Automatic Student Network Search for Knowledge Distillation

1 MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University, China

2 Ping An Technology (Shenzhen) Co. Ltd., Shenzhen, China

Introduction

Zhexi Zhang¹, Wei Zhu², Junchi Yan^{1*}, Peng Gao² and Guotong Xie^{2*}

- BERT is a commonly used pretrained language model (PLM) that obtains state-of-the-art results on 11 different natural language processing (NLP) tasks.
- However, BERT contains a large number of parameters and requires vast numbers of computational resources. Concretely, BERT_{BASE} has 110 million parameters while BERT_{LARGE} has 340 million.
- Knowledge distillation (KD) is a promising compression method and has achieved great success in compressing BERT. KD follows a student-teacher framework where the light-weight student network learns from the teacher.
- However, the student networks in previous KD studies are manually designed. Researchers have tried to compressing BERT into MLP, BiLSTM and network with less Transformer layers. These manually designed students are not optimal since they either still contain redundant parameters or have week representation ability.
- Motivated by the above observations, we propose to **automatically search for a compact student network for compressing BERT using neural architecture search (NAS)**.

NAS-KD: our Framework As illustrated in the right, NAS-KD is composed of NAS module and KD module. In NAS module, we select and combine the candidate operations in search space to generate a student cell. In KD module, multiple student cells are stacked as the student network in knowledge distillation. distillation objective The is formulated by three different loss functions. The total loss is returned to NAS module to update the parameters and structure of student cell.





KD Module

Except vanilla KD objective where the student learns from the ouput of teacher, we encourage each student cell to uniformly mimic a corresponding Transformer layer of BERT. The formulations are in the right.

11	$\sum_{i=1}$	$\sum_{j=1}$	$\ \ \mathbf{h}_{i,j}^s \ $	$\ _{2}^{2}$	$\ \mathbf{h}_{i,\frac{M^{t}}{M^{s}}}^{t}\ $
	$\mathcal{L} = c$	$\alpha \mathcal{L}_{ss}$	+(1 -	$\alpha)\mathcal{L}_{sh}$	$+\beta \mathcal{L}_{1}$

 $(\hat{y}_{i,j}^t \cdot \log(\hat{y}_{i,j}^s/t))$

Method	Par. (w/o Emb)	Par. (total)	Layers	SST-2	MRPC	QQP	MNLI-m	MNLI-mm	QNLI	RTE
BERT _{BASE} (Google)	85.2	109	12	93.5	88.9	71.2	84.6	83.4	90.5	66.4
BERT _{BASE} (Teacher)	85.2	109	12	93.1	87.7	71.2	83.5	82.8	90.3	66.1
Distilled BiLSTM _{SOFT}	0.96	10.1	1	90.7	-	68.2	73.0	72.6	-	-
TinyBERT (w/o DA)	4.8	9.7	4	-	82.4	-	80.5	81.0	-	-
BERTSMALL	4.8	9.7	4	87.6	83.2	66.5	75.4	74.9	84.8	62.6
DistilBERT	28.4	52.2	4	91.4	82.4	68.5	78.9	78.0	85.2	54.1
BERT ₃ -FT	23.9	45.7	3	86.4	80.5	65.8	74.8	74.3	84.3	55.2
BERT ₃ -KD	23.9	45.7	3	86.9	79.5	67.3	75.4	74.8	84.0	56.2
BERT ₃ -PKD	23.9	45.7	3	87.5	80.7	68.1	76.7	76.3	84.7	58.2
NAS-KD ₃	9.4±1.5	33.2 ± 1.5	3	86.9	79.3	67.5	76.1	75.5	83.9	58.9
BERT ₆ -FT	43.2	67.0	6	90.7	85.9	69.2	80.4	79.7	86.7	63.6
BERT ₆ -KD	43.2	67.0	6	91.5	86.2	70.1	80.2	79.8	88.3	64.7
BERT ₆ -PKD	43.2	67.0	6	92.0	85.0	70.7	81.5	81.0	89.0	65.5
NAS-KD ₆	$18.6 {\pm} 2.9$	42.4 ± 2.9	6	92.2	86.3	70.4	81.0	80.2	88.6	65.9