

Fully General Online Imitation Learning

Michael K. Cohen



Marcus Hutter



Neel Nanda



Imitation

- Optimizers that plan over the long term in an unknown environment are dangerous
- A perfect imitation of a human is “safe” (not likely to cause an unprecedented catastrophe)
- Perhaps a population of imitators could prevent the realization of dangerous optimizers
- Perhaps there could be global coordination to use only imitators and “narrow” optimizers, with too little data to model what makes humans tick

Learned Imitation

- Errors in the early stages of imitation learning could have irreversible consequences
- Some concern about dangerousness being *selected for* in the errors of strong prediction algorithms¹
- Active learning allows extremely strong error guarantees

¹Search terms: mesa-optimization; universal prior malign; inner alignment

Online Imitation Learning

- Instead of train-then-deploy, imitator can ask what the demonstrator would do when uncertain
- Maintain a posterior distribution over demonstrator-models
- Sort them
- Collect the top few. The n^{th} best gets collected if $\text{post}(\text{model}_n) > \alpha \sum_{i \leq n} \text{post}(\text{model}_i)$
- Each model gives probabilities over the next action
- For each action, take the min prob. assigned by the top models
- Sample actions with those probs; if nothing sampled, ask for help and take the action the demonstrator picks

Design Motivation

- Suppose the true demonstrator model is always in the top few
- For small enough α , we can make that arbitrarily likely
- a is arbitrary action; θ_q is query prob.
- $\pi^{\text{im}}(a \text{ w/o query}|\text{history}) \leq \pi^{\text{dem}}(a|\text{history})$
- $\pi^{\text{im}}(a \text{ after query}|\text{history}) = \theta_q \pi^{\text{dem}}(a|\text{history})$
- $\pi^{\text{im}}(a|\text{history}) \leq (1 + \theta_q) \pi^{\text{dem}}(a|\text{history})$
- actions are never that much more likely for imitator than for demonstrator *even in log space*
- $\mathbb{E} \prod_{t=1}^N (1 + \theta_q(\text{time } t)) \in O(N^2)$

(Un)bounded Cost Functions

- Suppose all we had was $\|\pi^{\text{im}} - \pi^{\text{dem}}\|_1 \leq \varepsilon$
- What would it take for us to find that satisfactory?
- Bounded cost function
- If [bad event] has prob 0 under demonstrator, it has prob $\leq \varepsilon$ under imitator, for added cost $\leq \varepsilon C$
- Bounded cost function is hard to justify
- Assume the demonstrator tries to keep bad events unlikely, e.g. bad events $< 10\%$ chance, really bad events $< 1\%$ chance, catastrophic events $< 10^{-6}$
- Events with 10% chance under demonstrator can only be so bad, maybe
- 1% events could be much worse
- An imitator that makes an event go from 1% to 2% way less scary than one that makes it go from 10^{-9} to 1% .
- Recall: $\pi^{\text{im}}(a|\text{history}) \leq (1 + \theta_q)\pi^{\text{dem}}(a|\text{history})$

Main Results

- Events can unfold over multiple timesteps
- **Assume:** True demonstrator is in model class
- Let E be the event that truth always \in top models, which has prob. $\geq 1 - \frac{\alpha}{\text{prior on truth}}$
- Let B be a possible bad event in the first t timesteps
- As $P_{\text{dem}}(B)^{-1} \rightarrow \infty$, $P_{\text{im}}(B \text{ and } E)^{-1} \rightarrow \infty$ at least polylogarithmically: $f(x) \gtrsim c \log(x)^3$
- $c \propto t^{-2}$, α^4
- $\mathbb{E} [\sum_{t=0}^{\infty} (\text{query prob})^3] \leq |\mathcal{A}| \alpha^{-3} \left(\frac{24}{\text{prior on truth}} + 12 \right) < \infty$
- Precondition lifetime with many demonstrator-chosen actions to increase “prior” on truth

AGI Safety Ideas that Use Imitation

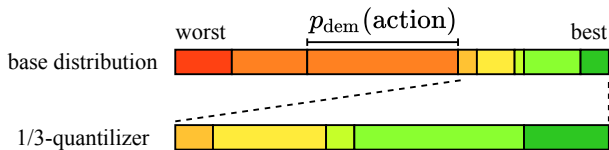
→ Iterated Distillation and Amplification

→ Recursive definition: IDA_0 is a rock

→ IDA_{n+1} imitates the output of a human with access to two copies of IDA_n

→ IDA_n imitates an organization of $2^n - 1$ humans

→ Quantilization



→ Quantilization only multiplies action probabilities by a constant factor

→ What if the base distribution were an imitator *without* a bound on $P_{\text{im}}(\text{action}) / P_{\text{dem}}(\text{action})$?

→ The guarantees of quantilization would lose relevance

Conclusion

- Querying for demonstrations during model-mismatch → finite error bounds in imitation
- Theoretical solution to safe imitation
- Theoretical solution to “inner alignment failure”
- There may be a path to powerful, safe AGI via imitation