

LORA: Low-Rank Adaptation of Large Language Models

Nguyễn Thị Ngọc Uyển Phạm Thị Cẩm

Khoa Toán - C - Tin học
Trường Đại học Khoa học Tự nhiên, HQGHN
Giảng viên hướng dẫn: PGS.TS Trần Trung Hieu

Hà Nội, 2025

Mc Ic

- 1 Tng quan
- 2 Phng phòp LoRA
- 3 Thc nghim vị ònh giò
- 4 Kt lun

Mc Ic

- 1 Tng quan
- 2 Phng phòp LoRA
- 3 Thc nghim vị ònh giò
- 4 Kt lun

Gii thiu

- Còc mĩ hnh ngũn ng ln (LLMs) nh GPT, BERT, LLaMA, T5 t nɦu thnh tu vt trı.
- Tuy nhiõn, vic tnh chnh toın b mĩ hnh (Full Fine-tuning) ùi hi:
 - Tị nguyõn tònھ toın ln (GPU, b nh)
 - Thı gian hun luyõn dı
 - Nguy c *overfitting*
- Gii phòp: **Parameter-Efficient Fine-Tuning (PEFT)** ch tnh chnh mt phn nh tham s.

ng c nghiên cứu

- **Thòch thc:**

- Full fine-tuning tiúu tn hịng trm GB b nh GPU.
- Mi tòc v cn mt bn tinh chnh riúng bit.

- **Gii phòp:** PEFT cho phỏp:

- Gi nguyên mủ hnh gc, ch hc vị tham s mi.
- D chia s, tài s dng mủ hnh.

- **LoRA (Low-Rank Adaptation):** gim ti **10.000** ln s tham s mị vn gi hiu nng tng ng.

Kin thc nn tng

- **Transformer**: kin trũc da trũn c ch *self-attention*, giúp mĩ hnh hc quan h gia còc t song song.
- **Full Fine-tuning**: cp nht toin b trng s mĩ hnh hiu nng cao nhng chi phò ln.
- **PEFT**: ch tinh chnh còc mĩ-un nh, giúp gim chi phò hun luyt mĩ vn duy trø hiu nng.

Còc phng phòp PEFT ph bin

- **Adapter Tuning:** thỏm tng ph nh gia còc tng Transformer.
- **Prefix/Prompt Tuning:** hc thỏm chui tham s c chỏn vọ u vọ.
- **BitFit:** ch tinh chnh còc tham s bias.
- **Compacter:** biu din hng thp cho adapter.
- **LoRA:** òp dng rịng buc hng thp lỏn ma trn trng s.

Mc Ic

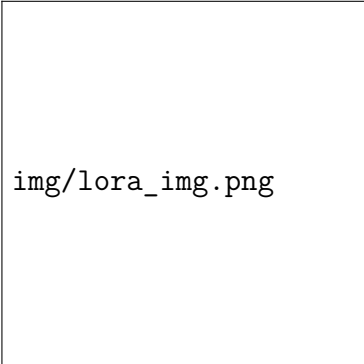
Nguyên lý hoạt động của LoRA

- Quan sát: thay đổi trọng số khi fine-tuning chỉ nằm trong không gian con hàng thấp.
- LoRA giới thiệu tham số hóa ma trận trọng số:

$$W = W_0 + BA, \quad B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}, \quad r \ll \min(d, k)$$

- Giới thiệu tham số huấn luyện $d \times k$ xuống $r(d + k)$.
- Với $r = 8$, giới thiệu 500 lần tham số cần tính chính xác.

Minh ha c ch LoRA



img/lora_img.png

Hình

minh ha c ch tài tham s hóa trng s trong LoRA.

- W_0 : trng s gc c gì nguỷn.
- A, B : hai ma trn hng thp c hun luyn.
- Kt qu: ch cn hun luyn rt ờt tham s.
- Khi hun luyn xong, cú th gp $(W_0 + BA)$ đứng cho suy lun mị khũng tng tr.

So sòn LoRA vi còc phng phòp khòc

Phng phòp	Thay i kin trữc	T l tham s	tr	c im
Adapter Tuning	Cú	13%	Cú	Hiu nng cao nhng tng sốu mĩ hnh
Prefix Tuning	Khũng	<1%	Khũng	n gin, ph thuc dji prefix
BitFit	Khũng	0.1%	Khũng	Rt nh, hiu nng gii hn
Compacter	Cú	0.5%	Nh	Adapter nỏn, hiu qu cao
LoRA	Khũng	<1%	Khũng	Rịng buc hng thp, d tồch hp

u vị nhc im ca LoRA

u im:

- Hiu qu tham s cao ($<1\%$ tham s hun luy).
- Khũng tng tr suy lun.
- Hun luy nhanh hn, tit kim b nh.
- D kt hp vi Adapter hoc Prefix Tuning.

Hn ch:

- Gi nh khũng gian hng thp cú th khũng ỡng cho mi tòc v.
- Giò tr r cn c chn phứ hp.
- Vi mũ hnh nh (nh BERT-base), li òch khũng òng k.

