## The Illustrated BERT, ELMo, and co. (How **NLP Cracked Transfer Learning)**

By: Quynh Nhi Tran

Week 8 - Personal Note

Learning Năm 2018 là thời điểm quan trọng đối với NLP. Đây là lúc mà cộng đồng NLP đã có một bước tiến lớn, đặc biệt trong việc tìm ra cách biểu diễn từ ngữ và câu sao cho

Điểm bùng phát trong NLP năm 2018 và vai trò của Transfer

- nắm bắt được ý nghĩa và mối quan hệ ngữ nghĩa. • Khái niệm **Transfer Learning** đã bùng nổ trong NLP, tạo ra một tình huống được ví
- như "NLP's ImageNet moment" (tương tự như sự phát triển của Machine Learning trong Computer Vision nhờ bộ dữ liệu ImageNet). • Transfer Learning giúp các mô hình có thể học từ các mô hình lớn đã được huấn
- luyện trên tập dữ liệu khổng lồ và sau đó tùy chỉnh (fine-tuning) để ứng dụng vào các bài toán cụ thể mà không phải huấn luyện từ đầu.
- Tác động của Học chuyển giao
  - Học chuyển giao Học truyền thống Cần huấn luyện nhiều Cách mạng hóa việc huấn luyện mô hình

**BERT:** Một cột mốc trong NLP • BERT (Bidirectional Encoder Representations from Transformers)

• Google đã cung cấp mã nguồn cùng với các phiên bản model đã được huấn luyện

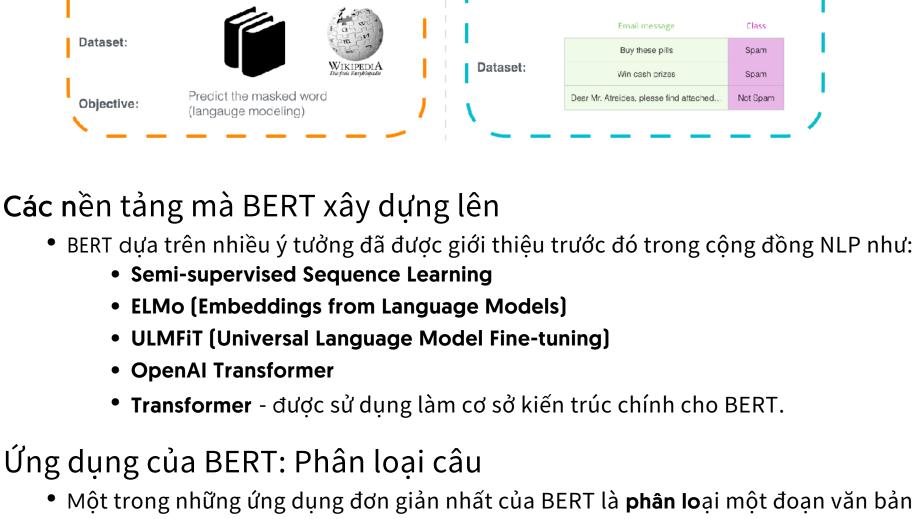
• Với BERT, ai cũng có thể tải xuống và sử dụng mô hình đã huấn luyện trên các tập

dữ liêu khổng lồ, tiết kiệm được nhiều thời gian và tài nguyên so với việc phải tư

## of text (books, wikipedia..etc). The model is trained on a certain task that enables it to grasp

huấn luyện mô hình từ đầu.

- 2 Supervised training on a specific task with a 1 - Semi-supervised training on large amounts labeled dataset. Supervised Learning Step
- patterns in language. By the end of the training process BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way 75% Classifier 25% Not Spam Semi-supervised Learning Step Model: Model: (pre-trained **BERT** in step #1)



Input **Features** 

. Ví dụ, mô hình có thể được huấn luyện để phân loại email là "spam" hay "not Output Prediction

• Quá trình huấn luyện mô hình này chủ yếu tập trung vào phần classifier (phân loại)

• Đây là một phương pháp thuộc supervised learning, đòi hỏi một labeled dataset (t

Classifier

15% Not Spam

Class

Not Spam

- (Feed-forward Help Prince Mayuko Transfer BERT Huge Inheritance softmax)

ở phần đuôi, chỉ tinh chỉnh một chút với BERT (Fine-Tuning).

