STraTA: Self-Training with Task Augmentation for Better Few-shot Learning

Tu Vu, Minh-Thang Luong , Quoc V. Le , Grady Simon , Mohit Iyyer EMNLP 2021

Motivation

Large generative model (e.g. GPT3) still lags behind well-finetuned model

This paper proposes STraTa to leverage unlabeled data:

- Self-Training
- Task Augmentation

Framework



Figure 2: An illustration of our Self-Training with Task Augmentation (STraTA) approach. In task augmentation, we train an NLI data generation model and use it to synthesize a large amount of *in-domain* NLI training data for each given target task, which is then used for auxiliary (intermediate) fine-tuning. Our self-training algorithm iteratively learns a better model using a concatenation of labeled and pseudo-labeled examples. At each iteration, we always start with the auxiliary-task model produced by task augmentation and train on a broad distribution of pseudo-labeled data.

Task Augmentation

- Goal:
 - Finetune LMs on auxiliary task before the target task
- Which task?
 - Domain mismatch: MNLI and SQuAD
- In this paper:
 - Finetune on NLI data

Task Augmentation (2)

- Synthetic NLI data generation
 - Generate from MNLI
 - (sentA, sentB)-> label
 - To NLI
 - (label, sentA) -> sentB
 - Finetune T5 model
- Benefit
 - Training label is free
 - Be able to overgeneration in-domain NLI training data
- Overgeneration
 - Generate 100 output samples per input (top-k=40)
- Filtering
 - Use BERT model finetuned on MNLI as NLI classifier
 - Filter synthetic data:
 - If produce the same label -> kept
 - Otherwise -> Remove

Self-training

- Goal
 - Improve the model using pseudo-labeled data
- A strong base model
 - Strong base model can produce quite good pseudo labeled data
 - Otherwise, errors can be propagated through incorrect data
- Steps

-

- Starts with the same strong base model
- Finetune all parameters using labeled and pseudo-labeled data

Self-training (2)

- Self-training on a broad distribution of pseudo-labeled data

- Confident data ~ narrow distribution
 - Overconfident teacher
 - Poorly calibrated
 - Harmful to the self-training
- Less confident data ~ broad distribution
 - Using the whole set of pseudo-labeled data
 - At each iteration, the pseudo-labels are regenerated as the teacher improves gradually

Experiment

- LMFT: Target task language model finetune
 - Finetune LM on in-domain unlabeled text using MLM signal
 - Finetune LM on the target task
- ITFT_MNLI: intermediate task finetuning on MNLI
 - Finetune LM on MNLI
 - Finetune LM on target task
- TA: task augmentation
- ST: self-training

