Machine Learning I 60629A

Attention and Transformers — Week #10



Transformers

- A deep learning model
 - Introduced in 2017 (Google researchers)
 - Uses the attention mechanism

* Most of the slides/figures/narrative are from David Berger (you might see them again in ML #2).

 Quickly adopted for modelling sequential data (text and images)—architecture behind LLMs

Today's plan

- ones
- Introduce attention
- Transformer block
- Examples of transformers in practice

Refresh our understanding of RNNs and bidirectional

- For concreteness, I will think of our data as a sequence of words
- And I use words and tokens synonymously
- - Tokenization is a topic beyond today's class

In practice, tokens are sub-words (e.g., a few letters)





RNNS

The hidden state at each timestep (h_t) Must contain all useful information up to time t.



Bidirectional RNNs







- The hidden state h_t must contain all useful information from to t
- The hidden state g_t must contain all useful information from the end (T) to t.
- Difficulties:
 - Exploding & vanishing gradients
 - Predictions will tend to use information from close neighbours



Encoder-Decoder (Also known as Sequence to Sequence — Seq2Seq)

 What if the output is a different length than the input? E.g., translation



- - hyperparameter

• h_4 must contain ("summarize") information from the input

• It's a bottleneck. Its size (num. neurons) is an important



three long narrative poems and a few other verses, some of uncertain authorship.



• Challenging to encode vast amounts of information

three long narrative poems and a few other verses, some of uncertain authorship.



His extant works, including collaborations, consist of some 39 plays, 154 sonnets, three long narrative poems and a few other verses, some of uncertain authorship.

- Challenging to encode vast amounts of information
- Could try to divide your input (e.g. sentence by sentence)



- Challenging to encode vast amounts of information
- Could try to divide your input (e.g. sentence by sentence)
 - Tasks (e.g., translation) have to local and global coherence

His extant works, including collaborations, consist of some 39 plays, 154 sonnets, three long narrative poems and a few other verses, some of uncertain authorship.



He is widely regarded as the greatest writer in the English language and the world's pre-eminent dramatist.[3][4][5] He is often called England's national poet and the "Bard of Avon" (or simply "the Bard"). His extant works, including collaborations, consist of some 39 plays, 154 sonnets, three long narrative poems and a few other verses, some of uncertain authorship.

- Challenging to encode vast amounts of information
- Could try to divide your input (e.g. sentence by sentence)
 - Tasks (e.g., translation) have to local and global coherence

• Instead, decode word by word by focusing on the relevant different parts of the input

(Soft) Attention





Encoder

Laurent Charlin — 60629

10

(Soft) Attention





Encoder

Laurent Charlin — 60629



exemple

de

RNN

(Soft) Attention





Encoder

Laurent Charlin — 60629



10

de

RNN



Advantages:

- No more bottleneck
- The decoder can consider different representations (information)
- The latent representation is now proportional to the length of the sequence.
- Shortcuts between the encoder and decoder
- Can model longer dependencies



Mathematical details:

$$\mathbb{P}(\mathbf{o}_t \mid \mathbf{o}_{1:t-1}, \mathbf{x}_{1:T}) = f(\mathbf{s}_t, \mathbf{c}_t),$$

where

$$\mathbf{c}_{t} = \sum_{j=1}^{T_{x}} \mathbf{a}_{t,j} \cdot [\mathbf{g}_{j}, \mathbf{h}_{j}]$$

Attention

- Timestep t (output)
- Hidden representation j (input)



Mathematical details:

The attention weights are obtained as follows:



Mathematical details:

The attention weights are obtained as follows:

$$\mathbf{a}_{t,j} = \left(\operatorname{softmax}(\mathbf{e}_{tj}) \right)_{j},$$

Where

 $\mathbf{e_{tj}} = \alpha(\mathbf{s_{t-1}}, \mathbf{h_j}, \mathbf{g_j}) \,.$

Function models the similarity between input and output representations



Mathematical details:

The attention weights are obtained as follows:

$$\mathbf{a_{t,j}} = \left(\mathsf{softmax}(\mathbf{e_{tj}})\right)_{\mathbf{j}},$$

Where

 $\mathbf{e}_{tj} = \alpha(\mathbf{s}_{t-1}, \mathbf{h}_j, \mathbf{g}_j) \,. \quad \blacktriangleleft$

Function models the similarity between input and output representations

The function $\alpha(\cdot)$ can be, for example, an MLP:



Mathematical details:

The attention weights are obtained as follows:

$$\mathbf{a_{t,j}} = \left(\mathsf{softmax}(\mathbf{e_{tj}})\right)_{\mathbf{j}},$$

Where

 $\mathbf{e}_{tj} = \alpha(\mathbf{s}_{t-1}, \mathbf{h}_j, \mathbf{g}_j) \,. \quad \blacktriangleleft$

Function models the similarity between input and output representations

The function $\alpha(\cdot)$ can be, for example, an MLP:

$$\alpha(\mathbf{s_{t-1}},\mathbf{h_j},\mathbf{g_j}) = \mathbf{v_a^{\top}} \mathsf{tanh}\left(\mathbf{W}_{\alpha}\mathbf{s_{t-1}} + \mathbf{U}_{\alpha}[\mathbf{h_j},\mathbf{g_j}]\right).$$

Translation task (EN -> FR)



Translation task (EN -> FR) (attention matrix)

Translation task (EN -> FR) (attention matrix)









end>

Translation task (EN -> FR) (attention matrix)





Soft Attention summary

• Dynamically weights the different contexts.

Laurent Charlin — 60629

• Empirically: Works for length of up to ~100 time steps

Toward Transformers



- A mix of sequential (h_i, g_j, S_t) and parallel (attention) processing
 - Expensive computationally
- Could a parallel architecture model sequential data?

Transformer Block



He is often called England's national poet and the "Bard of Avon" (or simply "the Bard"). His extant works, including collaborations, consist of some 39 plays, 154 sonnets, three long narrative poems and a few other verses, some of uncertain authorship.

- William Shakespeare (c. 23[a] April 1564 23 April 1616)[b] was an English playwright, poet and actor.
- He is widely regarded as the greatest writer in the English language and the world's pre-eminent dramatist.[3][4][5]

Word embedding

- Words are categorical and can be encoded using a 1of-K encoding (i.e., dummies)
- This fails to capture the semantics of words between words (e.g., foot/feet/toe)
- In RNNs/Transformers/etc. words are represented using a vector (of size D)
 - The representation is learned
 - Smaller-dimensional representation (D << K)

Word embedding

- Words are catego of-K encoding (i.e.
- This fails to capture
 words (e.g., foot/feedback
- In RNNs/Transforregime
 using a vector (of
- pot pan stirpoutern strain stew boil minutes_{stand} cook

pepper saltmix

oven

The representati

yolks beat whites wine thick cut roundlay thin pieces paper this pieces paper this pieces paper the pieces p

Smaller-dimensi

hour etc beef mrs side roashicken cake

per keep fryfat while fried potatoesup nix madethese his he corn good dramder madethese his he corn water into an bywereall or food their pint clean then from of the this they long wash dry until from of the this they long bake baking bake baking bake flavor cooking your juice who whiteup one when as so your juice who whiteup one when to it her she yollake flavor cooking till fire out at them but old see do hours cold ready let through againafter have first some hard just pot dish each hateetoone whole bethere another enough will parsley, nutmeg teaspoon fill lemon like little few other cough story only part tablespoonful story placetime without those ling who they cought ablespoonful tablespoonful tablesp

boil ounce servitable so hot very many its paste ten six five stand some should than about oil four twenty oil four twenty season fruit fruit new beaten half piece size

brown pudding

baked boiled sweet coffee

tea



The representation of each word (token) will be a combination of the representation of other words



• The procedure is unsupervised (self)



The representation of each word (token) will be a combination of the representation of other words

