

TRƯỜNG ĐẠI HỌC HỌC VĂN LANG  
KHOA CÔNG NGHỆ THÔNG TIN

**VAN LANG**  
UNIVERSITY



**BÁO CÁO ĐỒ ÁN MÔN HỌC:  
NHẬP MÔN PHÂN TÍCH DỮ LIỆU LỚN**

**Chủ đề:**

**ÁP DỤNG HADOOP CHO VIỆC XỬ LÝ DỮ LIỆU  
LỚN CỦA ỨNG DỤNG PHÂN TÍCH/ DỰ BÁO  
ĐIỂM HỌC SINH**

Nhóm sinh viên thực hiện (Họ tên - Mã SV):

1. Lê Minh Tâm - 207CT28471 (**Trưởng nhóm**)
2. Nguyễn Quốc An – 207CT27552
3. Trần Trung Kiên – 207CT10144
4. Võ Văn Thanh Nhân – 207CT10223

TP. Hồ Chí Minh – năm 2023

[illegible]

## MỤC LỤC

NHẬN XÉT CỦA GIẢNG VIÊN.....	2
PHỤ LỤC HÌNH ẢNH.....	4
LỜI NÓI ĐẦU.....	5
MỞ ĐẦU.....	6
1. LÝ DO CHỌN CHỦ ĐỀ.....	6
2. ĐỐI TƯỢNG, PHẠM VI TÌM HIỂU CỦA CHỦ ĐỀ.....	6
CHƯƠNG 1. GIỚI THIỆU CHỦ ĐỀ ĐỒ ÁN.....	7
1.1 KHÁI QUÁT MỘT SỐ NỘI DUNG CƠ BẢN VỀ HADOOP VÀ NGUỒN GỐC CỦA HADOOP: [1] [2].....	7
1.1.1 Hadoop là gì?.....	7
1.1.2 Nguồn gốc sự phát triển của Hadoop?.....	7
1.2 CÁC MODULE QUAN TRỌNG CỦA NỀN TẢNG HADOOP.....	7
1.2.1 Hadoop Distributed File System.....	8
1.2.2 MapReduce.....	8
1.3 SO SÁNH SỰ KHÁC BIỆT GIỮA HADOOP VÀ MONGODB [3] [4].....	9
CHƯƠNG 2. KẾT QUẢ THỰC HIỆN ĐỒ ÁN.....	11
2.1 THIẾT KẾ ĐƯỢC GIAO DIỆN CƠ BẢN CỦA ỨNG DỤNG QUẢN LÝ SINH VIÊN.....	11
2.2 TIẾN HÀNH TRIỂN KHAI CÁC CHỨC NĂNG CHO ỨNG DỤNG.....	11
KẾT LUẬN VÀ ĐỀ XUẤT.....	29
1. NHẬN XÉT VỀ HADOOP VÀ NÊU ƯU ĐIỂM VÀ NHƯỢC ĐIỂM CỦA NÓ.....	29
2. ĐỀ XUẤT TRONG TƯƠNG LAI.....	31
TÀI LIỆU THAM KHẢO.....	32
LỜI CẢM ƠN.....	33

## PHỤ LỤC HÌNH ẢNH

Hình 1: Giao diện hiển thị các chức năng hiển thị, hadoop ứng dụng, MapReduce.....	11
Hình 2: Tập dataset được hiển thị sử dụng MongoDB để lấy dữ liệu.....	11
Hình 3: Dữ liệu lấy từ MongoDB.....	12
Hình 4: Trực quan hóa dữ liệu.....	13
Hình 5: Dữ liệu điểm trung bình G1, G2, G3 của từng giới tính lấy từ MongoDB.....	14
Hình 6: Dữ liệu tổng học sinh từng giới tính lấy từ MongoDB.....	15
Hình 7: MrJob.....	15
Hình 8: Upload file.....	16
Hình 9: Kết quả MapReduce.....	16
Hình 10: Mã nguồn.....	17
Hình 11: Biểu Đồ Linear Regression.....	17
Hình 12: Tính toán.....	18
Hình 13: Mã nguồn.....	19
Hình 14: Mã nguồn.....	19
Hình 15: Mã nguồn.....	20
Hình 16: Mã nguồn.....	20
Hình 17: Mã nguồn.....	21
Hình 18: Mã nguồn.....	21
Hình 19: Mã nguồn.....	22
Hình 20: Kết quả.....	22
Hình 21: Kết quả.....	23
Hình 22: Kết quả.....	24
Hình 23: Kết quả.....	25
Hình 24: Kết quả.....	26
Hình 25: Kết quả.....	27
Hình 26: Kết quả.....	28

## LỜI NÓI ĐẦU

Viết một báo cáo đồ án môn học là một trong những việc được coi là khó mà chúng em phải hoàn thành trong quá trình học một môn học. Trong quá trình thực hiện đề tài chúng em đã gặp rất nhiều khó khăn và bế ngõ. Nếu không có những sự giúp đỡ và lời động viên chân thành của nhiều người có lẽ chúng em khó có thể hoàn thành tốt tiểu luận này. Đầu tiên chúng em xin gửi lời biết ơn chân thành đến thầy Hoàng Lê Minh, người trực tiếp hướng dẫn chúng em hoàn thành tiểu luận này.

Những ý kiến đóng góp của thầy là vô cùng hữu ích, nó giúp chúng em nhận ra các khuyết điểm của đồ án và rút kinh nghiệm cho những bài học lần sau. Cảm ơn thầy và các bạn trường Đại học Văn Lang là những người đã cùng nhóm em sát cánh và trải nghiệm để hoàn thành đồ án môn học.

Nhóm xin chân thành cảm ơn.

# MỞ ĐẦU

## 1. LÝ DO CHỌN CHỦ ĐỀ.

Xử lý lượng dữ liệu khổng lồ và Dữ liệu lớn ngày nay đã trở thành một trong những vấn đề quan trọng nhất trong khoa học máy tính. Những khó khăn với giáo dục trong lĩnh vực này phát sinh, các phương pháp và công cụ giảng dạy phù hợp là cần thiết. Việc xử lý một lượng lớn dữ liệu đến một cách nhanh chóng đòi hỏi phải lựa chọn và sắp xếp các nền tảng phần cứng mở rộng. Trong bài báo cáo, nhóm sẽ giả lập áp dụng Hadoop cho việc xử lý dữ liệu lớn của ứng dụng quản lý sinh viên. Việc sử dụng nền tảng học tập điện tử Moodle, một nền tảng dành riêng cho giảng dạy, có thể cho phép đội ngũ giảng viên và sinh viên tiếp xúc tốt hơn bằng cách tăng cường khả năng giao tiếp lẫn nhau. Giải pháp có thể cho phép đối phó với các vấn đề trong giáo dục sinh viên trong lĩnh vực Dữ liệu lớn.

## 2. ĐỐI TƯỢNG, PHẠM VI TÌM HIỂU CỦA CHỦ ĐỀ.

Đối tượng nghiên cứu của bài báo cáo là giáo viên, giảng viên, sinh viên, học sinh của các trường trung học và đại học của các tỉnh trên toàn nước. Chủ yếu lấy dữ liệu bao gồm các thông tin cơ bản của các học sinh và sinh viên từ hệ thống của các trường, từ đó nhóm thống kê thông tin và tạo dữ liệu chuẩn nhất.

Tạo ra một ứng dụng quản lý sinh viên bằng cách áp dụng Hadoop vào môi trường giáo dục, một nền tảng dành riêng cho các giáo viên, giảng viên từ đó dễ dàng lưu trữ và quản lý sinh viên, học sinh giúp đưa ra các quyết định và đề xuất cho tương lai.

# CHƯƠNG 1. GIỚI THIỆU CHỦ ĐỀ ĐỒ ÁN

## 1.1 KHÁI QUÁT MỘT SỐ NỘI DUNG CƠ BẢN VỀ HADOOP VÀ NGUỒN GỐC CỦA HADOOP: [1] [2]

### 1.1.1 Hadoop là gì?

Một câu hỏi quan trọng là loại cơ sở hạ tầng phần cứng và phần mềm nào nên được sử dụng để xử lý các tập dữ liệu lớn? Rõ ràng là nó sẽ là quá nhiều đối với một máy tính, chúng tôi cũng biết rằng ví dụ: Các hệ thống RDBM chỉ có thể được mở rộng đến một giới hạn cố định. Vì vậy, chúng tôi cần một hệ thống có thể được thu nhỏ tuyến tính và với chi phí hợp lý. Chúng ta cần tăng sức mạnh tính toán bằng cách thêm máy tính thay vì thay thế chúng.

Có một giải pháp như vậy - Hadoop, một nền tảng phần mềm có thể mở rộng cho điện toán phân tán. Hadoop có thể lưu trữ kích thước dữ liệu thực tế không giới hạn và có thể xử lý dữ liệu này trong môi trường phân tán. Đây là một giải pháp mã nguồn mở, miễn phí và tương đối đơn giản để mở rộng quy mô. Tất nhiên chúng ta phải ghi nhớ rằng một phần cứng mở rộng sẽ tạo ra chi phí.

