Recent Advances in Pre-trained Language Models

姜成翰

Cheng-Han Chiang

04.01.2022

1



更新日期:111年3月30日

本校防疫訊息--依據CDC 3月28日公告提醒本校教職員工生落實各項防疫措施

全校教職員工生大家好:

依據中央流行疫情指揮中心 3月28日發布的公告,請全校教職員工生在校內活動時,除用餐、室內 外運動、以及指揮中心規定的例外情形之外,其他活動均應全程配戴口罩。講課、以及室內外拍攝 個人/團體照,雖屬於指揮中心規定得免戴口罩的活動,但本校建議如無法與他人保持社交距離時, 仍應該配戴口罩。其他防疫措施提醒事項如下:

大家好: 疫情升溫,保健中心也接獲同學反應:有許多學生在上課時並未配戴口單,可能造成防疫漏洞。 尤其曾發生配戴口單的同學被迫與不戴口單的同學分在同一組討論,令配戴口單的同學感到很大的壓力。 保健中心在此籲請本校各單位再度對內加強宣導: 學生上課時應一律配戴口單,也請各單位提醒教師上課時應提醒學生全程配戴口單, 尊重自己也尊重他人,不應該讓遵守規定配戴口單的同學承受他人不配戴口單的防疫壓力,謝謝大家~



- •【機器學習2021】<u>自注意力機制</u> (Self-attention)(上)
- •【機器學習2021】<u>自注意力機制</u> (Self-attention)(下)





Highly Related Topics

•【機器學習2021】BERT簡介

•【機器學習2021】<u>GPT的野望</u>





Highly Related Topics

•【深度學習於人類語言處理 2020】<u>BERT and its family</u>



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
- Closing Remarks

Outline

- Background knowledge
 - Pre-trained Language Models
- The Problems of PLMs
- The Solutions of Those Problems
- Closing Remarks

• Neural Language Models: A neural network that defines the probability over sequences of words



- How are these language models trained?
 - Given an incomplete sentence, predict the rest of the sentence



- Autoregressive Language Models (ALMs): Complete the sentence given its prefix
 - Sentence completion



- Autoregressive Language Models (ALMs): Complete the sentence given its prefix
 - Sentence completion



- Autoregressive Language Models (ALMs): Complete the sentence given its prefix
 - Sentence completion



 Masked Language Models (MLMs): Use the unmasked words to predict the masked word



- Training a langauge model is self-supervised learning
- Self-supervised learning :Predicting any part of the input from any other part



• Transformer-based ALMs: Composed of stacked layers of transformer layers





- Training a langauge model is self-supervised learning
- Self-supervised learning :Predicting any part of the input from any other part



• Transformer-based PLMs: Composed of stacked layers of transformer layers





- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - Pre-training/預訓練: Using a large corpora to train a neural language model
 - Autoregressive pre-trained: GPT 系列 (GPT, GPT-2, GPT-3)
 - MLM-based pre-trained: BERT 系列 (BERT, RoBERTa, ALBERT)

The Free Encyclopedia



 We believe that after pre-training, the PLM learns some knowledge, encoded in its hidden representations, that can transfer to downstream tasks



- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - (Standard) fine-tuning/微調: Using the pre-trained weights of the PLM to initialize a model for a **downstream task**



- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - (Standard) fine-tuning/微調: Using the pre-trained weights of the PLM to initialize a model for a **downstream task**



- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - Fine-tuning PLMs on downstream tasks achieves exceptional performance on many kinds of downstream tasks

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	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - PLMs are widely applied to many different scenarios in different realms

KFU NLP Team at SMM4H 2019 Tasks. Want to Extract Adverse Drugs F A Simple Cross-Lingual Lemmatization and Morphology Tagging with Two-Stage Multilingual BERT Fine-Tuning

Zulfat Miftahu

S

Kazar TMU Transformer System Using BERT for Re-ranking at BEA 2019 Transfel Yasuh Grammatical Error Correction on Restricted Track

Incorporating medical knowledge in BERT for clinical relation extraction

machi

Arpita Roy and Shimei Pan Department of Information Systems University of Maryland, Baltimore County Maryland, USA {arpita2, shimei}@umbc.edu

- Pre-trained Language Models (PLMs)/ 預訓練語言模型
 - PLMs are every where



(a) The number of publications on "language models" and their citations in recent years.

- PLMs has shown great success on a variety of benchmark datasets in NLP
- The next goal is to make PLMs fit in real-life use case
 - How unrealistic is PLMs nowadays?

