

ĐỀ CƯƠNG CHI TIẾT

TÊN ĐỀ TÀI: XÂY DỰNG HỆ THỐNG LỌC RÁC CHO HỆ THỐNG PHÁT HIỆN TIN NÓNG TỪ CÁC TRANG TIN TỨC
Cán bộ hướng dẫn: TS. Huỳnh Ngọc Tín
Thời gian thực hiện: Từ ngày 04/09/2017 đến ngày 04/01/2018
Sinh viên thực hiện: Hoàng Anh Minh – 13520505 Lâm Tuấn Anh – 13520020
Nội dung đề tài: <ol style="list-style-type: none">Mục tiêu đề tài:<ul style="list-style-type: none">Tìm hiểu và đánh giá một số thuật toán phân lớp cho việc lọc rác dữ liệu tin tức.Xây dựng hệ thống lọc rác tin tức cho hệ thống phát hiện tin nóng áp dụng thuật toán phân lớp đã tìm hiểu.Phạm vi đề tài:<ul style="list-style-type: none">Nguồn dữ liệu: các bài viết từ báo chính thống Việt Nam.Ngôn ngữ: tiếng Việt.Các thuật toán tìm hiểu và áp dụng: Naive Bayes, SVM, J48Phương pháp thực hiện:<ul style="list-style-type: none">Tìm hiểu bài toán phân loại tin tức, tìm hiểu các phương pháp và các hướng tiếp cậnThử nghiệm đánh giá các phương pháp đã tìm hiểu: Thu thập dữ liệu tin tức từ cơ sở dữ liệu của công ty VCCorp. Tiến hành một số thống kê trên dữ liệu thu thập được. Huấn luyện và so sánh kết quả các thuật toán phân lớp để lọc nhiễu cho bài toán phân loại tin tức: Naive Bayes, SVM, J48.

4. Kết quả mong đợi:

- Hệ thống lọc rác cho hệ thống phát hiện tin nóng.
- Kiến thức về bài toán phân lớp văn bản
- Báo cáo đề tài

Kế hoạch thực hiện:

- Đặt vấn đề và phát biểu bài toán.
- Lên kế hoạch thực hiện đề tài.
- Tìm hiểu bài toán phân lớp văn bản.
- Tìm hiểu bài toán lọc rác.
- Tìm hiểu framework Struts2, thư viện React, Apache Lucene, công cụ học máy Weka.
- Thu thập và gán nhãn dữ liệu.
- Xử lý dữ liệu, huấn luyện và đánh giá một số mô hình phân lớp phổ biến.
- Cài đặt và hiện thực hệ thống lọc rác.
- Xây dựng và tích hợp hệ thống lọc rác cho hệ thống phát hiện tin nóng.
- Viết báo cáo.

Xác nhận của CBHD

(Ký tên và ghi rõ họ tên)

TP. HCM, ngày....thángnăm.....

Sinh viên

(Ký tên và ghi rõ họ tên)