CS388: Natural Language Processing Lecture 8: Pre-trained Encoders

Greg Durrett







P2 due Tuesday

Final project released, proposals due Feb 23



- Vectors: d_{model}
- Queries/keys: d_k , always smaller than d_{model}
- Values: separate dimension d_v , output is multiplied by W^o which is $d_v x d_{model}$ so we can get back to *d_{model}* before the residual
- FFN can explode the dimension with W_1 and collapse it back with W_2

 $FFN(x) = max(0, xW_1 + b_1)W_2 + b_2$

Recall: Transformers amodel & Norm Feed **d**internal Forward **a**model Add & Norm $d_v \rightarrow d_{model}$ Multi-Head Attention **d**_{model} Vaswani et al. (2017)



Recall: Training Transformer LMs





loss fcn = nn.NLLLoss() loss += loss_fcn(log_probs, ex.output_tensor) [seq len, num output classes] [seq len]

classes] to [batch * seq len, num classes]. You do not need to batch

P(w|context)

`*loss = -- log P(w*|context)

Total loss = sum of negative log likelihoods at each position

Batching is a little tricky with NLLLoss: need to collase [batch, seq len, num]



ELMo

► BERT

- Subword tokenization
- BERT results, BERT variants
- Applying BERT

Today

ELMo



- "Pre-train" a model on a large dataset for task X, then "fine-tune" it on a dataset for task Y
- Key idea: X is somewhat related to Y, so a model that can do X will have some good neural representations for Y as well
- ImageNet pre-training is huge in computer vision: learn generic visual features for recognizing objects
- GloVe can be seen as pre-training: learn vectors with the skip-gram objective on large data (task X), then fine-tune them as part of a neural network for sentiment/any other task (task Y)

What is pre-training?





- GloVe uses a lot of data but in a weak way
- GloVe gives a single embedding for each word is wrong

they swing the bats

- Identifying discrete word senses is hard, doesn't scale. Hard to identify how many senses each word has
- How can we make our word embeddings more context-dependent? Use language model pretraining!

GloVe is insufficient

they see the bats







- word embeddings
- useful word representations in the same way word2vec did

Context-dependent Embeddings



Train a neural language model to predict the next word given previous words in the sentence, use the hidden states (output) at each step as

This is the key idea behind ELMo: language models can allow us to form

Peters et al. (2018)







CNN over each word => RNN





next word

Representation of visited (plus vectors from another LM running backwards)

2048 CNN filters projected down to 512-dim







- Use the embeddings as a drop-in replacement for GloVe
- Huge gains across many high-profile tasks: NER, question answering, semantic role labeling (similar to parsing), etc.
- But what if the pre-training isn't just for the embeddings?



BERT



- Al2 made ELMo in spring 2018, GPT (transformer-based ELMo) was released in summer 2018, BERT came out October 2018
- Four major changes compared to ELMo:
 - Transformers instead of LSTMs

 - Bidirectional model with "Masked LM" objective instead of standard LM Fine-tune instead of freeze at test time
 - Operates over word pieces (byte pair encoding)

BERT





- ELMo is a unidirectional model (as is GPT): we can concatenate two unidirectional models, but is this the right thing to do?
- 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



John visited Madagascar yesterday

You could do this with a "onesided" transformer, but this "twosided" model can cheat







Masked Language Modeling

- BERT formula: take a chunk of text, mask out 15% of the tokens, and try to predict them
- Optimize

P(Madagascar | John visited [MASK] yesterday)

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



Devlin et al. (2019)





- 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)







What can BERT do?



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

- RTE, SWAG
- Artificial [CLS] token is used as the vector to do classification from
- Sentence pair tasks (entailment): feed both sentences into BERT
- BERT can also do tagging by predicting tags at each word piece

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

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

Devlin et al. (2019)





Natural Language Inference

Premise

A boy plays in the snow

A man inspects the uniform of a figure

An older and younger man smiling

- Long history of this task: "Recognizing Textual Entailment" challenge in 2006 (Dagan, Glickman, Magnini)
- Early datasets: small (hundreds of pairs), very ambitious (lots of world knowledge, temporal reasoning, etc.)

Hypothesis

- entails A boy is outside
- The man is sleeping contradicts Two men are smiling and neutral laughing at cats playing







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

What can BERT do?

RTE, SWAG



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-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
- Often requires tricky learning rate schedules ("triangular" learning rates with warmup periods)





Subword Tokenization



- Words are a difficult unit to work with. Why?
 - haven't seen...
- Character-level models were explored extensively in 2016-2018 but simply don't work well — becomes very expensive to represent sequences

Handling Rare Words

When you have 100,000+ words, the final matrix multiply and softmax start to dominate the computation, many params, still some words you





- Subword tokenization: wide range of schemes that use tokens that are between characters and words in terms of granularity
- These "word pieces" may be full words or parts of words
 - the _eco tax _port i co _in _Po nt de Bu is ...
- indicates the word piece starting a word (can think of it as the space) character).