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
- Closing Remarks

- Problem 1: Data scarcity in downstream tasks
- A large amount of labeled data is not easy to obtain for each downstream task



• Problem 2: The PLM is too big, and they are still getting bigger



(b) The model size and data size applied by recent NLP PTMs.

Han, Xu, et al. "Pre-trained models: Past, present and future." AI Open 2 (2021): 225-250.

- Problem 2: The PLM is too big
 - Need a copy for each downstream task



- Problem 2: The PLM is too big
 - Inference takes too long
 - Consume too much space



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic → Reducing the Number of Parameters
- Closing Remarks

- Problem 1: Data scarcity in downstream tasks
- A large amount of labeled data is not easy to obtain for each downstream task



Pre-trained Language Model (Fine-tuning)

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - Prompt Tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

- [CLS] Jack likes dog. [SEP] Jack loves ice cream. [SEP] >>> neutral
- [CLS] The spring break is coming soon. [SEP] The spring break was over. [SEP] >>> contradiction
- [CLS] I am going to have dinner. [SEP] I am going to eat something. [SEP] >>> entailment
- [CLS] Mary likes pie. [SEP] Mary hates pie. [SEP] >>> ?



Natural language inference

- [CLS] The spring break is coming soon. [SEP] The spring break was over. [SEP] >>> contradiction
- [CLS] I am going to have dinner. [SEP] I am going to eat something. [SEP] >>> entailment
- [CLS] Mary likes pie. [SEP] Mary hates pie. [SEP] >>> ?

?????

Natural language inference

Data-Efficient Fine-tuning: Prompt Tuning

- [CLS] The spring break is coming soon. Is it true that the spring break was over? >>> no
- [CLS] I am going to have dinner. Is it true that I am going to eat something? >>> yes
- [CLS] Mary likes pie. Is it true that Mary hates pie. [SEP]
 >>> ?



Natural language inference
- By converting the data points in the dataset into natural language prompts, the model may be easier to know what it should do
- [CLS] The spring break is coming soon.
 [SEP] The spring break was over. [SEP] >>>
 contradiction
- [CLS] I am going to have dinner. [SEP] I am going to eat something. [SEP] >>> entailment
- [CLS] Mary likes pie. [SEP] Mary hates pie.
 [SEP] >>> ?

- [CLS] The spring break is coming soon.
 Is it true that the spring break was over? >>> no
- [CLS] I am going to have dinner. Is it true that I am going to eat something?
 >> yes
- [CLS] Mary likes pie. Is it true that Mary hates pie. [SEP] >>> ?

• Format the downstream task as a language modelling task with predefined templates into natural language prompts



https://dictionary.cambridge.org/zht/%E8%A9%9E%E5%85%B8/%E8%8B%B1%E8%AA%9E-%E6%BC%A2%E8%AA%9E-%E7%B9%81%E9%AB%94/prompt



• What you need in prompt tuning

1. <u>A prompt template</u>: convert data points into a natural language prompt









• Prompt tuning





* I omit the [CLS] at the beginning and the [SEP] at the end

- Prompt tuning has better performance under data scarcity because
 - It incorporates human knowledge
 - It introduces no new parameters



Le Scao, Teven, and Alexander M. Rush. "How many data points is a prompt worth?." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

Lets see how prompts can help us under different level of data scarcity



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - Few-shot Learning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

- Few-shot learning: We have some labeled training data
 - "Some" ≈ 10 幾筆



Labeled Training data

- Good news: GPT-3 can be used for few-shot setting
- Bad news: GPT-3 is not freely available and contains 175B parameters

Few-shot

In addition to the task <u>description</u>, the model sees a few examples of the task. No gradient updates are performed.



Brown, Tom, et al. "Language models are few-shot learners." *Advances in neural information processing systems* 33 (2020): 1877-1901.

- Can we use smaller(?) PLMs and make them to perform well in fewshot learning?
- LM-BFF: <u>b</u>etter <u>f</u>ew-shot <u>f</u>ine-tuning of <u>l</u>anguage <u>m</u>odels

¹Alternatively, language models' <u>best friends forever</u>.

Core concept: prompt + demonstration



Gao, Tianyu, Adam Fisch, and Danqi Chen. "Making Pre-trained Language Models Better Few-shot Learners." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 2021.