ập dữ liệu có gắn nhãn) để mô hình học hỏi và đưa ra dự đoán.

Email message

Dear Mr. Atreides, please find attached...

• Sentiment Analysis: Phân tích đánh giá (tích cực hoặc tiêu cực).

Một số ví dụ khác về ứng dụng BERT trong NLP

BERTBASE

**ENCODER** 

**ENCODER** 

**ENCODER** 

• Ref: Transformer model (foundational concept for BERT)

https://jalammar.github.io/illustrated-transformer/

heads hơn so với Transformer ban đầu.

Buy these pills Spam Win cash prizes Spam



• BERT BASE: tương đương về kích thước với mô hình Transformer của OpenAI.

• BERT LARGE: lớn hơn rất nhiều và đạt kết quả tốt nhất vào thời điểm đó.

BERTLARGE

**ENCODER** 

**ENCODER** 

**ENCODER** 

**ENCODER** 

**ENCODER** 

BERTLARGE

BERTBASE • BERT là một Transformer Encoder Stack với 12 lớp Encoder trong phiên bản Base và

Input và Output của BERT

ở output.

[CLS] Help Prince Mayuko

12

Đầu vào và

được sử dụng làm input cho một classifier.

**ENCODER** 

**ENCODER** 

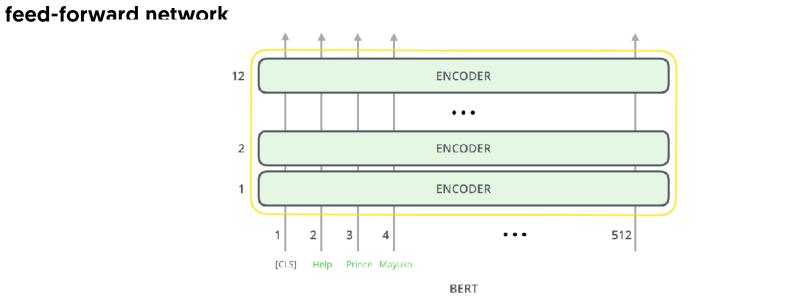
**ENCODER** 

Đầu ra

BERT có 2 kích thước mô hình:

→ BERT BASE → 12 lớp mã hóa Kiến trúc mô hình → BERT LARGE → 24 lớp mã hóa

24 lớp trong phiên bản Large, mỗi lớp có nhiều đơn vị mạng nơron và attention



• Mỗi vị trí trong chuỗi đầu vào trả về một vector có kích thước 768 (với BERT Base)

Đầu ra:

• Đối với bài toán phân loại câu, vector tại vị trí đầu tiên (tương ứng với token [CLS])

Kích thước

vector 768

85% Spam

Classifier

15% Not Spam

rd neural network + softmax)

512

Output

Prediction

-0.30 0.03

0.09

0.35

-0.28 -0

• Input của BERT: token đặc biệt [CLS] cho mục đích phân loại. ([CLS]: Classification)

• Cấu trúc hoạt động tương tự như Transformer Encoder: mô hình nhận sequence of

words làm input, và qua mỗi lớp oder, nó áp dụng self-attention và đưa qua

BERT 3 2 Help Parallels with Convolutional Nets

Khi làm việc với BERT, cấu trúc truyền tải thông tin giữa các thành phần của mạng giống

fully-connected để phân loại, BERT chuyển một vector đầu ra của mô hình Transformer

với cách mà mạng nơ-ron tích chập (CNN) hoạt động trong Computer Vision. Tương tự

