

LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

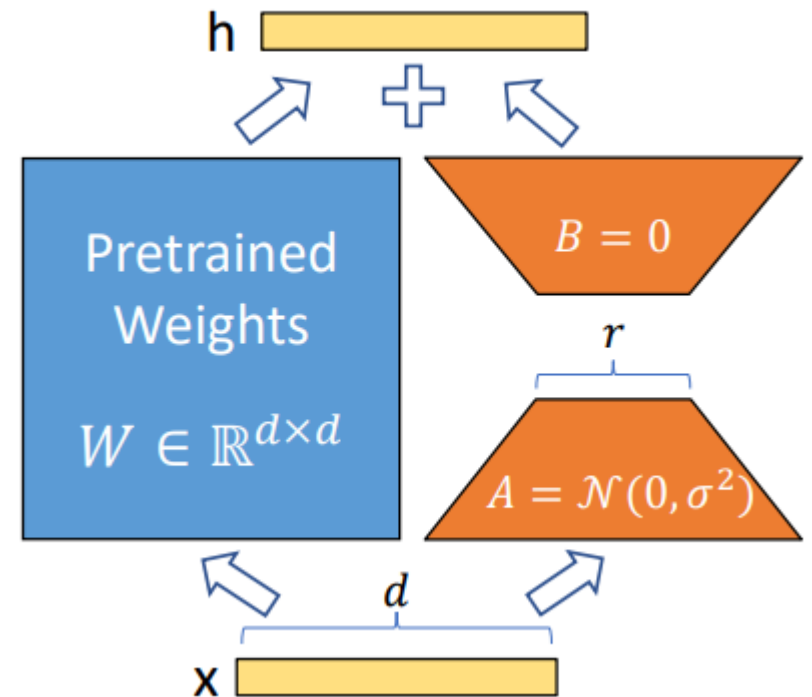
Zijiang Yang

Background

- Personalized model: heterogenous distribution/ multiple downstream tasks
- Full Finetuning: updating all parameters of pre-trained model
 - Pretrained model $P_{\Phi_0}(y|x)$ with parameter Φ_0 , update Φ_0 to $\Phi_0 + \Delta\Phi$
 - Dimension of $\Delta\Phi$ = Dimension of Φ , too many parameters to update
 - Learn a different $\Delta\Phi$ for each downstream task
- Existing solutions
 - Adding adapter layers: introduce inference latency
 - Prompt tuning: hard to directly optimize the prompt

LORA

- Inspiration: The change in weights ΔW has a low intrinsic rank. We can train some dense layers by optimizing rank decomposition matrices of the dense layers' change.
- Update parameters:
 - $W_0 + \Delta W = W_0 + BA$
 - $B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}, r \ll \min(d, k)$
 - $O(r(d + k)) \ll O(dk)$
- Training A and B
 - Initialization: A= random Gaussian, B=zeros
 - Forward pass: $h = W_0x + \Delta Wx = W_0x + BAx$
 - Use Adam Adapter to train A and B



Applying LORA to Transformer

- Only adapting the attentions weights (W_q, W_k, W_v, W_o) for downstream tasks, freeze MLP layers.
- Benefits:
 - Reduce memory and storage usage (VRAM, checkpoint size, GPUs)
 - Easy to switch between tasks by only changing LORA parameters instead of full parameters
- Limitations:
 - If we choose to absorb A and B into W to eliminate inference latency, it is not straight forward to batch inputs to different tasks with different A and B

Experiments

- Evaluate downstream task performance of Lora and other adaptation method on different models
- Lora trains fewer parameters without sacrificing too much accuracy