### 1.1.2 Nguồn gốc sự phát triển của Hadoop?

Hadoop là một Apache framework có mã nguồn mở được viết bằng Java. Hadoop cho phép người dùng phát triển các ứng dụng phân tán để lưu trữ, quản lý các tập dữ liệu. Một số ứng dụng có khả năng làm việc với hàng trăm node và hàng nghìn petabyte dữ liệu. Hadoop được phát triển dựa trên mô hình MapReduce, theo đó ứng dụng được chia thành nhiều phân đoạn khác nhau và chạy song song trên nhiều node khác nhau. Với cơ chế streaming, Hadoop còn cho phép phát triển các ứng dụng ở dạng phân tán dựa trên các ngôn ngữ lập trình C++, Python, Pearl, ...

Sự phát triển của Hadoop bắt đầu khi một số các kỹ sư phần mềm nhận ra rằng sẽ vô cùng hữu ích để có thể lưu trữ và phân tích các tập dữ liệu lớn hơn nhiều khả năng lưu trữ và truy cập thực tế trên một thiết bị lưu trữ vật lý (như đĩa cứng chẳng hạn). Khởi đầu của ý tưởng này có thể là do các thiết bị lưu trữ vật lý dần dần sẽ phải lớn hơn, cần nhiều thời gian hơn để thành phần đọc dữ liệu từ đĩa (nằm trong đĩa cứng, có thể là phần "head") di chuyển đến một phân đoạn cụ thể nào đó. Thay vào đó, nhiều thiết bị nhỏ hơn làm việc song song sẽ cho hiệu quả tốt hơn một thiết bị lớn.

## 1.2 CÁC MODULES QUAN TRỌNG CỦA NỀN TẢNG HADOOP

Nền tảng Hadoop chứa ba mô-đun quan trọng:

- HDFS – Hadoop Distributed File System;
- YARN – một khuôn khổ để lập lịch công việc và quản lý tài nguyên cụm;

- MapReduce – một hệ thống dựa trên YARN để xử lý song song các tập dữ liệu lớn.

Trong môi trường Hadoop phân tán, mỗi phần dữ liệu được lưu trữ trong một số bản sao (thường ít nhất là 3) trên các máy tính khác nhau (tức là các nút cụm). Một trong những giả định quan trọng nhất là dữ liệu được xử lý cục bộ. Điều đó có nghĩa là dữ liệu được xử lý ở nơi nó được lưu trữ (trên cùng một máy tính/nút) giúp giảm thiểu việc truyền mạng. Hơn nữa, hệ thống có khả năng chịu lỗi, tức là khi một trong các nút bị lỗi, kết quả từ nút này sẽ bị mất. Vì vậy, chỉ cần lặp lại các tính toán từ một nút bị hỏng. Hơn nữa, những tình huống như vậy được quản lý bởi chính hệ thống, vì vậy người dùng không cần thực hiện thêm hành động nào.

Bây giờ chúng tôi sẽ đưa ra mô tả ngắn gọn về các hạng mục chính của nền tảng Hadoop: HDFS và MapReduce.

### *1.2.1 Hadoop Distributed File System*

HDFS (Hadoop Distributed File System) là một hệ thống tệp phân tán, là một phần của hệ sinh thái Apache Hadoop. Nó được thiết kế để lưu trữ và quản lý khối lượng lớn dữ liệu trên nhiều nút trong cụm Hadoop.

Trong Hadoop, mỗi tệp được lưu trữ trong một hệ thống tệp được chia thành nhiều phần (block) và mỗi phần được lưu trữ trong một số bản sao trên các vị trí khác nhau. Kích thước 1 block thông thường là 128MB, kích thước này có thể thay đổi được bằng việc cấu hình.

### *1.2.2 MapReduce*

MapReduce là một framework (với các tool và method) để xử lý dữ liệu song song trong môi trường Hadoop.

- Một thao tác bản đồ trong đó đối với mỗi bản ghi, chúng tôi tính toán các cặp khóa-giá trị; tất cả các cặp có cùng khóa sẽ nằm trong cùng một nhóm, ví dụ: nếu chúng tôi xử lý dữ liệu thời tiết và đối với mỗi quan sát, chúng tôi tính toán một cặp nhiệt độ năm, thì tất cả các quan sát có cùng năm sẽ thuộc cùng một nhóm.
- Một hoạt động giảm trong đó đối với mỗi nhóm, chúng tôi tính toán một số tính năng từ giá trị (aggregate), ví dụ: đối với mọi nhóm quan sát được tính toán ở điểm trước, chúng tôi có thể tính toán nhiệt độ tối đa chẳng hạn. Sơ đồ này rất giống với việc nhóm và tổng hợp thông tin trong SQL.



### 1.3 SO SÁNH SỰ KHÁC BIỆT GIỮA HADOOP VÀ MONGODB [3] [4]

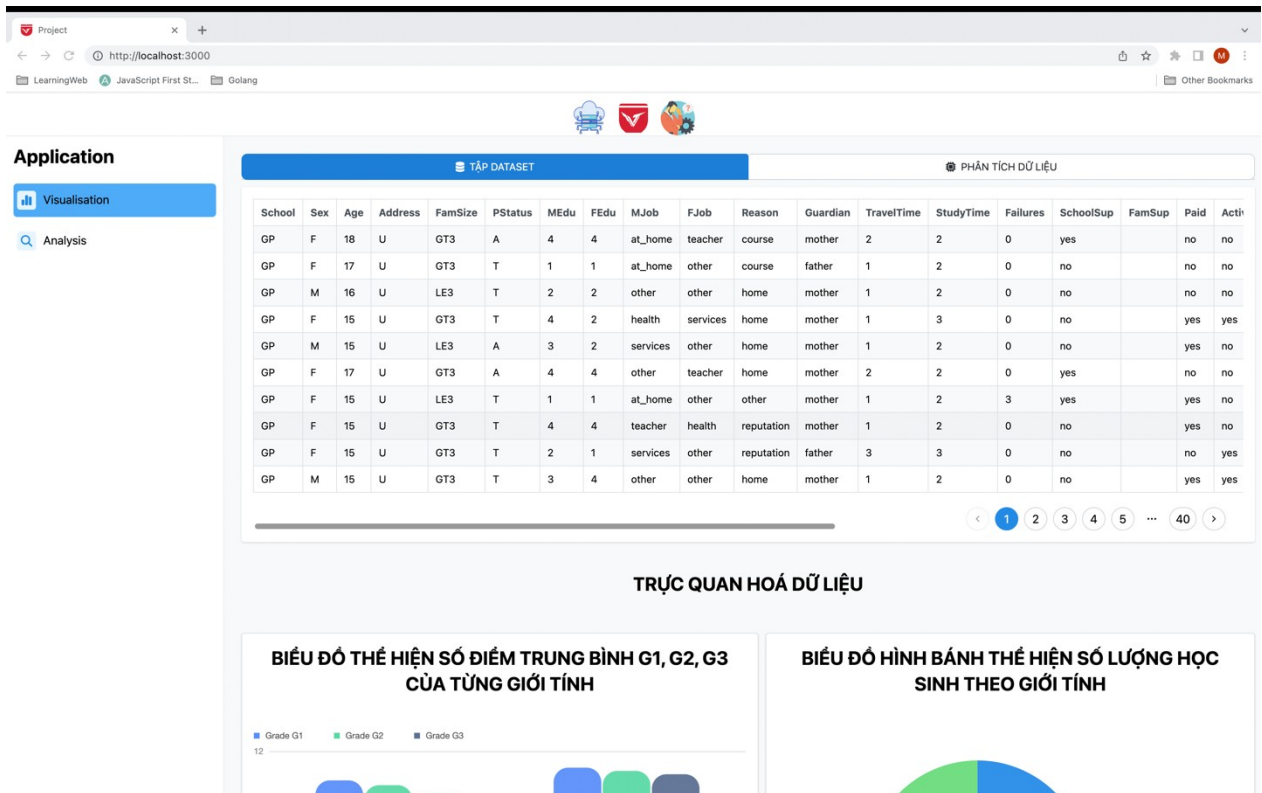
	Hadoop	MongoDB
<b>Data Model</b>	Hadoop là một khung xử lý dữ liệu phân tán chủ yếu tập trung vào xử lý hàng loạt các tập dữ liệu lớn. Nó sử dụng Hệ thống tệp phân tán Hadoop (HDFS) để lưu trữ và quản lý dữ liệu theo cách phân tán. Hadoop không có mô hình dữ liệu được xác định trước và có thể hoạt động với dữ liệu có cấu trúc, bán cấu trúc và phi cấu trúc.	MongoDB là một cơ sở dữ liệu định hướng tài liệu NoSQL cung cấp một lược đồ linh hoạt. Nó lưu trữ dữ liệu ở định dạng giống như JSON được gọi là BSON (JSON nhị phân). MongoDB được thiết kế để xử lý dữ liệu có cấu trúc và bán cấu trúc, đồng thời hỗ trợ các khả năng lập chỉ mục và truy vấn phong phú.
<b>Scalability</b>	Hadoop được thiết kế để mở rộng quy mô theo chiều ngang bằng cách thêm nhiều phần cứng hàng hóa hơn vào cụm. Nó có thể xử lý khối lượng dữ liệu lớn và phân phối khối lượng công việc xử lý trên nhiều nút, giúp nó có khả năng mở rộng cao.	MongoDB cũng có khả năng mở rộng và có thể xử lý các tập dữ liệu lớn. Nó hỗ trợ phân đoạn tự động, cho phép dữ liệu được phân phối trên nhiều nút để chia tỷ lệ theo chiều ngang.
<b>Data Processing Paradigm</b>	Hadoop tuân theo mô hình xử lý hàng loạt và được tối ưu hóa để xử lý khối lượng dữ liệu lớn bằng MapReduce hoặc các khung xử lý phân tán khác như Apache Spark. Nó vượt trội trong việc xử lý khối lượng công việc xử lý dữ liệu ngoại tuyến hoặc hàng loạt.	MongoDB hỗ trợ truy vấn thời gian thực và rất phù hợp với khối lượng công việc xử lý giao dịch trực tuyến (OLTP). Nó cung cấp khả năng truy vấn phong phú, lập chỉ mục và khung tổng hợp, làm cho nó phù hợp để truy cập và phân tích dữ liệu thời gian thực.
<b>Data Model Flexibility</b>	Hadoop có thể xử lý nhiều loại dữ liệu khác nhau, bao gồm dữ liệu có	Mô hình dữ liệu hướng tài liệu của MongoDB cung