512

như cách CNN truyền tải dữ liệu từ phần tích chập (convolution) sang phần

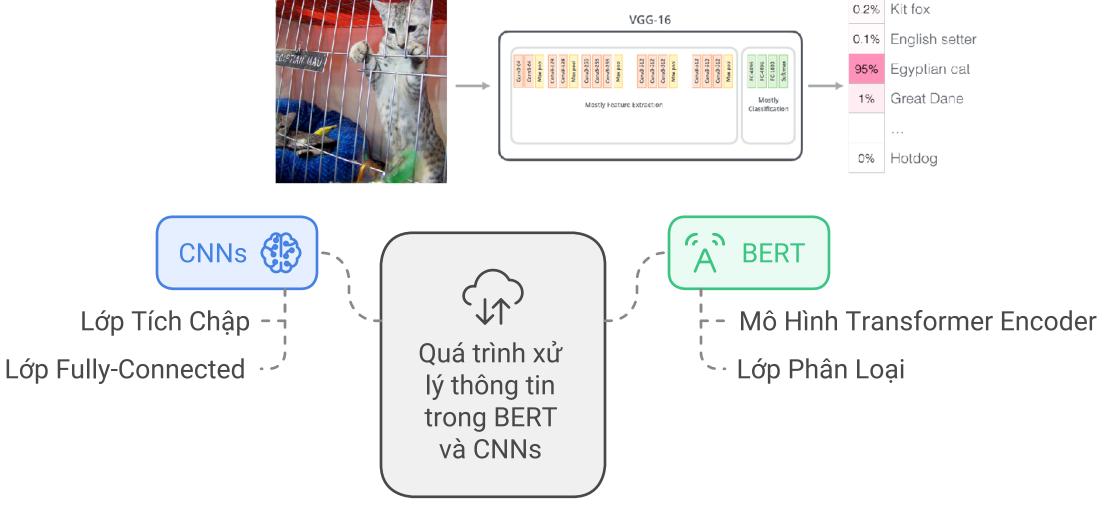
Encoder tới một lớp classifier để thực hiện các tác vụ xử lý ngôn ngữ.

Input

**Features** 

Kỷ Nguyên Mới của Word Embedding

văn bản, giúp giảm công sức và thời gian huấn luyện từ đầu.



Sự phát triển mới trong NLP đã thay đổi cách mã hóa từ vựng. Trước đây, các phương

pháp như Word2Vec và GloVe đã thống trị việc xử lý từ ngữ bằng cách tạo ra các vector

để biểu diễn từ ngữ theo các mối quan hệ ngữ nghĩa và ngữ pháp. Những mô hình này

cho phép sử dụng các embedding đã được tiền huấn luyện trên một lượng lớn dữ liệu

Ví dụ, GloVe embedding của từ "stick" là một vector bao gồm 200 giá trị, giúp biểu diễn

ý nghĩa từ vựng của "stick" nhưng vẫn cố định trong mọi ngữ cảnh. Tuy nhiên, vấn đề ở

0.06

The GloVe word embedding of the word "stick" - a vector of 200 floats (rounded to two decimals). It goes on for two hundred values.

0.30

0.33

-1.17

Nhạy cảm với

ngữ cảnh

đây là khi "stick" được sử dụng trong các ngữ cảnh khác nhau, ý nghĩa của nó có thể

المناع من من المناع الم

-0.26 | -0.12 | 0.23 | 1.04 | -0.16 | 0.31

Không phụ thuộc ngữ cảnh

ELMo: Ngữ Cảnh Là Quan Trọng

**Word Embedding Recap** 

-0.34 | -0.84 | 0.20

Vector từ cố Vector từ động định Embedding Truyên Embedding Ngữ Thống Cảnh So sánh các embedding từ truyền thống và

ngữ cảnh.

Để giải quyết vấn đề trên, ELMo (Embeddings from Language Models) được giới thiệu

câu cụ thể. Mô hình ELMo sử dụng mạng LSTM hai chiều (bi-directional LSTM) để gán

