

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

NGUYỄN THỊ THU HÒA

NGHIÊN CỨU MỘT SỐ KỸ THUẬT NLP
VÀ ỨNG DỤNG PHÂN LOẠI VĂN BẢN TIẾNG VIỆT

LUẬN VĂN THẠC SĨ
Chuyên ngành: Khoa học máy tính
Mã số: 8480101

Người hướng dẫn: TS Trần Quang Quý

LỜI CẢM ƠN

Em xin gửi lời cảm ơn sâu sắc tới thầy giáo hướng dẫn TS.Trần Quang Quý. Thầy đã giao đề tài và tận tình hướng dẫn em trong quá trình hoàn thành đề tài này.

Em xin gửi lời cảm ơn của mình tới các thầy cô giáo trong trường Đại học Công nghệ thông tin và Truyền Thông - ĐHTN, các thầy cô đã giảng dạy giúp đỡ em trong quá trình học tập tại trường.

Tôi xin chân thành cảm ơn các đồng nghiệp ở cơ quan, trường Cao đẳng Lào Cai. Cảm ơn gia đình bè bạn đã giúp đỡ tôi hoàn thành luận văn này.

Thái Nguyên, ngày 10 tháng 8 năm 2023
Sinh Viên

Nguyễn Thị Thu Hòa

LỜI CAM ĐOAN

Tôi xin cam đoan luận văn này do chính tôi thực hiện, dưới sự hướng dẫn của TS.Trần Quang Quý. Các kết quả lý thuyết được trình bày trong luận văn là sự tổng hợp từ các kết quả đã được công bố và có trích dẫn đầy đủ, kết quả của chương trình thực nghiệm trong luận văn này được tác giả thực hiện là hoàn toàn trung thực, nếu sai tôi hoàn toàn chịu trách nhiệm.

Thái Nguyên, ngày 10 tháng 8 năm 2023
Học viên

Nguyễn Thị Thu Hòa

Mục lục

DANH SÁCH BẢNG	v
DANH SÁCH HÌNH VẼ	v
DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT	vii
MỞ ĐẦU	1
Chương 1. GIỚI THIỆU VỀ XỬ LÝ NGÔN NGỮ TỰ NHIÊN	5
1.1. Các khái niệm cơ bản	5
1.1.1. Xử lý ngôn ngữ tự nhiên	5
1.1.2. Các đơn vị văn bản	5
1.2. Biểu diễn từ	6
1.2.1. onehot vector	7
1.2.2. Word2Vec	7
1.2.3. Glove	9
1.2.4. Biểu diễn t-SNE	10
1.2.5. SVD	11
1.3. Bài toán phân loại văn bản	12
1.3.1. Bài toán phân loại văn bản	12
1.3.2. Một số hướng tiếp cận	13
1.3.3. Một số độ đo mô hình phân loại	14
1.4. Một số thư viện hỗ trợ NLP	15
Chương 2. MỘT SỐ KỸ THUẬT HỌC SÂU TRONG NLP	17
2.1. Kiến trúc Transformer	18
2.1.1. Kiến trúc transformer	18
2.1.2. Mô hình chuỗi sang chuỗi	19
2.1.3. Kiến trúc tự tập trung	21
2.1.4. Các kỹ thuật trong transformer	22
2.1.5. Bộ mã hóa và giải mã trong transformer	25
2.1.6. Huấn luyện transformer	27
2.2. Mô hình bert	28
2.2.1. Tinh chỉnh bert	28
2.2.2. Mặt nạ ngôn ngữ	30
2.2.3. Các kiến trúc mô hình BERT	31
2.2.4. BERT trong Tiếng Việt	31
2.2.5. Một số kỹ thuật tokenize	34

2.3. Một số mô hình học sâu hiện đại khác	37
Chương 3. ỨNG DỤNG PHÂN LOẠI VĂN BẢN TIẾNG VIỆT	40
3.1. Chuẩn bị dữ liệu	41
3.1.1. Nguồn dữ liệu	41
3.1.2. Đọc và lưu dữ liệu	42
3.1.3. Tokenize nội dung	45
3.2. Thiết lập mô hình mạng	46
3.2.1. Cấu hình mô hình BERT	46
3.2.2. Kiến trúc mô hình	47
3.3. Huấn luyện mô hình	48
3.3.1. Thuật toán huấn luyện mô hình	48
3.3.2. Load mô hình BERT	52
3.3.3. Huấn luyện mô hình	53
Kết luận	55
Tài liệu tham khảo	56