Subword Tokenization

Sennrich et al. (2016)





- Subword tokenization: wide range of schemes that use tokens that are between characters and words in terms of granularity
- These "word pieces" may be full words or parts of words



Can achieve transliteration with this, subword structure makes some translations easier to achieve

Subword Tokenization

Sennrich et al. (2016)





- for i in range(num_merges): pairs = get_stats(vocab) cooccurrences best = max(pairs, key=pairs.get) vocab = merge_vocab(best, vocab)

- many whole words
- Most SOTA NMT systems use this on both source + target

Byte Pair Encoding (BPE)

Start with every individual byte (basically character) as its own symbol

- Count bigram character
- Merge the most frequent pair of adjacent characters

Doing 8k merges => vocabulary of around 8000 word pieces. Includes

Sennrich et al. (2016)



Byte Pair Encoding (BPE)



Original:	furiously					
BPE:	_fur	io	usl	у		
Unigram LM:	_fur	_fur ious				
Original:	Completely prepo]	
BPE:	_Com	ple	t	ely		
Unigram LM:		mplet	e	ly		

- BPE produces less linguistically plausible units than another technique based on a unigram language model: rather than greedily merge, find chunks which make the sequence look likely under a unigram LM
- Unigram LM tokenizer leads to slightly better BERT

- tricycles **Original: BPE:** $_t$ | ric | y | cles (b) **Unigram LM:** _tri | cycle | **' S** rous suggestions
- CompletelypreposteroussuggestionsCompletelypreposteroussuggestions

Bostrom and Durrett (2020)







What's in the token vocab?

...



@SoC_trilogy

I've just found out that several of the anomalous GPT tokens ("TheNitromeFan", "SolidGoldMagikarp", " davidjl", "Smartstocks", "RandomRedditorWithNo",) are handles of people who are (competitively? collaboratively?) counting to infinity on a Reddit forum. I kid you not.



Rank	User	Co
1	/u/davidjl123	16
2	/u/Smartstocks	11
3	/u/atomicimploder	10
4	/u/TheNitromeFan	84
5	/u/SolidGoldMagikarp	65
6	/u/RandomRedditorWithNo	63
7	/u/rideride	59
8	/u/Mooraell	57
9	/u/Removedpixel	55
10	/u/Adinida	48
11	/u/rschaosid	47





- tuned vocabulary; usually between 50k and 250k pieces (larger number of pieces for multilingual models)
- As a result, classical word embeddings like GloVe are not used. All Transformer models

Tokenization Today

All pre-trained models use some kind of subword tokenization with a

subword representations are randomly initialized and learned in the

BERT results, BERT variants



Evaluation: GLUE

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		

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	5
160GB of data instead of 16 GB	RoBERTa with BOOKS + WIKI + additional data (§3.2) + pretrain longer + pretrain even longer	16GB 160GB 160GB 160GB	8K 8K 8K 8K	100K 100K 300K 500K	93.6/87.3 94.0/87.7 94.4/88.7 94.6/89.4	89.0 89.3 90.0 90.2	
Dynamic masking: standard BERT uses the same MASK scheme for every epoch, RoBERTa recomputes them	BERT _{LARGE} with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	

New training + more data = better performance

RoBERTa

Liu et al. (2019)



93.7







- Discriminator to detect replaced tokens rather than a generator to actually *predict* what those tokens are
- More efficient, strong performance

ELECTRA

















Slightly better variant

$$\begin{split} A_{i,j} &= \{ \boldsymbol{H}_{i}, \boldsymbol{P}_{i|j} \} \times \{ \boldsymbol{H}_{j}, \boldsymbol{P}_{j|i} \}^{\mathsf{T}} \\ &= \boldsymbol{H}_{i} \boldsymbol{H}_{j}^{\mathsf{T}} + \boldsymbol{H}_{i} \boldsymbol{P}_{j|i}^{\mathsf{T}} + \boldsymbol{P}_{i|j} \boldsymbol{H}_{j}^{\mathsf{T}} + \boldsymbol{P}_{i|j} \boldsymbol{P}_{j|i}^{\mathsf{T}} \end{split}$$

That is, the attention weight of a word pair can be computed as a sum of four attention scores using disentangled matrices on their contents and positions as content-to-content, content-to-position, position-to-content, and position-to-position².