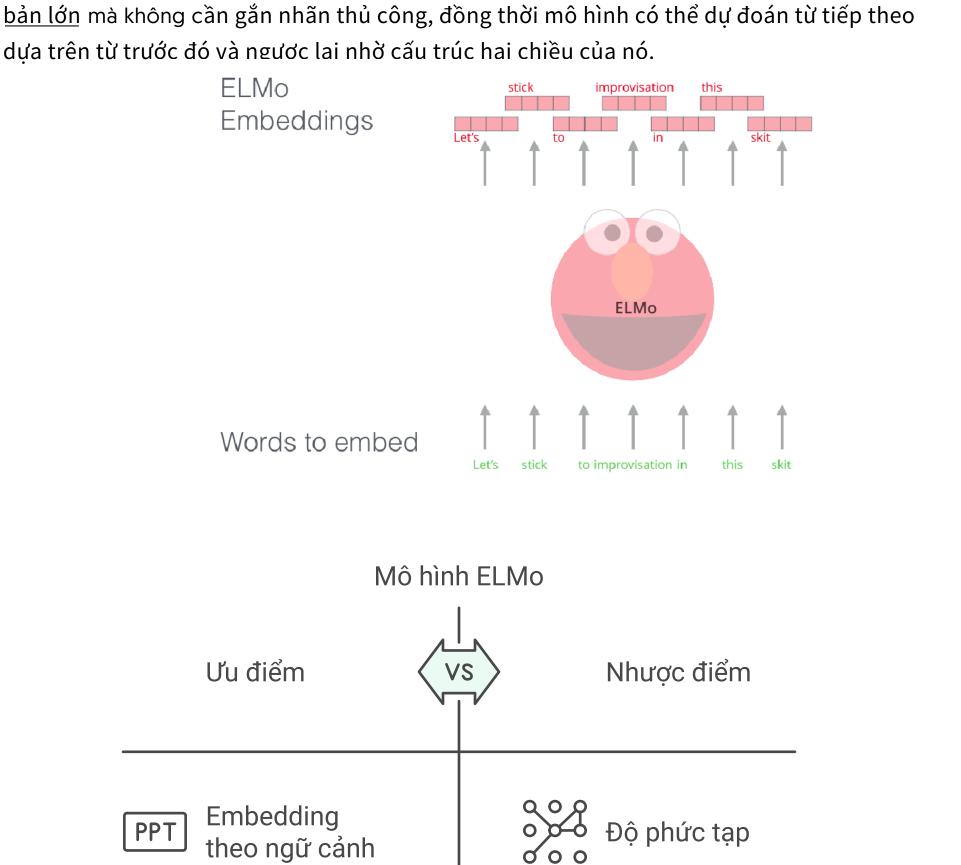
Hey ELMo, what's the embedding of the word "stick"?

There are multiple possible embeddings! Use it in a sentence.

nhằm tạo embedding dựa trên ngữ cảnh, từ đó nắm bắt ý nghĩa của từ ngữ trong từng

embedding cho từng từ dựa vào toàn bộ câu chứa nó. Điều này giúp ELMo hiểu được ngữ

cảnh tổng quát của từ, đồng thời cải thiện độ chính xác trong các tác vụ xử lý ngôn ngữ.





BERT: Từ Decoder sang Encoder với Bidirectional Context Mặc dù OpenAI Transformer đã thành công với mô hình unidirectional (chỉ nhìn về phía trước), nhưng vẫn còn một khoảng trống. BERT đã giải quyết vấn đề này bằng cách sử dụng bidirectional conditioning để mô hình nhìn được cả ngữ cảnh phía trước và phía sau của một từ. Thay vì dùng decoder, BERT dùng các lớp encoder được xếp chồng và giới thiệu phương pháp Masked Language Model (MLM). Với MLM, BERT chọn một số từ trong câu và mask chúng để mô hình dự đoán từ đó, nhờ vậy mô hình tận dụng ngữ cảnh hai BER LUTA MAN TOUR BER LE LINGTH THOU THOU TO BER LEVEL sơ bộ. Nhiệm vụ này giúp mô hình học cách xác định liệu một câu có khả năng là câu tiếp theo của một câu khác không, giúp BERT xử lý tốt hơn các tác vụ yêu cầu nhận biết mối quan hệ giữa các câu, như kiểm tra tính tương đồng và trả lời câu hỏi. Mô hình Chuyên Biệt và Tính năng Feature Extraction với BERT Ngoài fine-tuning, BERT còn có thể dùng để tạo contextualized embeddings cho các từ trong văn bản. Bằng cách tận dụng embeddings từ BERT đã được huấn luyện sẵn, các mô hình không cần fine-tuning cho từng tác vụ vẫn có thể cải thiện đáng kể hiệu quả,

embeddings cho các tác vụ khác nhau trong NLP. Việc sử dụng BERT cho feature extraction, cùng với fine-tuning, là cách tiếp cận mạnh

mẽ cho nhiều ứng dụng NLP, từ nhận diện tên thực thể đến đo lường tính tương đồng

