

$$1. \text{ a) } J_{MLE}(\theta) = E_{p(x)} E_{\pi_{\text{teacher}}} [\log \pi_{\theta}(y|x)] \\ = \int_x p(x) \int_y \pi_{\theta}(y|x) \log \pi_{\theta}(y|x) dy dx$$

$$\boxed{J'_{MLE}(\theta) = \int_x p(x) \int_y \pi_{\theta}(y|x) \frac{\partial}{\partial \theta} [\log \pi_{\theta}(y|x)] dy dx}$$

Stochastic gradient algorithm

1. Sample N datapoints (x_i, y_i)

$$\{x_i\}^N \sim P_{\text{data}}(x) \text{ and } \{y_i\}^N \sim \pi_{\theta}(y_i|x_i)$$

2. Approximate $J'_{MLE}(\theta)$ from finite sample

$$J'_{MLE}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \theta} \log \pi_{\theta}(y_i|x_i)$$

3. Update parameters

$$\theta^{(t+1)} = \theta^{(t)} + \eta J'(\theta)$$

$$b) J_{RL} = E_{p(x)} E_{\pi_{\theta}(y|x)} [r(x, y)] \quad \frac{\partial}{\partial \theta} \log \pi_{\theta} = \frac{1}{\pi_{\theta}} \frac{\partial}{\partial \theta} \pi_{\theta}$$

$$J'_{RL} = \int_x p(x) \int_y \frac{\partial}{\partial \theta} [\pi_{\theta}(y|x)] r(x, y) dy dx \quad \pi_{\theta} \frac{\partial}{\partial \theta} \log \pi_{\theta} = \frac{\partial}{\partial \theta} \pi_{\theta}$$

$$= \int_x p(x) \int_y \pi_{\theta}(y|x) \frac{\partial}{\partial \theta} [\log \pi_{\theta}(y|x)] r(x, y) dy dx$$

$$= \boxed{E_{p(x)} E_{\pi_{\theta}(y|x)} [r(x, y) \frac{\partial}{\partial \theta} \log \pi_{\theta}(y|x)]}$$

Stochastic gradient algorithm

1. Sample N datapoints (x_i, y_i)

$$\{x_i\}^N \sim P_{\text{data}}(x) \text{ and } \{y_i\}^N \sim \pi_{\theta}(y_i|x_i)$$

2. Approximate $J'_{MLE}(\theta)$ from finite sample

$$J'_{MLE}(\theta) \approx \frac{1}{N} \sum_{i=1}^N [r(x_i, y_i) \cdot \frac{\partial}{\partial \theta} \log \pi_{\theta}(y_i|x_i)]$$

3. Update parameters

$$\theta^{(t+1)} = \theta^{(t)} + \eta J'(\theta)$$

current probability distribution

What if we substitute $r(x, y) \rightarrow r(x, y) - b(x)$

$$\begin{aligned}
 J'_{RL, \text{new}} &= \int_x p(x) \int_y \frac{\partial}{\partial \theta} [\pi_\theta(y|x)] [r(x, y) - b(x)] dy dx \\
 &= J'_{RL, \text{old}} - \int_x p(x) \int_y \frac{\partial}{\partial \theta} [\pi_\theta(y|x)] b(x) dy dx \\
 &= J'_{RL, \text{old}} - \int_x p(x) b(x) \int_y \frac{\partial}{\partial \theta} [\pi_\theta(y|x)] dy dx \\
 &= J'_{RL, \text{old}} - \int_x p(x) b(x) \cancel{\frac{\partial}{\partial \theta} [1]} dx = 1 \\
 &= J'_{RL, \text{old}} \quad \boxed{\therefore J'_{RL} \text{ does not change}}
 \end{aligned}$$

If we set $r(x, y) = 1$ for J_{RL} and J'_{RL}

$$\begin{aligned}
 E_{\pi_\theta} \left[\frac{\partial}{\partial \theta} \log \pi_\theta \right] &= \int_y \pi_\theta(y|x) \cdot \frac{1}{\pi_\theta} \frac{\partial}{\partial \theta} \pi_\theta dy \\
 &= \int_y \frac{\partial}{\partial \theta} \pi_\theta dy \\
 &= \frac{\partial}{\partial \theta} \int_y \pi_\theta dy \\
 &= \frac{\partial}{\partial \theta} [1] \\
 &= \boxed{0}
 \end{aligned}$$

c) Stochastic Gradient Algorithms for MLE and RL

- Similarities
 - draw (x_i, y_i) from a joint distribution to estimate $J'(\theta)$
- Differences
 - RL includes a reward factor when calculating $J'(\theta)$
 - RL model learns from itself, so the distribution $\pi_\theta(y|x)$ from which we sample y_i changes each update

2a) Explanation of Code

We first create a custom dataset using some predefined prompts and answers generated using GPT-2. We (the user) rank the predictions generated, and this labeled data is used as the training data for training the reward model. For the preprocessing step: the token sequence inputs use padding to make them same length, and an attention map to distinguish between actual tokens and padding. GPT-2 is used as the base for the reward model, and the first 70% of hidden layers are frozen (only latter 30% are updated).