• LM-BFF

K=

16

- Prompt tuning: No new parameters are introduced during fine-tuning
- Automatic template searching

	MNLI (acc)	MNLI-mm (acc)	SNLI (acc)	QNLI (acc)	RTE (acc)	MRPC (F1)	QQP (F1)	STS-B (Pear.)
Majority [†]	32.7	33.0	33.8	49.5	52.7	81.2	0.0	-
Prompt-based zero-shot [‡]	50.8	51.7	49.5	50.8	51.3	61.9	49.7	-3.2
Fine-tuning	45.8 (6.4)	47.8 (6.8)	48.4 (4.8)	60.2 (6.5)	54.4 (3.9)	76.6 (2.5)	60.7 (4.3)	53.5 (8.5)
Prompt-based FT (man)	68.3 (2.3)	70.5 (1.9)	77.2 (3.7)	64.5 (4.2)	69.1 (3.6)	74.5 (5.3)	65.5 (5.3)	71.0 (7.0)
+ demonstrations	70.7 (1.3)	72.0 (1.2)	79.7 (1.5)	69.2 (1.9)	68.7 (2.3)	77.8 (2.0)	69.8 (1.8)	73.5 (5.1)
Prompt-based FT (auto)	68.3 (2.5)	70.1 (2.6)	77.1 (2.1)	68.3 (7.4)	73.9 (2.2)	76.2 (2.3)	67.0 (3.0)	75.0 (3.3)
+ demonstrations	70.0 (3.6)	72.0 (3.1)	77.5 (3.5)	68.5 (5.4)	71.1 (5.3)	78.1 (3.4)	67.7 (5.8)	76.4 (6.2)
Fine-tuning (full) [†]	89.8	89.5	92.6	93.3	80.9	91.4	81.7	91.9

Gao, Tianyu, Adam Fisch, and Danqi Chen. "Making Pre-trained Language Models Better Few-shot Learners." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 2021.

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - Semi-supervised Learning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

• Semi-Supervised learning: We have some labeled training data and a large amount of unlabeled data



Pre-trained Language Model

• Semi-Supervised learning: We have some labeled training data and a large amount of unlabeled data

It's Not Just Size That <u>Matters:</u> Small Language Models Are Also Few-Shot Learners

Timo Schick 1,2 and Hinrich Schütze 1

¹ Center for Information and Language Processing, LMU Munich, Germany ² Sulzer GmbH, Munich, Germany

timo.schick@sulzer.de

Schick, Timo, and Hinrich Schütze. "It's Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.* 2021.

- Pattern-Exploiting Training (PET)
 - Step 1: Use different prompts and verbalizer to prompt-tune different PLMs on the labeled dataset



- Pattern-Exploiting Training (PET)
 - Step 2: Predict the unlabeled dataset and combine the predictions from different models



- Pattern-Exploiting Training (PET)
 - Step 3: Use a PLM with classifier head to train on the soft-labeled data set



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - Zero-shot Learning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

- Zero-shot inference: inference on the downstream task without any training data
- If you don't have training data, then we need a model that can zeroshot inference on downstream tasks





- GPT-3 shows that zero-shot (with task description) is possible
 - Only if your model is large enough

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.





- Where does this zero-shot ability spring from?
 - Hypothesis: during pre-training, the training datasets implicitly contains a mixture of different tasks
 - QA

Q: I got 4 papers. Should I expect this load in the future?

A: The average monthly load for reviewers should be much closer to 2, but in certain periods (close to large conferences), it's possible that the load is higher.

• Summarization

Finetuned Language Models are Zero-Shot Learners 🛛 🔤

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, Quoc V Le

29 Sept 2021 (modified: 10 Feb 2022) ICLR 2022 Oral Readers: 🚱 Everyone Show Bibtex Show Revisions

Keywords: natural language processing, zero-shot learning, language models

Abstract: This paper explores a simple method for improving the zero-shot learning abilities of language models. We show that instruction tuning—finetuning language models on a collection of datasets described via instructions—substantially improves zero-shot performance on unseen tasks. We take a 137B parameter pretrained language model and instruction tune it on over 60 NLP datasets verbalized via natural language instruction templates. We evaluate this instruction-tuned model, which we call FLAN, on unseen task types. FLAN substantially improves the performance of its unmodified counterpart and surpasses zero-shot 175B GPT-3 on 20 of 25 datasets that we evaluate. FLAN even outperforms few-shot GPT-3 by a large margin on ANLI, RTE, BoolQ, AI2-ARC, OpenbookQA, and StoryCloze. Ablation studies reveal that number of finetuning datasets, model scale, and natural language instructions are key to the success of instruction tuning.

