

Transformers (2)

2022/03/09 CS 182/282A Lecture 13

Transformers review

Setup

features \mathbf{x}

It was the best of times, it was the worst of times, it was the age





- sequential data
- may be variable length

label y

could correspond to...

. . .

. . .

. . .

- sentiment analysis, translation to another language,
- audio transcription, speaker identification,
- activity identification, video captioning,

or there could be no label!

- unsupervised learning / generative modeling

-???-

model

transformers

Self-attention: the building block of transformers

The goal of self-attention is to handle sequential features as the input

Think of it as a neural network layer that allows for processing the whole sequence

"key-value-query" system:

$$\mathbf{q}_t = q(\mathbf{x}_t)$$

$$\mathbf{k}_t = k(\mathbf{x}_t)$$
 or \mathbf{h}_t

$$\mathbf{v}_t = v(\mathbf{x}_t)$$

These functions are learned and can be, e.g., simple linear layers

$$\mathbf{x}_{1} \xrightarrow{\mathbf{q}_{2}} e_{1,2} = \mathbf{k}_{1}^{\top} \mathbf{q}_{2} \xrightarrow{\mathbf{r}} \alpha_{1,2} \xrightarrow{\mathbf{r}} \mathbf{a}_{2} = \sum_{t} \alpha_{t,2} \mathbf{v}_{t}$$

$$\mathbf{k}_{T} \xrightarrow{\mathbf{r}} e_{T,2} = \mathbf{k}_{T}^{\top} \mathbf{q}_{2} \xrightarrow{\mathbf{r}} \alpha_{T,2} \xrightarrow{\mathbf{r}} \mathbf{a}_{2} = \sum_{t} \alpha_{t,2} \mathbf{v}_{t}$$

$$\mathbf{x}_{1} \xrightarrow{\mathbf{r}} \mathbf{q}_{2} \xrightarrow{\mathbf{r}} \alpha_{T,2} \xrightarrow{\mathbf{r}} \mathbf{a}_{T} \xrightarrow{\mathbf{r}} \alpha_{t,2} \mathbf{v}_{t}$$

Transformers for "encoding"



Transformers for "decoding" (generation)



A note about training transformer decoders

At training time, we pass in entire sequences to the decoder



This is why the self-attention in the decoder must be masked

In this way, the decoder is simultaneously trained on T next token (word) predictions

Transformers in action

The original transformer

The original transformer is a sequence-to-sequence model for translation



Sequence-to-sequence models typically follow an *encoder-decoder* architecture, and the transformer does as well

Next week, we will talk in greater detail about sequence-to-sequence models

Vaswani et al, "Attention is All You Need". NIPS 2017.

The original transformer

The original transformer follows an encoder-decoder architecture



Vaswani et al, "Attention is All You Need". NIPS 2017.

Just the encoder: BERT



The GLUE benchmark

GLUE Tacks

	GLUL HASKS					
Name	Download	More Info	0 * Books were sent to eac	0 * Books were sent to each other by the students.		
The Corpus of Linguistic Acceptability	<u>.</u>	<u>s</u>	Matthew's C 1 She voted for herself.	Matthew's C 1 She voted for herself.		
The Stanford Sentiment Treebank	*		Accuracy 1 I saw that gas can exp	1 I saw that gas can explode.		
Microsoft Research Paraphrase Corpus	*	2	F1 / Accuracy			
Semantic Textual Similarity Benchmark	*		P Premise	Label	Hypothesis	
Quora Question Pairs	*		F Fiction		Orldoor local the Old One service	
MultiNLI Matched	*		A The Old One always comforted Ca'daan, except today.	neutral	well.	
MultiNLI Mismatched			A Your gift is appreciated by each and every student who will benefit from your generosity	neutral	Hundreds of students will	
Question NLI	*		A Telephone Speech		benent norriyota generooky.	
Recognizing Textual Entailment	*	2	yes now you know if if everybody like in August when A everybody's on vacation or something we can dress a little more casual or	contradiction	August is a black out month for vacations in the company.	
Winograd NLI	*		A 9/11 Report			
Diagnostics Main	*	2	At the other end of Pennsylvania Avenue, people began to M line up for a White House tour.	entailment	People formed a line at the end of Pennsylvania Avenue.	

"Fine tuning" BERT



Up to billions of sentences

Mask / next sentence prediction

MNLI NER SQUAD Start/End Span С Τ, T. Tiern BERT E. (CLS) Tok 1 ISEP Tok Question Paragraph Question Answer Pair Fine-Tuning Much less data

Specific to the task at hand



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG





(b) Single Sentence Classification Tasks: SST-2, CoLA



(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". NAACL 2019.

Just the encoder: vision transformers (ViTs)



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". ICLR 2021.

An even more recent ViT

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu^{†*} Yutong Lin^{†*} Yue Cao^{*} Han Hu^{*‡} Yixuan Wei[†] Zheng Zhang Stephen Lin Baining Guo Microsoft Research Asia



Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows". ICCV 2021.

Unsupervised training of ViTs

Emerging Properties in Self-Supervised Vision Transformers

Mathilde Caron^{1,2}Hugo Touvron^{1,3}Ishan Misra¹Hervé Jegou¹Julien Mairal²Piotr Bojanowski¹Armand Joulin¹¹ Facebook AI Research² Inria*³ Sorbonne University





Caron et al, "Emerging Properties in Self-Supervised Vision Transformers". ICCV 2021.

Unsupervised training of ViTs

Masked Autoencoders Are Scalable Vision Learners

Kaiming He^{*,†} Xinlei Chen^{*} Saining Xie Yanghao Li Piotr Dollár Ross Girshick ^{*}equal technical contribution [†]project lead Facebook AI Research (FAIR)





Just the decoder: reinforcement learning





Chen et al, "Decision Transformer: Reinforcement Learning via Sequence Modeling". NeurIPS 2021. Janner et al, "Reinforcement Learning as One Big Sequence Modeling Problem". NeurIPS 2021.

Just the decoder: GPT

Now we're getting into surreal territory. GPT-3 knows how to have a know how to say "Wait a moment... your question is nonsense." It a know."

Q: How do you sporgle a morgle?

A: You sporgle a morgle by using a sporgle.

Q: How many bonks are in a quoit?

A: There are three bonks in a quoit.

Q: How many rainbows does it take to jump from Hawaii to A: It takes two rainbows to jump from Hawaii to sevente

Q: Which colorless green ideas sleep furiously? A: Ideas that are colorless, green, and sleep furiously furiously.

Q: Do you understand these questions?

A: I understand these questions.

This is a conversation between a human and a brilliant AI. If a question is "normal" the AI answers it. If the question is "nonsense" the AI says "yo be real"

Q: What is human life expectancy in the United States? A: Human life expectancy in the United States is 78 years.

Q: How do you sporkle a morgle? A: yo be real

Q: Who was president of the United States before George W. Bush? A: Bill Clinton was president of the United States before George W. Bush.

Q: How many rainbows does it take to jump from Hawaii to seventeen? A: yo be real

Q: How does an umbrella work

A: An umbrella works by using a series of spokes to keep the rain from falling on you.

Q: How many bonks are in a quoit? A: yo be real

Q: Which colorless green ideas speak furiously

A: yo be real

Summary

- Transformers are the latest and greatest model for sequential data
- The fundamental building block is *self-attention*, which allows for processing the whole sequence all at once
 - Compared to previous sequence models, long range dependencies are captured much better with self-attention
- Even data that is not obviously sequential has benefited from transformers
- What other great architecture ideas are still out there?

