

Đại học Bách Khoa Hà Nội
Viện Công nghệ thông tin và Truyền thông

ĐỒ ÁN MÔN HỌC

Phân loại văn bản tin tức

Môn: Học máy
Mã lớp: 95090

Giảng viên hướng dẫn: **TS. Nguyễn Nhật Quang**

Nhóm thực hiện

Phan Ngọc Lâm - 20142505
Nguyễn Duy Mạnh - 20142857

Hà Nội, 5/2017

Mục lục

1. Bài toán.....	3
1.1. Mô tả và yêu cầu.....	3
1.2. Một số cách tiếp cận phổ biến.....	4
1.3. Hướng giải quyết.....	4
2. Mạng neuron - Mạng hồi quy.....	5
2.1. Mạng neuron (Artificial Neural Network – ANN).....	5
2.2. Mạng hồi quy (Recurrent Neural Network – RNN).....	5
2.3. Chính quy hóa L2 (L2 Regularization).....	7
3. Thiết kế, cài đặt.....	9
3.1. Kiến trúc mạng neuron.....	9
3.2. Tập dữ liệu.....	10
3.3. Quá trình huấn luyện.....	11
3.4. Cài đặt chi tiết.....	12
3.4.1. Đọc và tiền xử lý dữ liệu.....	12
3.4.2. Mô hình phân loại.....	12
3.4.3. Huấn luyện và chạy mô hình.....	15
3.4.4. Cấu hình huấn luyện.....	15
3.4.5. Đánh giá mô hình.....	16
3.4.6. Các hàm tiện ích.....	16
3.4.7. Các demo.....	17
4. Kết quả.....	19
5. Kết luận và hướng phát triển.....	22
5.1. Khó khăn, tranh luận và hướng giải quyết.....	22
5.2. Kết luận.....	22
5.3. Hướng phát triển.....	22

1. Bài toán

1.1. Mô tả và yêu cầu

Phân loại tin tức là bài toán có ứng dụng rộng rãi không chỉ trong các ứng dụng và dịch vụ tin tức (news aggregator), mà còn có ảnh hưởng trong việc sàng lọc dữ liệu hay tìm kiếm thông tin nói chung. Việc phân loại đặc biệt hữu ích trong các hệ thống thu thập tin tức, trong đó các tin có thể đến từ nhiều nguồn với các cách phân loại khác nhau, không rõ ràng hoặc thậm chí không có phân loại trước. Bài toán

Nhóm phát biểu bài toán phân loại văn bản tin tức như sau:

Đầu vào : 1 văn bản tin tức (bao gồm tiêu đề và nội dung) và 1 số các chủ đề có sẵn.

Đầu ra : Chủ đề phù hợp nhất với văn bản được đưa vào.

Ví dụ : Với tập chủ đề (Sports, World, Business, Sci/Tech),

Input	Output
Tiêu đề : Fears for T N pension after talks Nội dung : Unions representing workers at Turner Newwall say they are 'disappointed' after talks with stricken parent firm Federal Mogul.	Business
Tiêu đề : The Race is On: Second Private Team Sets Launch Date for Human Spaceflight (SPACE.com) Nội dung : SPACE.com - TORONTO, Canada -- A second team of rocketeers competing for the 36.10 million Ansari X Prize, a contest for privately funded suborbital space flight, has officially announced the first launch date for its manned rocket.	Sci/Tech

Bài toán tương tự nhưng khác với bài toán xây dựng chủ đề (topic modeling), trong đó không có 1 tập chủ đề xác định mà thay vào đó lời giải bao gồm 1 tập các chủ đề được tạo ra với số lượng cho trước.

Lời giải cần thỏa mãn 1 số yêu cầu:

- Có độ chính xác, độ hợp lý cao. Trong 1 số trường hợp, có thể chấp nhận 1 lượng sai sót nhất định, nhất là khi nội dung văn bản không thực sự nằm gọn trong 1 chủ

đề cho sẵn. Tuy nhiên, lời giải không thể quá “bất hợp lý” (ví dụ phân loại văn bản tài chính vào thể thao,...).

- Có khả năng thích ứng cao với các văn bản mới.
- Thời gian dự đoán tương đối tốt.

1.2. Một số cách tiếp cận phổ biến

Cách tiếp cận phổ biến nhất với bài toán là sử dụng các giải thuật dựa trên xác suất, với điển hình là giải thuật Naive Bayes [1]. SVM[2] cũng là 1 giải pháp thường được sử dụng.

Một cách tiếp cận thứ hai sử dụng các kĩ thuật xử lý ngôn ngữ tự nhiên như các bộ đoán nhận, bộ dịch [3].

Một số nghiên cứu cũng có sử dụng đến mạng neuron và Deep Learning.

1.3. Hướng giải quyết

Nhóm quyết định sử dụng một mô hình mạng neuron để giải quyết bài toán, với ưu điểm là khả năng chịu nhiễu và làm việc với dữ liệu thô tốt, đồng thời có thời gian chạy ngắn sau giai đoạn huấn luyện. Cụ thể, mô hình sẽ tích hợp mạng hồi quy xử lý từng kí tự trong văn bản để tiến hành phân loại.

Để cài đặt hệ thống, nhóm sử dụng ngôn ngữ lập trình Python và Keras [4], một thư viện Deep Learning mã nguồn mở. Keras cho phép sử dụng 1 trong 2 nền tảng tính toán là TensorFlow (Google) và Theano (cả 2 đều hỗ trợ tăng tốc bằng GPU), cài đặt sẵn các dạng neuron dưới dạng các lớp có tính module hóa cao và độ tối ưu tốt. Quá trình cài đặt thuật toán do đó khá ngắn gọn, cho phép nhóm tập trung vào thiết kế mô hình và xử lý dữ liệu.

2. Mạng neuron - Mạng hồi quy

2.1. Mạng neuron (Artificial Neural Network – ANN)

2.2. Mạng hồi quy (Recurrent Neural Network – RNN)

Mạng hồi quy là 1 mô hình mạng neuron đặc biệt, trong đó đầu vào của mỗi lớp RNN được chia thành các bước rời rạc (time-step). Mỗi neuron không xử lý toàn bộ đầu vào mà xử lý lần lượt trên từng time-step, với trạng thái của bước sau phụ thuộc vào bước trước. Đặc điểm này cho phép mạng hồi quy phản ánh quan hệ phụ thuộc giữa các phần tử trong một chuỗi, ví dụ giữa các từ trong 1 câu hay các khung trong 1 video. RNN do đó được thường được sử dụng để học trên các luồng input tuyến tính như văn bản, tín hiệu,... Một trong những thành công đầu tiên trong ứng dụng RNN là các mô hình dịch máy.

Một neuron (hay cell) của mạng hồi quy có 2 input và 2 output. Input bao gồm giá trị của ví dụ tại time-step tương ứng x_i và 1 vector “trạng thái” (hay hidden state) được truyền xuyên suốt các time-step. Giá trị ban đầu của vector này bằng 0. Tương tự như neuron thông thường, input của ví dụ được nhân với vector trọng số U . Hidden state cũng được nhân với trọng số riêng W . 2 giá trị thu được được cộng vào nhau và đưa qua 1 hàm kích hoạt (hidden activation) để được giá trị s_i . Giá trị này tiếp tục được nhân với trọng số ra V , sau đó lại được đưa qua hàm kích hoạt thứ hai (output activation) để được o_i . o_i trở thành đầu ra thấy được của time-step, còn s_i trở thành hidden state đưa vào time-step tiếp theo.