Danh sách bảng

3.1	Bảng mô tả dữ liệu huấn luyện.	42
3.2	Bảng mô tả dữ liệu test	43

Danh sách hình vẽ

1.1	So sánh giữa CBOW và Skip gram	9
2.1	Kiến trúc transformer	20
2.2	Kiến trúc seq2seq	20
2.3	Các tầng trong bộ mã hóa và giải mã	21
2.4	Kiến trúc tự tập trung	22
2.5	Kỹ thuật tập trung đa đầu	23
2.6	Kỹ thuật biểu diễn vị trí trong transformer	25
2.7	Minh họa dự đoán ở bước thời gian t của tầng tự tập trung	27
2.8	Tiền trình huấn luyện trước và tinh chỉnh của BERT	29
2.9	Kiến trúc bert base và bert large	32
3.1	Kiến trúc tinh chỉnh của BERT cho tác vụ phân loại.	48

DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT

NLP	Xử lý ngôn ngữ tự nhiên
Word Representation	Biểu diễn từ.
Tokenize	Tách từ.
$ \mathbf{x} _2$	Chuẩn 2 của véc tơ, chuẩn Euclide
\mathbb{R}^n	Không gian véc tơ thực n chiều.
Sklearn	Thư viện học máy cơ bản trong python.
Word embedding	Nhúng từ.
SVD	Phân tích ma trận bằng phương pháp suy biến.
PCA	Một phương pháp giảm chiều, phân tích thành phần chính.
Word2vec	Một phương pháp đưa từ thành véc tơ.
Layers	Các lớp, tầng trong mạng nơ ron.
Skip-Grams	Một phương pháp nhúng từ.
$P(x y)$	Xác suất xuất hiện x với điều kiện y .
Layer	Tầng, lớp mạng.
CBOV	Một mô hình nhúng từ.
TensorFlow	Một framework hỗ trợ thực hành deeplearning bằng python.
epochs	Số lần huấn luyện mô hình.
t-SNE	Một thuật toán giảm chiều dùng để biểu diễn từ.
gensim	Một thư viện xử lý ngôn ngữ tự nhiên.
Bert	Mô hình biểu diễn từ hai chiều ứng dụng Transformer
Attention	Cơ chế chú ý
fine tuning	Tinh chỉnh mô hình (tham số)
MLM	Mặt nạ LM - tác vụ tinh chỉnh biểu diễn từ
Fully connect	Kết nối đầy đủ
SOTA	State-of-the-art Kết quả tốt nhất
Roberta	Một kiến trúc, thuật toán của facebook trên pytorch.
BPE	Mã hóa Byte Pair Encoding

MỞ ĐẦU

Xử lý ngôn ngữ tự nhiên (NLP) đại diện cho một phân nhánh quan trọng trong lĩnh vực Trí tuệ nhân tạo, đặc biệt chú trọng vào việc khám phá sự tương tác giữa máy tính và ngôn ngữ tự nhiên của con người, bất kể đó là qua hình thức tiếng nói (speech) hay văn bản (text) [1], [2]. Mục tiêu vượt qua trong lĩnh vực này là giúp máy tính thấu hiểu và thực hiện các nhiệm vụ liên quan đến ngôn ngữ của con người một cách hiệu quả, bao gồm việc tương tác giữa con người và máy, cải thiện chất lượng giao tiếp giữa con người và con người, cũng như tối ưu hóa quá trình xử lý cả văn bản và lời nói.

Xử lý ngôn ngữ tự nhiên có nguồn gốc từ những năm 1940 của thế kỷ 20, trải qua một hành trình phát triển đa dạng với một loạt phương pháp và mô hình xử lý khác nhau. Các giai đoạn quan trọng bao gồm việc áp dụng các phương pháp ô-tô-mát và mô hình xác suất vào những năm 1950, việc khai thác ký hiệu và các phương pháp ngẫu nhiên trong những năm 1970, tiến bộ sử dụng học máy truyền thống vào đầu thế kỷ 21, và đặc biệt là sự bùng nổ của học sâu trong thập kỷ gần đây [3], [4].

