
Deconfounded Imitation Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Standard imitation learning can fail when the expert demonstrators have different
2 sensory inputs than the imitating agent. This partial observability gives rise to
3 hidden confounders in the causal graph, which lead to the failure to imitate. We
4 break down the space of confounded imitation learning problems and identify three
5 settings with different data requirements in which the correct imitation policy can
6 be identified. We then introduce an algorithm for deconfounded imitation learning,
7 which trains an inference model jointly with a latent-conditional policy. At test
8 time, the agent alternates between updating its belief over the latent and acting
9 under the belief. We show in theory and practice that this algorithm converges
10 to the correct interventional policy, solves the confounding issue, and can under
11 certain assumptions achieve an asymptotically optimal imitation performance.

12 1 Introduction

13 In imitation learning (IL), an agent learns a policy directly from expert demonstrations without
14 requiring the specification of a reward function. This paradigm could be essential for solving real-
15 world problems in autonomous driving and robotics where reward functions can be difficult to shape
16 and online learning may be dangerous. However, standard IL requires that the conditions under which
17 the agent operates exactly match those encountered by the expert. In particular, they assume that
18 there are no *latent confounders*—variables that affect the expert behavior, but that are not observed by
19 the agent. This assumption is often unrealistic. Consider a human driver who is aware of the weather
20 forecast and lowers its speed in icy conditions, even if those are not visible from observations. An
21 imitator agent without access to the weather forecast will not be able to adapt to such conditions.

22 In such a situation, an imitating agent may take their own past actions as evidence for the values of
23 the confounder. A self-driving car, for instance, could conclude that it is driving fast, thus there can be
24 no ice on the road. This issue of *causal delusion* was first pointed out in Ortega and Braun [2010a,b]
25 and studied in more depth by Ortega et al. [2021]. The authors analyze the causal structure of this
26 problem and argue that an imitator needs to learn a policy that corresponds to a certain interventional
27 distribution. They then show that the classic DAgger algorithm [Ross et al., 2011], which requires
28 querying experts at each time step, solves this problem and converges to the interventional policy.

29 In this paper, we present a solution to a confounded IL problem, where both the expert policy and
30 the environment dynamics are Markovian. The solution does not require querying experts. We first
31 present a characterization of confounded IL problems depending on properties of the environment
32 and expert policy (section 3). We then show theoretically that an imitating agent can learn behaviors
33 that approach optimality when the above Markov assumptions and a recurrence property hold.

34 We then introduce a practical algorithm for deconfounded imitation learning that does not re-
35 quire expert queries (section 4). An agent jointly learns an inference network for the value
36 of latent variables that explain the environment dynamics as well as a latent-conditional pol-
37 icy. At test time, the agent iteratively samples latents from its belief, acts in the environ-

38 ment, and updates the belief based on the environment dynamics. An imitator steering a self-
 39 driving car, for instance, would learn how to infer the weather condition from the dynamics of
 40 the car on the road. This inference model can be applied both to its own online experience
 41 as well as to expert trajectories, allowing it to imitate the behavior adequate for the weather.
 42

43 Finally, our deconfounded imitation learning algorithm is
 44 demonstrated in a multi-armed bandit problem. We show
 45 that the agent quickly adapts to the unobserved properties
 46 of the environment and then behaves optimally (section 5).

47 2 Imitation learning and latent confounders

48 We begin by introducing the problem of confounded imi-
 49 tation learning. Following Ortega et al. [2021], we discuss
 50 how behavioral cloning fails in the presence of latent con-
 51 founders. We then define the interventional policy, which
 52 solves the problem of confounded imitation learning.

53 2.1 Imitation learning

54 Imitation learning learns a policy from a dataset of expert demonstrations via supervised learning.
 55 The expert is a policy that acts in a (reward-free) Markov decision process (MDP) defined by a tuple
 56 $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P(s' | s, a), P(s_0))$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $P(s' | s, a)$ is
 57 the transition probability, and $P(s_0)$ is a distribution over initial states. The expert’s interaction with
 58 the environment produces a trajectory $\tau = (s_0, a_0, \dots, a_{T-1}, s_T)$. The expert may maximize the
 59 expectation over some reward function, but this is not necessary (and some tasks cannot be expressed
 60 through Markov rewards Abel et al. [2021]). In the simplest form of imitation learning, a behavioral
 61 cloning policy $\pi_\eta(a | s)$ parametrized by η is learned by minimizing the loss $-\sum_{s,a \in \mathcal{D}} \log \pi_\eta(a | s)$,
 62 where \mathcal{D} is the dataset of state-action pairs collected by the expert’s policy.