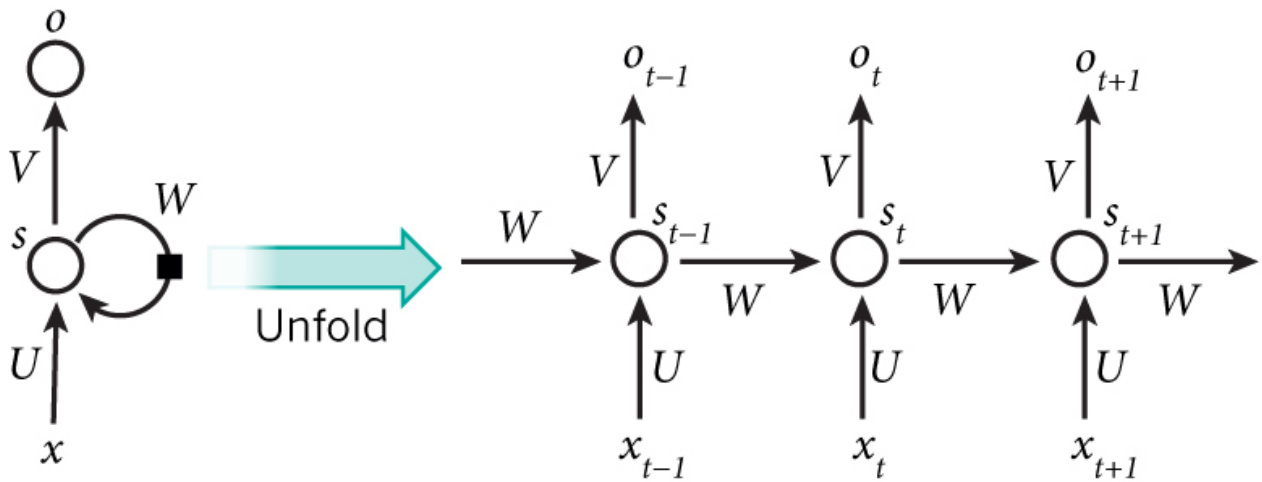
$$s_i = \sigma_g(Ux_i + Ws_{i-1} + b_s)$$

$$o_i = \sigma_h(Vs_i + b_o)$$

trong đó σ_g là hàm kích hoạt ẩn, σ_h là hàm kích hoạt đầu ra, b là các trọng số bias.

Trong cài đặt, σ_g thường được chọn là hàm tanh() hoặc sigmoid(). σ_h được chọn dựa trên output mong muốn của neuron.

Trên lý thuyết, một neuron có thể nhận vào các ví dụ có độ dài bất kỳ. Tuy nhiên, để tăng tốc độ tính toán, ta thường muốn gộp các ví dụ thành 1 ma trận và xử lý theo batch. Do đó khi cài đặt các ví dụ thường được thêm (pad) các kí tự đặc biệt để đạt cùng độ dài.

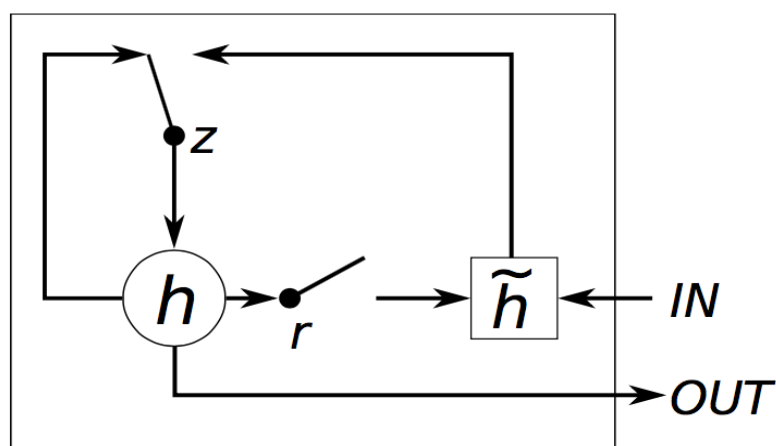


Hình 1: Mô hình 1 RNN Cell [0]

Các vector trọng số U, W, V là mục tiêu học của mỗi neuron và là như nhau với mỗi time-step trong quá trình dự đoán.

Để học các tham số của U, W, V , RNN sử dụng 1 biến thể của thuật toán lan truyền ngược là lan truyền ngược qua thời gian (Back-propagation Through Time – BPTT). BPTT coi các xử lý ở các time-step là các neuron có chung trọng số, hay nói cách khác, ta trải các time-step thành 1 mạng neuron. Sau đó, tương tự như back-propagation, ta áp dụng quy tắc chuỗi đạo hàm trên từng time-step, với lỗi truyền từ time-step cuối đến time-step đầu. Tuy nhiên, vì các trọng số được chia sẻ giữa các bước, nên ta lấy tổng gradient ở mỗi bước và update trọng số sau khi đi hết time-step đầu.

Trong Deep Learning, ta thường sử dụng 2 biến thể của mạng hồi quy là Long Short-Term Memory (LSTM) và Gated Recurrent Unit (GRU). Phần sau sẽ tập trung giải thích cơ chế hoạt động của GRU - mô hình được sử dụng trong báo cáo.



Hình 2: Mô hình 1 GRU Cell

Khác với 1 neuron hồi quy thông thường, neuron của GRU có output bằng hidden state (được kí hiệu h). Giá trị này được tính dựa trên input đầu vào của time-step, giá trị đầu ra của time-step trước và trọng số của 2 “cổng” update z và reset r . Công thức tính của h như sau:

$$\begin{aligned} z_i &= \sigma_g(W_z x_i + U_z h_{i-1} + b_z) \\ r_i &= \sigma_g(W_r x_i + U_r h_{i-1} + b_r) \\ h_i &= z_i \circ h_{i-1} + (1 - z_i) \circ \sigma_h(W_h x_i + U_h (r_i \circ h_{i-1}) + b_h) \end{aligned}$$

Cổng update z quyết định lượng thông tin trong hidden state được lưu giữ lại (z tăng $\rightarrow z_i \circ h_{i-1}$ tăng). Ngược lại, cổng reset quyết định lượng thông tin trong hidden state được dùng để kết hợp với input. Nếu z có giá trị lớn và xấp xỉ 1 (σ_g thường mặc định là sigmoid), ta gần như giữ nguyên hidden state và bỏ qua input đầu vào. Đặc tính này sẽ ứng với các neuron học được các phụ thuộc “dài”, giữa các từ đứng xa nhau trong câu. Ngược lại, với z nhỏ và r nhỏ, neuron gần như reset lại hidden state cũ, tuy nhiên lại kết hợp nó với input của time-step. Đây là các neuron học được các phụ thuộc “ngắn”, giữa các từ đứng gần nhau trong câu.

Trong một mạng GRU, mỗi neuron sẽ học được một dạng dependency nhất định trong các trọng số W và U . Các thuộc tính được sinh ra do đó phản ánh rất tốt cả quan hệ ngữ nghĩa ngắn (phủ định, khẳng định,...) và dài (chủ đề, ngữ cảnh,...).

2.3. Chính quy hóa L2 (L2 Regularization)

Một vấn đề thường gặp trong các thuật toán học máy là overfit, trong đó mô hình học quá sát với tập học và mất khả năng khái quát hóa, thể hiện ở độ chính xác tập học rất cao nhưng độ chính xác trên các ví dụ chưa học lại rất thấp.