We then fine-tune a GPT-2 model using PPO and our pre-trained reward model to generate summaries from online posts (using Reddit TL;DR dataset for training). Only the last 8 layers will be refined. We define a reward function to be the difference in scores (calculated using our reward model) between the original sample data and the predicted sample data using the GPT model. This serves as the loss function for which we backpropagate and fine-tune the model.

2b) Tested Code Notebook

I was able to use Label Studio to annotate the custom dataset.

The screenshot shows the Label Studio interface with a dataset grid. The columns are: ID, Completed, Order, Annotated by, prompt, answer1, and answer2. The data rows are:

ID	Completed	Order	Annotated by	prompt	answer1	answer2
1	Nov 17 2024, 14:13:01	1	0	EN	What is the latest news on the stock market?	Let the spotlight shine on something big, something that matters. If you haven't
2	Nov 17 2024, 14:13:35	1	0	EN	What is the current state of the economy?	I'm seeing some of the data back on here, about how much we need to increase
3	Nov 17 2024, 14:14:56	1	0	EN	What are the latest developments in technology?	Will it ever be able to be incorporated into existing devices? This survey
4	Nov 17 2024, 14:15:12	1	0	EN	What is the political situation in the Middle East?	The Arab League, of course, is very concerned about what's going on in the
5	Nov 17 2024, 14:15:17	1	0	EN	What are the latest trends	We believe the new trend

I then was able to train the reward model based on the labeled data.

```

1 training_args = TrainingArguments(
2     output_dir="rm_checkpoint/",
3     num_train_epochs=100,
4     logging_steps=10,
5     gradient_accumulation_steps=4,
6     save_strategy="steps",
7     evaluation_strategy="steps",
8     per_device_train_batch_size=4,
9     per_device_eval_batch_size=4,
10    eval_accumulation_steps=1,
11    eval_steps=10,
12    save_steps=10,
13    warmup_steps=100,
14    logging_dir="./logs",
15    fp16=True,
16    bf16=False,
17    learning_rate=1e-4,
18    # deepspeed="ds_config_gpt_j.json"
19    save_total_limit=1
20 )
21
22 Trainer(
23     model=model,
24     args=training_args,
25     train_dataset=train_dataset,
26     compute_metrics=compute_metrics,
27     eval_dataset=val_dataset,
28     data_collator=data_collator,
29 ).train()

```

Step	Training Loss	Validation Loss	Accuracy
10	0.575200	0.778556	1.000000
20	0.545700	0.774925	0.857143
30	0.484300	0.773462	0.857143
40	0.396500	0.785950	0.857143
50	0.269000	0.835993	0.857143
60	0.149300	0.940576	0.857143
70	0.057200	1.131552	0.857143
80	0.012300	1.351765	0.571429
90	0.001700	1.518373	0.571429
100	0.000700	1.607360	0.571429

TrainOutput(global_step=100, training_loss=0.249198, 'train_steps_per_second': 0.283, 'total_flos': 0.0,

Unfortunately, I was not able to train the language model because I could not get trlx to install properly. (I was using Google Colab since I do not have a GPU). Attempting to train on my local machine's CPU was extremely slow. Note: the authors of the jupyter notebook commented that trlx is not compatible with Google Colab and cannot be trained if using Google Colab.

```

▶ 1 # !pip install -U git+https://github.com/CarperAI/trlx.git
2
3 import os
4 from typing import List
5
6 import torch
7 from datasets import load_dataset
8 from reward_model import GPTRewardModel
9 from tqdm import tqdm
10 from transformers import AutoTokenizer
11
12 import trlx
13 from trlx.data.configs import (
14     ModelConfig,
15     OptimizerConfig,
16     SchedulerConfig,
17     TokenizerConfig,
18     TrainConfig,
19     TRLConfig,
20 )
21 from trlx.models.modeling_ppo import PPOConfig
--INSERT--
→ -----
ModuleNotFoundError                         Traceback (most recent call last)
<ipython-input-22-6754258bbef> in <cell line: 13>()
    11
    12 import trlx
--> 13 from trlx.data.configs import (
    14     ModelConfig,
    15     OptimizerConfig,

```

ModuleNotFoundError: No module named 'trlx.data'