Mc Ic

Mc tiều thc nghiim

- ònh giò LoRA so vi Full Fine-Tuning trổn bị toòn sinh ngũn ng t nhĩn.
- D liu: **E2E NLG Challenge** 50.000 cp (MR, cóu mũ t).
- Cóu hi chørnh:
 - LoRA cú t hieu nng tng ng Full Fine-Tuning?
 - Mc tit kim tị nguyổn ra sao?

Thit lp thc nghim

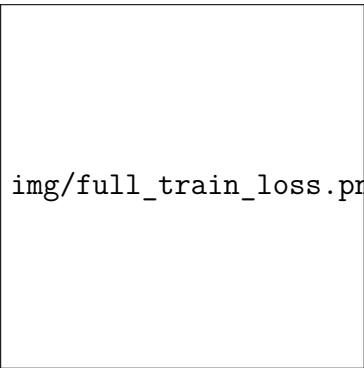
Mũ hnh: GPT-2 Medium (355M tham s).

Cu hnh:

- Optimizer: AdamW, $lr = 2e-4$, batch size = 4
- Epochs = 5, dropout = 0.1
- LoRA rank = 4, scaling factor = 32
- Ch tinh chnh **0.3% tham s mũ hnh**

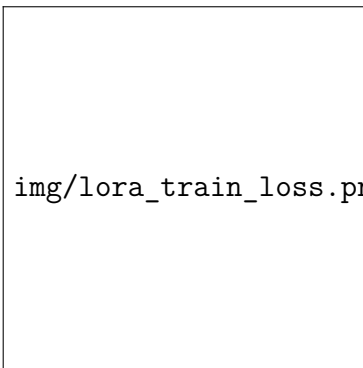
ònh giò: BLEU, NIST, METEOR, ROUGE-L, CIDEr.

So sòn Train Loss gia còc mĩ hnh



img/full_train_loss.png

Train



img/lora_train_loss.png

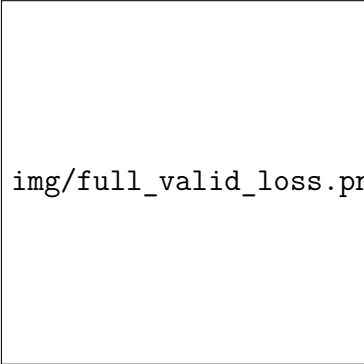
Train

Loss ca mĩ hnh Full Fine-tuning

Loss ca mĩ hnh LoRA

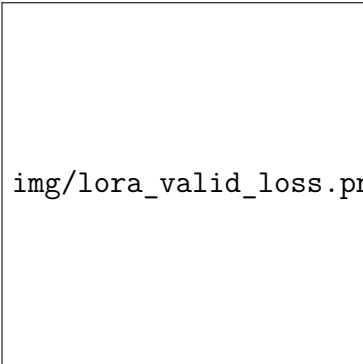
- C hai mĩ hnh u hi t sau khong 34 epoch.
- LoRA cú ng loss mt hn, gim nhanh hn giai on u.
- iu nìy cho thy vic gii hn khũng gian hng thp giúp quò trnh ti u n nh

So sánh Validation Loss giữa các mô hình



img/full_valid_loss.png

Validation Loss của mô hình Full Fine-tuning



img/lora_valid_loss.png

Validation Loss của mô hình LoRA

- Validation Loss của LoRA giảm nhanh và ổn định hơn.
- Full fine-tuning có xu hướng dao động và dễ xảy ra overfitting.
- Kết luận: LoRA giúp mô hình tăng khả năng tổng quát hóa tốt hơn trên tập kiểm tra.

Kết quả thí nghiệm

Tiêu chí	Full Fine-tuning	LoRA
BLEU	0.6734	0.6803
NIST	8.5452	8.6419
METEOR	0.4622	0.4556
ROUGE-L	0.7033	0.7018
CIDEr	2.3843	2.4465

Nhận xét: LoRA đạt hiệu suất tổng thể cao hơn Full Fine-Tuning chỉ cần huấn luyện 0.3% tham số.

Phón tồch vị nhn xố

- LoRA hi t nhanh hn, gim dao ng gradient.
- Giúp trờnh overfitting vị n nh trổn tp validation.
- Kt qu chng minh: thay i trng s quan trng ch nm trong khũng gian hng thp.

Mc Ic

Kt lun vị hng phòt trin

- LoRA lị phng phòp tinh chnh hieu qu, tit kim tị nguổn mị vn gi hieu nng cao.
- Phứ hp vi h thng cú hn ch phn cng (GPU nh).
- Hng m:
 - ng dng trổn tòc v túm tt, i thoi, truy xut thũng tin.
 - Kt hp LoRA vi Adapter hoc Prefix Tuning.
 - Tì u giò tr rank r vị h s α .

Xin chón thịnh cm n!

Cm n thy vị còc anh/ch ố lng nghe.