	cấu trúc, bán cấu trúc và phi cấu trúc. Nó cung cấp sự linh hoạt trong việc xử lý và phân tích các định dạng dữ liệu khác nhau.	cấp tính linh hoạt trong việc xử lý các cấu trúc dữ liệu khác nhau. Nó cho phép lưu trữ dữ liệu không có lược đồ, làm cho nó có thể thích ứng với các yêu cầu dữ liệu thay đổi.
<b>Data Storage</b>	HDFS của Hadoop được tối ưu hóa để lưu trữ khối lượng lớn dữ liệu trên nhiều nút theo cách phân tán. Nó cung cấp khả năng chịu lỗi và tính sẵn sàng dữ liệu cao.	MongoDB lưu trữ dữ liệu theo cách phân tán trên nhiều nút. Nó sử dụng một định dạng lưu trữ linh hoạt có thể xử lý một loạt các kích thước và cấu trúc dữ liệu.

Tóm lại, Hadoop là khung xử lý dữ liệu phân tán được tối ưu hóa để xử lý hàng loạt và phân tích dữ liệu quy mô lớn, trong khi MongoDB là cơ sở dữ liệu định hướng tài liệu NoSQL được thiết kế để truy vấn thời gian thực và lưu trữ dữ liệu linh hoạt. Lựa chọn giữa Hadoop và MongoDB phụ thuộc vào các trường hợp sử dụng cụ thể, yêu cầu xử lý dữ liệu, nhu cầu về khả năng mở rộng và tính linh hoạt của mô hình dữ liệu mong muốn. Chúng cũng có thể được sử dụng cùng nhau trong một số tình huống nhất định, trong đó Hadoop được sử dụng để xử lý hàng loạt và MongoDB đóng vai trò là kho lưu trữ dữ liệu thời gian thực.

## CHƯƠNG 2. KẾT QUẢ THỰC HIỆN ĐỒ ÁN

### 2.1 THIẾT KẾ ĐƯỢC GIAO DIỆN CƠ BẢN CỦA ỨNG DỤNG QUẢN LÝ SINH VIÊN



Hình 1: Giao diện hiển thị các chức năng hiển thị, hadoop ứng dụng, MapReduce

## 2.2 TIẾN HÀNH TRIỂN KHAI CÁC CHỨC NĂNG CHO ỨNG DỤNG

TẬP DATASET

PHÂN TÍCH DỮ LIỆU

School	Sex	Age	Address	FamSize	PStatus	MEdu	FEdu	MJob	FJob	Reason	Guardian	TravelTime	StudyTime	Failures	SchoolSup	FamSup	Paid	Acti
GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother	2	2	0	yes		no	no
GP	F	17	U	GT3	T	1	1	at_home	other	course	father	1	2	0	no		no	no
GP	M	16	U	LE3	T	2	2	other	other	home	mother	1	2	0	no		no	no
GP	F	15	U	GT3	T	4	2	health	services	home	mother	1	3	0	no		yes	yes
GP	M	15	U	LE3	A	3	2	services	other	home	mother	1	2	0	no		yes	no
GP	F	17	U	GT3	A	4	4	other	teacher	home	mother	2	2	0	yes		no	no
GP	F	15	U	LE3	T	1	1	at_home	other	other	mother	1	2	3	yes		yes	no
GP	F	15	U	GT3	T	4	4	teacher	health	reputation	mother	1	2	0	no		yes	no
GP	F	15	U	GT3	T	2	1	services	other	reputation	father	3	3	0	no		no	yes
GP	M	15	U	GT3	T	3	4	other	other	home	mother	1	2	0	no		yes	yes

1

2

3

4

5

...

40

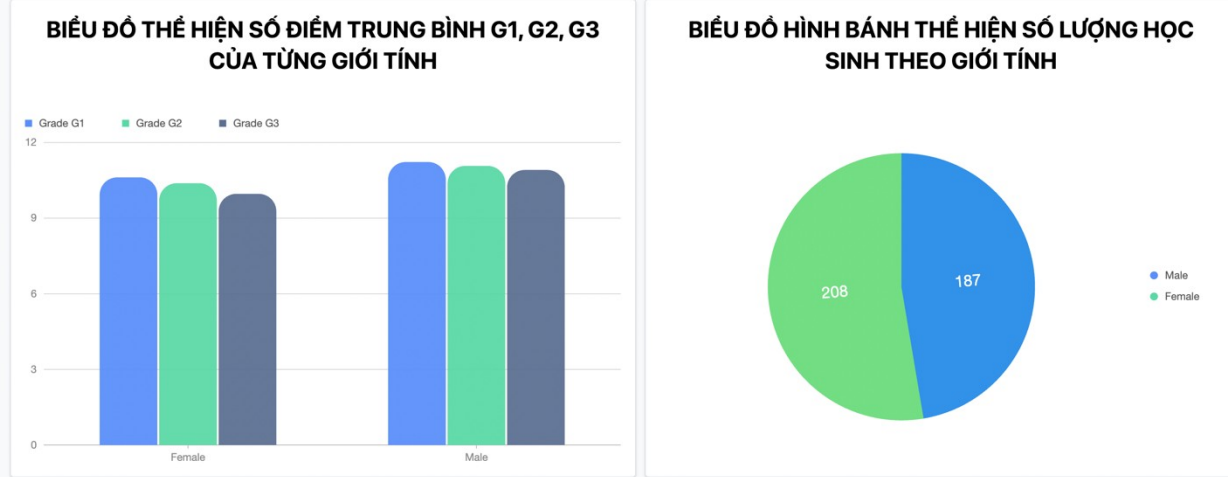
Hình 2: Tập dataset được hiển thị sử dụng MongoDB để lấy dữ liệu

```
mongosh mongodb://<credentials>@127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000 (com.docker.cli)
defaultdb> db.student.find();
[
  {
    _id: ObjectId("64aff51a663ed762f8df0edd"),
    school: 'GP',
    sex: 'F',
    age: 18,
    address: 'U',
    famsize: 'GT3',
    Pstatus: 'A',
    Medu: 4,
    Fedu: 4,
    Mjob: 'at_home',
    Fjob: 'teacher',
    reason: 'course',
    guardian: 'mother',
    traveltime: 2,
    studytime: 2,
    failures: 0,
    schoolsup: 'yes',
    famsup: 'no',
    paid: 'no',
    activities: 'no',
    nursery: 'yes',
    higher: 'yes',
    internet: 'no',
    romantic: 'no',
    famrel: 4,
    freetime: 3,
    goout: 4,
    Dalc: 1,
    Walc: 1,
    health: 3,
    absences: 6,
    G1: 5,
    G2: 6,
    G3: 6
  },
  {
    _id: ObjectId("64aff51a663ed762f8df0ede"),
    school: 'GP',
    sex: 'F',
    age: 17,
    address: 'U',
    famsize: 'GT3',
    Pstatus: 'T',
    Medu: 1,
    Fedu: 1,
    Mjob: 'at_home',
    Fjob: 'other',
    reason: 'course',
    guardian: 'father',
    traveltime: 1,
    studytime: 2,
    failures: 0,
    schoolsup: 'no',
    famsup: 'yes',
    paid: 'no',
    activities: 'no',
    nursery: 'no',
    higher: 'yes',
    internet: 'yes',
```

*Hình 3: Dữ liệu lấy từ MongoDB*

Từ dữ liệu dataset trên chúng ta có thể trực quan hóa dữ liệu để thể hiện xu hướng, tìm ra insightful của data.

