông số quan trọng như nhiệt độ trung bình, cao nhất và thấp nhất theo từng năm.
- Chương trình cho thấy tiềm năng ứng dụng lớn của Hadoop trong các bài toán xử lý dữ liệu khí tượng và khoa học môi trường.

5.2 Hướng phát triển

- **Tích hợp thêm yếu tố phân tích:** Bao gồm độ ẩm, lượng mưa, tốc độ gió để mở rộng phân tích khí hậu.
- **Visualization:** Sử dụng các công cụ như Tableau hoặc Power BI để trực quan hóa kết quả.
- **Tối ưu hóa hiệu suất:** Sử dụng các kỹ thuật caching và nén dữ liệu để giảm thời gian xử lý.
- **Ứng dụng thực tế:** Mở rộng hệ thống để xử lý dữ liệu khí tượng theo thời gian thực, hỗ trợ các cơ quan nghiên cứu khí hậu và dự báo thời tiết.

BẢNG PHÂN CHIA CÔNG VIỆC CỦA TỪNG THÀNH VIÊN

| Họ và tên | Công việc |
|------------------------|--|
| Vũ Minh Đức | <ul style="list-style-type: none"> • Tìm hiểu tổng quan về Hadoop • Tìm hiểu về Map Reduce • Xây dựng lớp Mapper, Reducer, Driver • Làm slide thuyết trình |
| Quách Đắc Chính | <ul style="list-style-type: none"> • Tìm hiểu tổng quan về Hadoop • Lên ý tưởng về việc lấy dữ liệu từ weather api • Xây dựng lớp get_api để cào và xây dựng dữ liệu từ weather api |
| Thái Nguyễn Hoàng Bách | <ul style="list-style-type: none"> • Tìm hiểu tổng quan về Hadoop • Phân tích dữ liệu thu thập được • Chạy demo chương trình • Viết báo cáo • Xây dựng phần code python |
| Đào Duy Hưng | <ul style="list-style-type: none"> • Tìm hiểu tổng quan về Hadoop • Hỗ trợ các phần code bị lỗi • Tìm hiểu về các hướng phát triển của đề tài • Xây dựng phần code python |