Xử lý ngôn ngữ tự nhiên có khả năng được phân chia thành hai hướng phát triển quan trọng, không hoàn toàn độc lập, gồm xử lý tiếng nói (speech processing) và xử lý văn bản (text processing). Xử lý tiếng nói tập trung vào việc nghiên cứu và phát triển các thuật toán cũng như chương trình máy tính nhằm xử lý ngôn ngữ con người ở dạng tiếng nói (dữ liệu âm thanh). Trong lĩnh vực này, nhận dạng tiếng nói và tổng hợp tiếng nói đứng ra như những ứng dụng quan trọng. Nhận dạng tiếng nói liên quan đến việc chuyển đổi ngôn ngữ từ dạng tiếng nói sang dạng văn bản, trong khi tổng hợp tiếng nói thực hiện chuyển ngôn ngữ từ dạng văn bản thành tiếng nói.

Xử lý văn bản tập trung vào việc phân tích dữ liệu văn bản. Các ứng dụng quan trọng trong lĩnh vực này bao gồm tìm kiếm và truy xuất thông tin, dịch máy, tóm tắt văn bản tự động và kiểm tra chính tả tự động. Thậm chí, xử lý văn bản có thể tiếp tục được chia thành hai phân nhánh nhỏ hơn, bao gồm hiểu văn bản và sinh văn bản. Hiểu văn bản liên quan đến các nhiệm vụ phân tích văn bản, trong khi sinh văn bản liên quan đến việc tạo ra văn bản mới, như trong các ứng dụng dịch máy hoặc tóm tắt văn bản tự động [4].

Bài toán phân loại văn bản là một trong những bài toán quan trọng trong lĩnh vực xử lý ngôn ngữ tự nhiên. Quy trình cơ bản để xây dựng mô hình phân loại văn bản sử dụng các mô hình học máy, học sâu là [4]:

- Chuẩn bị dữ liệu: Thu thập và chuẩn bị dữ liệu cho bài toán phân loại văn bản. Dữ liệu này gồm các văn bản đã được gán nhãn (label) cho từng loại phân loại. Ví dụ: nếu muốn phân loại email là spam hay không, dữ liệu sẽ bao gồm các email đã được gán nhãn là "spam" hoặc "không phải spam".
- Tiền xử lý dữ liệu: Dữ liệu văn bản thường cần được tiền xử lý để loại bỏ các ký tự đặc biệt, chuyển đổi về dạng chuẩn (ví dụ: chuyển đổi thành chữ thường), loại bỏ từ dừng (stop words) và thực hiện các quy trình như tokenization, stemming hoặc lemmatization.
- Trích xuất đặc trưng: Văn bản cần được biểu diễn thành các đặc trưng số học để máy học có thể xử lý. Một số phương pháp phổ biến để trích xuất đặc trưng từ văn bản là sử dụng bag-of-words, TF-IDF (Term Frequency-Inverse Document Frequency), Word2Vec, hoặc BERT (Bidirectional Encoder Representations from Transformers).
- Xây dựng mô hình: Chọn một thuật toán phù hợp để huấn luyện mô hình phân loại. Các thuật toán phổ biến dùng trong phân loại văn bản bao gồm Naive Bayes, Logistic Regression, Support Vector Machines (SVM), Random Forest, hoặc các mô hình học sâu như Recurrent Neural Networks (RNNs) và Convolutional Neural Networks (CNNs).
- Huấn luyện và đánh giá mô hình: Sử dụng dữ liệu đã được gán nhãn, huấn luyện mô hình trên tập huấn luyện và đánh giá hiệu suất của mô hình trên tập kiểm tra. Các phép đo đánh giá thông thường bao gồm độ chính xác (accuracy), độ phủ (recall), độ chính xác trung bình trọng số (weighted precision), và F1-score.
- Tinh chỉnh mô hình

Bài toán phân loại văn bản tiếng Việt có một số điểm khác biệt so với phân loại văn bản nói chung do đặc thù ngôn ngữ và ngữ cảnh văn hóa. Một số điểm khác biệt quan trọng [9]:

- Ngôn ngữ và Từ vựng: Ngôn ngữ tiếng Việt có các đặc điểm riêng, bao gồm cấu trúc ngữ pháp, từ vựng và cách ngữ âm. Điều này tạo ra một số thách thức trong việc biểu diễn từ vựng, xử lý ngữ pháp và thậm chí việc phát hiện từ ngữ mới hay viết sai.
- Từ viết tắt và biểu ngữ cụ thể: Ngôn ngữ tiếng Việt thường sử dụng rất nhiều từ viết tắt, biểu ngữ cụ thể và ngôn ngữ thông tin (ví dụ: "e", "mình", "mk" thay cho "tôi", "mình"). Điều này có thể làm cho việc biểu diễn ngôn ngữ và phân loại trở nên phức tạp hơn.