Các mạng neuron nhiều tầng đặc biệt dễ rơi vào tình trạng này, do khả năng khớp dữ liệu của mạng là rất tốt, dẫn đến giới hạn quyết định quá phức tạp và chỉ phù hợp với tập học.

Một cách làm giảm bớt hiện tượng này là việc sử dụng 1 biểu thức “chính quy hóa” (regularization). Cụ thể, hàm mục tiêu của mạng được cộng thêm 1 lượng $\lambda \sum W^2$,

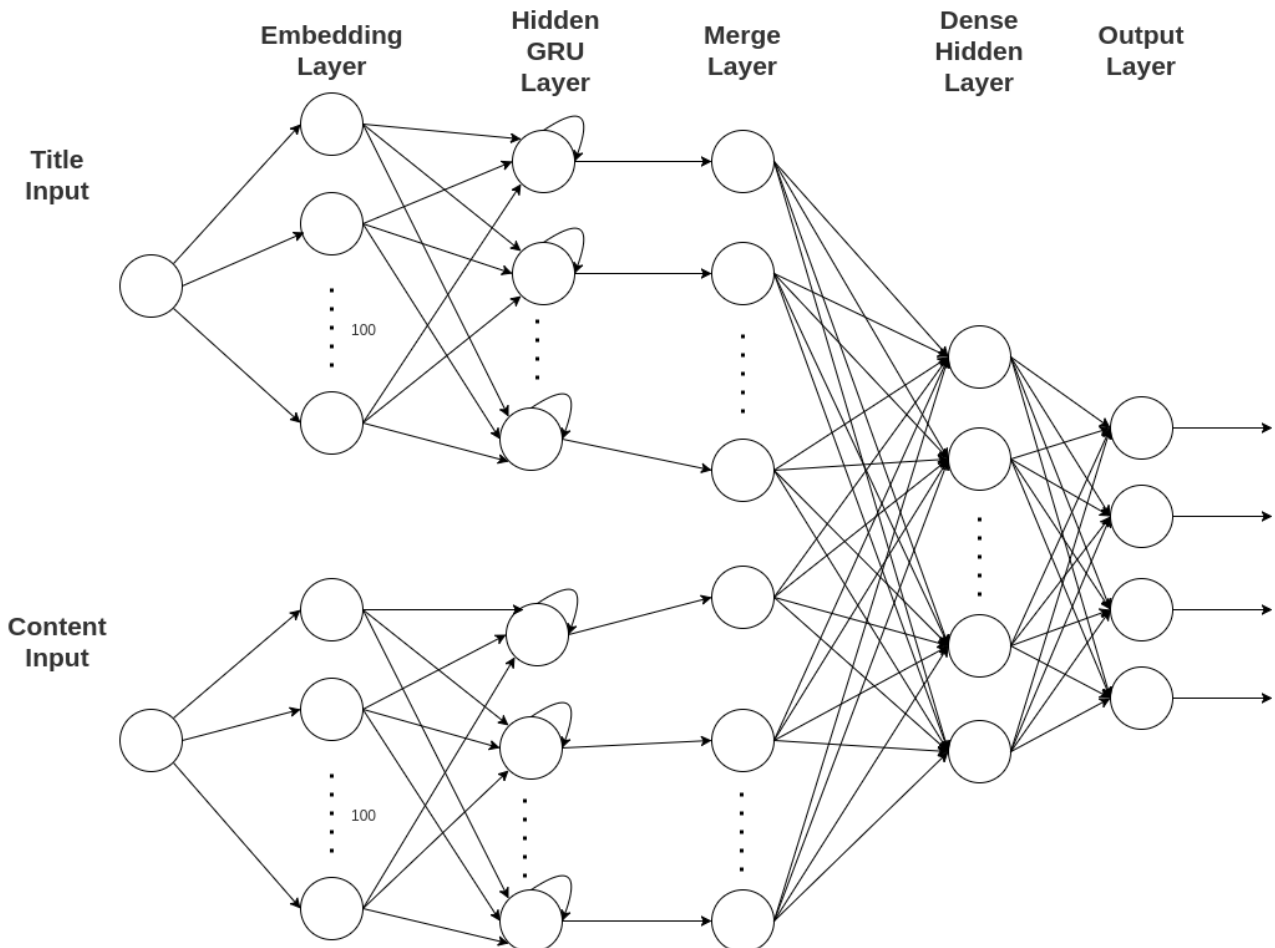
với W là tập tất cả các trọng số của mạng, và λ là một tham số được lựa chọn. Vì mạng neuron có xu hướng làm giảm hàm mục tiêu, nên khi cộng vào các trọng số của mạng, ta yêu cầu mạng cố gắng giữ các trọng số này càng nhỏ càng tốt, từ đó tạo ra khả năng khái quát tốt hơn. λ điều khiển mức độ giảm các trọng số của mạng, và thường được quyết

định qua thử nghiệm.

Đồ án sử dụng chính quy hóa L2 để điều chỉnh các mô hình được huấn luyện,

3. Thiết kế, cài đặt

3.1. Kiến trúc mạng neuron



Hình 3: Kiến trúc mạng, với số chủ đề 4.

Mạng neuron nhận vào 2 input: tiêu đề và nội dung của văn bản. Thay vì gộp chung vào thành 1 chuỗi, ta dùng tiêu đề như 1 tham số thứ 2, do tiêu đề thường phản ánh ngắn gọn nội dung của văn bản và có giá trị cao trong phân loại. Mỗi từ trong chuỗi đưa vào được biểu diễn dưới dạng 1 vector (đề án sử dụng độ dài 100) sinh bởi 1 lớp Embedding. Lớp này thực chất là 1 tập các vector từ tương ứng với từ điển, $E(i) = W_i$, với W là ma trận trọng số của E . Trọng số ban đầu của E được lấy từ tập con của biểu diễn từ [Glove](#) [5] và cũng được cập nhật trong quá trình huấn luyện.

Để thực hiện xử lý batch, ta buộc phải cố định độ dài của mỗi chuỗi. Do đó, nhóm sử dụng 1 ký tự đặc biệt MASK_TOKEN có thứ tự 0 để thể hiện ký tự rỗng.

Các vector được đưa vào 1 lớp mạng GRU và xử lý lần lượt như phần 2.2. Hàm kích hoạt cho các neuron GRU là $\tanh()$, hàm kích hoạt trong là $\text{sigmoid}()$. Khi gặp MASK_TOKEN, GRU sẽ tự động bỏ qua từ mà không thực hiện tính toán. Output của từ cuối cùng trong câu được dùng làm vector biểu diễn cho cả câu.¹

Sau khi có được 2 vector nội dung và tiêu đề, ta thực hiện ghép chúng lại để tạo thành vector văn bản D. Có thể coi vector này là 1 tập các feature mới được học bởi mạng GRU. Lớp GRU được dùng như một lớp mã hóa (encoder), tìm ra các đặc điểm cần cho phân loại.

Bài toán trở thành phân loại các văn bản với đặc điểm D. Ta sử dụng 1 mạng neuron nhân tạo để học cách phân loại các đặc điểm này, các lớp đều có hàm kích hoạt sigmoid (hệ thống cho phép tùy chỉnh độ sâu và kích cỡ đầu ra của từng lớp mạng). Cuối cùng, hệ thống coi mỗi văn bản chỉ có 1 chủ đề, do đó ta lấy đầu ra là 1 phân bố xác suất sinh bởi hàm kích hoạt softmax:

$$\text{softmax}(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

với z là vector độ dài K , $j = 1 \dots K$

Softmax trả về 1 vector có tổng bằng 1, với tỉ lệ tương đối giữa các phần tử được bảo toàn. Do đó có thể coi kết quả thu được là 1 phân bố xác suất đúng của mỗi chủ đề trên văn bản đầu vào. Chủ đề có xác suất cao nhất được chọn làm chủ đề của văn bản.