63 2.2 Confounded imitation learning

64 We now extend the imitation learning setup to allow for some variables $\theta \in \Theta$ that are observed by
 65 the expert, but not the imitator. We define a family of Markov Decision processes as a latent space Θ ,
 66 a distribution $P(\theta)$, and for each $\theta \in \Theta$, a reward-free MDP $\mathcal{M}_\theta = (\mathcal{S}, \mathcal{A}, P(s' | s, a, \theta), P(s_0 | \theta))$.
 67 We assume there exists an expert policy $\pi_{\text{exp}}(a | s, \theta)$ for each MDP. When it interacts with the
 68 environment, it generates the following distribution over trajectories τ :

$$P_{\text{exp}}(\tau | \theta) = P(s_0 | \theta) \prod_{t=0}^{T-1} P(s_{t+1} | s_t, a_t; \theta) \pi_{\text{exp}}(a_t | s_t; \theta)$$

69 The imitator does not observe the latent θ . It may thus need to implicitly infer it from the past
 70 transitions, so we take it to be a non-Markovian policy $\pi_\eta(a_t | s_1, a_1, \dots, s_t)$, parameterized by η .
 71 The imitator generates the following distribution over trajectories:

$$P_\eta(\tau | \theta) = P(s_0 | \theta) \prod_{t=0}^{T-1} P(s_{t+1} | s_t, a_t; \theta) \pi_\eta(a_t | s_0, a_0, \dots, s_t)$$

72 The Bayesian networks associated to these distributions are shown in figure 1.

73 The goal of imitation learning in this setting is to learn imitator parameters η such that when the
 74 imitator is executed in the environment, the imitator agrees with the expert’s decisions, meaning we
 75 wish to maximise

$$\mathbb{E}_{\theta \sim P(\theta)} \mathbb{E}_{\tau \sim P_\eta(\tau; \theta)} \left[\sum_{s_t, a_t \in \tau} -\log \pi_{\text{exp}}(a_t | s_t, \theta) \right]. \quad (1)$$

76 If the expert solves some task (e. g. maximizes some reward function), this amounts to solving the
 77 same task.

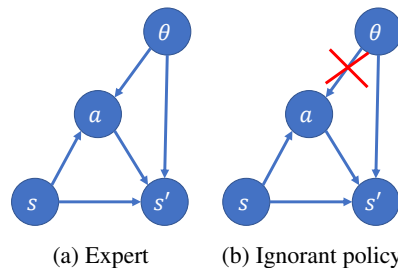


Figure 1: Bayes nets for (a) a confounded transition and (b) non-confounded transition.

78 **2.3 Naive behavioral cloning**

79 If we have access to a data set of expert demonstrations, one can learn an imitator via behavioral
80 cloning on the expert’s demonstrations. At optimality, this learns the *conditional policy*:

$$\pi_{\text{cond}}(a_t | s_1, a_1, \dots, s_t) = \mathbb{E}_{\theta \sim p_{\text{cond}}(\theta | \tau)} \pi_{\text{exp}}(a_t | s_t, \theta) \quad (2)$$

$$p_{\text{cond}}(\theta | \tau) \propto p(\theta) \prod_t p(s_{t+1} | s_t, a_t, \theta) \pi_{\text{exp}}(a_t | s_t, \theta) \quad (3)$$

81 Following Ortega et al. [2021], consider the following example of a confounded multi-armed bandit
82 with $\mathcal{A} = \Theta = \{1, \dots, 5\}$ and $\mathcal{S} = \{0, 1\}$:

$$p(\theta) = \frac{1}{5}, \quad \pi_{\text{exp}}(a_t | s_t, \theta) = \begin{cases} \frac{6}{10} & \text{if } a_t = \theta \\ \frac{1}{10} & \text{if } a_t \neq \theta \end{cases}, \quad P(s_{t+1} = 1 | s_t, a_t, \theta) = \begin{cases} \frac{3}{4} & \text{if } a_t = \theta \\ \frac{1}{4} & \text{if } a_t \neq \theta. \end{cases} \quad (4)$$

83 The expert knows which bandit arm is special (and labeled by θ) and pulls it with high probability,
84 while the imitating agent does not have access to this information.