_

Model	SNLI	QQP	QNLI	SST-2	SciTail	SST-5	STS-B	SICK-E	SICK-R	CR	MRPC	RTE
Full												
BERTLARGE	91.1	88.4	91.9	92.4	95.3	53.7 _{0.9}	89.6 _{0.2}	$87.9_{0.6}$	$84.4_{0.4}$	$91.7_{0.6}$	$89.0_{0.8}$	68.6 _{7.2}
+ LMFT	91.0	88.1	90.4	93.5	95.3	$54.0_{0.4}$	89.5 _{0.2}	87.7 _{0.5}	$84.0_{0.5}$	91.6 _{0.8}	$89.5_{1.0}$	66.5 _{7.3}
+ $ITFT_{MNLI}$	91.1	88.2	91.6	93.5	96.5	$54.0_{0.8}$	90.3 _{0.3}	$89.9_{0.2}$	86.3 _{0.3}	$92.0_{0.6}$	89.7 _{0.9}	82.3 _{1.4}
+ TA	91.9	88.5	92.5	94.7	96.9	55.7 _{0.8}	90.9 _{0.2}	90.7 _{0.3}	87.0 _{0.3}	93.3 _{0.6}	90.8 _{0.7}	83.8 _{1.1}
LIMITED (1024 total training examples)												
BERTLARGE	$77.4_{0.6}$	$74.1_{1.0}$	81.7 _{0.9}	89.8 _{0.6}	90.9 _{0.7}	49.1 _{1.3}	88.20.4	$84.8_{0.7}$	$80.2_{0.4}$	$91.2_{0.6}$	85.7 _{1.7}	66.8 _{2.7}
+ LMFT	75.8 _{1.5}	71.60.5	80.52.0	88.9 _{0.8}	87.7 _{2.3}	49.2 _{3.1}	$88.4_{0.4}$	$83.2_{0.6}$	$78.5_{0.6}$	90.9 _{0.7}	$84.9_{1.1}$	65.2 _{3.4}
+ $ITFT_{MNLI}$	$85.2_{0.4}$	$74.0_{0.5}$	83.5 _{0.5}	$90.0_{0.8}$	92.1 _{1.1}	$49.4_{1.2}$	$87.8_{0.8}$	88.8 _{0.5}	83.20.7	$91.3_{0.7}$	86.40.9	81.1 _{1.3}
+ TA	87.3 _{0.3}	75.7 _{0.5}	85.0 _{0.5}	91.7 _{0.7}	92.3 _{1.1}	$51.4_{1.0}$	89.0 _{0.6}	89.4 _{0.4}	84.3 _{0.4}	92.6 _{0.6}	88.0 _{0.8}	82.9 _{1.8}
FEW-SHOT (8 training examples per class)												
BERT _{BASE}	43.72.2	55.9 _{6.5}	59.0 _{10.9}	59.1 _{8.4}	67.1 _{6.6}	30.5 _{2.0}	73.6 _{4.5}	61.3 _{4.1}	59.7 _{2.7}	65.2 _{8.2}	$72.4_{10.2}$	51.4 _{2.5}
+ LMFT	45.2 _{3.9}	57.2 _{6.2}	57.6 _{9.1}	64.9 _{8.7}	$64.0_{8.0}$	33.4 _{1.9}	75.44.4	59.3 _{4.0}	58.3 _{2.0}	$72.4_{6.0}$	$73.9_{8.6}$	50.9 _{3.9}
+ ITFT _{MNLI}	75.2 _{5.7}	63.7 _{7.0}	62.8 _{5.1}	76.8 _{7.2}	75.8 _{5.6}	35.0 _{2.6}	80.21.1	80.41.9	73.5 _{2.7}	79.2 _{3.6}	$74.3_{8.0}$	62.213.5
+ TA	83.3 _{0.8}	68.7 _{1.5}	70.1 _{3.4}	80.3 _{6.6}	78.5 _{3.2}	37.4 _{3.0}	80.71.5	81.1 _{2.4}	$75.9_{1.8}$	86.52.2	74.5 _{6.5}	67.6 _{7.1}
+ ST	65.0 _{5.8}	69.9 _{5.9}	71.6 _{11.3}	$62.7_{10.4}$	$68.6_{8.3}$	33.9 _{3.5}	80.52.2	68.1 _{4.5}	64.0 _{2.4}	$78.2_{6.3}$	$80.5_{1.8}$	50.7 _{3.1}
+ $ITFT_{MNLI}$ + ST	83.20.3	70.7 _{5.9}	81.5 _{1.2}	88.0 _{2.1}	83.7 _{4.4}	39.5 _{2.0}	84.20.8	81.8 _{2.6}	75.8 _{2.2}	85.6 _{2.3}	80.6 _{1.2}	62.5 _{12.0}
+ STraTA	85.7 _{0.2}	$74.5_{0.4}$	82.1 _{0.5}	90.1 _{0.8}	86.3 _{3.5}	$41.3_{1.5}$	84.7 _{0.5}	84.9 _{1.2}	77.6 _{1.6}	90.5 _{0.8}	$81.0_{0.8}$	70.6 _{2.4}
BERTLARGE	43.1 _{4.4}	58.5 _{4.7}	64.4 _{6.1}	66.1 _{8.7}	68.8 _{9.5}	35.2 _{1.3}	74.6 _{3.8}	66.5 _{4.5}	66.6 _{3.3}	72.0 _{6.0}	79.9 _{2.0}	53.1 _{3.3}
+ LMFT	39.6 _{2.6}	52.7 _{4.7}	52.21.6	66.3 _{9.3}	$66.4_{10.6}$	36.8 _{2.9}	$75.4_{9.4}$	58.8 _{6.9}	51.6 _{7.0}	75.6 _{5.9}	80.52.4	52.8 _{4.8}
+ $ITFT_{MNLI}$	$79.9_{3.1}$	62.6 _{9.0}	64.5 _{4.4}	80.7 _{5.0}	$72.3_{11.2}$	36.42.1	75.54.0	77.8 _{3.8}	73.5 _{2.8}	82.6 _{3.0}	72.87.9	69.7 _{14.6}
+ TA	$84.8_{0.7}$	64.6 _{6.3}	71.54.0	85.5 _{1.4}	79.0 _{4.5}	38.5 _{3.0}	78.9 _{2.4}	81.2 _{3.9}	77.5 _{1.4}	88.6 _{1.3}	$78.2_{6.6}$	77.0 _{6.3}
+ ST	69.3 _{9.2}	$74.3_{1.2}$	85.4 _{1.7}	$81.9_{12.2}$	$79.9_{4.8}$	$42.0_{1.5}$	82.82.3	77.3 _{3.1}	73.1 _{2.3}	88.1 _{1.3}	$81.2_{0.5}$	53.9 _{4.3}
+ $ITFT_{MNLI}$ + ST	85.4 _{0.3}	$74.8_{0.7}$	86.1 _{1.1}	89.7 _{0.7}	86.24.2	$42.2_{2.0}$	84.1 _{1.7}	84.3 _{2.0}	78.4 _{1.3}	$89.3_{1.0}$	$81.4_{1.2}$	72.75.4
+ STraTA	87.3 _{0.3}	75.1 _{0.2}	86.4 _{0.8}	91.7 _{0.7}	87.3 _{2.9}	43.0 _{2.3}	$84.5_{1.6}$	86.3 _{1.8}	79.0 _{1.0}	90.0 _{0.6}	81.5 _{0.7}	77.1 _{5.4}
Prompt-be	ased (LM	-BFF; G	ao et al., 2	2021) and	l entailme	nt-based	(EFL; Wa	ung et al., 2	021) fine-ti	uning app	proaches	
RoBERTa LARGE	$38.4_{1.3}$	58.8 _{9.9}	52.7 _{1.8}	60.5 _{3.1}	_	-	24.5 _{8.4}	-	-	$61.9_{5.1}$	76.1 _{3.9}	55.0 _{1.3}
+ LM-BFF	$52.0_{1.7}$	68.2 _{1.2}	61.8 _{3.2}	$79.9_{6.0}$	-	-	66.0 _{3.2}	-	-	88.62.3	78.5 _{2.3}	63.3 _{2.1}
+ EFL	81.0 _{1,1}	$67.3_{2.6}$	68.0 _{3.4}	90.8 _{1.0}	_	-	71.0 _{1.3}	_	_	92.3 _{0.4}	$76.2_{1.3}$	85.8 _{0.9}