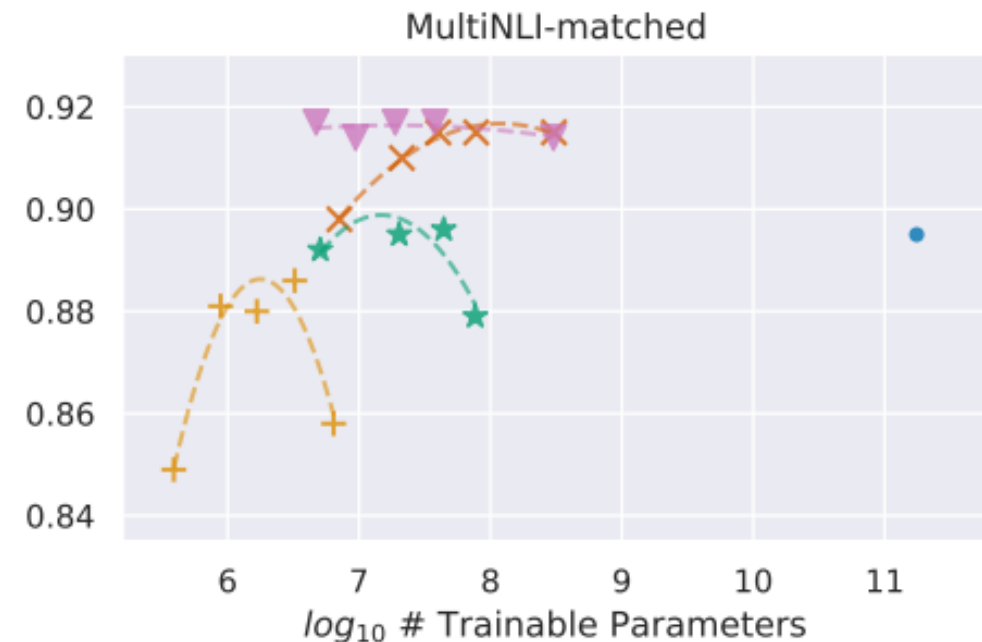
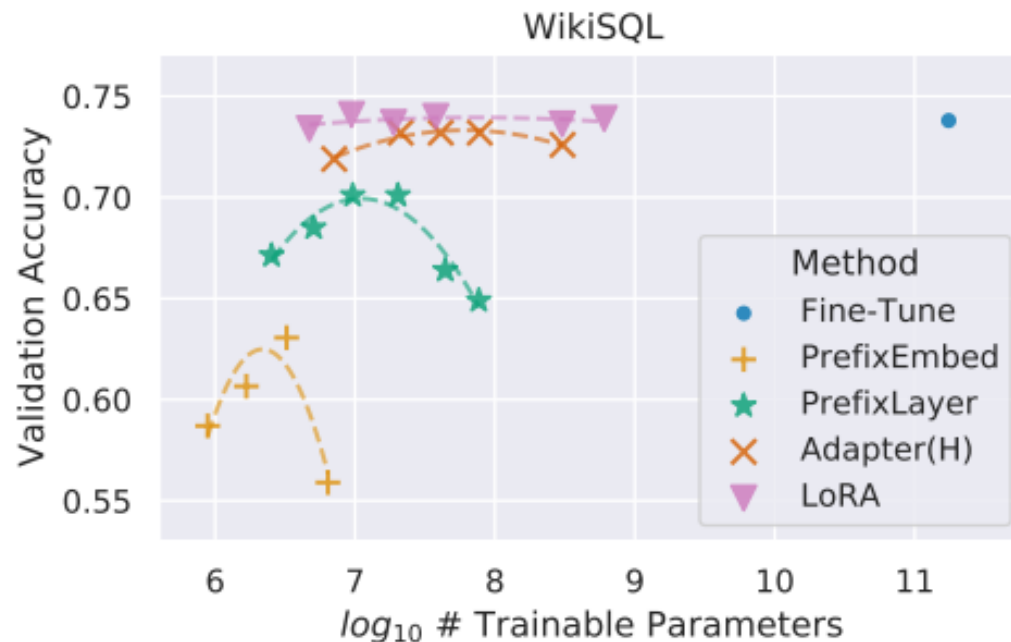
Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB _{base} (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
RoB _{base} (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
RoB _{base} (Adpt ^D)*	0.3M	87.1 \pm .0	94.2 \pm .1	88.5 \pm 1.1	60.8 \pm .4	93.1 \pm .1	90.2 \pm .0	71.5 \pm 2.7	89.7 \pm .3	84.4
RoB _{base} (Adpt ^D)*	0.9M	87.3 \pm .1	94.7 \pm .3	88.4 \pm .1	62.6 \pm .9	93.0 \pm .2	90.6 \pm .0	75.9 \pm 2.2	90.3 \pm .1	85.4
RoB _{base} (LoRA)	0.3M	87.5 \pm .3	95.1\pm.2	89.7 \pm .7	63.4 \pm 1.2	93.3\pm.3	90.8 \pm .1	86.6\pm.7	91.5\pm.2	87.2
RoB _{large} (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB _{large} (LoRA)	0.8M	90.6\pm.2	96.2 \pm .5	90.9\pm1.2	68.2\pm1.9	94.9\pm.3	91.6 \pm .1	87.4\pm2.5	92.6\pm.2	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2 \pm .3	96.1 \pm .3	90.2 \pm .7	68.3\pm1.0	94.8\pm.2	91.9\pm.1	83.8 \pm 2.9	92.1 \pm .7	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5\pm.3	96.6\pm.2	89.7 \pm 1.2	67.8 \pm 2.5	94.8\pm.3	91.7 \pm .2	80.1 \pm 2.9	91.9 \pm .4	87.9
RoB _{large} (Adpt ^H)†	6.0M	89.9 \pm .5	96.2 \pm .3	88.7 \pm 2.9	66.5 \pm 4.4	94.7 \pm .2	92.1 \pm .1	83.4 \pm 1.1	91.0 \pm 1.7	87.8
RoB _{large} (Adpt ^H)†	0.8M	90.3 \pm .3	96.3 \pm .5	87.7 \pm 1.7	66.3 \pm 2.0	94.7 \pm .2	91.5 \pm .1	72.9 \pm 2.9	91.5 \pm .5	86.4
RoB _{large} (LoRA)†	0.8M	90.6\pm.2	96.2 \pm .5	90.2\pm1.0	68.2 \pm 1.9	94.8\pm.3	91.6 \pm .2	85.2\pm1.1	92.3\pm.5	88.6
DeB _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB _{XXL} (LoRA)	4.7M	91.9\pm.2	96.9 \pm .2	92.6\pm.6	72.4\pm1.1	96.0\pm.1	92.9\pm.1	94.9\pm.4	93.0\pm.2	91.3