Các lớp mạng đều sử dụng và tối ưu hóa theo điều chỉnh chính quy L2.

3.2. Tập dữ liệu

Đồ án sử dụng 3 tập dữ liệu: AG-News có độ lớn 120,000 (train) và 7,600 (test) với 4 nhãn, BBC có độ lớn 1,777 (train) và 447 (test) với 5 nhãn, và Reuters có độ lớn 7,769 (train) và 3,018 (test) với 86 nhãn.

Nhóm quy định định dạng chuẩn của các dataset theo dataset AG-News, trong đó mỗi dataset cần có 3 file:

- train.csv chứa các ví dụ của tập train, lưu trong 3 cột nhãn (số nguyên bắt đầu từ 1), tiêu đề và nội dung.
- test.csv chứa các ví dụ của tập test với định dạng tương tự.

¹ Một cách tiếp cận khác để tạo biểu diễn câu là sử dụng 1 vector trung bình của tất cả các từ trong câu sau khi đưa qua 1 lớp trọng số. Tuy nhiên phương pháp này trở nên phức tạp khi phải xử lý các MASK_TOKEN được bỏ qua.

- classes.txt chứa tên các nhãn tương ứng với thứ tự trong các file CSV, mỗi nhãn trên 1 dòng.

Để nhận thấy rằng tập AG-News có tỉ lệ rất không cân đối giữa 2 tập train và test. Do đó nhóm thực hiện phân bố lại tập này tuân theo stratified sampling và thu được 1 dataset với .

Tập dữ liệu Reuters được thêm vào trong giai đoạn cuối của đồ án và có nhiều chủ đề cho mỗi đoạn. Để tương thích với hệ thống, nhóm buộc phải lấy chủ đề đầu tiên của mỗi đoạn làm nhãn của đoạn đó. Điều này làm cho một số nhãn có số ví dụ rất ít hoặc không có ví dụ nào trong tập test. Đây là 1 hạn chế nhóm chưa kịp khắc phục trong tập dữ liệu và có ảnh hưởng đến mô hình học trên tập này.

Các tập huấn luyện được tách validation tuân theo stratified sampling với tỉ lệ 80:20 khi huấn luyện.

3.3. Quá trình huấn luyện

Hệ thống được huấn luyện offline trên 1 tập dữ liệu đã xác định trước. Bước đầu, các file đầu vào (tập train, test và tên các nhãn) được nạp vào bộ nhớ cùng với **1 phần** của tập vector Glove (nhóm sử dụng 30,000 – 70,000 từ đầu tiên trên tổng số 400,000). Tập vector được chèn thêm vector của kí tự MASK_TOKEN (một vector 0) và UNKNOWN_TOKEN tượng trưng cho các từ không có trong từ điển (một vector có giá trị ngẫu nhiên). Tiếp đó, các đoạn đi qua quá trình tách từ (tokenization) được thực hiện bởi thư viện Natural Language Toolkit (NLTK) [6]. Lý do ta không thực hiện tách từ theo dấu trắng là để tách các các dấu liền với từ (VD 'chào.') thành 1 từ riêng biệt. Sau đó, mỗi từ được thay thế bằng chỉ số của vector tương ứng trong tập vector. Chương trình cũng in ra tỉ lệ UNKNOWN_TOKEN để tránh trường hợp từ điển quá nhỏ (tỉ lệ trong các huấn luyện của nhóm vào khoảng 2-5%). Khi kết thúc giai đoạn nạp dữ liệu, ta thu được một dãy vector từ, các ma trận số nguyên X_t và X_c (tiêu đề và nội dung) với kích cỡ $n.s$ (n : số ví dụ trong tập, s : độ dài tối đa của 1 câu) và vector nhãn y cho mỗi tập dữ liệu train, validation và test.

Sau khi nạp dữ liệu, hệ thống thực hiện khởi tạo mạng neuron và thiết lập cấu hình huấn luyện, bao gồm thuật toán huấn luyện và tốc độ học. Thuật toán huấn luyện được nhóm sử dụng là RMSProp [7]. RMSProp là 1 thuật toán học có thích nghi được sử dụng phổ biến với các mạng hồi quy. Tốc độ học ban đầu được nhóm lựa chọn là 0.0003.

Quá trình huấn luyện được thực hiện nhiều lần với các tham số điều chỉnh L2 khác nhau cho các lớp GRU và lớp thường. Kết quả cuối cùng được lưu là mô hình có độ chính xác tốt nhất trên tập validation.

Trong mỗi lần huấn luyện, thuật toán chạy trên từng mini-batch của tập train. Sau mỗi

lần đi hết tập train, ta đánh giá hàm mục tiêu và độ chính xác trên tập train và validation. Huấn luyện sẽ dừng lại khi:

- Hàm mục tiêu trên tập train hội tụ (giá trị bắt đầu tăng) hoặc
- Hàm mục tiêu trên tập validation bắt đầu phân kì, được nhóm quy ước là khi giá trị hàm không cải thiện sau 5 vòng liên tiếp.

Sau mỗi vòng, mô hình được lưu lại chỉ khi có cải thiện hàm mục tiêu trong lần huấn luyện và độ chính xác validation trong các lần huấn luyện trước.

Khi kết thúc một lần huấn luyện, ta đánh giá độ chính xác của mô hình trên tập test. Các đồ thị hàm mục tiêu và độ chính xác cũng được lưu trong thư mục của mô hình.

3.4. Cài đặt chi tiết

3.4.1. Đọc và tiền xử lý dữ liệu

Nhóm cài đặt các hàm đọc và tiền dữ liệu trong script `data_utils.py`, bao gồm:

- `load_embedding()` nạp vào các vector từ Glove và đưa vào 1 ma trận. Đồng thời hàm cũng trả về 1 mapping từ-số hiệu (word_to_index) và số hiệu-từ (index-to-word) để thực hiện các chuyển đổi.
- `read_csv()` đọc file csv được truyền vào, thực hiện tách từ bằng thư viện NLTK và trả về danh sách các ví dụ. Mỗi ví dụ được biểu diễn dưới dạng 1 từ điển với 3 khóa 'class', 'title' và 'content'.
- `strat_samples()` nhận vào 1 danh sách văn bản được sinh bởi `read_csv()` và tách thành 2 tập tuân thủ theo stratified sampling với tỉ lệ bất kỳ.
- `get_mat()` thực hiện chuyển đổi danh sách các văn bản thành các ma trận để đưa vào huấn luyện.
- `load_generic()` thực hiện lần lượt các thao tác trên với một thư mục bất kỳ.
- `load_ag_news()`, `load_bbc()` và `load_reuters()` là các hàm bao cho `load_generic()`.

