

Lecture 2 Transformers

CSC485/2501

Announcements

- A1 release today (later in the afternoon).
- Essay 1 released Due Sep 16.
- A1 Tutorials:
 - T2: Fri, Sep 20.
 - T3: Wed, Sep 25.

Language Model

Models that estimate:

Probability (Some sentence over here.)

Neural Models









tensor([[0.7913, 0.4032, 0.4254, ..., 0.4260, 0.2358, 0.4319], [0.5912, 0.8657, 0.9194, ..., 0.1741, 0.1045, 0.0873], [0.7583, 0.4037, 0.6144, ..., 0.8484, 0.5134, 0.0949], ..., [0.5109, 0.5798, 0.1756, ..., 0.9258, 0.2265, 0.8081], [0.0229, 0.0628, 0.7615, ..., 0.5301, 0.6956, 0.0535], [0.6283, 0.3913, 0.7945, ..., 0.3634, 0.2575, 0.0871]])

Image: 863 x 625 pixels

Vectorize

Predict

Image vs. Natural Language

Vision (Image)

- Fixed input size (e.g., 960x720 pixels).
- Defined and continuous input range (RGB).



Natural Language

- Dynamic input sequence length.
- Discrete and large input space.
 - GPT-2 tokenizer: 50257 tokens.
 - LLaMA-3.1 tokenizer: 128000 tokens.





Probability (Some sentence over here.)?



Training Task 1: Next Token Prediction







Probability($\overset{\sim}{\bullet} \overset{\leftarrow}{\bullet} \overset{\sim}{\bullet} \overset{\circ}{\bullet}$) =P($\overset{\sim}{\bullet}$)*P($\overset{\leftarrow}{\bullet}$ | $\overset{\leftarrow}{\bullet}$)*P($\overset{\sim}{\bullet}$ | $\overset{\leftarrow}{\bullet}$ | $\overset{\leftarrow}{\bullet}$

Training Task 1: Next Token Prediction

Training Task 2: Masked Language Modelling

• MASK 15% of the tokens.

The Vanishing Gradient Problem

Formula Source: Razvan Pascanu et al. (2013)

© SuperDataScience

Transformers

Attention Mechanism

Attention Mechanism

Self Attention

What is o(like)? I.e., the attention score from like to the sentence.

*: made-up numbers, not real.

*: made-up numbers, not real.

Multi-Head Self Attention

Transformers

BERT: Pretrain-Finetune Paradigm

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

Class Label Τ, TN С T_2 ... BERT E_{N} E_[CLS] E₂ E [CLS] Tok 1 Tok N Tok 2 ... Single Sentence

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

Position Encoding

Input	[CLS] my	dog is	cute [SEP]	he likes	play	##ing	[SEP]
Token Embeddings	E _[CLS] E _{my}	E _{dog} E _{is}	E _{cute} E _[SEP]	E _{he} E _{likes}	E _{play}	E _{##ing}	E _[SEP]
Segment Embeddings	+ + E _A E _A	+ + E _A E _A	+ + E _A E _A	E _B E _B	+ E _B	+ E _B	+ E _B
	+ +	+ +	+ +	+ +	+	+	+
Position Embeddings	E ₀ E ₁	E ₂ E ₃	E ₄ E ₅	E ₆ E ₇	E ₈	E ₉	E ₁₀

$$ec{p}_t = egin{bmatrix} \sin(\omega_1.t) \ \cos(\omega_1.t) \ \sin(\omega_2.t) \ \cos(\omega_2.t) \ ec{p}_t = f(t)^{(i)} \coloneqq \begin{cases} \sin(\omega_k.t), & ext{if } i = 2k \ \cos(\omega_k.t), & ext{if } i = 2k+1 \end{cases}$$
 where $\omega_k = rac{1}{10000^{2k/d}}$

Position Encoding

- Encodings of any two distinct positions are distinct
- Each position maps to only one encoding
- Test sentences may be longer than training
- Distance between two positions should be constant across sentences (of varying lengths).

Review

- Embeddings:
 - Token + Sentence + Position
- Multi-Head Attention
- Feed forward module (MLP)

Layers

More Transformer in This Course

- Use Transformer Models to build all kinds of applications.
 - Parser (A1!), Word sense, QA, LLMs...
- Why do Transformers work so well?
 - Interpretability.
- How can we control them.
 - Interpretability, model editing.

