

A primer in BERTology

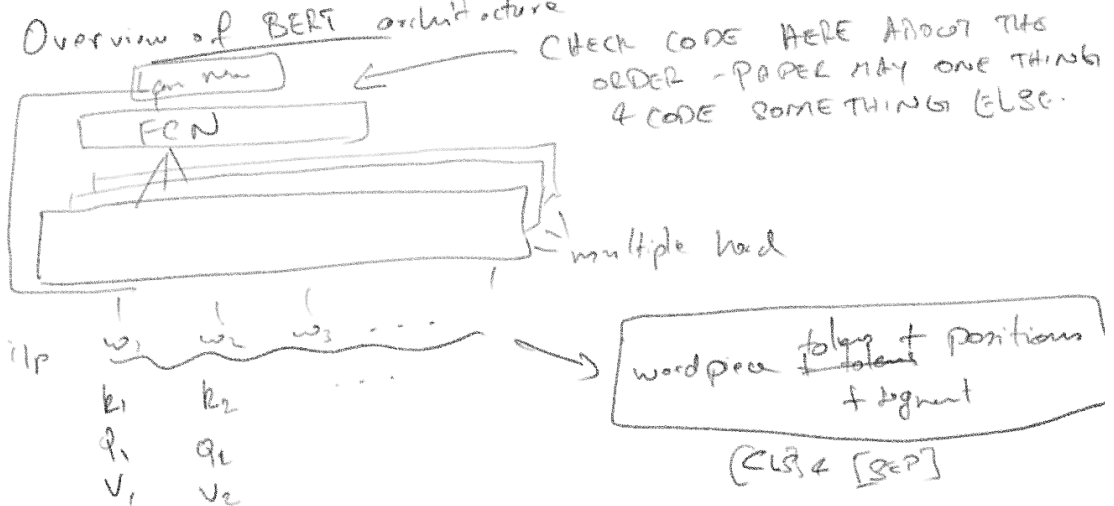
Feb 2020 - Umav Lowell

Abstract

Transformer used a lot - but do not understand inner workings
This paper describes what is known to date
Propose modifications for further research

- ① limits hypothesis derived development b/c of lack of understanding of inner workings
 - Unlike CNN transformer have little cognitive motivation
 - Size hinders even ablation studies

② Overview of BERT architecture



2 stages

- a) Pre-training
- b) Fine-tuning

2 tasks

- MLM
- NSP (next sentence prediction)

2 versions

- a) Base
- b) Large

③ BERT EMBEDDINGS

word2vec & Glove are static, BERT repr is contextualized.

- ① BERT embeddings form clear clusters corresponding to

word sense

- ⑥ Late layers in BERT produce more content specific repr
[How did they measure this?]

④ What knowledge does BERT have?

Popular approaches:

fill-in-the-gap probes

Analysis of Attn weights

probing classifiers

④.1 Syntactic knowledge

- a) BERT representations are hierarchical rather than linear. in order to apart from word order information, there is also something like a syntactic tree - Lin 2019 (getting inside BERT's linguistic knowledge)
- b) BERT embeddings also encode POS, syntactic chunks and roles [what do you learn from context? - 2019]
- c) Hewitt & Manning learned transformation matrices that recover the dependency
- * d) Sawaher 2019 - Approximate BERT repr with Tensor Product Decomposition Networks.
- * e) Predictions were not altered with shuffled word order, truncated sentences, removed subjects & objects
- d) BERT encodes syntactic structure, but does not rely on that knowledge.

④.2 Semantic knowledge

- a) tip a robin < tip a chef < tip a waiter
↑
incorrect filter for semantic roles that are semantically related rather than not related

[How was this preference measured?]

- ⑥ BERT struggles with representations of numbers - probes
1 billion? why

due to word piece tokenizer [what are word piece tokens?
not straight away stored with w2Vec or GloVe?
[Also, how are numbers encoded in w2Vec?]

4.3 World knowledge

a) "Dante ~~was~~ born in [MASK]."

knowledge induction by filling in the blanks.

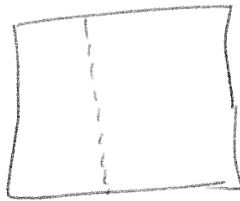
b) Common Sense - [BERT knows that people can walk into
houses, and houses are big, but it doesn't know that
houses are bigger than people - [How did they
evaluate this? Does BERT "know" what is "big" / "small"?]
If it can learn it, then do you say it has
common sense?

5) Localizing Linguistic knowledge.

5.1 Self-attention Heads

a) - Estimating dark secrets of BERT. (look into this - they've
classified attention types using CNNs).

Much of the model encodes [SEP] and [CLS] and
productions of vertical attention



← vertical pattern in attn matrix

This redundancy must be related to overparameterization.