### TRỰC QUAN HOÁ DỮ LIỆU



*Hình 4: Trực quan hóa dữ liệu*

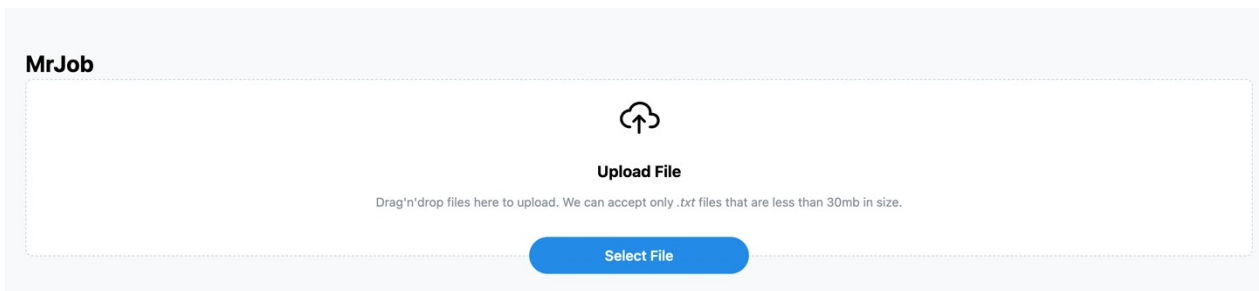
```
mongosh mongodb://<credentials>@127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000 (com.docker.cli)
defaultdb> db.student.aggregate([
...     {
...         "$group": {
...             "_id": "$sex",
...             "avgGradeG1": {"$avg": "$G1"},
...             "avgGradeG2": {"$avg": "$G2"},
...             "avgGradeG3": {"$avg": "$G3"}
...         }
...     }
... ])
[
  {
    _id: 'M',
    avgGradeG1: 11.229946524064172,
    avgGradeG2: 11.074866310160427,
    avgGradeG3: 10.914438502673796
  },
  {
    _id: 'F',
    avgGradeG1: 10.620192307692308,
    avgGradeG2: 10.389423076923077,
    avgGradeG3: 9.966346153846153
  }
]
defaultdb> █
```

Hình 5: Dữ liệu điểm trung bình G1, G2, G3 của từng giới tính lấy từ MongoDB

```
mongosh mongodb://<credentials>@127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000 (com.docker.cli)
defaultdb> db.student.aggregate([
...     {
...         "$group": {
...             "_id": "$sex",
...             "quantity": {"$sum": 1}
...         }
...     }
... ])
[ { _id: 'M', quantity: 187 }, { _id: 'F', quantity: 208 } ]
defaultdb>
```

Hình 6: Dữ liệu tổng học sinh từng giới tính lấy từ MongoDB

Mô tả: MrJob là hàm thư viện sử dụng trong ngôn ngữ Python để có thể sử dụng map reduce khi không cần phải cài đặt hadoop trên máy. Ở ứng dụng này chúng ta có thể bỏ một file text có đuôi là .txt để có thể mapping và reduce dữ liệu.



Hình 7: MrJob

Mô tả: Sau khi tải file thành công, UI sẽ hiển thị thông tin chờ người dùng bấm nút chạy MrJob





```

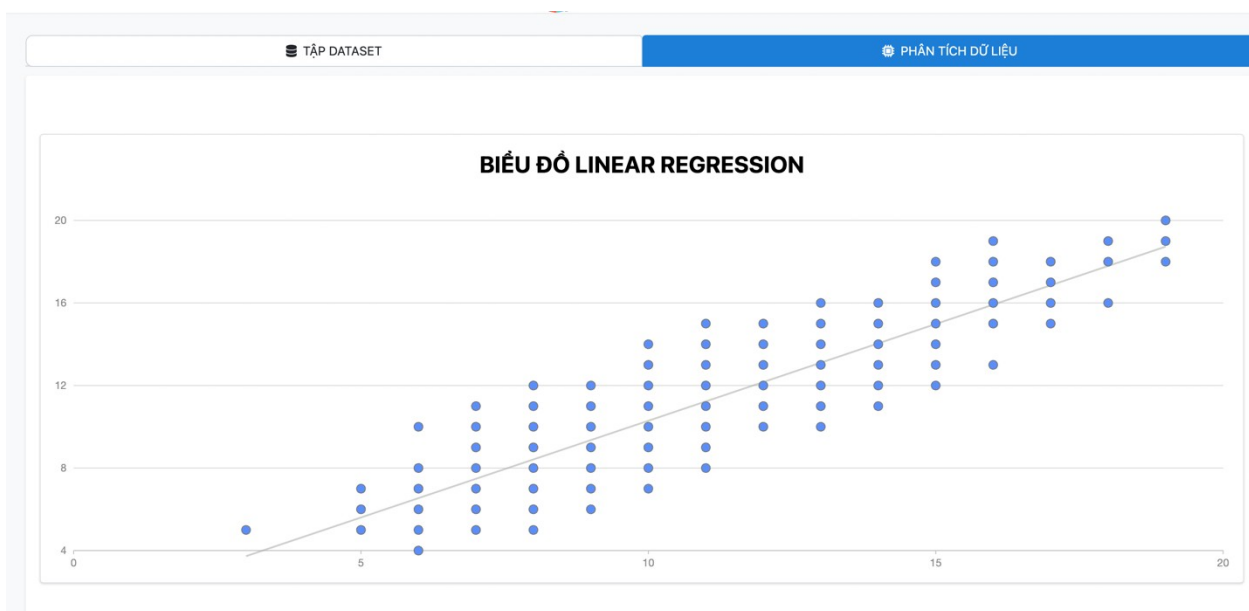
mapping.py > ...
1  from mrjob.job import MRJob
2
3  class Count(MRJob):
4      """ The below mapper() function defines the mapper for MapReduce and takes
5          key value argument and generates the output in tuple format .
6          The mapper below is splitting the line and generating a word with its own
7          count i.e. 1 """
8      def mapper(self, _, line):
9          for word in line.split():
10             yield(word, 1)
11      """ The below reducer() is aggregating the result according to their key and
12          producing the output in a key-value format with its total count"""
13      def reducer(self, word, counts):
14          yield(word, sum(counts))
15
16  if __name__ == "__main__":
17      Count.run()

```

Hình 10: Mã nguồn

Từ tập dataset sau khi đã được sàng lọc và trực quan hóa thì chúng ta có thể áp dụng dữ liệu cho mô hình học máy (Linear Regression).

Mô tả: Trong biểu đồ scatter plot dưới đây, chúng ta có một đường hồi quy (regression line) được sử dụng để biểu diễn dữ liệu và cho thấy khả năng áp dụng Linear Regression để dự đoán điểm số của học sinh.



Hình 11: Biểu Đồ Linear Regression

Sau khi tính toán kết quả dự đoán như sau:

Công thức hồi quy tuyến tính	
Singular Linear Regression	$Y = \beta_0 + \beta_1 X$
Multiple Linear Regression	$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$
Trong đó:	
	X: là biến độc lập
	Y: là biến phụ thuộc
	$\beta_0$ : là intercept (biến chặn)
	$\beta_n$ : là slope or coefficient (hệ số tương quan)
Tính toán	
Tổng G1:	3455.0
Tổng G2:	3366.0
Tổng G3:	3263.0
Tổng G1 square:	3259.604430379746
Tổng G2 square:	4443.708860759492
Tổng (G1 + G2) square:	3171.689873417723
Tổng (G1 + G3) square:	3681.8449367088615
Tổng (G2 + G3) square:	4930.854430379746
Coefficient $\beta_1$ :	0.1631382846956674
Coefficient $\beta_2$ :	0.9931862151944502
Intercept $\beta_0$ :	-2.0370492847090187
Dự đoán 20% từ tập dữ liệu	
SSR (Sum Squared Regression):	332.7375322315393
SSR (total sum of squares):	1619.8987341772147
Mean Absolute Error (MAE):	1.2622306979018147
Mean Squared Error (MSE):	4.211867496601764
R-squared ( $R^2$ ) Score:	0.7945936216805893

Hình 12: Tính toán

Mô tả: Sàng lọc dữ liệu, lấy những cột features quan trọng để áp dụng mô hình học máy. Ở trường hợp này lấy những cột G1, G2, G3 và tách train data và test data.

```

1 extractData.py > ...
2 #!/usr/bin/python3
3
4 # -*- coding: utf-8 -*-
5 #####
6
7 import sys
8 import csv
9 from sklearn.model_selection import train_test_split
10 from io import StringIO
11 import pandas as pd
12
13
14 # Read CSV data from sys.stdin
15 csv_data = sys.stdin.read()
16
17 # Create a StringIO object to simulate a file-like object
18 csv_file = StringIO(csv_data)
19
20 # Use pandas read_csv to read the CSV data
21 csv_reader = pd.read_csv(csv_file, delimiter=',')
22
23 g_df = csv_reader.loc[:, ["G1", "G2", "G3"]]
24 x_filter = g_df.dropna()
25
26 x_filter['G1'] = x_filter['G1'].apply(lambda x: str(x).replace(",", ".").astype(float))
27 x_filter['G2'] = x_filter['G2'].apply(lambda x: str(x).replace(",", ".").astype(float))
28 x_filter['G3'] = x_filter['G3'].apply(lambda x: str(x).replace(",", ".").astype(float))
29 X = x_filter.iloc[:, :2]
30 y = x_filter['G3']
31 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
32
33 X_G1_train = X_train["G1"].array
34 X_G2_train = X_train["G2"].array
35 Y_train_arr = y_train.array
36
37 X_G1_test = X_test["G1"].array
38 X_G2_test = X_test["G2"].array
39 Y_test_arr = y_test.array
40
41 with open("training_data.txt", "w") as f:
42     for row in range(len(X_G1_train)):
43         f.write("{}{}{}{}\n".format(X_G1_train[row], X_G2_train[row], Y_train_arr[row], Y_train_arr[row]))
44     for row in range(len(X_G1_test)):
45         f.write("{}{}{}{}\n".format(X_G1_test[row], X_G2_test[row], Y_test_arr[row], Y_test_arr[row]))
46