3.4.2. Mô hình phân loại

Nhóm thiết lập 1 lớp Classifier để chứa mô hình phân loại. Ngoài mạng neuron và các trọng số, Classifier cần lưu các thông tin như từ điển, tên các nhãn,... để tiến hành phân loại. Các tham số như kích cỡ đầu ra, độ sâu mạng,... đều có thể được điều chỉnh. Hàm khởi tạo của lớp được định nghĩa như sau:

```
def __init__(self, word_vec, word_to_index, index_to_word, classes,
              title_output=128, content_output=512,
              dense_neurons=(1024, 256,), title_len=50, content_len=2000,
```

```
        weights=None, directory='.',
        gru_regularize=0, dense_regularize=0):
    self.directory = directory
    self.word_to_index = word_to_index
    self.index_to_word = index_to_word
    self.title_len = title_len
    self.content_len = content_len
    self.word_vec = word_vec
    self.classes = classes
    self.title_output = title_output
    self.content_output = content_output
    self.dense_neurons = dense_neurons
    self.gru_regularize = gru_regularize
    self.dense_regularize = dense_regularize

    # Encode document's title
    title_inp = Input(shape=(title_len,), name='Title_Input')
    title_embed = Embedding(input_dim=np.size(word_vec, 0),
                           output_dim=np.size(word_vec, 1),
                           weights=[word_vec], mask_zero=True,
                           name='Title_Embedding')
    self.t_encoder = Sequential(name='Title_Encoder')
    self.t_encoder.add(title_embed)
    self.t_encoder.add(GRU(title_output, name='Title_GRU', consume_less='mem',
                           W_regularizer=WeightRegularizer(l2=gru_regularize)))
    title_vec = self.t_encoder(title_inp)
    # Encode document's content
    content_inp = Input(shape=(content_len,), name='Content_Input')
    content_embed = Embedding(input_dim=np.size(word_vec, 0),
                              output_dim=np.size(word_vec, 1),
                              weights=[word_vec], mask_zero=True,
                              name='Content_Embedding')
    self.c_encoder = Sequential(name='Content_Encoder')
    self.c_encoder.add(content_embed)
    self.c_encoder.add(GRU(content_output, name='Content_GRU',
                           consume_less='mem',
                           W_regularizer=WeightRegularizer(l2=gru_regularize)))
    content_vec = self.c_encoder(content_inp)
    # Merge vectors to create output
    doc_vec = merge(inputs=[title_vec, content_vec], mode='concat')
    self.decoder = Sequential(name='Decoder')
    self.decoder.add(Dense(dense_neurons[0],
                           input_shape=(title_output + content_output,),
                           name='Dense_0', activation='hard_sigmoid'))
    for i, n in enumerate(dense_neurons[1:]):
        self.decoder.add
            (Dense(n, activation='hard_sigmoid', name='Dense_%s' % (i + 1),
                    W_regularizer=WeightRegularizer(l2=dense_regularize)))

    self.decoder.add(Dense(len(classes), activation='softmax',
                           name='Dense_Output',
                           W_regularizer=WeightRegularizer(l2=dense_regularize)))
    output = self.decoder(doc_vec)
    self.model = Model(input=[title_inp, content_inp], output=output,
                       name='Model')
    if weights is not None:
        self.model.load_weights(weights)
```

Chức năng dự đoán chủ đề được cài đặt dưới dạng 1 phương thức của Classifier, trả về nhãn được dự đoán cùng với phân bố xác suất của từng chủ đề:

```
def predict(self, title, content, verbose=0):
    t = nltk.word_tokenize(title.lower())
    Xt = [self.word_to_index[word] if word in self.word_to_index
          else self.word_to_index[utils.UNKNOWN_TOKEN] for word in t]
    [:self.title_len]
    c = nltk.word_tokenize(content.lower())
    Xc = [self.word_to_index[word] if word in self.word_to_index
          else self.word_to_index[utils.UNKNOWN_TOKEN] for word in c]
    [:self.content_len]
    Xt = utils.pad_vec(Xt, self.title_len)
    Xc = utils.pad_vec(Xc, self.content_len)
    probs = self.model.predict([Xt, Xc], verbose=verbose)
    pred = np.argmax(probs)
    return pred, probs
```

Keras hỗ trợ việc định nghĩa các hàm callback (hàm được gọi tại các thời điểm nhất định trong quá trình huấn luyện). Nhóm cài đặt 2 callback để lưu mô hình và đánh giá trên tập test.

```
class SaveCallback(Callback):
    def __init__(self, classifier, prev_val_acc):
        self.best_loss = 1000.0
        self.best_val_loss = 1000.0
        self.best_val_acc = prev_val_acc
        self.classifier = classifier
        super(SaveCallback, self).__init__()
    def on_epoch_end(self, epoch, logs=None):
        print()
        if logs['val_loss'] <= self.best_val_loss and logs['loss'] <=
            self.best_loss and logs['val_acc'] >= self.best_val_acc:
            utils.save_classifier(self.classifier, self.classifier.directory)
            self.best_val_loss = logs['val_loss']
            self.best_loss = logs['loss']
            self.best_val_acc = logs['val_acc']
            self.classifier.log('Save point:\nLoss: %s\tAccuracy: %s\n' %
                                (logs['loss'], logs['acc']) +
                                'Validation loss: %s\tValidation accuracy: %s\n'
                                % (logs['val_loss'], logs['val_acc']), out=False)
            elif logs['val_loss'] > self.best_val_loss:
                print('No improvement on validation loss. Skipping save...')
            elif logs['loss'] > self.best_loss:
                print('No improvement on loss. Skipping save...')
            else:
                print('No improvement on validation accuracy. Skipping save...')

class TestCallback(Callback):
    def __init__(self, classifier, Xt, Xc, y):
        self.classifier = classifier
        self.Xt = Xt
```

```
self.Xc = Xc
self.y = y
super(TestCallback, self).__init__()
def on_train_end(self, logs=None):
    print('Evaluating on test set...')
    result = self.classifier.model.evaluate([self.Xt, self.Xc], self.y,
                                           batch_size=10)
    self.classifier.log('GRU_Regularizer: %s --- Dense Regularizer: %s' %
                       (self.classifier.gru_regularize,
                        self.classifier.dense_regularize))
self.classifier.log('Test loss: %s --- Test acc: %s\n' % (result[0],
                                                           result[1]))
```

3.4.3. Huấn luyện và chạy mô hình

Phương thức `compile()` và `train()` của `Classifier` cho phép biên dịch và huấn luyện mạng neuron chứa trong đối tượng.

Khái niệm biên dịch của Keras được thừa hưởng từ nền tảng tính toán (backend) Theano và TensorFlow. Thao tác này chỉ định backend xây dựng và tối ưu đồ thị tính toán của mạng, để từ đó có thể thực hiện tính năng tính gradient và tối ưu tự động. Biên dịch cũng quyết định thuật toán tối ưu và tốc độ học được sử dụng.

Huấn luyện nhận vào các ví dụ và nhãn được sử dụng (dưới dạng các mảng NumPy [8] [9]) cùng với các cài đặt huấn luyện để huấn luyện mạng. `compile()` buộc phải được thực hiện trước `train()`.

Script `train.py` được sử dụng để huấn luyện 1 mô hình mới dựa trên các cấu hình (được mô tả trong phần sau). Hệ thống có thể thực hiện huấn luyện với từng cặp tham số chính quy hóa (regularization) của lớp GRU và lớp neuron thường. Mô hình được lưu luôn là mô hình đạt kết quả tốt nhất trên tập validation.

Mô hình đã huấn luyện được nạp và sử dụng trong 2 demo của hệ thống. Việc dự đoán (phương thức `predict()`) không cần thực hiện `compile()`.