One-sentence Summary: "Instruction tuning", which finetunes language models on a collection of tasks described via instructions, substantially boosts zero-shot performance on unseen tasks.

- Hypothesis: multi-task training enables zero-shot generalization
 - Why not train a model with multi-task learning on a bunch of dataset?





 Fine-tuning with some types of tasks and zero-shot inference on other types of tasks



• Sometimes achieves performance better than GPT-3 (175B parameters) with *only 11B* parameters



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - Summary
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

Data-Efficient Fine-tuning: Summary

• Use natural language prompts and add scenario-specific designs



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
- Closing Remarks

The Problems of PLMs

- Problem 2: The PLM is too big
 - Need a copy for each downstream task



The Problems of PLMs

- Problem 2: The PLM is too big
 - Inference takes too long
 - Consume too much space



Reducing the Number of Parameters

- Problem: PLM is too large (in terms of numbers of parameters, model size, and the storage needed to store the model)
- Solution: Reduce the number of parameters
 - Smaller pre-trained model?

Reducing the Number of Parameters

• Pre-train a large model, but use a smaller model for the downstream tasks


Reducing the Number of Parameters

• Share the parameters among the transformer layers



Lan, Zhenzhong, et al. "ALBERT: A Lite BERT for Self-supervised Learning of Language Representations." *International Conference on Learning Representations*. 2019.

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Parameter-Efficient Fine-tuning
- Closing Remarks

• Use a small amount of parameters for each downstream task



• Use a small amount of parameters for each downstream task



- What is standard fine-tuning really doing?
 - Modify the <u>hidden representations</u> (h) of the PLM such that it can perform well on downstream task



He, Junxian, et al. "Towards a unified view of parameter-efficient transfer learning." *arXiv preprint arXiv:2110.04366* (2021).

- What is standard fine-tuning really doing?
 - Modify the <u>hidden representations</u> (h) of the PLM such that it can perform well on downstream task



He, Junxian, et al. "Towards a unified view of parameter-efficient transfer learning." *arXiv preprint arXiv:2110.04366* (2021).

preprint arXiv:2110.04366 (2021).

Fine-tuning = modifying the hidden representation based on a PLM
 Before Fine-tuning
 After Fine-tuning



85

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Parameter-Efficient Fine-tuning
 - Adapter
- Closing Remarks

preprint arXiv:2110.04366 (2021).

Use special submodules to modify hidden representations!
 Before Fine-tuning
 After Fine-tuning



Parameter-Efficient Fine-tuning: Adapter

PMLR, 2019.

• Adapters: small trainable submodules inserted in transformers





• Adapters: During fine-tuning, only update the adpaters and the classifier head



Houlsby, Neil, et al. "Parameter-efficient transfer learning for NLP." *International Conference on Machine Learning*. PMLR, 2019.

• Adapters: All downstream tasks share the PLM; the adapters in each layer and the classifier heads are the task-specific modules



Houlsby, Neil, et al. "Parameter-efficient transfer learning for NLP." *International Conference on Machine Learning*. PMLR, 2019.

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Parameter-Efficient Fine-tuning
 - LoRA
- Closing Remarks

• Use special submodules to modify hidden representations!



He, Junxian, et al. "Towards a unified view of parameter-efficient transfer learning." *arXiv preprint arXiv:2110.04366* (2021).

Parameter-Efficient Fine-tuning: LoRA

• LoRA: Low-Rank Adaptation of Large Language Models



• LoRA





Hu, Edward J., et al. "Lora: Low-rank adaptation of large language models." arXiv preprint arXiv:2106.09685 (2021).

• LoRA: All downstream tasks share the PLM; the LoRA in each layer and the classifier heads are the task-specific modules



Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Parameter-Efficient Fine-tuning
 - Prefix Tuning
- Closing Remarks

• Use special submodules to modify hidden representations!



He, Junxian, et al. "Towards a unified view of parameter-efficient transfer learning." *arXiv preprint arXiv:2110.04366* (2021).

• What is *prefix*

prefix
<i>noun</i> [C]
UK ◀》 /ˈpriː.fiks/ US ◀》 /ˈpriː.fiks/
prefix noun [C] (GRAMMAR) Guide word: helps you find the right meaning when a word has more than one meaning
B2 LANGUAGE
a letter or group of letters added to the boginning of a word to make a new word
beginning of a word to make a new word
前綴

• 放在某個東西前面的東西

• Prefix Tuning: Insert trainable prefix in each layer



Li, Xiang Lisa, and Percy Liang. "Prefix-Tuning: Optimizing Continuous Prompts for Generation." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).* 2021.



