CS388: Natural Language Processing

Lecture 19: Pretrained Transformers





Credit: ???



Project 2 due Tuesday

Presentation day announcements next week

Administrivia



Each word forms a "query" which then computes attention over each word

$$\alpha_{i,j} = \operatorname{softmax}(x_i^\top x_j) \quad \text{scalar}$$

$$x_i' = \sum_{j=1}^n lpha_{i,j} x_j$$
 vector = sum of scalar

Multiple "heads" analogous to different convolutional filters. Use

$$\alpha_{k,i,j} = \operatorname{softmax}(x_i^\top W_k x_j) \quad x'_{k,i} = \sum_{j=1}^n \alpha_{k,i,j} V_k x_j$$

Vaswani et al. (2017)

Recall: Self-Attention



parameters W_k and V_k to get different attention values + transform vectors







Recall: Transformers







- a one-hot vector

Augment word embedding with position embeddings, each dim is a sine/cosine wave of a different frequency. Closer points = higher dot products

Works essentially as well as just encoding position as Vaswani et al. (2017)







BERT

• GPT/GPT2

Analysis/Visualization

This Lecture

BERT



- Al2 made ELMo in spring 2018, GPT was released in summer 2018, BERT came out October 2018
- Three major changes compared to ELMo:
 - Transformers instead of LSTMs (transformers in GPT as well)
 - Bidirectional <=> Masked LM objective instead of standard LM
 - Fine-tune instead of freeze at test time

BERT





ELMo reprs look at each direction in isolation; BERT looks at them jointly



A stunning ballet dancer, Copeland is one of the best performers to see live.



BERT

ELMo

- "ballet dancer"
- "ballet dancer/performer"









John

visited Madagascar yesterday

BERT

How to learn a "deeply bidirectional" model? What happens if we just



visited Madagascar yesterday John

Transformer LMs have to be "onesided" (only attend to previous tokens), not what we want





Masked Language Modeling

- BERT formula: take a chunk of text, predict 15% of the tokens
 - For 80% (of the 15%), replace the input token with [MASK]
 - For 10%, replace w/random
 - For 10%, keep same

How to prevent cheating? Next word prediction fundamentally doesn't work for bidirectional models, instead do masked language modeling







- Input: [CLS] Text chunk 1 [SEP] Text chunk 2
- 50% of the time, take the true next chunk of text, 50% of the time take a random other chunk. Predict whether the next chunk is the "true" next
- BERT objective: masked LM + next sentence prediction



Next "Sentence" Prediction



BERT Architecture

- BERT Base: 12 layers, 768-dim per wordpiece token, 12 heads. Total params = 110M
- BERT Large: 24 layers, 1024-dim per wordpiece token, 16 heads. Total params = 340M
- Positional embeddings and segment embeddings, 30k word pieces
- This is the model that gets pre-trained on a large corpus

Input

Token

Segment

Position



Devlin et al. (2019)









- CLS token is used to provide classification decisions
- Sentence pair tasks (entailment): feed both sentences into BERT
- BERT can also do tagging by predicting tags at each word piece

What can BERT do?





Entails

Transformer

Transformer

. . .

[CLS] A boy plays in the snow [SEP] A boy is outside

- How does BERT model this sentence pair stuff?
- Transformers can capture interactions between the two sentences, even though the NSP objective doesn't really cause this to happen

What can BERT do?





What can BERT NOT do?

BERT cannot generate text (at least not in an obvious way)

- Not an autoregressive model, can do weird things like stick a [MASK] at the end of a string, fill in the mask, and repeat
- Masked language models are intended to be used primarily for "analysis" tasks



Fine-tune for 1-3 epochs, batch size 2-32, learning rate 2e-5 - 5e-5



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

Fine-tuning BERT

- Large changes to weights up here (particularly in last layer to route the right information to [CLS])
- Smaller changes to weights lower down in the transformer
- Small LR and short fine-tuning schedule mean weights don't change much
- More complex "triangular learning rate" schemes exist