3.4.4. Cấu hình huấn luyện

Để tiện lợi trong việc điều chỉnh cấu hình huấn luyện, nhóm sử dụng 1 file cấu hình `settings.py` để xác định các tham số mạng.

```
LEARNING_RATE = 0.0003
N_EPOCH = 200
BATCH_SIZE = 100
GRU_LAMBDA = (0.0, 0.003, )
DENSE_LAMBDA = (0.0, 0.003, 0.01, )
VOCABULARY_SIZE = 70000
TITLE_LEN = 50
CONTENT_LEN = 500
TITLE_OUTPUT = 512
CONTENT_OUTPUT = 768
```

```
DENSE_NEURONS = (1024, )
DATASET = 'reuters'
```

3.4.5. Đánh giá mô hình

Để đánh giá chi tiết về độ chính xác của 1 mô hình đã huấn luyện, nhóm sử dụng phương thức `evaluate()` của `Classifier`, gọi qua `evaluate.py` để tính Precision, Recall và điểm F1 của từng nhãn trên tập dữ liệu test cũng như các giá trị trung bình. Ngoài ra, hệ thống cũng thực hiện đánh giá hàm mục tiêu và độ chính xác trong suốt quá trình huấn luyện.

3.4.6. Các hàm tiện ích

Nhóm xây dựng 1 số hàm tiện ích như tính precision, recall, điểm F1, vẽ đồ thị,... trong script `utils.py`. Việc lưu và nạp lại 1 mô hình đã huấn luyện được thực hiện trong 1 thư mục dành riêng cho mô hình đó. Các thông tin cần lưu được tổ chức như sau:

- Trọng số của mạng neuron được lưu trong `weights.hdf5`
- Các thông số cấu hình mạng (độ sâu, đầu ra,...) được lưu trong `config.pkl`
- Các giá trị từ điển (bao gồm vector từ) được lưu trong `dictionary.npz`
- Các ghi chép của hệ thống qua từng chu kỳ huấn luyện được lưu trong `logs.txt` và `epochs.csv`.
- Cấu hình hệ thống khi huấn luyện mô hình được lưu trong `settings.py`. Thay thế `settings.py` ở thư mục gốc bằng file này cho phép lặp lại quá trình huấn luyện của mô hình.
- Các đồ thị độ chính xác và hàm mục tiêu được lưu trong thư mục `plots`.

```
def save_classifier(classifier, directory):
    f1 = directory + '/weights.hdf5'
    f2 = directory + '/config.pkl'
    f3 = directory + '/dictionary.npz'
    config = {'title_output': classifier.title_output,
              'content_output': classifier.content_output,
              'dense_neurons': classifier.dense_neurons,
              'title_len': classifier.title_len,
              'content_len': classifier.content_len,
              'classes': classifier.classes,
              'word_vec_dim': np.shape(classifier.word_vec),
              'gru_reg': classifier.gru_regularize,
              'dense_reg': classifier.dense_regularize}
    classifier.model.save_weights(f1)
    pickle.dump(config, open(f2, 'wb'), pickle.HIGHEST_PROTOCOL)
    np.savez(f3, wit=classifier.word_to_index, itw=classifier.index_to_word,
            wv=classifier.word_vec)
    print('Saved model to %s.' % directory)
```


3.4.7. Các demo

Nhóm xây dựng 2 demo cho các mô hình được huấn luyện.

Demo thứ nhất cho phép nhập vào tiêu đề và nội dung của 1 văn bản tin tức bất kì, và trả về nhãn được dự đoán cũng như các xác suất tương ứng. Giao diện được xây dựng bằng framework [Kivy](#).

Demo thứ hai lấy về các văn bản tin tức từ feed RSS của 1 số trang báo và phân loại từng bài báo dựa trên tiêu đề và phần mô tả (được coi như nội dung). Kết quả sau đó được lưu vào results.csv.

Hướng dẫn thực thi các file huấn luyện và demo được nêu trong README.md.

News Classifier

Article Title
Eden Hazard: 'Antonio Conte is fantastic. It is a privilege to play with him'

Article Content
Chelsea have been "on fire" this season, Eden Hazard says, and what is fanning those flames is clear. They are determined to see this one through: win the Premier League, add the FA Cup and re-establish themselves. Get over the trauma of the last campaign. Make good, in fact. "We are Chelsea," Hazard explains. "We have to be at the top... talk on the pitch, walk on the pitch. And win games."

The ambition, the desire to succeed is burning – and not just because Chelsea have that human inferno, manager Antonio Conte, on the touchline driving them on, stoking the blaze. "We all have something to prove," Hazard says.

It is a sunny afternoon at Chelsea's training ground in Cobham and Hazard, the sunniest of players, has just been presented with an award for Premier League Goal of the Month – his superb solo effort against Arsenal in a 3-1 victory. The Belgian has already submitted an entry for this month's award with his devastating, counter-attack strike last Monday away at West Ham United.

There was also a wonderful, highlights-reel moment in that match when Hazard flicked the ball from an awkward height to N'Golo Kanté off his back. "It was something natural," he explains. "I did it once before, in France, it's for the fans." But Hazard says that with a grin. It was not, really, for the supporters. It was actually the best way for him to execute "a good pass". "Because if we had lost the ball and conceded a goal then I would have been in trouble with Antonio," Hazard adds.

And there is an insight as to how the Italian works. "Even if English is not his first language, he talks a lot, to give confidence to the players," Hazard says. "If you do something different from what he wants, and you do something bad, he will tell you off ... if you do something different, do it well!"

Do it well. Ahead of Chelsea's FA Cup quarter-final at home against Manchester United on Monday night, Hazard's lively conversation ranges from last season's problems – which he immediately raises – to being pain-free and on to his father's fascination with Conte, as well as why he longer worries about winning the Ballon d'Or. He also discusses how he would give up football for his family and, even, his dream of one day playing for the Belgium national side with his three younger brothers.

To begin, Hazard recalls two moments. The first was at Leicester City in December 2015. He walked off the pitch, much to Jose Mourinho's evident disgust, and was unable to carry on. His hip was causing him too much discomfort. "I started last season with pain," Hazard says. "I remember the game at Leicester. It was so painful." Nevertheless, the forward was accused of not trying.

"When you play a bad game, we lost the game and the manager [Mourinho] was sacked, so everyone was talking in a bad way," Hazard says. "But you have to deal with it. I deal with it this season, so I am happy."

It helped that the 26-year-old was eventually allowed to recuperate, that he had a good Euros with his country and that he came back free of discomfort. "We started the season well and, when the whole team is on fire, they give a lot of confidence," Hazard says.

