

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Lê Đình Hợp

**Thuật toán đánh chỉ mục ngược với MapReduce
và ứng dụng trong việc đánh giá ý kiến
của học sinh Hòa Bình trên mạng xã hội**

Chuyên ngành: Khoa học máy tính

Mã số : 60 48 01 01

Người hướng dẫn khoa học: PGS TS Đỗ Trung Tuấn

Thái Nguyên, 12 - 2016

Lời cam đoan

Tôi xin cam đoan:

Những kết quả nghiên cứu được trình bày trong luận văn là hoàn toàn trung thực, của tôi, không vi phạm bất cứ điều gì trong luật sở hữu trí tuệ và pháp luật Việt Nam. Nếu sai, tôi hoàn toàn chịu trách nhiệm trước pháp luật.

TÁC GIẢ LUẬN VĂN

Lê Đình Hợp

Lời cảm ơn

Tôi xin chân thành cảm ơn Trường Đại học Công nghệ thông tin và Truyền thông - Đại học Thái Nguyên đã tạo điều kiện thuận lợi cho tôi hoàn thành khóa học này.

Tôi xin chân thành cảm ơn các Thầy Cô giáo – Các nhà khoa học đã trực tiếp giảng dạy truyền đạt những kiến thức chuyên ngành Khoa học máy tính cho tôi trong những tháng năm học tập tại trường.

Đặc biệt tôi xin bày tỏ lòng biết ơn chân thành và sâu sắc tới PGS TS Đỗ Trung Tuấn đã tận tình hướng dẫn, dìu dắt và chỉ bảo cho tôi những kiến thức về chuyên môn thiết thực và những chỉ dẫn khoa học quý báu để tôi hoàn thành bản luận văn này.

Luận văn này còn nhiều thiếu sót, rất mong được các thầy cô giáo trong hội đồng chấm luận văn xem xét, góp ý kiến để luận văn được hoàn thiện hơn.

Tôi xin chân thành cảm ơn!

Thái Nguyên, tháng 12 năm 2016

Mục lục

| | |
|---|-----|
| Lời cam đoan | i |
| Lời cảm ơn..... | iii |
| Mục lục | iv |
| Danh sách các từ viết tắt..... | vi |
| Danh mục các hình vẽ, bảng biểu..... | vii |
| Chương mở đầu | 9 |
| Đặt vấn đề..... | 9 |
| Đối tượng và phạm vi nghiên cứu | 11 |
| Hướng thực hiện đề tài | 11 |
| Những nội dung nghiên cứu chính | 11 |
| CHƯƠNG 1 MÔ HÌNH MapReduce | 12 |
| 1.1. Tổng quan về MapReduce | 12 |
| 1.1.1. Sự quan trọng của MapReduce..... | 12 |
| 1.1.2. Các ý tưởng của MapReduce | 13 |
| 1.1.3. Cấu trúc dữ liệu trong MapReduce | 15 |
| 1.1.4. Mapper và Reducer | 15 |
| 1.1.5. Partitioner và Combiner | 17 |
| 1.2. Bộ khung thực thi | 19 |
| 1.2.1. Lập lịch..... | 19 |
| 1.2.2. Di chuyển dữ liệu và mã lệnh..... | 19 |
| 1.2.3. Đồng bộ hóa | 20 |
| 1.2.4. Xử lý lỗi..... | 20 |
| 1.3. Hệ thống file phân tán..... | 20 |
| 1.3.1. Kiến trúc của HDFS | 21 |
| 1.3.2. Nhiệm vụ của NameNode | 21 |
| 1.3.3. Nhiệm vụ của DataNode | 22 |
| 1.3.4. Nhiệm vụ của Secondary NameNode | 23 |
| CHƯƠNG 2 THUẬT TOÁN XỬ LÝ DỮ LIỆU VĂN BẢN VỚI MapReduce | 25 |
| 2.1. Thiết kế thuật toán MapReduce cơ bản | 25 |