Pretraining	Adaptation	NER CoNLL 2003	SA SST-2	Nat. lang MNLI	g. inference SICK-E	Semantic SICK-R	textual si MRPC	milarity STS-B
Skip-thoughts		_	81.8	62.9	_	86.6	75.8	71.8
ELMo		91.7	91.8	79.6	86.3	86.1	76.0	75.9
	٨	91.9	91.2	76.4	83.3	83.3	74.7	75.5
	$\Delta = -$	0.2	-0.6	-3.2	-3.3	-2.8	-1.3	-0.4
		92.2	93.0	84.6	84.8	86.4	78.1	82.9
BERT-base	٨	92.4	93.5	84.6	85.8	88.7	84.8	87.1
	$\Delta = -$	0.2	0.5	0.0	1.0	2.3	6.7	4.2

Fine-tuning BERT

BERT is typically better if the whole network is fine-tuned, unlike ELMo

Peters, Ruder, Smith (2019)



Corpus	Train	Test	Task	Metrics	Domain					
	Single-Sentence Tasks									
CoLA SST-2	8.5k 67k	1k 1.8k	acceptability sentiment	Matthews corr. acc.	misc. movie reviews					
	Similarity and Paraphrase Tasks									
MRPC STS-B QQP	3.7k 7k 364k	1.7k 1.4k 391k	paraphrase sentence similarity paraphrase	acc./F1 Pearson/Spearman corr. acc./F1	news misc. social QA questions					
			Infere	ence Tasks						
MNLI QNLI RTE WNLI	393k 105k 2.5k 634	20k 5.4k 3k 146	NLI QA/NLI NLI coreference/NLI	matched acc./mismatched acc. acc. acc. acc.	misc. Wikipedia news, Wikipedia fiction books					

Evaluation: GLUE

Wang et al. (2019)





System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Ave
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79
BERTLARGE	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81

- Huge improvements over prior work (even compared to ELMo)
- imply sentence B), paraphrase detection

Results

Effective at "sentence pair" tasks: textual entailment (does sentence A

Devlin et al. (2018)







	"Robustly optimized BERT"	Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	(
		RoBERTa						
	160GB of data instead of	with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	
16 GB		+ additional data ($\S3.2$)	160GB	8K	100K	94.0/87.7	89.3	
	16 GB	+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	
		+ pretrain even longer	160GB	8K	500K	94.6/89.4	90.2	
		BERT _{LARGE}						
	Dynamic masking: standard	with BOOKS + WIKI	13GB	256	1 M	90.9/81.8	86.6	
	BERT USES THE SAME WASK							
	scheme for every epoch,							
	RoBFRTa recomputes them							

New training + more data = better performance

RoBERTa

Liu et al. (2019)



93.7



GPT/GPT2



- "ELMo with transformers" (works better than ELMo)
- Train a single unidirectional transformer LM on long contexts
- ► GPT2: trained on 40GB of text collected from upvoted links from reddit
- ▶ 1.5B parameters by far the largest of these models trained as of March 2019
- Because it's a language model, we can generate from it

OpenAl GPT/GPT2

Parameters	Layers	d_{model}
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

Radford et al. (2019)



OpenAl GPT2



SYSTEM PROMPT (HUMAN-WRITTEN)

MODEL COMPLETION (MACHINE-WRITTEN, SECOND TRY) Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back

The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.

She was carrying a pair of black and white striped gloves and a small black bag.

slide credit: OpenAl



1) How novel is the stuff being generated? (Is it just doing nearest neighbors on a large corpus?)

2) How do we understand and distill what is learned in this model?