```

Hình 13: Mã nguồn

Mô tả: Mã nguồn để mapping dữ liệu trong mapred.

```

1 mapper.py > ...
2 #!/usr/bin/python3
3
4 import sys
5
6 for line in sys.stdin:
7     line = line.strip()
8     Gn = line.split(',')
9     if len(Gn) == 4:
10         print("{}{}{}{}" .format(Gn[0], Gn[1], Gn[2], Gn[3]))

```

Hình 14: Mã nguồn

Mô tả: Mã nguồn để reduce và sử dụng Linear Regression để tính toán mô hình học máy và dự đoán điểm số của học sinh.

```

reducer-linearRegression.py > ...
1  #!/usr/bin/python3
2
3  from operator import itemgetter
4  import sys
5
6  current_word = None
7  current_count = 0
8  total_g2 = 0
9  word = None
10 total_g1 = 0
11 total_g2 = 0
12 total_g3 = 0
13
14 total_g1_square = 0
15 total_g2_square = 0
16 total_g3_square = 0
17
18 total_g1_g3 = 0
19 total_g2_g3 = 0
20 total_g1_g2 = 0
21
22 total_g1_test = 0
23 total_g2_test = 0
24 total_g3_test = 0
25
26 total_g1_test_square = 0
27 total_g2_test_square = 0
28 total_g3_test_square = 0
29
30 total_g1_g3_test = 0
31 total_g2_g3_test = 0
32 total_g1_g2_test = 0
33
34 length_grade = 0
35
36 length_grade_test = 0
37
38 testData = []
39
40 actualData = []
41

```

Hình 15: Mã nguồn

```

reducer-linearRegression.py > ...
42
43 for line in sys.stdin:
44     line = line.strip()
45     # g1, g2, g3 = line.split('; ', 2)
46     Gn = line.split(';')
47
48     if (len(Gn) == 4 and Gn[len(Gn) - 1] == "Train"):
49         try:
50             total_g1 += float(Gn[0])
51             total_g2 += float(Gn[1])
52             total_g3 += float(Gn[2])
53
54             total_g1_square += float(Gn[0]) ** 2
55             total_g2_square += float(Gn[1]) ** 2
56             total_g3_square += float(Gn[2]) ** 2
57
58             total_g1_g2 += float(Gn[0]) * float(Gn[1])
59             total_g1_g3 += float(Gn[0]) * float(Gn[2])
60             total_g2_g3 += float(Gn[1]) * float(Gn[2])
61
62             length_grade += 1
63
64         except ValueError:
65             continue
66     elif len(Gn) == 4 and Gn[len(Gn) - 1] == "Test":
67         try:
68             actualData.append(float(Gn[2]))
69
70             testData.append([float(Gn[0]), float(Gn[1])])
71
72             total_g1_test += float(Gn[0])
73             total_g2_test += float(Gn[1])
74             total_g3_test += float(Gn[2])
75
76             total_g1_test_square += float(Gn[0]) ** 2
77             total_g2_test_square += float(Gn[1]) ** 2
78             total_g3_test_square += float(Gn[2]) ** 2
79
80             total_g1_g2_test += float(Gn[0]) * float(Gn[1])
81             total_g1_g3_test += float(Gn[0]) * float(Gn[2])
82             total_g2_g3_test += float(Gn[1]) * float(Gn[2])
83
84             length_grade_test += 1
85
86         except ValueError:
87             continue
88

```

Hình 16: Mã nguồn

```

reducer-linearRegression.py > ...
90 if length_grade > 0:
91
92     sum_g1 = total_g1
93     print("Sum of G1:", sum_g1)
94
95     sum_g2 = total_g2
96     print("Sum of G2:", sum_g2)
97
98     sum_g3 = total_g3
99     print("Sum of G3:", sum_g3)
100
101     sum_g1_square = total_g1_square
102     regression_sum_g1_square = sum_g1_square - ((sum_g1 ** 2) / length_grade)
103     #  $\sum x_1^2 = \sum x_1^2 - (\sum x_1)^2 / n$ 
104     print("Regression of Sum G1 square:", regression_sum_g1_square)
105
106     sum_g2_square = total_g2_square
107     regression_sum_g2_square = sum_g2_square - ((sum_g2 ** 2) / length_grade)
108     #  $\sum x_2^2 = \sum x_2^2 - (\sum x_2)^2 / n$ 
109     print("Regression of Sum G2 square:", regression_sum_g2_square)
110
111     sum_g1_g2 = total_g1_g2
112     #  $\sum x_1x_2 = \sum x_1x_2 - (\sum x_1\sum x_2) / n$ 
113     regression_sum_g1_g2 = sum_g1_g2 - ((sum_g1*sum_g2) / length_grade)
114     print("Regression of Sum (G1 + G2) square:", regression_sum_g1_g2)
115
116     sum_g1_g3 = total_g1_g3
117     #  $\sum x_1y = \sum x_1y - (\sum x_1\sum y) / n$ 
118     regression_sum_g1_g3 = sum_g1_g3 - ((sum_g1+sum_g3) / length_grade)
119     print("Regression of Sum (G1 + G3) square:", regression_sum_g1_g3)
120
121     sum_g2_g3 = total_g2_g3
122     #  $\sum x_2y = \sum x_2y - (\sum x_2\sum y) / n$ 
123     regression_sum_g2_g3 = sum_g2_g3 - ((sum_g2 + sum_g3) / length_grade)
124     print("Regression of Sum (G2 + G3) square:", regression_sum_g2_g3)
125
126     b1 = ((regression_sum_g2_square * regression_sum_g1_g3) - (regression_sum_g1_g2 * regression_sum_g2_g3)) / (regression_sum_g1_square*regression_sum_g2_square - regression_sum_g1_g2 ** 2) # slope
127     b2 = ((regression_sum_g1_square * regression_sum_g2_g3) - (regression_sum_g1_g2 * regression_sum_g1_g3)) / (regression_sum_g1_square*regression_sum_g2_square - regression_sum_g1_g2 ** 2)
128
129     print("Coefficient b1:", b1)
130     print("Coefficient b2:", b2)
131 else:
132     # Print all the values that declared above
133     print("res", total_g1, total_g2, total_g3)
134