85 If we roll out the naive behavioral cloning policy in this environment, shown in Figure 2, we see the
86 causal delusion at work. At time t , the latent that is inferred by p_{cond} takes past actions as evidence
87 for the latent variable. This makes sense on the expert demonstrations, as the expert is cognizant
88 of the latent variable. However, during an imitator roll-out, the past actions are not evidence of the
89 latent, as the imitator is blind to it. Concretely, the imitator will take its first action uniformly and
90 later tends to repeat that action, as it mistakenly takes the first action to be evidence for the latent.

91 **2.4 Interventional policy**

92 A solution to this issue is to only take as evidence the data that was actually informed by the latent,
93 which are the transitions. This defines the following imitator policy:

$$\pi_{\text{int}}(a_t | s_1, a_1, \dots, s_t) = \mathbb{E}_{\theta \sim p_{\text{int}}(\theta | \tau)} \pi_{\text{exp}}(a_t | s_t, \theta), \quad p_{\text{int}}(\theta | \tau) \propto p(\theta) \prod_t p(s_{t+1} | s_t, a_t, \theta). \quad (5)$$

94 In a causal framework, that corresponds to treating the choice of past actions as interventions. In
95 the notation of the do-calculus [Pearl, 2009], this equals $p(a_t | s_1, \text{do}(a_1), s_2, \text{do}(a_2), \dots, s_t)$. The
96 policy in equation (5) is therefore known as *interventional policy* [Ortega et al., 2021].

97 **3 Deconfounding imitation learning**

98 We now present our theoretical results on how imitation learning can be deconfounded. We first show
99 that the interventional policy is optimal in some sense, before analyzing in which settings it can be
100 learned.

101 **3.1 Optimality of the interventional policy**

102 Under some reasonable assumptions, the interventional policy approaches the expert’s policy, as we
103 prove in the appendix 2.

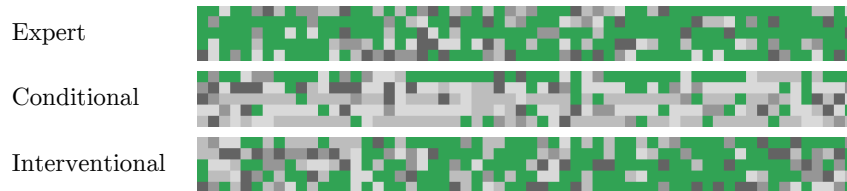


Figure 2: Actions from rollouts from bandit environment (4). The x -axis is episode time. In the y -axis five roll-outs are shown with from the expert and policies (3) and (5). Colors denote actions, with the correct arm labelled green. The interventional imitator tends to the expert policy, while the conditional policy tends to repeat itself.

104 **Theorem 3.1** (Informal). *If the interventional inference $p_{\text{int}}(\theta \mid \tau_{<t})$ approaches the true latent of the*
105 *environment as $t \rightarrow \infty$ on the rollouts of π_{int} , and if the expert maximises some reward that is fixed*
106 *across all environments, then as $t \rightarrow \infty$, the imitator policy $\pi_{\text{int}}(a_t \mid s)$ approaches the expert policy.*

107 *Proof.* See lemma 2.1 in the appendix.

108 The requirement here means that the transition dynamics must be informative about the latent —
109 we consider latent confounders that manifest in the dynamics, not those that affect only the agent
110 behavior.¹ In this case, the interventional policy thus presents a solution to the confounding problem.
111 In the rest of this paper we focus on the question if and how it can be learned from data.

112 Note that the interventional policy only guarantees *asymptotic* optimal performance. However, it is not
113 guaranteed to be the policy that adapts to a given test environment in the fastest possible way. Efficient
114 exploration may require learning behavior that the expert never demonstrated Zhou et al. [2019].

115 3.2 Identification

116 The rest of this section concerns the question: under which assumptions on the model and data avail-
117 ability can we identify the interventional policy. Under the typical rules of do-calculus, identification
118 should be impossible without observing the latent. However, we discern three tiers of decreasing
119 strength of model assumptions and increasing data requirements in which identification is neverthe-
120 less possible without observing the latent. In this way, we split the confounding problem discussed
121 by Ortega and Braun [2010a,b], Ortega et al. [2021] into three separate classes, provide conditions
122 for when each class applies, and weaken the requirements for identifiability for two of these classes.