Standard Self-Attention



 Prefix Tuning: Only the prefix (key and value) are updated during finetuning





Li, Xiang Lisa, and Percy Liang. "Prefix-Tuning: Optimizing Continuous Prompts for Generation." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).* 2021.

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Parameter-Efficient Fine-tuning
 - Soft Prompting
- Closing Remarks

Parameter-Efficient Fine-tuning: Soft Prompting

- Soft Prompting
 - Prepend the prefix embedding at the input layer



Tuning." Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. 2021.

Parameter-Efficient Fine-tuning: Soft Prompting

- Soft Prompting can be considered as the *soften* version of prompting
 - (Hard) prompting: add words in the input sentence (fine-tune the model while fixing the prompts)



Lester, Brian, Rami Al-Rfou, and Noah Constant. "The Power of Scale for Parameter-Efficient Prompt Tuning." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

• Soft Prompts: vectors (can be initialized from some word embeddings)



• Hard Prompts: words (that are originally in the vocabulary)



Lester, Brian, Rami Al-Rfou, and Noah Constant. "The Power of Scale for Parameter-Efficient Prompt Tuning." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

• Benefit 1: Drastically decreases the task-specific parameters

	Adapter	LoRA	Prefix Tuning	Soft Prompt	
Task-specific parameters*	$\Theta(d_{model}rL)$	$\Theta(d_{model}rL)$	Θ(<mark>d_{model}nL</mark>)	$\Theta(d_{model}n)$	
Percent Trainable	<5%	<0.1%	<0.1%	<0.05%	
Illustration	+ r Nonlinearity r	r r	<i>n</i> :Prefix length $\boldsymbol{k}_{p_1} \boldsymbol{v}_{p_1} \boldsymbol{k}_{p_n} \boldsymbol{v}_{p_n}$	n:Prefix length	

*not including the classifier head

 Benefit 2: Less easier to overfit on training data; better out-of-domain performance

In domain	Dataset	Domain	Model	Soft Prompt	Δ
dataset	SQuAD	Wiki	94.9 ± 0.2	94.8 ± 0.1	-0.1
	TextbookQA BioASO	Book Bio	54.3 ± 3.7	66.8 ±2.9 79 1 ±0 3	+12.5
OOD test dataset	RACE	Exam	59.8 ± 0.6	60.7 ±0.5	+1.2+0.9
	RE	Wiki Movie	88.4 ± 0.1	88.8 ±0.2	+0.4
	DROP	Wiki	68.9 ±0.7 68.9 ±1.7	67.1 ± 1.9	-1.2 -1.8

• Benefit 3: Fewer parameters to fine-tune; a good candidate when training with small dataset

	low-resource				high-resource			
Model	CHEMPROT	ACL-ARC	SCIERC	HYP.	RCT	AGNEWS	HELPFUL.	IMDB
	(4169)	(1688)	(3219)	(515)	(180k)	(115k)	(115k)	(20k)
$RoBaft^{\dagger}$	81.9 _{1.0}	63.0 _{5.8}	77.3 _{1.9}	86.6 _{0.9}	87.2 _{0.1}	93.9 _{0.2}	65.1 _{3.4}	95.0 _{0.2}
RoBaft*	81.7 _{0.8}	$65.0_{3.6}$	$78.5_{1.8}$	88.9 _{3.3}	$87.0_{0.1}$	93.7 _{0.2}	69.1 _{0.6}	$95.2_{0.1}$
RoBaadapter ₂₅₆	82.9 _{0.6}	$67.5_{4.3}$	$80.8_{0.7}$	$90.4_{4.2}$	87.10.1	$93.8_{0.1}$	$69.0_{0.4}$	$95.7_{0.1}$

He, Ruidan, et al. "On the Effectiveness of Adapter-based Tuning for Pretrained Language Model Adaptation." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).* 2021.

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Early Exit
- Closing Remarks

Early Exit

- Problem 1: The PLM is too big
 - Inference takes too long



Early Exit

- Inference using the whole model takes too long
- Simpler data may require lesser effort to obtain the answer
- Reduce the number of layers used during inference


Early Exit

• Add a classifier at each layer



Xin, Ji, et al. "BERxiT: Early Exiting for BERT with Better fine-tuning and extension to regression." *Proceedings of the 16th conference of the European chapter of the association for computational linguistics: Main Volume*. 2021

Early Exit

• How do we know which classifier to use?