Model	CoLA Mcc	QQP Acc	MNLI-m/mm Acc	SST-2 Acc	STS-B Corr	QNLI Acc	RTE Acc	MRPC Acc	Avg.
BERT _{large}	60.6	91.3	86.6/-	93.2	90.0	92.3	70.4	88.0	84.05
RoBERTa _{large}	68.0	92.2	90.2/90.2	96.4	92.4	93.9	86.6	90.9	88.82
XLNet _{large}	69.0	92.3	90.8/90.8	97.0	92.5	94.9	85.9	90.8	89.15
ELECTRA _{large}	69.1	92.4	90.9/-	96.9	92.6	95.0	88.0	90.8	89.46
DeBERTa _{large}	70.5	92.3	91.1/91.1	96.8	92.8	95.3	88.3	91.9	90.00

He et al. (2021)





Using BERT

HuggingFace Transformers: big open-source library with most pre-trained architectures implemented, weights available

Lots of standard models...

Model architectures

Transformers currently provides the following NLU/NLG architectures:

- 1. **BERT** (from Google) released with the paper **BERT**: **Pre-training of Deer** Understanding by Jacob Devlin, Ming-Wei Chang, Kenton Lee and Krist
- 2. GPT (from OpenAI) released with the paper Improving Language Under Radford, Karthik Narasimhan, Tim Salimans and Ilya Sutskever.
- 3. GPT-2 (from OpenAI) released with the paper Language Models are Un Jeffrey Wu*, Rewon Child, David Luan, Dario Amodei** and Ilya Sutskev
- 4. Transformer-XL (from Google/CMU) released with the paper Transform Fixed-Length Context by Zihang Dai*, Zhilin Yang*, Yiming Yang, Jaime
- 5. XLNet (from Google/CMU) released with the paper XLNet: Generalized Understanding by Zhilin Yang*, Zihang Dai*, Yiming Yang, Jaime Carbon
- 6. XLM (from Facebook) released together with the paper Cross-lingual Li and Alexis Conneau.
- 7. RoBERTa (from Facebook), released together with the paper a Robustly

and "community models"







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)









Applying BERT



- Compared to ELMo, BERT is very good at sentence-pair tasks
 - Paraphrase detection
 - Semantic textual similarity
 - Textual entailment

Two Tasks

Question answering (not really a sentence pair, but it's a pair of inputs)



- Show people captions for (unseen) images and solicit entailed / neural / contradictory statements
- >500,000 sentence pairs
- One possible architecture:

300D BiLSTM: 83% accuracy (Liu et al., 2016)

One of the first big successes of LSTMbased classifiers (sentiment results were more marginal)



Bowman et al. (2015)





Drawn from multiple genres of text

Premise

Fiction

The Old One always comforted Ca'daan, except t

Letters

Your gift is appreciated by each and every studer will benefit from your generosity.

Telephone Speech

yes now you know if if everybody like in August w everybody's on vacation or something we can dread more casual or

9/11 Report

At the other end of Pennsylvania Avenue, people line up for a White House tour.

MNLI Dataset

	Label	Hypothesis
today.	neutral	Ca'daan knew the Old One very well.
nt who	neutral	Hundreds of students will benefit from your generosity.
vhen ess a little	contradiction	August is a black out month for vacations in the company.
began to	entailment	People formed a line at the end of Pennsylvania Avenue.
		Williams et al. (2018)





How do models do it?





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

- transformed
- **But**, models are often overly sensitive to lexical overlap

Transformers can easily learn to spot words or short phrases that are

Williams et al. (2018)





Question Answering

- Many types of QA:
- We'll focus on factoid questions being answered from text
 - E.g., "What was Marie Curie the first female recipient of?" unlikely you would have this answer in a database
 - Not appropriate: "When was Marie Curie born?" probably answered in a DB Not appropriate: "Why did World War II start?" — no simple answer



Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the Nobel Prize. Famous musicians include Władysław Szpilman and Frédéric Chopin. Though Chopin was born in the village of Żelazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.

Answer = Nobel Prize

how to retrieve these effectively (will discuss when we get to QA)

SQuAD

Assume we know a passage that contains the answer. More recent work has shown





Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the **Nobel Prize**....

Predict answer as a pair of (start, end) indices given question q and passage p; compute a score for each word and softmax those

$$P(\text{start} | q, p) =$$

$$f$$

$$recipient$$

SQuAD

0.010.010.850.01t of the **Nobel Prize** .

P(end | q, p) = same computation but different params







What was Marie Curie the first female recipient of ? [SEP] One of the most famous people born in Warsaw was Marie ...

In a couple lectures, we will look at what BERT learns when trained on this kind of data

QA with BERT

Devlin et al. (2019)





Pre-trained models and BERT are very powerful for a range of NLP tasks

These models have enabled big advances in NLI and QA specifically

Next time: pre-trained decoders (GPT-3) and encoder-decoder models (T5)