[Did I notice this? No, I don't think so, my
patterns were very different, but they all had
a vertical "look" to it].

Periods, commas, are as frequent as [CLS] & [SEP], so
the model learns to rely on them a lot more.

[SEP] gets increased attn starting layer 5, but its
importance in prediction drops [How was this
importance measured?]

5.2 BERT Layers

- a) Decrease in linear word order information around layers. Accompanied by increased knowledge of Hierarchical sentence structure [Lin 2019]
- b) Syntactic information is the most prominent in the middle BERT layers. Hewitt & Manning: 6-9 layers - reconstructed the middle BERT layers
Goldberg 2019: Subject-verb agreement
- c) Middle layers are the most transferable across tasks
- d) Final layers of BERT are more task specific.
- in pre-training - MLM task
if you after fine-tuning, if you restore the early layers to its original weight, it does not dramatic hurt the model performance.

6 Training

6.1 Pretraining

a) Removing NSP does not hurt (sometimes improve performance)

b) Diverse masks during MLM pre-training within an epoch

c) XLNet

6.2 Model Architecture Choices

a) No. of heads not as important as no. of layers

b) Large batch training (8k examples) improves.

c) $bsz = 32k$, training time reduces with no post-decay

d) Normalise [CLS]

e) Model pre-training by stacking - "worm stack"

6.3 Fine-tuning BERT

④ ~~Self~~ attention on [SEP] basically tells BERT what to ignore. [SEP] attn increases on finetuning tasks.

⑤ Propose using weighted repr. of all layers instead of only the last layer.

⑥ Adaptor Modules - reduces the ~~training~~ ^{finetuning} cost to a fraction of original

⑦ How big should BERT be?

⑦.1 Overparameterization

- [Pay less Attn with Lightweight & Dynamic Convolutions.] Feb 2019 - Good paper.

⑧ Current models do not make good use of the parameters they already have.

Specialized heads do the heavy lifting - the rest can be pruned - 2019

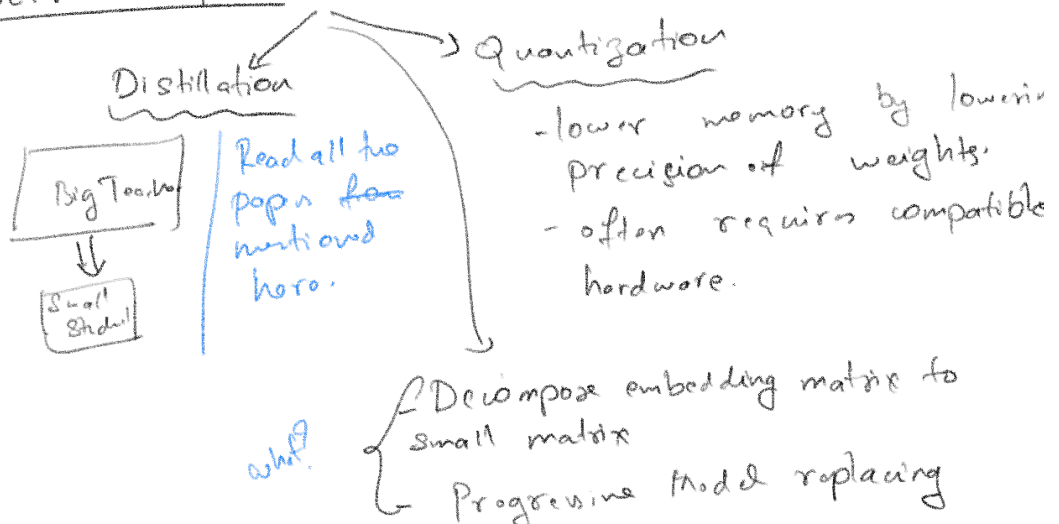
What does BERT look at? Analysis of BERT's attention - 2019

Most layers could be reduced to a single head.

⑨ Some heads are harmful for downstream tasks.

⑩ Attn dropout could be the reason why some heads are redundant.

⑦.2 BERT Compression.



② Multi-lingual BEK1



⑨ Discussions

⑨.1 Limitations

- ② "the fact that a linguistic pattern is not there, doesn't mean that it is absent or the presence of one doesn't mean that it is used"

→ more complex probe might be able to recover more information. - BUT IT BECOMES LESS CLEAR WHETHER WE ARE TALKING ABOUT THE ORIGINAL MODEL.

- ⑤ Ongoing debate about the merits of attention as a tool for interpreting deep learning models.
(WHAT? - On identifiability in Transformers, Attention is not explanation / Attention is not not explanation, Is Attention Interpretable?)

⑨.2 Incentivize dataset development