Experiments

Model & Method	# Trainable Parameters	E2E NLG Challenge				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 M (FT)*	354.92M	68.2	8.62	46.2	71.0	2.47
GPT-2 M (Adapter ^L)*	0.37M	66.3	8.41	45.0	69.8	2.40
GPT-2 M (Adapter ^L)*	11.09M	68.9	8.71	46.1	71.3	2.47
GPT-2 M (Adapter ^H)	11.09M	67.3 \pm .6	8.50 \pm .07	46.0 \pm .2	70.7 \pm .2	2.44 \pm .01
GPT-2 M (FT ^{Top2})*	25.19M	68.1	8.59	46.0	70.8	2.41
GPT-2 M (PreLayer)*	0.35M	69.7	8.81	46.1	71.4	2.49
GPT-2 M (LoRA)	0.35M	70.4\pm.1	8.85\pm.02	46.8\pm.2	71.8\pm.1	2.53\pm.02
GPT-2 L (FT)*	774.03M	68.5	8.78	46.0	69.9	2.45
GPT-2 L (Adapter ^L)	0.88M	69.1 \pm .1	8.68 \pm .03	46.3 \pm .0	71.4 \pm .2	2.49\pm.0
GPT-2 L (Adapter ^L)	23.00M	68.9 \pm .3	8.70 \pm .04	46.1 \pm .1	71.3 \pm .2	2.45 \pm .02
GPT-2 L (PreLayer)*	0.77M	70.3	8.85	46.2	71.7	2.47
GPT-2 L (LoRA)	0.77M	70.4\pm.1	8.89\pm.02	46.8\pm.2	72.0\pm.2	2.47 \pm .02

Experiments

- Scaling up to GPT-3
- LoRA matches or exceeds the fine-tuning baseline on all three datasets
- Not all methods benefit monotonically from having more trainable parameters



Which weight matrices in Transformer should we apply LORA to?

- Even we only adapt Lora to attention heads, there are still 4 kinds of weights to be finetuned. Which one is more appropriate?