Xin, Ji, et al. "BERxiT: Early Exiting for BERT with Better fine-tuning and extension to regression." *Proceedings of the 16th conference of the European chapter of the association for computational linguistics: Main Volume*. 2021



conference of the European chapter of the association for computational linguistics: Main Volume. 2021

117

Early Exit

• Early exit reduces the inference time while keeping the performance

	RTE	MRPC	SST-2	QNLI	QQP	MNLI-(m/mm)	STS-B
	Score Layer	Score Layer	Score Layer	Score Layer	Score Layer	Score Layer	Score Layer
				BERT _{BAS}	Е		
RAW	66.4 12	88.9 12	93.5 12	90.5 12	71.2 12	84.6/83.4 12	85.8 12
ALT	$\begin{array}{r} 101\% -44\% \\ 99\% -54\% \\ 96\% -64\% \end{array}$	99% -30% 97% -56% 94% -74%	98% -65% 96% -79% 94% -87%	99% -42% 98% -63% 95% -71%	99% -56% 97% -75% 93% -84%	99%/99% -37% 97%/97% -57% 93%/92% -72%	$\begin{array}{rrrr} 95\% & -50\% \\ 91\% & -67\% \\ 85\% & -75\% \end{array}$

	BERTLARGE							
RAW	70.1 24	89.3 24	94.9 24	92.7 24	72.1 24	86.7/85.9 24	86.5	24
	95% -33%	99% -32%	100% -32%	97% -62%	98% -74%	99%/99% -36%	97%	-39%
ALT	94% - 46%	98% - 46%	99% -61%	95% -73%	96% - 82%	96%/97% -57%	90%	-62%
	88% -62%	94% -71%	96% -78%	91% -83%	91% -89%	90%/90% -75%	76%	-80%

Xin, Ji, et al. "BERxiT: Early Exiting for BERT with Better fine-tuning and extension to regression." *Proceedings of the 16th conference of the European chapter of the association for computational linguistics: Main Volume*. 2021

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
 - Labeled Data Scarcity \rightarrow Data-Efficient Fine-tuning
 - PLMs Are Gigantic \rightarrow Reducing the Number of Parameters
 - Summary
- Closing Remarks

Reducing the Number of Parameters: Summary

- Parameter-efficient fine-tuning: Reduce the task-specific parameters in downstream task
- Early exit: Reduce the models that are involved during inference

Outline

- Background knowledge
- The Problems of PLMs
- The Solutions of Those Problems
- Closing Remarks

Closing Remarks

- What we address in this lecture
 - Making PLM smaller, faster, and more parameter-efficient
 - Deploying PLMs when the labeled data in the downstream task is scarce

Data-Efficient Fine-tuning: Prompt Tuning

• Prompt tuning





* I omit the [CLS] at the beginning and the [SEP] at the end

Schick, Timo, and Hinrich Schütze. "Exploiting Cloze-Questions for Few-Shot Text Classification and Natural Language Inference." *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 2021.

Parameter-Efficient Fine-tuning

• Benefit 1: Drastically decreases the task-specific parameters

	Adapter	LoRA	Prefix Tuning	Soft Prompt
Task-specific parameters*	$\Theta(d_{model}rL)$	$\Theta(d_{model}rL)$	Θ(<mark>d_{model}nL</mark>)	$\Theta(d_{model}n)$
Percent Trainable	<5%	<0.1%	<0.1%	<0.05%
Illustration	+ r Nonlinearity r	r r	<i>n</i> :Prefix length $\boldsymbol{k}_{p_1} \boldsymbol{v}_{p_1} \boldsymbol{k}_{p_n} \boldsymbol{v}_{p_n}$	n:Prefix length

*not including the classifier head

Closing Remarks

- What we address in this lecture
 - Making PLM smaller, faster, and more parameter-efficient
 - Deploying PLMs when the labeled data in the downstream task is scarce
- The problems are not completely solved yet
- The problems we discuss are just a small part of problems of PLMs
 - Why does self-supervised pre-training work
 - Interpretability of the model's prediction
 - Domain adaptation
 - Continual learning/lifelong learning
 - Security and privacy

To Learn More

- AACL-IJCNLP 2022 Tutorial (11.24.2022)
 - Recent Advances in Pre-trained Language Models: Why Do They Work and How Do They Work.

Cheng-Han Chiang National Taiwan University dcml0714@gmail.com Yung-Sung Chuang CSAIL, MIT yungsung@mit.edu Hung-yi Lee National Taiwan University hungyilee@ntu.edu.tw