123 Tier 1: identifiability from demonstrations

124 In the first tier, we make the strongest assumption and show that the interventional policy is identifiable
125 only from the expert’s demonstrations.

126 **Theorem 3.2** (Informal). *Let the expert’s demonstrations be recurrent, meaning that each trajectory*
127 *contains every state and action pair infinitely often. Then the interventional policy is identifiable only*
128 *from expert demonstrations.*

129 *Proof.* See theorem 1.2 in the appendix.

130 The assumption of recurrence is for instance satisfied if for any state, the expert policy has full support
131 over actions. In that case, on one trajectory, with one value of the latent θ , we are able to accurately
132 estimate the distributions over both the dynamics $p(s' \mid s, a, \theta)$ and the expert $\pi_{\text{exp}}(a \mid s, \theta)$. Then
133 from the distribution over trajectories, we estimate the distribution over distributions and thus identify
134 $p(\theta)$ and the θ -conditional dynamics and expert distributions. The interventional policy can then be
135 constructed from these models via equation (5).

136 In practice, this method suffices when the demonstrations are of sufficient length and the expert and
137 dynamics are sufficiently noisy to sufficiently cover the state and action space in each trajectory.

138 Tier 2: identifiability from demonstrations and explorations

139 In case the expert’s demonstrations are not recurrent, but we are able to execute an exploratory policy
140 in the environment, the interventional policy may still be identifiable:

141 **Theorem 3.3** (Informal). *Consider the setting in which in addition to the expert demonstration data,*
142 *we can collect trajectories by rolling out an exploration policy. Let the exploration data be recurrent.*
143 *Let the true latents be identifiable from the expert trajectories given the true state-transition model*
144 *$p(s' \mid s, a, \theta)$. Then the interventional policy is identifiable from expert demonstrations and exploration*
145 *data.*

146 *Proof.* See theorem 1.4 in the appendix.

147 How can we identify the interventional policy from expert data and exploration data? Under the

¹For instance, our results do not apply to latents that specify the task the agent solves but do not affect the environment dynamics at all. That setting, a staple of the meta-RL literature, generally requires a single-shot or few-shot setting where some demonstrations are required for each value of the latent at test time.

148 assumption of recurrence of exploration trajectories, which is satisfied if the exploration policy has
 149 full support and any state is reachable from any state in the environment, the interventional policy can
 150 be identified in the following way. First, we estimate the distribution over the state transitions from
 151 the exploratory data. This allows us to construct the interventional inference model $p_{\text{int}}(\theta | \tau)$ from
 152 equation (5). With that model, we can infer the latent for any expert trajectory and learn a policy
 153 $\pi_\eta(a | s, \theta)$ with behavioral cloning on the demonstrations.

154 We need to make the crucial assumption that the inference collapses to the true latent of the dynamics —
 155 in other words, the state transitions encountered by the expert are informative about the latent
 156 variables. (If the expert only visits states where the latent variables do not influence the transition
 157 dynamics, it stands to reason that we cannot deconfound the imitator.) Under this assumption, this
 158 method identifies the interventional policy, as we prove in theorem 1.4 in the appendix.

159 We will demonstrate the identifiability of the interventional policy in this tier in practice in the
 160 following section.

161 Tier 3: identifiability from expert queries

162 If the dynamics of the environment are such that no exploratory policy leads to recurrent trajectories,
 163 then the interventional policy can only be identified if we can interact in the environment and query the
 164 expert for its decisions. The procedure is to learn a recurrent policy $\pi_\eta(a_t | s_1, a_1, \dots, s_t)$, execute it
 165 in the environment, query the expert — aware of the true latent — for its actions, and train the imitator
 166 with maximum likelihood to match those actions. This is similar to the setup in DAgger [Ross et al.,
 167 2011]. In Ortega et al. [2021], it is shown that at optimality this converges to the interventional policy.

168 Examples for problems that fall into tier 3 include cases where the environment dynamics is not
 169 Markov given the latent, where the expert policy is not Markov given the latent, or where the expert
 170 only visits regions of the state space in which the dynamics are insensitive to the latent.

