

Reinforcement Learning with Human Feedback (lec10) (CS224N)

Language Models can act as a way of indirectly world models. They model agents, beliefs, and actions.

It may not be to predict next text, you'd also want to predict agents' beliefs, background.

LMs as multitask assistants

zero-shot learning: no examples, no pretrain updates.

GPT2 did better than SOTA supervised +
finetuned
in LAMBADA task.

GPT3:

in-context learning: examples of the task before your
(no pretrain update) example.

Zero-shot:

Translate English to French

cheese =

One-shot: ~ 10% increase

Translate English to French

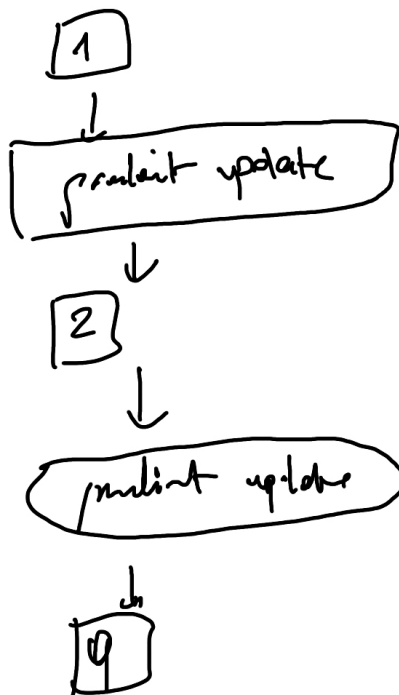
sea otter \Rightarrow loutre de mer

cheese =

Few-shot

couple of shots.

Traditional fine-tuning:



Fine-tune training:

Frozen model,
no prompt updates.

Chain of thought prompting:

Student

Q
Answer
Q
Answer

CoT:

Q
A with
steps of
reasoning.

Q
should reason
in steps +
answer

Zero-shot CoT prompting:

Q:

Answer: Let's think step by step.
followed by step by step
answering.

← still manual CoT is better though.

Language models aren't aligned with user intent.

Pretraining on language modeling.



to train finetuning (on many tasks)

to allow the model to adapt to tasks.

Instruction finetuning:

Collect examples of (ins, output) pairs across many tasks and finetune a Language Model.

evaluate on unseen tasks.

Multitask LMs.

Messine Multitask Language Understanding
benchmark

That isn't quite possible on creative generation tasks.

LM penalizes all token-level mistakes equally.
Same one wrong way than others.

for each sample from LM s , let's assume we've

$$R(s) \in \mathbb{R}$$

maximize reward.

$E_{s \sim P_{\theta}(s)}[R(s)]$ how do we choose LM params θ to maximize?

Policy gradient methods for maximizing parameters

Monte Carlo sampling:

The best way to approximate any expectation is to take samples and average them.

$$\theta_{t+1}: \theta_t + \alpha \cdot \frac{1}{m} \cdot \sum_{i=1}^m R(s_i) \nabla \theta_t \cdot \log P_{\theta_t}(s_i)$$

if R is positively increasing, take gradient steps to maximize $P_{\theta}(s_i)$

if R is negative and decreasing, take steps to minimize $P_{\theta}(s_i)$

turn-in-the-loop is expensive
so modeling preferences
is nice.

LM LM $\phi(s)$ to predict human preferences
from a dataset and optimize for BMDP

Pairwise comparisons of deciding which option is more preferable also works. (Rate ascent)

Reinl Model does come close to the single human preference if there's enough data.

$LM_{P^{PT}}(s) + RM_{\phi}(s) + \text{Policy Gradient}$

$$R(s) = RM_{\phi}(s) - \beta \cdot \underbrace{\log \left(\frac{P_{\theta}^{RL}(s)}{P^{PT}(s)} \right)}_{\text{penalty}}$$

Penalty when

$$P_{\theta}^{RL}(s) > P^{PT}(s)$$

penalty.

KL divergence
between $P_{\theta}^{RL}(s)$
and $P^{PT}(s)$

Instruct GPT, scaling up RCHF to thousands of tasks. (30k tasks)