| | |
|--|-----------|
| 2.1.1. Gộp lớn cục bộ | 26 |
| 2.1.2. Bộ hai và bộ ba..... | 30 |
| 2.1.3. Tính toán tần số tương đối | 33 |
| 2.1.4. Sắp xếp thứ cấp | 36 |
| 2.2 Thuật toán tính chỉ mục ngược để tìm kiếm dữ liệu văn bản | 36 |
| 2.2.1. Dò tìm Web | 37 |
| 2.2.2 Thuật toán chỉ mục ngược | 39 |
| 2.2.3. Cài đặt theo cơ bản | 41 |
| 2.2.4. Cài đặt thuật toán cải tiến | 43 |
| 2.2.5. Nén chỉ mục..... | 45 |
| 2.3. Về tìm kiếm | 52 |
| CHƯƠNG 3 THỬ NGHIỆM THUẬT TOÁN ĐÁNH GIÁ Ý KIẾN TRÊN | |
| MẠNG XÃ HỘI | 56 |
| 3.1 Mã nguồn mở Solr | 56 |
| 3.1.1. Giới thiệu..... | 56 |
| 3.1.2. Các tính năng chính của Solr: | 56 |
| 3.2 Mã nguồn mở Nutch..... | 56 |
| 3.2.1. Các lý do để tự xây dựng một Search Engine | 56 |
| 3.2.2. Các tính năng chính của Nutch | 57 |
| 3.3. API biểu đồ Facebook | 58 |
| 3.4. Solr trên Hadoop và tìm kiếm thử nghiệm | 60 |
| 3.4.1. Sơ đồ..... | 60 |
| 3.4.1. Cài đặt cụm máy Hadoop | 62 |
| 3.4.2. Cài đặt Nutch tích hợp với Solr..... | 67 |
| 3.4.3. Thu thập dữ liệu | 69 |
| 3.5. Thực hiện tìm kiếm thử nghiệm trên tập chỉ mục đã thu thập được. | 72 |
| Kết luận | 75 |

Danh sách các từ viết tắt

| | |
|------|---------------------------------|
| CNTT | Công nghệ Thông tin |
| HDFS | Hadoop Distributed File System |
| URL | Uniform Resource Locator |
| HTML | HyperText Markup Language |
| LISP | LISt Processing |
| ML | Markup Language |
| HPC | High-Performance Computing |
| NAS | Network-Attach Storage |
| SAN | Storage Area Network |
| GFS | Google File System |
| SPOF | Single Point Of Failure |
| SNN | Secondary NameNode |
| APW | Associated Press Wordstream |
| REST | Representational State Transfer |
| PRAM | Parallel Random Access Machine |
| BSP | Bulk Synchronous Parallel |

Danh mục các hình vẽ, bảng biểu

| | |
|---|----|
| Hình 1.1. Mô hình chia để trị | 14 |
| Hình 1.2. Hàm Map và Fold trong Functional Programming | 15 |
| Hình 1.3. Hai pha Map và Reduce của một MapReduce job | 16 |
| Hình 1.4. Mô hình MapReduce đầy đủ các thành phần | 19 |
| Hình 1.5. Kiến trúc của HDFS | 21 |
| Hình 1.6. Vai trò của NameNode và DataNode trong HDFS | 23 |
| Hình 1.7. Kiến trúc HDFS đầy đủ | 23 |
| Hình 2.1. Bảo toàn trạng thái trong Hadoop | 26 |
| Hình 2.2. Tiến trình hoạt động của chương trình WordCount | 27 |
| Hình 2.3. Thời gian chạy của thuật toán "pairs" và "stripes" | 32 |
| Hình 2.4. Ví dụ minh họa cặp giá trị | 35 |
| Hình 2.5. Minh họa đơn giản của một chỉ mục ngược | 40 |
| Hình 2.6: Minh họa đơn giản cơ sở thuật toán lập chỉ mục ngược trong MapReduce với ba mapper và hai reducer | 43 |
| Hình 2.7. Mười số nguyên dương đầu tiên trong nguyên phân, γ , và mã Golomb ($b = 5, 10$) | 49 |
| Hình 2.8. Ma trận Term-document | 53 |
| Hình 3.1. Sơ đồ hoạt động của Nutch khi sử dụng như một Crawler | 57 |
| Hình 3.2. Sơ đồ đầy đủ của Nutch khi sử dụng như một Search Engine | 58 |
| Hình 3.3. Facebook | 58 |
| Hình 3.4. Trao đổi qua API | 59 |
| Hình 3.5: Mô hình tổng quan của hệ thống khảo sát | 60 |
| Hình 3.6: Sơ đồ giai đoạn đánh chỉ mục | 61 |
| Hình 3.7: đánh chỉ mục với MapRedece trên Solr | 61 |
| Hình 3.8: Giao diện làm việc của Solr | 68 |
| Hình 3.9: Giao diện làm việc của Facebook Graph API | 69 |
| Hình 3.10: Access Token của một trình Facebook Graph API | 70 |