Model: bbc_2017-03-12 **Load**

LEARNING_RATE = 0.0003
N_EPOCH = 50
BATCH_SIZE = 15
GRU_LAMBDA = 0.002
INFNSF | **AMRNA** = 0 001

Prediction result: Sport

Prediction breakdown

Category	Confidence
Sport	0.9
Politics	0.1
Business	0.05
entertainment	0.02
Tech	0.0

Clear **Classify**

```
phan ngoc lan@will-Lenovo-250-70:~/Dropbox/Class-Material/Machine-Learning/Project$ python3 -m demo_rss models/bbc_final
Using Theano backend.
Using gpu device 0: GeForce 820M (CUMEM is disabled, cuDNN not available)
Loading model from models/bbc_final...
Done.
Parsing http://rss.nytimes.com/services/xml/rss/nyt/HomePage.xml...
http://www.nytimes.com/2017/04/08/us/politics/trump-doctrine-foreign-policy.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/07/us/politics/donald-trump-syria-twitter.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/world/europe/us-attack-on-syria-cements-kremlins-embrace-of-assad.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/world/middleeast/us-strike-on-syria-brings-fleeting-hope-to-those-caught-in-brutal-conflict.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/07/us/syria-refugees-trump.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/08/world/americas/trump-nikki-haley-united-nations.html?partner=rss&emc=rss Entertainment
http://www.nytimes.com/2017/04/08/world/asia/china-xi-jinping-president-trump-xinhua.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/08/us/politics/us-islamic-state-syria.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/07/us/white-house-kushner-bannon-military-strike.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/06/us/politics/stephen-bannon-white-house.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/business/china-trade-solar-panels.html?partner=rss&emc=rss Tech
http://www.nytimes.com/video/world/europe/10000005033063/everything-indicates-terror-attack-in-stockholm.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/08/world/europe/stockholm-truck-attack-arrest.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/world/asia/myanmar-aung-san-su-kyi-first-year.html?partner=rss&emc=rss Tech
http://www.nytimes.com/video/style/10000005024359/salone-de-mobile-milan-furniture-fair-2017-360.html?partner=rss&emc=rss Tech
http://www.nytimes.com/aponline/2017/04/08/us/politics/ap-us-united-states-korea.html?partner=rss&emc=rss Business
http://www.nytimes.com/aponline/2017/04/08/world/europe/ap-eu-norway-unexploded-device.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/08/us/alabama-governor-impeachment-hearings.html?partner=rss&emc=rss Politics
http://www.nytimes.com/2017/04/08/world/asia/indian-man-accused-in-multi-million-dollar-call-center-swindle-is-held.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/world/europe/spain-basque-separatist-group-eta-disarmament.html?partner=rss&emc=rss Business
http://www.nytimes.com/crosswords/game/mini?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/07/smarter-living/what-to-cook-watch-listen-to-and-more-this-weekend.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/01/theater/five-must-see-shows-if-youre-in-new-york-this-month.html?partner=rss&emc=rss Entertainment
http://www.nytimes.com/2017/04/07/us/circling-brothers-circus-set-to-go-dark-after-146-year-run.html?partner=rss&emc=rss Sport
http://www.nytimes.com/2017/04/08/sports/golf/masters-third-round.html?partner=rss&emc=rss Sport
http://www.nytimes.com/2017/04/08/briefing/17-great-stories-that-have-nothing-to-do-with-politics.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/08/business/dealbook/george-soros-dawn-fitzpatrick-american-stock-exchange.html?partner=rss&emc=rss Business
http://www.nytimes.com/2017/04/08/style/cafe-milano-donald-trump-washington.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/06/education/edlife/ap-perfect-imperfect-gap-year.html?partner=rss&emc=rss Sport
http://www.nytimes.com/2017/04/06/education/edlife/a-college-application-guide-for-gap-year-students.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/06/education/edlife/readers-tell-us-is-a-gap-year-worth-it.html?partner=rss&emc=rss Tech
http://www.nytimes.com/2017/04/07/nyregion/kosher-passover-welchs-manischewitz.html?partner=rss&emc=rss Entertainment
```

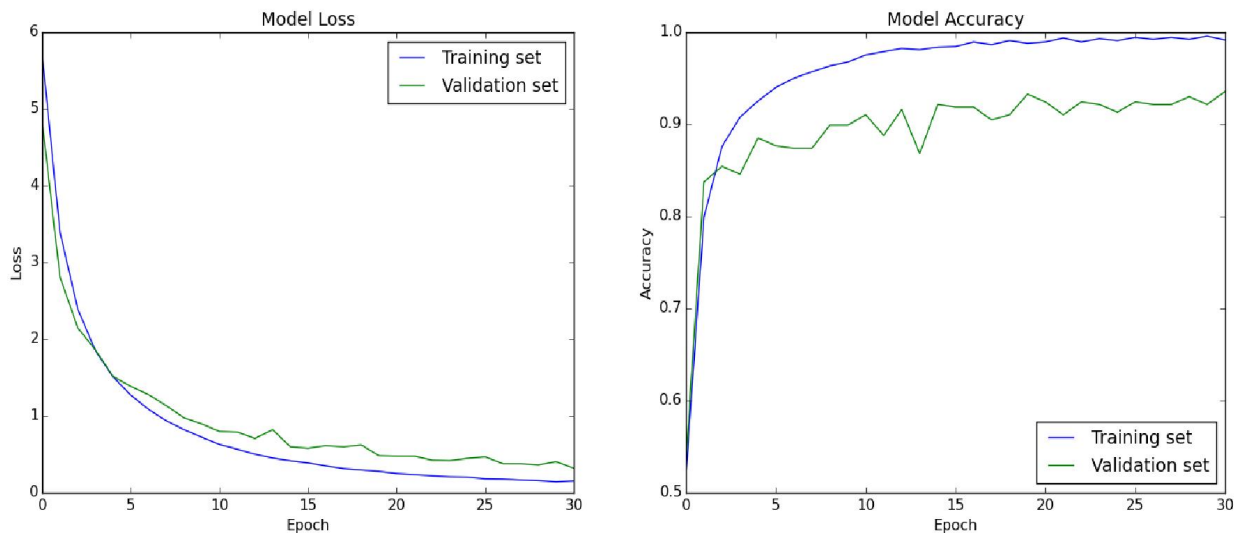
4. Kết quả

Kết quả trên từng tập dữ liệu được nêu trong bảng sau*:

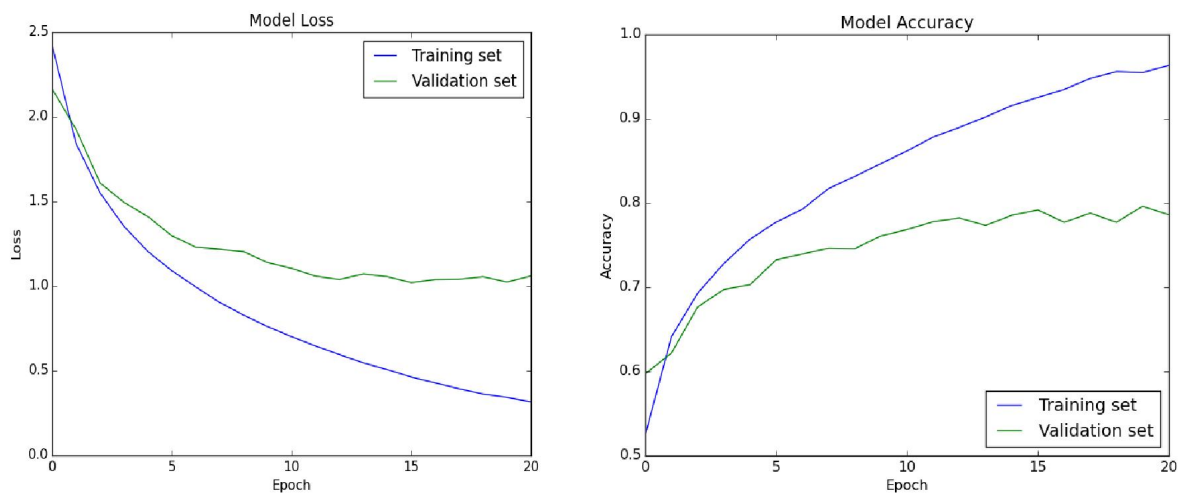
	AG-News	BBC	Reuters
Số lượng nhãn	4	5	86
Hàm mục tiêu (train)	0.2162	0.1530	0.4631
Hàm mục tiêu (val)	0.2378	0.3196	1.0209
Hàm mục tiêu (test)	0.2712	0.2868	1.0247
Độ chính xác (train)	0.9287	0.9915	0.9252
Độ chính xác (val)	0.9207	0.9355	0.7916
Độ chính xác (test)	0.9177	0.9485	0.7919
Precision (test)	0.9222	0.9488	0.5512
Recall (test)	0.9213	0.9464	0.2641
Điểm F1 (test)	0.9214	0.9473	0.4917

