Announcements

- eCIS surveys released take a snapshot of the "done" page for your final project submission
- FP check-ins due Friday
- ▶ A4, A5 grading

Today: recap of BERT, applying BERT pre-trained models for generation: GPT-2, dialogue, summarization

Recall: Self-Attention > Each word forms a "query" which then computes attention over each word

 $\alpha_{i,j} = \operatorname{softmax}(x_i^\top W x_j))$ scalar $x_i' = \sum lpha_{i,j} V x_j$ vector = sum of scalar *

۲



> Multiple "heads" analogous to different convolutional filters. Use parameters W_k and V_k to get different attention values + transform vectors

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

Vaswani et al. (2017)



Recall	Next	"Sentence"	Prediction
--------	------	------------	------------

Input: [CLS] Text chunk 1 [SEP] Text chunk 2

۲

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

NotNext		Madagas †	car			enjoyed †		like ↑	
				Transf	ormer)
				Transf	ormer)
[CLS] Joh	n visited	[MASK]	yesterday	and	really	[MASK]	it [SE	P] / [mask	[] Madonna
								Devlin e	t al. (2019

Recall: BERT Architecture

 BERT Base: 12 layers, 768-dim per wordpiece token, 12 heads.
 Total params = 110M

۲

- 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







 Transformers can capture interactions between the two sentences, even though the NSP objective doesn't really cause this to happen

What can BERT NOT do?

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

- Can fill in MASK tokens, but can't generate left-to-right (well, you could put MASK at the end repeatedly, but this is slow)
- Masked language models are intended to be used primarily for "analysis" tasks

Fine-tuning BERT

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



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

- 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

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

Evaluation: GLUE

Wang et al. (2019)

		R	esu	lts					
System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	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	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERTBASE	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Huge improvements over prior work (even compared to ELMo)

 Effective at "sentence pair" tasks: textual entailment (does sentence A imply sentence B), paraphrase detection

Devlin et al. (2018)





OpenAl G	PT/GPT2				Oper
 "ELMo with transformers" (works b Train a single unidirectional transformation in the second sec	etter than ELM rmer LM on lon	o) g contexts		SYSTEM PROM (HUMAN-WRITTE	Miley Cyrus was caught shopli N Hollywood Boulevard today.
 GPT2: trained on 40GB of text collected from upvoted links 	Parameters	Layers	d _{model}	(MACHINE-WRITTE SECOND TR	N, the store by security guard Y) The singer was wearing a bl Lines' on the front and 'Fa
 from reddit 1.5B parameters — by far the 	117M 345M 762M	12 24 36	768 1024 1280		Scroll down for video Shoplifting: Miley Cyrus wa and Fitch on Hollywood Boul
when it came out in March 2019	1542M	48	1600		The singer was wearing a bl Lines' on the front and 'Fa
Because it's a language model, we observe the second se	can generate fro	om it Radford	et al. (2019)		The singer was also wearing black jacket, black jeans a She was carrying a pair of small black baa.

Open Questions	Pre-Training Cost (with Google/AWS)
 How novel is the stuff being generated? (Is it just doing nearest neighbors on a large corpus?) 	▶ BERT: Base \$500, Large \$7000
2) How do we understand and distill what is learned in this model?	 GPT-2 (as reported in other work): \$25,000
 How do we harness these priors for conditional generation tasks (summarization, generate a report of a basketball game, etc.) 	 XLNet (BERT variant): \$30,000 — \$60,000 (unclear)
4) Is this technology dangerous? (OpenAI pursued a "staged release" strategy and didn't release biggest model)	 This is for a single pre-training rundeveloping 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)
	https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-model

	OpenAl GPT2
SYSTEM PROMPT (HUMAN-WRITTEN)	Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.
MODEL COMPLETION (MACHINE-WRITTEN, SECOND TRY)	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.





	ELIZA	ELIZA
 Created 1964-1966 at MIT, heavily scripted DOCTOR script was most successful: repeats user's input, asks inane questions 	You are like my father in some ways. WHAT RESEMBLANCE DO YOU SEE You are not very aggressive but I think you don't want me to notice that. WHAT MAKES YOU THINK I AM NOT VERY AGGRESSIVE You don't argue with me. WHY DO YOU THINK I DON'T ARGUE WITH YOU You are afraid of me. DOES IT PLEASE YOU TO BELIEVE I AM AFRAID OF YOU My father is afraid of everybody. WHAT ELSE COMES TO MIND WHEN YOU THINK OF YOUR FATHER Bullies. DOES THAT HAVE ANYTHING TO DO WITH THE FACT THAT YOUR BOYFRIEND MADE YOU COME HERE	 Identify keyword, identify context, apply transformation rule (.*) you (.*) me My (.) (.*) Why do you think I \$2 you? What else comes to mind when you think of your \$1? You don't argue with me. WHY DO YOU THINK I DON'T ARGUE WITH YOU Very little need to generate new content, but can only have one type of conversation
	Weizenbaum (1966)	