| | |
|---|----|
| Hình 3.11: Thu thập dữ liệu từ trang mạng của trường THPT Hoàng Văn Thụ | 70 |
| Hình 3.12: Giao diện theo dõi quá trình làm việc của MapReduce | 71 |
| Bảng 3.2: Kết quả thu thập dữ liệu ở 2 chế độ | 72 |
| Hình 3.13: Giao diện trang web tìm kiếm trên Solr | 73 |
| Bảng 3.3: Một số kết quả truy vấn theo chủ đề | 73 |

Chương mở đầu

Đặt vấn đề

Trong thời đại hiện nay, công nghệ thông tin được ứng dụng tại mọi lĩnh vực trong cuộc sống, với một hệ thống máy tính người ta có thể làm được rất nhiều công việc, tiết kiệm được thời gian, công sức và tiền bạc. Với sự phát triển vượt bậc của Internet hiện nay, lượng thông tin ngày càng nhiều, sự tăng trưởng có thể nói là được tính bằng cấp số nhân, theo một nghiên cứu thì cứ khoảng 5 năm thì lượng tri thức của nhân loại sẽ tăng gấp đôi, với lượng thông tin đồ sộ trên mạng hiện nay thì việc tìm kiếm và khai thác thông tin là một công việc hết sức quan trọng, mang lại nhiều lợi ích về khoa học và kinh tế.

Cùng với sự ra đời của Internet, sự xuất hiện và phát triển không ngừng của lĩnh vực thương mại điện tử, các lĩnh vực nghiên cứu xã hội khiến cho việc xúc tiến các hoạt động kinh doanh hoặc nghiên cứu, quảng bá sản phẩm dịch vụ diễn ra trên khắp các kênh thông tin xã hội, đặc biệt là trên Internet.

Như chúng ta đã biết ngày nay mọi thông tin đều được đưa lên các trang mạng xã hội dưới dạng các Posts và rất nhiều người dùng để lại các nhận xét (comments) về các thông tin được đưa lên, ta thấy đó chính là một kho thông tin vô cùng hữu ích, nếu ta có thể tìm kiếm và phân loại dữ liệu ấy, chúng ta có thể thu được các kết quả khảo sát cần thiết phục vụ cho các hoạt động nghiên cứu hoặc các hoạt động sản xuất kinh doanh. Kết quả khảo sát ấy có thể là những tỉ lệ như "thích" (like) hay là không có ý kiến gì đối với một vấn đề được đưa ra. Việc tìm kiếm và xử lý và tổng hợp các thông tin hữu ích đó cần phải có một mô hình đáp ứng được nhu cầu về việc có thể làm việc trên một lượng dữ liệu lớn và tốc độ cao.

Mô hình MapReduce là một mô hình lập trình giúp các ứng dụng có thể xử lý nhanh hơn một lượng dữ liệu lớn dữ liệu trên các máy phân tán song song, độc lập với nhau từ đó giúp rút ngắn thời gian xử lý toàn bộ dữ liệu. MapReduce có thể chạy trên các phần cứng thông thường (commodity hardware), không đòi hỏi các server chạy MapReduce phải là các máy tính có cấu hình đặc biệt mạnh mẽ. Do vậy chi phí triển khai Mapreduce sẽ rẻ hơn.