```

Hình 17: Mã nguồn

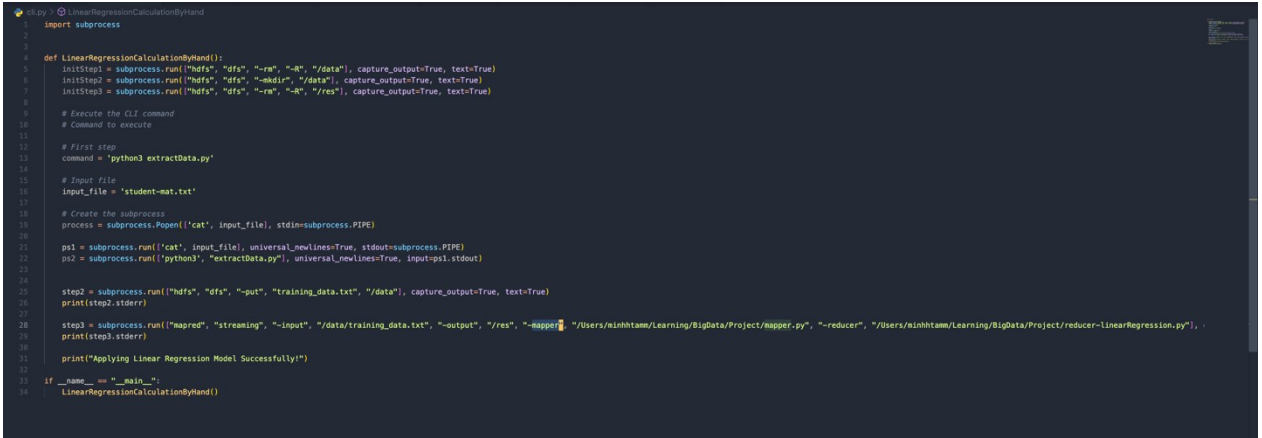
```

reducer-linearRegression.py > ...
135
136 # The formula to calculate b0 is: y_hat - b1x1_hat - b2x2_hat (hat mean average(mean in statistic))
137
138 avg_sum_g3 = sum_g3 / length_grade
139 avg_sum_g1 = sum_g1 / length_grade
140 avg_sum_g2 = sum_g2 / length_grade
141
142 b0 = avg_sum_g3 - b1*avg_sum_g1 - b2*avg_sum_g2
143 print("Intercept b0:", b0)
144
145
146 mean_actual_value = total_g3_test / length_grade_test
147 #  $\hat{y} = b0 + b1x_1 + b2x_2$ 
148 SSR = 0 # SSR => Sum of Squared Residuals
149 SST = 0 # SST => Total Sum of Squares
150 abs_residual = 0
151 for i in range(len(testData)):
152     y_pred = b0 + b1*testData[i][0] + b2*testData[i][1]
153     residual = actualData[i] - y_pred
154     print(f"Predicted values: {y_pred} and actual values: {actualData[i]} => Residual: {residual}")
155     SSR += residual**2
156     abs_residual += abs(residual)
157     SST += ((actualData[i] - mean_actual_value)**2)
158
159 # Calculate the Mean Squared Error(MSE)
160 MSE = SSR / len(testData)
161
162 # Calculate the Total Sum of Squares (SST)
163 SST = SST
164
165 MAE = abs_residual / len(testData)
166
167 print("SSR:", SSR)
168 print("SSR:", SST)
169 R_Square = 1 - (SSR / SST)
170 print("Mean Absolute Error (MAE):", MAE)
171 print("Mean Squared Error (MSE):", MSE)
172 print("R-squared (R^2) Score:", R_Square)
173

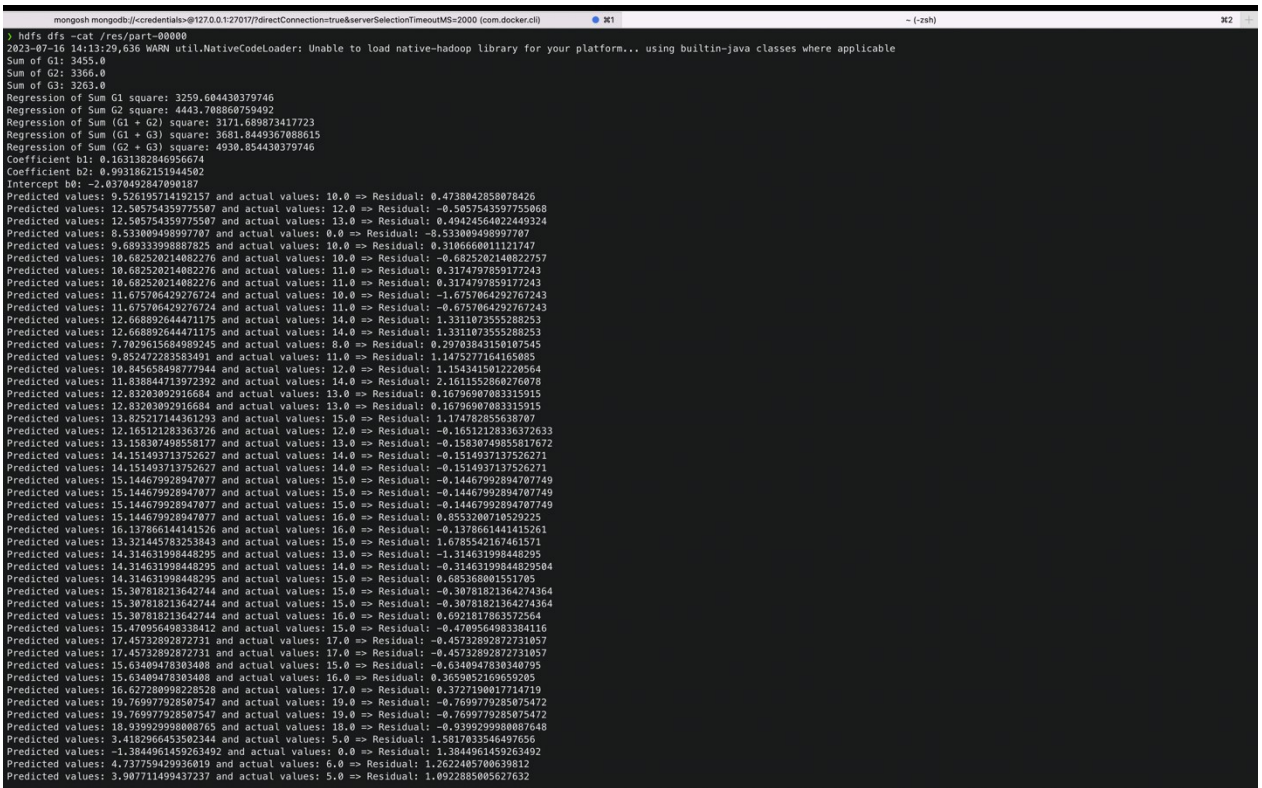
```

Hình 18: Mã nguồn

Mô tả: sử dụng CLI (command line interface) để sàng lọc dữ liệu và sử dụng map reduce để tính toán mô hình học máy.