Để thử nghiệm BERT, bạn có thể bắt đầu với notebook **BERT FineTuning v**ới Cloud TPUs

trên Google Colab, hỗ trợ TPUs, CPUs và GPUs. Các tài nguyên chính bao gồm:

đạt gần đến hiệu quả của các mô hình fine-tuning đầy đủ. BERT cung cấp nhiều lựa chọn

Thử nghiệm với BERT

ừ nguyên bản.

Start

Start

Start

Start

Start

Similarity

Multiple Choice

Text 1

Text 2

Context

Context

Context

Delim

Delim

Delim

Delim

Delim

Text 2

Text 1

Answer 1

Answer 2

Answer N

Extract

Extract

Extract

Extract

Extract

• modeling.py (class BertModel) chứa cấu trúc chính của BERT. • run\_classifier.py là ví dụ về quá trình fine-tuning và tạo lớp phân loại. • tokenization.py là tokenizer của BERT, chia nhỏ từ thành WordPieces thay vì dùng t Hiện có nhiều mô hình BERT được huấn luyện sẵn, hỗ trợ nhiều ngôn ngữ khác nhau từ

Transformer

Transformer

Transformer

Transformer

Transformer

+)→ Linear

Linear

Linear

tiếng Anh đến tiếng Trung, cũng như một mô hình đa ngôn ngữ được huấn luyện trên Wikipedia với 102 ngôn ngữ.

Source: https://jalammar.github.io/illustrated-bert/?authuser=0

Độ chính xác Tốn tài nguyên được cải thiện tính toán Không cần gán Cần dữ liệu lớn IS0 nhãn thủ công Khả năng giải √ Học hai chiều thích hạn chế Đào tạo trên văn bản lớn **ULMFiT: Transfer Learning cho NLP** Trước đây, các mô hình NLP gặp khó khăn khi muốn tận dụng transfer learning giống như trong lĩnh vực thị giác máy tính (Computer Vision). ULMFiT đã cải thiện vấn đề này bằng cách cung cấp một mô hình ngôn ngữ (language model) và một quy trình tinh chỉnh (fine-tuning) để tận dụng tối đa những gì mô hình học được từ giai đoạn tiền huấn luyện. Transformer: Vượt Ra Khỏi LSTM Sự ra đời của mô hình Transformer, với cấu trúc Encoder-Decoder, đã tạo ra bước ngoặt trong NLP, đặc biệt trong các tác vụ yêu cầu phụ thuộc dài hạn như dịch máy. Transformer giúp mô hình hóa các phụ thuộc dài hạn tốt hơn LSTM và phù hợp cho các tác vụ như dịch máy. Tuy nhiên, để sử dụng nó trong phân loại câu hoặc các tác vụ khác, việc tinh chỉnh thêm là cần thiết. OpenAl Transformer: Huấn luyện sơ bộ với Decoder Layers Trong lĩnh vực NLP, để thực hiện transfer learning, chúng ta không cần toàn bộ Transformer mà chỉ cần phần decoder. Decoder là lựa chọn tự nhiên cho language

Oh, okay. Here: "Let's stick to improvisation in this Oh in that case, the embedding is: -0.02, -0.16, 0.12, -0.1 ....etc Điểm đặc biệt của ELMo là được huấn luyện thông qua nhiệm vụ Language Modeling, tức là <u>dự đoán từ tiếp theo trong một chuỗi từ.</u> Việc này cho phép ELMo <u>học từ các mẫu văn</u>