

Week 6 Thursday Paper 1: [Direct Preference Optimization](#)

Supplemental: [Preference Tuning LLMs with Direct Preference Optimization Methods](#)

Questions:

1. How is DPO different from the RLHF algorithm?

Ans: DPO offers a new parameterization of the reward model in RLHF, which allows deriving the corresponding optimal policy directly, simplifying the RLHF problem into a classification task. Unlike RLHF, which involves complex and unstable procedures, DPO is more stable, computationally lighter, and doesn't require sampling from the language model during fine-tuning.

2. What is meant by "Your Language Model Is Secretly a Reward Model"?

Ans: This concept suggests that the language model, through the DPO process, implicitly acts as a reward model. Instead of separately learning a reward model and then optimizing a policy based on it, DPO integrates these steps, allowing the LM to directly reflect human preferences.

3. Explain the reward model parameterization used in DPO.

Ans: In DPO, the reward model is implicitly parameterized through a change of variables, transforming the optimization problem into one that can be solved directly without separate reward learning. This results in a system where the LM's updates inherently optimize preference satisfaction as can be seen in equation 7. We start with the same RL objective as prior work, equation 3 and it is straightforward to show that the optimal solution to the KL-constrained reward maximization objective in equation 3 becomes equation 4. Then, taking logarithm on both sides and with some algebra, we obtain equation 5. Substituting the reparameterization in Eq. 5 for $r^*(x, y)$ into the preference model Eq. 1, the partition function cancels, and we can express the human preference probability in terms of only the optimal policy π^* and reference policy π_{ref} . Now that we have the probability of human preference data in terms of the optimal policy rather than the reward model, we can formulate a maximum likelihood objective for a parametrized policy π_θ .

4. How does DPO address instability issues associated with actor-critic algorithms like PPO?

Ans: DPO mitigates instability associated with actor-critic methods like PPO by eliminating the need for the reinforcement learning loop and instead using a simple classification loss. This approach reduces the complexity and computational demands, leading to a more stable training process.

5. What are some of the advantages or drawbacks of DPO?

Ans: Advantages: DPO is computationally efficient, simpler to implement, and provides stable performance without the need for extensive hyperparameter tuning.

Drawbacks: As with any model, there may be concerns regarding generalization to new or unseen distributions, and there is less control over explicit reward function dynamics compared to traditional RL approaches.

6. Do you have any questions?

Week 6 Thursday Paper 2: [Behavioral Cloning from Observation](#)

Supplemental: [Supervised Learning of Behaviors](#)

Questions:

1. How is Behavioral Cloning different from Inverse Reinforcement Learning?

Ans: Behavioral Cloning directly copies the actions of an expert given the states they were taken in, effectively treating the imitation learning problem as a supervised learning problem. Inverse Reinforcement Learning, however, attempts to deduce the underlying reward function that the expert seems to be optimizing, and then uses this inferred reward to learn the policy.

2. Explain the two phase approach used in the BCO algorithm.

Ans: The BCO algorithm utilizes a two-phase approach as follows,

1. Inverse Dynamics Model Learning: Before any demonstration, the agent learns an agent-specific inverse dynamics model through exploration. During this phase, the agent collects interactions from the environment, storing state transitions and the corresponding actions. This phase is designed to infer missing action information from observed state transitions, constructing a model that maps state transitions to actions based on maximum likelihood estimation. This enables the agent to understand the dynamics of its actions within the environment without needing prior knowledge of the expert's specific actions
2. Behavioral Cloning: After acquiring the inverse dynamics model, the agent observes state-only demonstrations and uses the previously learned model to infer the missing actions of the expert. With these inferred actions, the agent applies a modified version of behavioral cloning to learn the imitation policy. This approach allows the agent to mimic the expert's behavior based on state transitions alone, bypassing the need for explicit action information from the expert. If allowed, post-demonstration environment interaction can further refine the learned model and imitation policy, balancing between the quality of imitation and the amount of interaction needed for improvement

3. How does BCO leverage model-based learning, and what advantages does it offer over model-free methods?

Ans: BCO leverages model-based learning by using the learned inverse dynamics model to infer actions, which can then be used for behavioral cloning. This approach is more sample-efficient than model-free methods since it relies on understanding the dynamics of the environment rather than learning a policy through trial and error.

4. Why is it important for the agent to acquire experience in a self-supervised fashion before observing expert demonstrations?

Ans: It's important for the agent to acquire experience in a self-supervised fashion before observing expert demonstrations to build a foundational understanding of the environment dynamics, which helps in better interpreting and mimicking the expert's behavior.

5. What are the advantages and disadvantages of doing Behavioral Cloning from Observation?

Ans: Advantages of BCO include not needing explicit action information from the expert, potentially reducing the amount of data needed for training, and enabling learning from state-only demonstrations. Disadvantages include possible inaccuracies in action inference and the general challenges of behavioral cloning, such as compounding errors and overfitting to the expert's trajectories.

6. Do you have any questions about the paper?