171 4 Practical algorithm

172 We now introduce a practical algorithm for training an agent from expert data in the presence of
 173 latent confounders. For simplicity, we focus on tier 2 of the settings discussed in section 3, i. e.
 174 assume access to the expert data as well as the ability to gather more data interactively. In appendix 3,
 175 we describe an algorithm that applies to tier 1 (i. e. does not require the ability to gather more data
 176 interactively), but faces a more difficult learning problem in practice.

177 As outlined in section 3, we use the access to the family of MDPs defined by the latent θ to learn an
 178 inference model for the latent variable, infer the latents on the expert trajectories using the learned
 179 model, and learn a policy conditioned on the inferred latent, which imitates the expert based on the
 180 demonstrations. At test time, the agent uses posterior sampling (or Thompson sampling) [Thompson,
 181 1933]: it alternates between updating its belief about which MDP it is in and acting optimally under
 182 its current belief.

183 **Components** The agent consists of: (1) an inference model q_ϕ , which maps trajectories $\tau =$
 184 (s_0, a_0, \dots, s_n) to a belief over a latent variable $q_\phi(\hat{\theta} | \tau)$; (2) a dynamics model p_ψ , which maps
 185 latent, state, and action to a distribution over the next state $p_\psi(s' | s, a, \hat{\theta})$; (3) a prior over latents
 186 $p(\hat{\theta})$; and (4) a latent-conditional policy $\pi_\eta(a | s, \hat{\theta})$.

187 **Training the inference model** The inference model q_ϕ , dynamics model p_ψ , and prior $p(\hat{\theta})$ form a
 188 latent variable model for trajectory data. Similarly to a variational autoencoder (VAE) [Kingma and
 189 Welling, 2013], they are trained by minimizing a variational inference objective² given by

$$\mathcal{L}_i = \mathbb{E}_{\hat{\theta} \sim q_\phi(\hat{\theta} | \tau_{:t}^i)} \left[\log p_\psi(s_{t+1} | s_t, a_t, \hat{\theta}) \right] - \beta D_{KL} \left(q_\phi(\hat{\theta} | \tau_{:t}^i) \parallel p(\hat{\theta}) \right). \quad (6)$$

²Variational inference is not a necessary component of the algorithm, it is just a particular design choice we make for convenience. Another option would be exact inference over a latent variable model, which is in some cases tractable. We have verified empirically that that leads to similar results.

190 At timestep t , the encoder q_ϕ takes as input the trajectory observed until timestep t , $\tau_{:t} =$
 191 $(s_0, a_0, \dots, a_{t-1}, s_t)$, and predicts a distribution of the belief over the latent $\hat{\theta}$. The decoder p_ψ is a
 192 dynamics model, which predicts the next state s_{t+1} given the state s_t , action a_t , and a sample from
 193 the belief distribution. Unlike a VAE, the proposed latent variable model is not an autoencoder, since
 194 instead of decoding the input data distribution, the dynamics model p_ψ decodes the future states. To
 195 minimize the loss in equation (6), the inference model needs to encode in the learned latent $\hat{\theta}$ all of
 196 the information in the ground-truth latent θ that is important to predict the state transition.

197 Crucially, the training data for the inference model is collected from from the interactions of an
 198 exploration policy with the environment. As exploration policy, we use the latent-conditional policy
 199 π_η , sampling $\hat{\theta}$ from the prior $p(\hat{\theta})$. Other exploration policies may be used as long as they explore
 200 sufficiently diverse trajectories and do not depend on the true latent θ . However, we cannot simply
 201 train the inference model on expert data, as those trajectories were collected from the θ -dependent
 202 expert policy. An inference model trained on that data would exhibit the problem of causal delusions
 203 that we are trying to solve.

204 **Training the imitation policy** The learned inference model is used for inferring the latents in the
 205 expert data distribution. These inferences are valid despite the distribution shift from exploratory data
 206 to expert data, as the transition model is the same.

207 Standard imitation learning can then be used to learn a policy conditional on the inferred latents from
 208 the expert demonstrations. We use behavioral cloning as the imitation learning method because it is
 209 simple and reasonably effective. We show the full training algorithm in algorithm 1.