MapReduce làm đơn giản hóa các giải thuật tính toán phân tán. Với MapReduce, bạn chỉ cần cung cấp hai hàm Map và Reduce cùng với một số thành phần xử lý dữ liệu đầu vào. Do vậy, các nhà phát triển ứng dụng phân tán có thể tập trung nhiều hơn cho phần logic của ứng dụng, bỏ qua các chi tiết phức

tạp của việc phân tán xử lý. Sự ra đời của MapReduce đã mở ra cho các doanh nghiệp và các trung tâm nghiên cứu cơ hội xử lý các nguồn dữ liệu đồ sộ với chi phí thấp và thời gian nhanh hơn. Hiện nay, đã có nhiều công ty lớn triển khai sử dụng mô hình MapReduce trong việc kinh doanh và khảo sát.

Công ty Amazon sử dụng MapReduce để xử lý các file log trong quá trình mua hàng của khách hàng để dự đoán được xu hướng mua hàng...

Facebook có thể xử lý được khối lượng hơn 10 tỷ hình ảnh mà họ đang lưu trữ để thu thập các thông tin về hình ảnh, và thu thập 15 terabyte dữ liệu mỗi ngày vào một kho dữ liệu quy mô Petabyte để thực hiện việc khảo sát và đánh giá xu hướng người dùng.

Việc nghiên cứu về xu hướng, đánh giá khảo sát một vấn đề trên quy mô lớn luôn là 1 vấn đề gặp nhiều khó khăn. Trước đây các nhà khảo sát, đánh giá ý kiến trên các đối tượng nghiên cứu thường sử dụng phương pháp thủ công rất tốn kém và mất rất nhiều thời gian để tổng hợp tin tức, chẳng hạn như muốn khảo sát ý kiến của học sinh đối với một số thay đổi trong chương trình học, người ta không thể lựa chọn hỏi ý kiến của tất cả các học sinh mà chỉ có thể lựa chọn một số địa điểm đặc trưng để thực hiện khảo sát, và đôi khi, kết quả của những khảo sát này không mang được tính khách quan vì tâm lý e ngại của các em học sinh. Và những cuộc khảo sát này, đôi khi phải thực hiện trong vòng một vài năm mới có thể có kết quả tổng hợp. Như vậy là mất rất nhiều công sức, của cải và thời gian. Với việc thực trạng hiện nay hầu hết rất cả các em trong lứa tuổi học sinh, sinh viên đều biết sử dụng và thích tham gia các mạng xã hội trên Internet (đặc biệt là Facebook) thì việc tìm kiếm một từ khóa có tần suất xuất hiện cao sẽ phản ánh được những xu hướng, những ý kiến của người dùng hơn là việc khảo sát thủ công rất nhiều và việc nhận về những kết quả khảo sát ý kiến. Tổng hợp các thông tin trên máy tính với sự hỗ trợ của mô hình MapReduce sẽ giúp chúng ta có thể thực hiện quá trình đánh giá, khảo sát ý kiến hết sức nhanh chóng và mang lại hiệu quả, cũng như tiết kiệm được rất nhiều thời gian và tiền bạc.

Với những nhu cầu cấp thiết trên, học viên thực hiện nghiên cứu kỹ thuật chỉ mục ngược (Inverted Indexing) đó là phương pháp thực hiện quét một lần trên văn bản sau đó lập danh sách các thuật ngữ (từ, cụm từ) trong file đó và bao gồm cả những thông tin đi kèm với mỗi thuật ngữ (term) (vị trí, tần suất, độ quan trọng...). Các thông tin này sẽ được tổ chức theo một cấu trúc dữ liệu riêng và được gọi là chỉ mục. Với phương pháp đánh chỉ mục ngược kết hợp với mô hình MapReduce sẽ giải quyết được những hạn chế trước đây trong phương pháp thông kê, đánh giá ý kiến trên một quy mô lớn, và đó là lý do học viên lựa chọn