Example:



• The procedure is unsupervised (self)

The animal didn't cross the street because it was too tired

Mathematical Details

(Recall) Soft attention

$$a_{t,j} = \left(\mathsf{softmax}(e_{ij})\right)_{j},$$



Mathematical Details

(Recall) Soft attention

$$a_{t,j} = \left(\mathsf{softmax}(e_{ij})\right)_{j},$$



Self attention

$$a_{i,j} = \left(\mathsf{softmax}(e_{ij})\right)_j,$$

where

$$\mathbf{e_{ij}} = \left((\mathbf{x_i}\mathbf{W}^k)(\mathbf{x_j}\mathbf{W}^q) \right)$$

Mathematical Details

(Recall) Soft attention

$$a_{t,j} = \left(\mathsf{softmax}(e_{ij})\right)_j,$$



<u>Calculate its output</u>

$$\mathbf{c}_{t} = \sum_{j=1}^{T_{x}} \mathbf{a}_{t,j} \cdot [\mathbf{g}_{j}, \mathbf{h}_{j}]$$

Self attention

$$\mathbf{a}_{i,j} = \left(\text{softmax}(\mathbf{e}_{ij}) \right)_i,$$

where

$$\mathbf{e_{ij}} = \left((\mathbf{x_i}\mathbf{W}^k)(\mathbf{x_j}\mathbf{W}^q) \right)$$

Mathematical Details

(Recall) Soft attention

$$a_{t,j} = \left(\mathsf{softmax}(e_{ij})\right)_j,$$



<u>Calculate its output</u>

$$\mathbf{c}_{t} = \sum_{j=1}^{T_{x}} \mathbf{a}_{t,j} \cdot [\mathbf{g}_{j}, \mathbf{h}_{j}]$$

Self attention

$$\mathbf{a}_{i,j} = \left(\text{softmax}(\mathbf{e}_{ij}) \right)_{j},$$

where

$$\mathbf{e_{ij}} = \left((\mathbf{x_i}\mathbf{W}^k)(\mathbf{x_j}\mathbf{W}^q) \right)$$

X: la représentation en entrée des mots (NxD) z_i : la représentation en sortie du j'ème mot (1xD)

<u>Calculate its output</u>

Often referred to as "Key-Query-Value" attention

$$\label{eq:softmax} \begin{split} \textbf{z}_{j} &= (x_{i}W^{k})\, softmax((x_{i}W^{k})\, (x_{j}W^{q})) \\ & \overbrace{\text{Value}}^{} \qquad \overbrace{\text{Key}}^{} \overbrace{\text{Query}}^{} \end{split}$$





Visual representation N=2, D=4





*: in practice the value of attention pre-softmax (QK^{\top}) is often normalized by \sqrt{D} for optimisation stability

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



Multi-Head Self-Attention (MHSA)

- and sharks live in the water)
- head:

$$\mathbf{x'_j} = [\mathbf{z_j^1 z_j^2 \cdots z_j^H}]\mathbf{W^0}$$

Attention captures the similarity between two words

 One may wish to capture several different similarities (e.g., whale and humans are both mammals, whales

 You can learn the parameters for multiple attention mechanisms. Each mechanism is called an attention

 x'_i : the representation of word j after self-attention (D)

 z_i^h : the representation of word j from attention head h (D') H: the number of heads

 W^0 : weight matrix (D'H x D)

Multi-Head Self-Attention (MHSA)

Example:

The animal didn't cross the street because it was too tired







Two heads

<u>https://poloclub.github.io/transformer-explainer/</u>

- Shows attention
- Sub-word tokens

Laurent Charlin — 60629

Demo

Transformer Block



He is often called England's national poet and the "Bard of Avon" (or simply "the Bard"). His extant works, including collaborations, consist of some 39 plays, 154 sonnets, three long narrative poems and a few other verses, some of uncertain authorship.

Transformers... transform?

