

BÀI TẬP

Môn X lý ngôn ngữ tự nhiên

Giáo viên: Lê Thanh Hoàng
Email: huonglt@soict.hust.edu.vn

Đây là các bài tập lý thuyết cho sinh viên. Sinh viên có thể xuất tài liệu khác nhau từ các cơ sở dữ liệu có sẵn để giao bài.

1. Tìm hiểu cấu trúc dữ liệu tìm kiếm thông tin Google hiện tại và các kỹ thuật xử lý trong tìm kiếm thông tin của Google

2. Khai phá dữ liệu văn bản: quy tắc nhúng trang web có phải là trang web cá nhân (home page) hay không.

3. Các tính năng pháp xác định biên giới câu.

4. Phân tích cú pháp

5. Phân tích ngữ nghĩa: gì là quy tắc ngữ pháp tham chiếu trong các câu của PTCP ngữ nghĩa

6. Xây dựng chương trình cho phép chuyển đổi các tài liệu văn bản về mặt ngữ nghĩa sang CSDL và các truy vấn dữ liệu để xác định sự liên quan (bảng liên kết CSDL). CSDL có thể bằng tiếng Việt hoặc tiếng Anh. Hãy tìm kiếm các công cụ có sẵn như Gate hay Lucence.

Ví dụ :

a. Thu thập các thông tin liên hệ của các tổ chức có thông tin trên mạng và lưu vào 1 file XML hoặc 1 CSDL gồm có: tên, địa chỉ, số điện thoại, số fax, email. Tiêu chí tìm kiếm chủ yếu là phân loại, ví dụ, tìm các trường đại học và cao đẳng VN, hoặc tìm các công ty tin học Hà Nội.

b. Thu thập thông tin về các cửa hàng bán lẻ điện thoại di động có thông tin trên mạng và lưu vào 1 file XML hoặc 1 CSDL gồm có: tên điện thoại, hãng, tính năng, giá tiền, nơi bán, địa chỉ, điện thoại liên hệ, email liên hệ.

c. Thu thập thông tin về các hội thảo công nghệ thông tin và lưu vào 1 file XML hoặc 1 CSDL gồm có: tên hội thảo, phạm vi hội thảo (trong nước, quốc tế, châu á,...), địa điểm, thời gian diễn ra hội thảo, địa chỉ trang Web, deadline abstract, deadline fullpaper, acceptance time. Tiêu chí tìm kiếm hội thảo chủ yếu là phân loại dựa trên năm diễn ra, ví dụ, call for papers, 2007, 2008, natural language processing.

d. Trích rút tên riêng từ các bài báo tiếng Việt

e. Nhận diện tên thực thể

7. Tóm tắt văn bản, tóm tắt văn bản

8. Phân cụm văn bản

9. Phân loại văn bản:

- phân loại thư, lịch rác
- phân loại trang web

10. Dịch máy từ ngữ kê hoạch ngôn ngữ tự nhiên.

11. Tìm kiếm thông tin:

- xu hướng phát triển công nghệ tìm kiếm siêu dữ liệu và cài đặt

Yêu cầu:

Mỗi nhóm có tối đa 4 người. Nghiên cứu lý thuyết từ 1-2 người. Các tài liệu cài đặt chương trình (có thể tìm kiếm mã nguồn mở có sẵn) từ 1-4 người. Tất cả các nhóm nộp báo cáo và demo chương trình (nếu có). Mỗi người trong nhóm nộp tham gia báo cáo phần kết quả của mình.

Vấn đề báo cáo:

- Báo cáo có $n \geq 8$ trang
- Nội dung liên quan đến cài đặt chương trình, báo cáo về độ hiệu quả tài liệu kiểm thử có phân tích đánh giá mức độ liên quan, phân tích phần cài đặt chương trình (các cấu trúc dữ liệu, thuật toán), mức độ kết quả đạt được, đánh giá chính xác và hiệu quả phát triển.
- Tất cả các báo cáo nộp phải rõ ràng đóng góp của từng thành viên trong nhóm thể hiện tài liệu. Báo cáo có phần tài liệu tham khảo.

Một số mã nguồn mở và X lý ngôn ngữ tự nhiên:

Stanford's Core NLP Suite (viết bằng Java): <http://stanfordnlp.github.io/CoreNLP/>

Natural Language Toolkit (viết bằng Python): <http://www.nltk.org/>

Apache Lucene and Solr: <http://lucene.apache.org/>

Apache OpenNLP (viết bằng Java): <http://opennlp.apache.org/>

Apache UIMA: <https://uima.apache.org/>

GATE (General architecture for text engineering, viết bằng Java): <https://gate.ac.uk/>

Các tài liệu tham khảo:

1. Grant Ingersoll, Thomas Morton, Drew Farris. [*Taming Text*](#) : cho người lập trình biết về NLP và Search. Mỗi chương có ví dụ sử dụng các mã nguồn mở.
2. Steven Bird, Ewan Klein, and Edward Loper. [*Natural Language Processing with Python*](#) : hướng dẫn sử dụng NLTK qua các công cụ phân loại văn bản, trích rút thông tin, ...