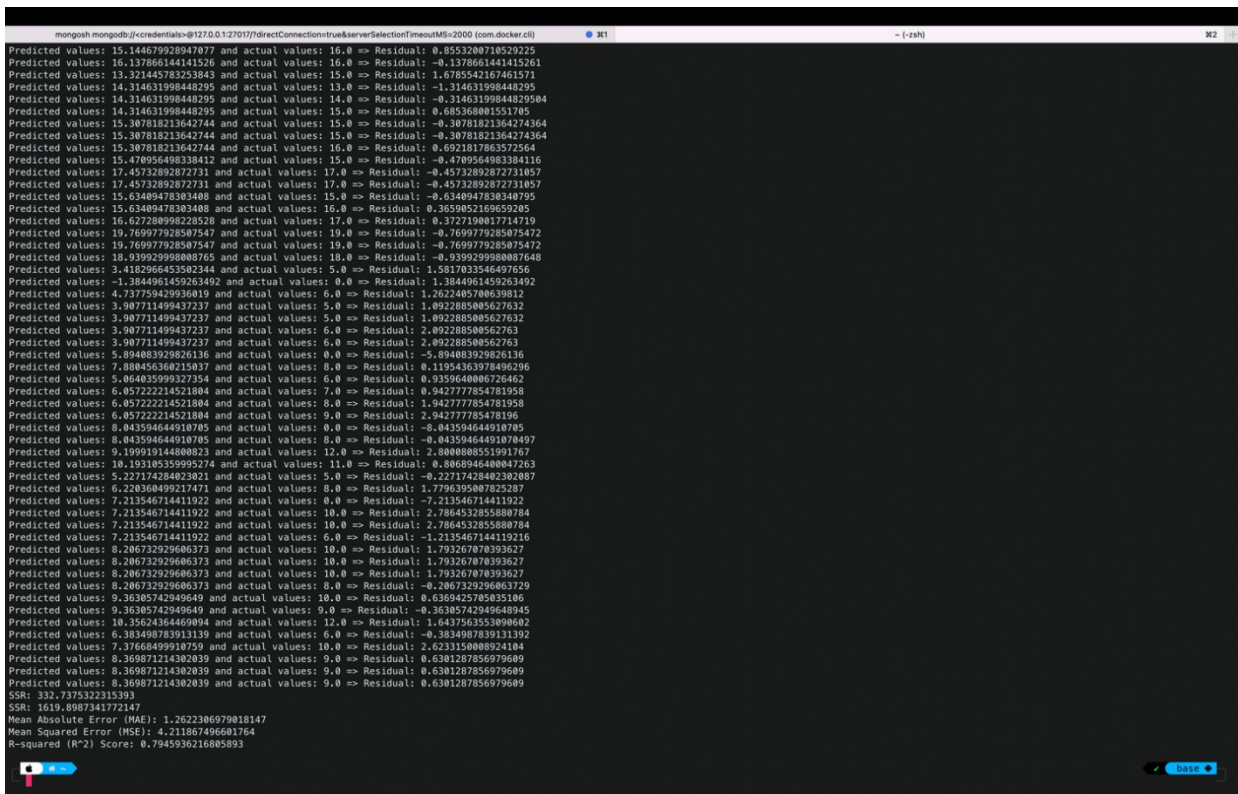


Hình 19: Mã nguồn



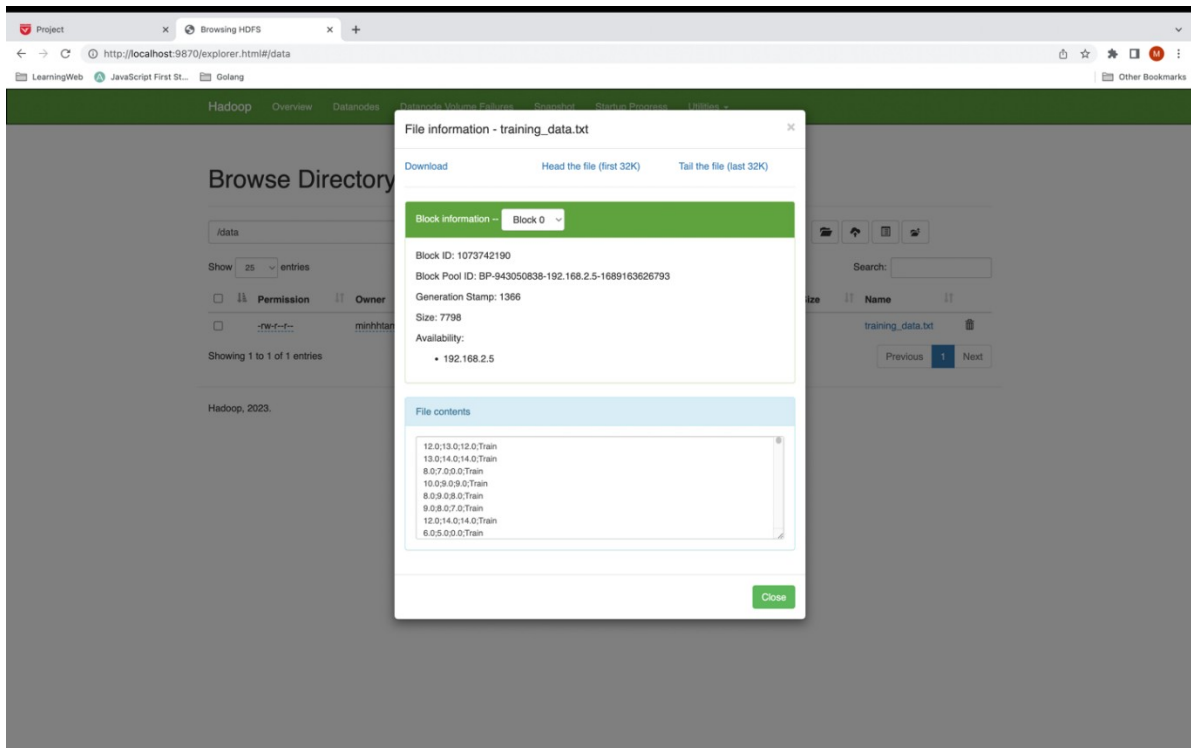
Hình 20: Kết quả





Hình 21: Kết quả

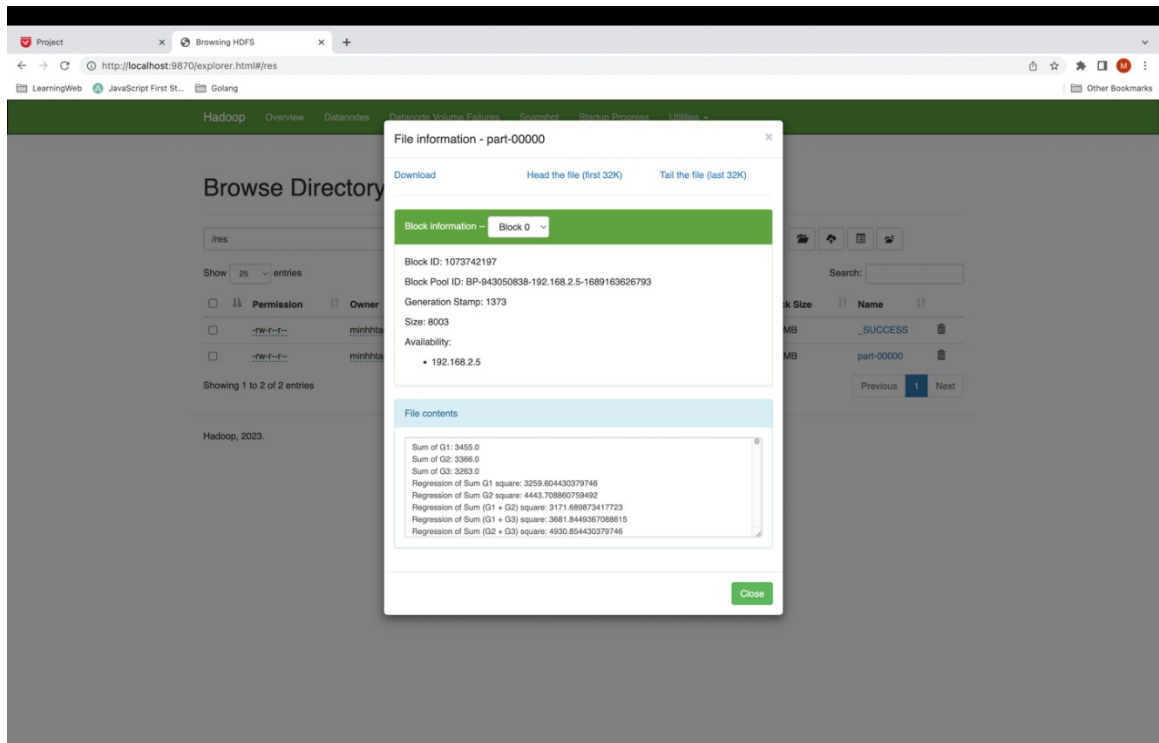
Mô tả: Kết quả trên localhost:9870 của tập dữ liệu được tách ra.



Hình 22: Kết quả

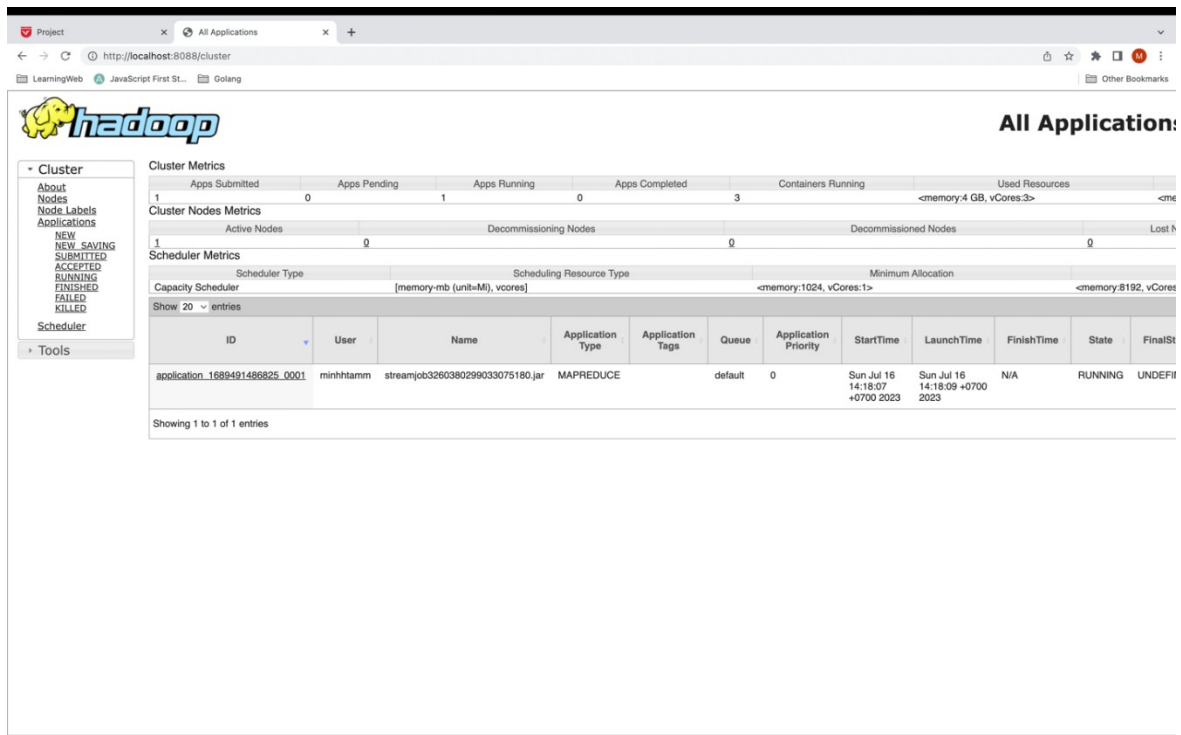


Mô tả: Kết quả được hiện thì trên localhost.



Hình 23: Kết quả

Mô tả: Job công việc của MapReduce trên localhost:8088



Hình 24: Kết quả

Mô tả: Theo dõi công việc trên hadoop

The screenshot shows the Hadoop web interface for monitoring application\_1689491486825\_0001. The interface includes a sidebar with navigation links like Cluster, About, Nodes, Node Labels, Applications, NEW, NEW SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED, Scheduler, and Tools. The main content area displays the application details, including the user (minhtamm), application name (streamjob3260380299033075180.jar), application type (MAPREDUCE), application priority (0), and application state (RUNNING). It also shows the application's progress, including the number of containers preempted and the aggregate resource allocation. A table at the bottom lists the application's attempts, with the first attempt (attempt\_1689491486825\_0001\_000001) shown as completed on Sun Jul 16 14:18:07 +0700 2023.

Attempt ID	Started	Node	Logs	Nodes blacklisted by the app	Nodes blacklisted by the system
attempt_1689491486825_0001_000001	Sun Jul 16 14:18:07 +0700 2023	http://192.168.2.5:8042	Logs	0	0

Hình 25: Kết quả

Mô tả: Logs để tìm ra lỗi hoặc kết quả đầu ra của Job

The screenshot shows the Hadoop web interface for viewing the logs of container\_1689491486825\_0001\_01\_000001. The interface includes a sidebar with navigation links like Resource Manager, RM Home, Node Manager, and Tools. The main content area displays the local logs for the container, including the directory info, launch, prelaunch, and stdout logs. The logs show the container's startup process and the execution of the application.