- Transformer blocks rely o MLP
- How can we use them for language tasks (e.g., translation)?
- Transformer blocks can be combined and refined
 - Transformer blocks can be used as encoders & decoders
 - For decoding, a few changes are required

Transformer blocks rely on attention heads followed by an

A few other characteristics

- Word embeddings
 - Words are categorical and can be encoded using a 1-of-K encoding (i.e., dummies)
 - This fails to capture the semantics of words between words (e.g., foot/feet/toe)
 - In RNNs/Transformers/etc. words are represented using a vector (of size D)

A few other characteristics

- Positional embedding
 - Attention looses the position of words into account
 - In many tasks, position is essential
 - "cat eats plant" is different from "plant eats cat"
 - As a remedy, transformers also learn positional representations
 - Vectors that are specific to the position (can be absolute or relative)

Word + Positional embeddings



A few other characteristics

- Normalization
 - Layer normalization is often applied. It helps to learn by normalizing each token's dimensions to have 0-mean and 1 standard deviation.
- Residual connections
 - $z_i = \text{Self-Attention}(X) + x_i$
 - Biases learned functions to be "simple", helps to learn
 - Can be used for self-attention & MLP

Laurent Charlin — 60629

What about decoding?

<EOS>

What about decoding?

- Transformers transform word representations
 - Decoding requires using these words representations to obtain the following word
 - Add a "softmax"-layer at the end:

$$P(\mathbf{o}_{\mathbf{n}} \mid \mathbf{o}_{1:\mathbf{n}-1}) = -$$

- You can only attend to previous words
 - Attention matrix is constrained to be (upper) triangular

- $\frac{\exp(\mathbf{g}_{\mathbf{w}}^{\top}\mathbf{x}_{n-1})}{\sum_{\mathbf{w}}^{W}\exp(\mathbf{g}_{\mathbf{w}}^{\top}\mathbf{x}_{n-1})}$
- $g_{\scriptscriptstyle W}\!\!:\!$ parameters W: number of words in the vocabulary

Encoder-Decoder Transformer

Encoder

- Self-Attention
- Positional embedding
- Encoder
- Decoder
- Softmax-layer
- Multiple transformer blocks (Nx)

- Self-Attention
- Positional embedding
- Encoder
- Decoder
- Softmax-layer
- Multiple transformer blocks (Nx)

- Self-Attention
- Positional embedding
- Encoder
- Decoder
- Softmax-layer
- Multiple transformer blocks (Nx)

- Self-Attention
- Positional embedding
- Encoder
- Decoder
- Softmax-layer
- Multiple transformer blocks (Nx)

Objective

- Often in two (or more) stages:
 - First stage
 - Next-word prediction

 - Other tasks:
 - Next-sentence prediction
 - Second stage
 - Preference and/or downstream task fine-tuning

Random-word prediction (mask some words in the input/output)

One of many transformers

- Partial list: <u>https://huggingface.co/docs/transformers/en/index</u>
- Encoder-only models:
 - BERT (RoBERTa), ALBERT
- Encoder-Decoder models:
 - BART
- Decoder-only models:
 - Generative pre-trained transformer (GPT), BLOOM (for code), Llama
- (Llama-7B -> 7 billion parameters)
- Also used for image modelling, audio, multi-modalities, etc.

• *Most of the above transformers refer to a system that was trained not only an architecture. Many of these models now have tens of billions of parameters

Summary

- Self-attention is the key ingredient
 - Efficient and provides good results
 - Recent LLMs can attend to very long contexts (some refinements of attention)
 - Enables training transformers on much larger-scale datasets (internet size)
- Transformers has become the standard model often surpassing the performance of RNNs/CNNs when trained on enough data

References

- 2024, <u>https://arxiv.org/abs/2304.10557</u>
- <u>arxiv.org/abs/1706.03762</u>

<u>https://jalammar.github.io/illustrated-transformer/</u>

An Introduction to Transformers, Richard E. Turner,

Attention Is All You Need, Vaswani et al., 2017, <u>http://</u>