Seq2seq mod	dels			Lack of Di	versity	
What are you doing <s></s>	n goinghome [STOP]	 Training to maximi responses: 	ze likelihood gives a Input: What are you doi -0.86 I don't know. -1.03 I don't know! -1.06 Nothing. -1.09 Get out of the way	ng? -1.09 Get out of here. -1.09 I'm going home. -1.09 Oh my god! -1.10 I'm talking to you.	common
Can we just train seq2seq models to "trans response?	ate" from utterance	to		Input: what is your nam -0.91 I don't know. -0.92 I don't know!	e? -1.55 My name is Robert.	
	System	BLEU		-0.92 I don't know, sir.	-1.58 My name is John.	
Hard to evaluate with automatic metrics:	RANDOM MT HUMAN	0.33 3.21 6.08		-0.97 Oh, my god! Input: How old are you -0.79 I don't know. -1.06 I'm fine. -1.17 I'm all right. -1.17 I'm not sure.	-1.59 My name's John. -1.64 Twenty-five. -1.66 Five. -1.71 Eight.	Li et al. (2016)

	Pers	onaChat	Meena		
Persona 1Persona 2I like to skiI am an artistMy wife does not like me anymoreI have four childrenI have went to Mexico 4 times this yearI recently got a catI hate Mexican foodI enjoy walking for exerciseI like to eat cheetosI love watching Game of Thrones		Persona 2 I am an artist I have four children I recently got a cat I enjoy walking for exercise I love watching Game of Thrones	 2.6B-parameter seq2seq model (larger than GPT-2) Trained on 341GB of online conversations scraped from public social media Sample responses: 		
	[PERSON 1:] In [PERSON 2:] Hello ! How are you today [PERSON 1:] I am good thank you, how	/ ? / are vou	Responses for "Why do you like the ocean?"	Score	
	[PERSON 2:] Great, thanks ! My childre [PERSON 1:] Nice ! How old are your cl [PERSON 2:] I have four that range in ag [PERSON 1:] I do not have children at th [PERSON 2:] That just means you get to [PERSON 1:] And Cheetos at the momen	n and I were just about to watch Game of Thrones. hildren? ge from 10 to 21. You? he moment. keep all the popcorn for yourself. nt!	 I love the beach. I like that it doesn't have anything in it. Just open water. It's dark and mysterious. I discovered it when I was really young and it was a dream of mine since. The ocean is full of wonders, mysteries and great things to discover. I love snorkeling. I moved to Hawaii a year ago and I'm going to college for marine biology. So there's that. The sea sounds cool and it's really satisfying to swim in. I dunno. It's heantiful I guess. It reflects my upbringing. I reckon 	-12.98 -14.03 -14.62 -14.87	
► Effo	rts to imbue seq2seq models	s with "personality" Zhang et al. (2018)	5. Well I'm from the pacific northwest and always somehow end up back there. I've been to Hawaii but I'm not one for the beach, it's too hot and the air is so dry. Adiwardana et	-15.36	













.

Neural Abstractive Systems

Actually this model ends up copying most of the time... 40% of the generated summary sentences are direct copies from article



BART: Pre-trained seq2seq model

	BART Summary
The researchers examined three types of coral in reefs off the coast of Fiji The researchers found when fish were plentiful, they would eat algae and seaweed off the corals, which appeared to leave them more resistant to the bacterium Vibrio corallillyticus, a bacterium associated with bleaching. The researchers suggested the algae, like warming temperatures, might render the corals' chemical defenses less effective, and the fish were protecting the coral by removing the algae.	Fisheries off the coast of Fiji are protect- ing coral reefs from the effects of global warming, according to a study in the jour- nal Science.
Sacoolas, who has immunity as a diplomat's wife, was involved in a traffic collision Prime Minister Johnson was questioned about the case while speaking to the press at a hospital in Wat- ford. He said, "I hope that Anne Sacoolas will come back if we can't resolve it then of course I will be raising it myself personally with the White House."	Boris Johnson has said he will raise the is- sue of US diplomat Anne Sacoolas' diplo- matic immunity with the White House.

	Takeaways
Þ	Pre-trained models are remarkably good at generating text
•	Story generation, dialogue systems, summarization, etc. have gotten way better in the past few years
•	Still much more to do: these systems usually don't have anything to say . Goal-oriented dialogue and grounded/embodied systems (e.g., a dialogue system on a robot are much tougher to get working
Þ	Next time: other languages

•

•