	# of Trainable Parameters = 18M						
Weight Type Rank r	W_q 8	W_k 8	W_v 8	W_o 8	W_q, W_k 4	W_q, W_v 4	W_q, W_k, W_v, W_o 2
WikiSQL ($\pm 0.5\%$)	70.4	70.0	73.0	73.2	71.4	73.7	73.7
MultiNLI ($\pm 0.1\%$)	91.0	90.8	91.0	91.3	91.3	91.3	91.7

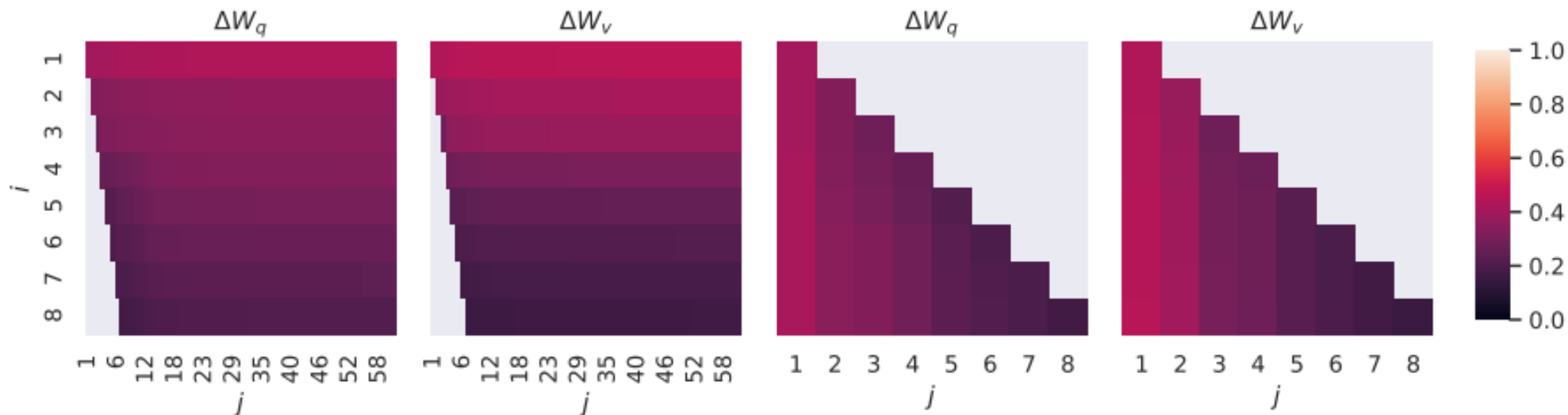
What is the optimal rank r for LORA?

- Adapting W_q and W_v can achieve very high accuracy with rank $r=1$, but training W_q alone needs much larger r .
- Increasing r does not cover a more meaningful subspace
 - > Updating matrix ΔW could have a very small “intrinsic rank”

	Weight Type	$r = 1$	$r = 2$	$r = 4$	$r = 8$	$r = 64$
WikiSQL($\pm 0.5\%$)	W_q	68.8	69.6	70.5	70.4	70.0
	W_q, W_v	73.4	73.3	73.7	73.8	73.5
	W_q, W_k, W_v, W_o	74.1	73.7	74.0	74.0	73.9
MultiNLI ($\pm 0.1\%$)	W_q	90.7	90.9	91.1	90.7	90.7
	W_q, W_v	91.3	91.4	91.3	91.6	91.4
	W_q, W_k, W_v, W_o	91.2	91.7	91.7	91.5	91.4

Subspace similarity between different r

- We do singular value decomposition to $A_{r=8}$ and $A_{r=64}$, and get $U_{r=8}$, $U_{r=64}$. Then we measure the normalized subspace similarity based on Grassmann distance.
- The top singular-vector directions of $A_{r=8}$ and $A_{r=64}$ are the most useful



How does the adaptation matrix ΔW compare to W ?

- We project W onto the r -dimensional subspace of ΔW by computing $U^T W V^T$
- ΔW has a stronger correlation with W compared to a random matrix
 - $\rightarrow \Delta W$ enhance some features of W
- ΔW only amplifies directions that are not emphasized in W instead of repeating the top singular directions
- The amplification factor is huge

	$r = 4$			$r = 64$		
	ΔW_q	W_q	Random	ΔW_q	W_q	Random
$\ U^T W_q V^T\ _F =$	0.32	21.67	0.02	1.90	37.71	0.33
$\ W_q\ _F = 61.95$	$\ \Delta W_q\ _F = 6.91$			$\ \Delta W_q\ _F = 3.57$		

Discussion

- Are LoRA adapters more robust to malicious client updates compared to full-model FL?
- If we apply Lora to federated learning, will the heterogeneity of data among clients undermine the "low-rank assumption" of LoRA?
- Does LoRA, by only transmitting low-rank parameters, provide better privacy protection?