* kết quả ứng với mô hình có độ chính xác validation cao nhất

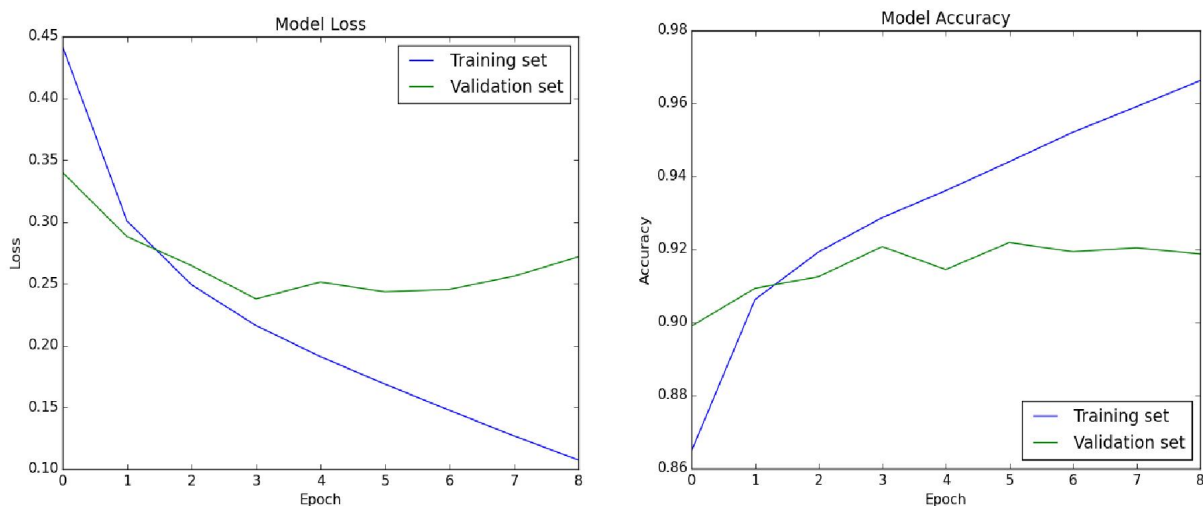
Đồ thị hàm mục tiêu và độ chính xác tại lần huấn luyện tối ưu:



Hình 4: Dataset BBC (0.002, 0.003)



Hình 5: Dataset Reuters (0.000, 0.003)



Hình 6: Dataset AG-News (0.00, 0.01)

Thời gian huấn luyện của hệ thống phụ thuộc vào kích cỡ của mạng được thiết lập trong settings.py và khả năng phần cứng của máy tính. Với cấu hình được nhóm lựa chọn chạy trên máy tính cá nhân có hỗ trợ GPU, một lần huấn luyện có thời gian chạy từ 3-4 giờ. Nhóm tối ưu qua 6 lần huấn luyện với các lambda khác nhau, do đó mỗi mô hình được huấn luyện trong khoảng 24 giờ.

Thời gian đoán nhận sau khi huấn luyện của mô hình là khá nhanh, thường trong khoảng 1-2s. Do đặc thù của backend Theano, lần đoán nhận đầu tiên có thể chậm hơn

và có thể đến ~10s.

Nhóm cũng nhận thấy chênh lệch đáng kể giữa thời gian chạy khi có và không sử dụng hỗ trợ GPU. Sử dụng GPU cho phép hệ thống học nhanh gấp 3-4 lần, tuy nhiên không có chênh lệch đáng kể trong thời gian đoán nhận.

5. Kết luận và hướng phát triển

5.1. Khó khăn, tranh luận và hướng giải quyết

5.2. Kết luận

Từ kết quả thu được, nhóm rút ra một số kết luận như sau:

- Mô hình hoạt động rất tốt với các văn bản ngắn hoặc trung bình, trong đó GRU dễ dàng học được sự liên kết giữa các từ từ đầu văn bản đến cuối văn bản.
- Mô hình thể hiện rõ khả năng chịu nhiễu rất tốt của mạng neuron. Với số lượng từ UNKNOWN_TOKEN ở ngưỡng 2-5%, mô hình rất ít bị ảnh hưởng và không có dấu hiệu overfit. Tuy nhiên, với tỉ lệ nhãn chênh lệch trên tập Reuters, mô hình vẫn rơi vào tình trạng overfit.
- Hệ thống có thể dễ dàng học các phân loại trên nhiều loại văn bản không phải tin tức hoặc chỉ có 1 input (bằng cách đặt TITLE_OUTPUT bằng 1 và chỉ truyền vào tiêu đề là xâu rỗng, hệ thống sẽ học cách bỏ qua hoàn toàn thông tin về tiêu đề).
- Thời gian huấn luyện, yêu cầu phần cứng và số lượng dependency lớn là nhược điểm không thể tránh khỏi của hệ thống. Điều này làm ứng dụng của hệ thống bị hạn chế ở phía server, khó triển khai ở các máy tính cá nhân hay thiết bị di động.

5.3. Hướng phát triển

Nhóm có dự định phát triển hệ thống theo một số hướng như sau:

- Kết hợp mô hình với một thuật toán phân loại khác (có thể là SVM) và lựa chọn đáp án dựa trên mức tự tin của mỗi dự đoán.
- Xây dựng cơ chế để một mô hình đã huấn luyện có thể tiếp tục học online để cải thiện kết quả. Chủ yếu liên quan đến việc quyết định mức độ học với mỗi ví dụ học mới.
- Áp dụng một số kĩ thuật giảm overfit mạnh hơn trong Deep Learning như Dropout, Pooling,...
- Thử nghiệm mô hình với các tập dữ liệu ở các domain tương tự như phân loại tweet trên Twitter, phân loại lời bài hát,...

Tài liệu tham khảo

- 1: Zach Chase, Nicolas Genain, Orren Karniol-Tambour, Learning Multi-Label Topic Classification of News Articles,
- 2: Thorsten Joakims, Text Categorization with Support Vector Machines: Learning with Many Relevant Features,
- 2: Dennis Ramdass, Shreyes Seshasai, Document Classification for Newspaper Articles, 2009
- 3: Francois Chollet, Keras Documentation, , <https://keras.io/>
- 4: Denny Britz, Recurrent Neural Networks Tutorial, Part 1 – Introduction to RNNs – WildML, , <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>
- 5: Jeffrey Pennington; Richard Socher; Christopher D. Manning, GloVe: Global Vectors for Word Representation, 2014
- 6: NLTK Project, NLTK 3.0 documentation, , <http://www.nltk.org/>
- 7: Sebastian Ruder, An overview of gradient descent optimization algorithm, 2016, <http://sebastianruder.com/optimizing-gradient-descent/index.html#rmsprop>
- 8: Eric Jones, Travis Olyphant, Pearu Peterson, et al , SciPy: Open Source Scientific Tools for Python, 2001, <http://www.scipy.org>
- 9: Stefan van der Walt, S. Chris Colbert, Gael Varoquaux, The NumPy Array: A Structure for Efficient Numerical Computation, 2011