3) How do we harness these priors for conditional generation tasks (summarization, generate a report of a basketball game, etc.)

model, not 1.5B yet)

- 4) Is this technology dangerous? (OpenAI has only released 774M param



- authors, and headline
- Humans rank Grover-generated propaganda as more realistic than real "fake news"
- Fine-tuned Grover can detect Grover propaganda easily authors argue for releasing it for this reason
- NOTE: Not a GAN, discriminator trained separately from the generator

Grover

Sample from a large language model conditioned on a domain, date,

			Un	paired A	Accuracy	y Paired Accuracy	
			(Generate	or size	Generator size	
			1.5B	355M	124M	1.5B 355M 124M	
		Chance		50.0		50.0	
Se	1.5B	Grover-Mega	92.0	98.5	99.8	97.4 100.0 100.0	
· Siz		Grover-Large	80.8	91.2	98.4	89.0 96.9 100.0	
itor	355M	BERT-Large	73.1	75.9	97.5	84.1 91.5 99.9	
mina		GPT2	70.1	78.0	90.3	78.8 87.0 96.8	
scri		GROVER-Base	70.1	80.0	89.2	77.5 88.2 95.7	
Dis	124M	BERT-Base	67.2	76.6	84.1	80.0 89.5 96.2	
		GPT2	66.2	71.9	83.5	72.5 79.6 89.6	
	11 M	FastText	63.8	65.6	69.7	65.9 69.0 74.4	

Zellers et al. (2019)





- BERT: Base \$500, Large \$7000
- Grover-MEGA: \$25,000
- XLNet (BERT variant): \$30,000 \$60,000 (unclear)
- This is for a single pre-training run...developing new pre-training techniques may require many runs
- Fine-tuning these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets)

<u>https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/</u>

Pre-Training Cost (with Google/AWS)





- NVIDIA: trained 8.3B parameter GPT model (5.6x the size of GPT-2)
- Arguable these models are still underfit: larger models still get better held-out perplexities

Pushing the Limits



NVIDIA blog (Narasimhan, August 2019)

Epoch





Number of tokens	Repeats	GLUE	CNNDM	SQuAD	SGLUE	EnDe	\mathbf{EnFr}	EnRo
\star Full dataset	0	83.28	19.24	80.88	71.36	26.98	39.82	27.65
2^{29}	64	82.87	19.19	80.97	72.03	26.83	39.74	27.63
2^{27}	256	82.62	19.20	79.78	69.97	27.02	39.71	27.33
2^{25}	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
2^{23}	4,096	76.34	18.33	70.92	59.29	26.37	38.84	25.81

Colossal Cleaned Common Crawl: 750 GB of text

We still haven't hit the limit of bigger data being useful

Google T5

Raffel et al. (October 23, 2019)







For downstream tasks: feed document into both encoder + decoder, use decoder hidden state as output

Good results on dialogue, summarization tasks

BART



Lewis et al. (October 30, 2019)





Analysis



What does BERT learn?





Heads on transformers learn interesting and diverse things: content heads (attend based on content), positional heads (based on position), etc.

Clark et al. (2019)



What does BERT learn?



Head 8-10



Still way worse than what supervised systems can do, but interesting that this is learned organically

Head 8-11

Head 5-4

Clark et al. (2019)











Probing BERT

Try to predict POS, etc. from each layer. Learn mixing weights

$$\mathbf{h}_{i,\tau} = \gamma_{\tau} \sum_{\ell=0}^{L} s_{\tau}^{(\ell)} \mathbf{h}_{i}^{(\ell)}$$

representation of wordpiece *i* for task τ

- Plot shows s weights (blue) and performance deltas when an additional layer is incorporated (purple)
- BERT "rediscovers the classical NLP pipeline": first syntactic tasks then semantic ones







Compressing BERT

- Remove 60+% of BERT's heads with minimal drop in performance
- DistilBERT (Sanh et al., 2019): nearly as good with half the parameters of BERT (via knowledge distillation)



(b) Evolution of accuracy on the MultiNLI-matched validation set when heads are pruned from BERT according to I_h (solid blue) and accuracy difference (dashed green).

Michel et al. (2019)





- analysis task
- These techniques are here to stay, unclear what form will win out
- purely from academia
- this cost should come down)

BERT-based systems are state-of-the-art for nearly every major text

Role of academia vs. industry: no major pretrained model has come

Cost/carbon footprint: a single model costs \$10,000+ to train (though