Local Logs:

```
directory.info : Total file length is 2684 bytes.  
launch_container.sh : Total file length is 5385 bytes.  
prelaunch.err : Total file length is 0 bytes.  
prelaunch.out : Total file length is 100 bytes.  
stdout : Total file length is 1484 bytes.  
stderr : Total file length is 0 bytes.  
syslog : Total file length is 67885 bytes.
```

Hình 26: Kết quả

## KẾT LUẬN VÀ ĐỀ XUẤT

### 1. NHẬN XÉT VỀ HADOOP VÀ NÊU ƯU ĐIỂM VÀ NHƯỢC ĐIỂM CỦA NÓ.

Hadoop là một khung điện toán phân tán mạnh mẽ được thiết kế để xử lý và phân tích khối lượng dữ liệu lớn. Nó cung cấp khả năng mở rộng, khả năng chịu lỗi, hiệu quả chi phí và tính linh hoạt, làm cho nó trở thành lựa chọn phổ biến để xử lý dữ liệu lớn. Với hệ thống tệp phân tán (HDFS), Hadoop cung cấp khả năng lưu trữ có khả năng mở rộng và chịu lỗi, trong khi hệ sinh thái các công cụ và khung mở rộng khả năng của nó cho các khối lượng công việc xử lý dữ liệu khác nhau.

Tuy nhiên, việc thiết lập và quản lý Hadoop có thể phức tạp và có thể không phù hợp để xử lý thời gian thực hoặc độ trễ thấp. Ngoài ra, nó yêu cầu đầu tư cơ sở hạ tầng và có những cân nhắc về đường cong học tập và quản lý dữ liệu. Để dạy học sinh trong lĩnh vực Dữ liệu lớn, các công cụ và phương pháp phù hợp là cần thiết.

Nhìn chung, điểm mạnh của Hadoop nằm ở khả năng mở rộng, khả năng chịu lỗi và khả năng xử lý các tác vụ xử lý dữ liệu đa dạng, khiến nó trở thành một khung có giá trị để xử lý và phân tích dữ liệu lớn.

Ưu điểm	Nhược điểm
<b>Scalability:</b> Hadoop có khả năng mở rộng cao và có thể xử lý khối lượng dữ liệu lớn bằng cách phân phối khối lượng công việc trên nhiều nút trong một cụm. Nó cho phép khả năng mở rộng theo chiều ngang, cho phép các tổ chức thêm nhiều nút hơn vào cụm khi dữ liệu tăng lên, đảm bảo xử lý và lưu trữ hiệu quả.	<b>Complexity:</b> Việc thiết lập và quản lý cụm Hadoop có thể phức tạp và đòi hỏi kiến thức cũng như chuyên môn chuyên môn. Nó liên quan đến việc định cấu hình nhiều thành phần, quản lý tài nguyên cụm và đảm bảo hiệu suất và độ tin cậy tối ưu.
<b>Fault Tolerance:</b> Hadoop cung cấp các cơ chế chịu lỗi tích hợp. Nó sao chép dữ liệu trên nhiều nút, đảm bảo tính khả dụng của dữ liệu ngay cả trong trường hợp lỗi nút. Hadoop tự động phân phối lại dữ liệu và tác vụ cho các nút khỏe mạnh, giảm thiểu tác động của lỗi đối với quá trình xử lý dữ liệu.	<b>Latency:</b> Hadoop được thiết kế chủ yếu cho khối lượng công việc xử lý hàng loạt, có nghĩa là nó có thể không phù hợp cho các yêu cầu xử lý thời gian thực hoặc độ trễ thấp. Chi phí sao chép dữ liệu, xáo trộn dữ liệu và tính toán phân tán có thể gây ra độ trễ trong kết quả xử lý.

<p><b>Cost-Effective:</b></p> <p>Hadoop được thiết kế để chạy trên phần cứng hàng hóa, làm cho nó trở thành một giải pháp hiệu quả về chi phí để xử lý dữ liệu lớn. Nó loại bỏ nhu cầu về phần cứng chuyên dụng đắt tiền và cho phép các tổ chức tận dụng tài nguyên phần cứng giá cả phải chăng để xây dựng cụm Hadoop.</p>	<p><b>Steep Learning Curve:</b></p> <p>Hadoop và hệ sinh thái của nó có đường cong học tập dốc cho các nhà phát triển và quản trị viên. Việc hiểu các khái niệm, mô hình lập trình và công cụ liên quan đến Hadoop đòi hỏi thời gian và công sức.</p>
<p><b>Flexibility in Data Processing:</b></p> <p>Hadoop hỗ trợ nhiều tác vụ xử lý dữ liệu, bao gồm xử lý hàng loạt, xử lý thời gian thực, học máy, v.v. Hệ sinh thái của nó bao gồm các công cụ như MapReduce, Apache Spark, Hive và Pig, mang lại tính linh hoạt và khả năng mở rộng cho các yêu cầu xử lý khác nhau.</p>	<p><b>Hardware and Infrastructure Requirements:</b></p> <p>Hadoop yêu cầu một cụm máy, có thể liên quan đến các khoản đầu tư cơ sở hạ tầng và phần cứng trả trước. Ngoài ra, việc quản lý cụm Hadoop quy mô lớn có thể yêu cầu tài nguyên lưu trữ và tính toán đáng kể.</p>
<p><b>Distributed Storage:</b></p> <p>Hệ thống tệp phân tán Hadoop (HDFS) cung cấp một hệ thống lưu trữ phân tán và có khả năng chịu lỗi. Nó cho phép lưu trữ và xử lý các tập dữ liệu lớn trên nhiều nút, mang lại độ sẵn sàng và độ bền dữ liệu cao.</p>	<p><b>Data Management:</b></p> <p>Mặc dù Hadoop cung cấp một hệ thống lưu trữ phân tán và có thể mở rộng, nhưng nó có thể không phù hợp nhất với mọi loại dữ liệu. Hadoop phù hợp hơn với dữ liệu phi cấu trúc và bán cấu trúc, nhưng nó có thể không tối ưu cho dữ liệu có cấu trúc cao hoặc xử lý giao dịch.</p>
<p><b>Data Locality:</b></p> <p>Nguyên tắc định vị dữ liệu của Hadoop nhằm mục đích xử lý dữ liệu trên cùng một nút nơi dữ liệu được lưu trữ. Điều này làm giảm truyền dữ liệu qua mạng, giảm thiểu chi phí mạng và cải thiện hiệu suất xử lý tổng thể.</p>	

Hadoop vượt trội trong việc xử lý hàng loạt quy mô lớn, lưu trữ phân tán và khả năng chịu lỗi, nhưng nó có thể yêu cầu lập kế hoạch cẩn thận, cân nhắc cơ sở hạ tầng và chuyên môn để tận dụng hết tiềm năng của nó.

## **2. ĐỀ XUẤT TRONG TƯƠNG LAI**

Tóm lại, tương lai của Hadoop nằm ở khả năng tích hợp với điện toán đám mây, áp dụng kiến trúc kết hợp, tiến bộ xử lý luồng, vùng chứa và điều phối, khả năng phân tích nâng cao và học máy, cải thiện quản trị và bảo mật dữ liệu cũng như vai trò của nó trong việc kích hoạt kiến trúc hồ dữ liệu. Những tiến bộ và xu hướng này sẽ nâng cao hơn nữa các khả năng và ứng dụng của Hadoop trong việc giải quyết các yêu cầu xử lý dữ liệu lớn phức tạp và thúc đẩy sự đổi mới trong quá trình ra quyết định dựa trên dữ liệu.

## TÀI LIỆU THAM KHẢO

- [1] bigdataviet, "BIGDATAVIET," BIGDATAVIET, 8 August 2015. [Online]. Available: <https://bigdataviet.wordpress.com/2015/08/08/hadoop-la-gi/>.
- [2] P. N. Nghia, "viblo," Viblo, 22 February 2020. [Online]. Available: <https://viblo.asia/p/tim-hieu-ve-hadoop-bJzKmOBXl9N>.
- [3] Simplilearn, "simplilearn," Simplilearn Solutions, 30 March 2023. [Online]. Available: <https://www.simplilearn.com/hadoop-vs-mongodb-article#:~:text=MongoDB%20is%20a%20C%2B%2B%20based,optimizes%20space%20better%20than%20MongoDB..>
- [4] mongodb, "mongodb," mongodb, [Online]. Available: <https://www.mongodb.com/compare/hadoop-vs-mongodb>.

## LỜI CẢM ƠN

Để hoàn thành tiểu luận này, em xin gửi lời cảm ơn chân thành đến:

Xin cảm ơn giảng viên bộ môn - Thầy Hoàng Lê Minh đã giảng dạy tận tình, chi tiết để em có đủ kiến thức và vận dụng chúng vào bài tiểu luận này.

Do chưa có nhiều kinh nghiệm làm đề tài cũng như những hạn chế về kiến thức, trong bài tiểu luận chắc chắn sẽ không tránh khỏi những thiếu sót. Rất mong nhận được sự nhận xét, ý kiến đóng góp, phê bình từ phía Thầy để bài tiểu luận được hoàn thiện hơn.

Lời cuối cùng, em xin kính chúc thầy nhiều sức khỏe, thành công.