210 **Test time** At test time, the agent faces an environment with an unknown latent and needs to adapt
 211 to the correct expert behavior. We solve this problem by alternating between updating a posterior
 212 belief over the latent and acting under the current belief, i. e. sampling a latent from the current belief
 213 and acting like the corresponding expert.

214 Concretely, the agent initially samples a latent from the prior $\hat{\theta} \sim p(\hat{\theta})$ and an action $a \sim \pi_\eta(a|s, \hat{\theta})$
 215 to imitate the expert corresponding to that latent. It observes the state transition and computes the
 216 posterior belief with the inference network. Another latent is sampled from the updated belief, and so
 217 on. Once the inference has converged to match the true latent for the environment, the true expert for
 218 the environment will be imitated consistently. We summarize the test-time behavior in algorithm 2.

219 In practice, using posterior sampling successfully requires that the imitator conditioned on a randomly
 220 sampled latent does not end up in states the true expert would not have visited. For example, in a
 221 robotic manipulation task, after accidentally pushing the object outside the reach of the robot, none
 222 of the expert policies could continue the task. A simple way of avoiding this problem is allowing
 223 multiple episodes of interaction with the task, where the environment is reset to a state sampled
 224 from the initial state distribution after each episode. Solving this problem for arbitrary environments
 225 without resets requires methods beyond behavioral cloning; we leave this for future work.

226 5 Experiments

227 To test our method in practice, we conduct an experiment in the multi-armed bandit problem proposed
 228 by Ortega et al. [2021] and described in section 2. With this experiment, we aim to answer the
 229 following questions. First, how big is the effect of confounding on naive imitation learning — large
 230 enough to justify the use of specialized methods? Second, is our algorithm capable of identifying the
 231 interventional policy? Finally, does the interventional policy converge to the expert policy?

232 **Setup** In the experiments, the expert policy is defined in equation (4). We consider episodes of
 233 length 100. We are only interested in the effects of confounding on imitation learning, and not
 234 issues arising from optimization. Therefore, in order to avoid overfitting to a finite expert dataset,
 235 we resample new data from the expert for each update of the learning algorithms. Each learning
 236 algorithm is run for ten independent seeds and the results are averaged. The hyperparameters for the
 237 algorithms are provided in appendix 4.

Algorithm 1 Behavioral cloning with latent inference

Require: Initial parameters of the imitation policy η , inference model ϕ , and dynamics model ψ , an expert dataset $\{\tau_e^j\}$, an MDP \mathcal{M} , true latent distribution $p(\theta)$, learning rates α_1 , α_2 , and α_3 .

while not done **do**

$\theta \sim p(\theta)$

$s_0 \sim p_0(s_0)$

$\tau = s_0, t = 0$

for $t \leq H$ **do**

$a_t \sim \pi(a_t | s_t)$

▷ Sample action from a Markov policy, e.g. a random policy

$s_{t+1} = p(s_{t+1} | s_t, a_t, \theta)$

▷ True dynamics

Append (a_t, s_{t+1}) to τ .

$t = t + 1$

end for

$\phi = \phi - \alpha_1 \nabla_{\phi} \hat{\mathcal{L}}(\tau)$ with $\hat{\mathcal{L}}(\tau)$ being a sample estimate of equation 6.

$\psi = \psi - \alpha_2 \nabla_{\psi} \hat{\mathcal{L}}(\tau)$ with $\hat{\mathcal{L}}(\tau)$ being a sample estimate of equation 6.

Infer latents $\hat{\theta}_H^j \sim q_{\phi}(\hat{\theta}_H^j | \tau_e^j)$ for expert trajectories $\{\tau_e^j\}$.

$\eta \leftarrow \eta - \alpha_3 \nabla_{\eta} \sum_j \sum_{s, a \in \tau_e^j} \log \pi_{\eta}(a | s, \hat{\theta}_H^j)$.

end while

Algorithm 2 Trained agent imitating an expert

Require: Trained parameters of the imitation policy η , inference model ϕ an MDP \mathcal{M} , and a latent θ .

$s_0 \sim p_0(s_0)$

$\tau = s_0, t = 0$

for $t \leq H$ **do**

$\hat{\theta}_t \sim q_{\phi}(\hat{\theta}_t | \tau)$

▷ Infer the latent for trajectory

$a_t \sim \pi_{\eta}(a_t | s_t, \hat{\theta}_t)$

▷ Condition on the inferred latent

$s_{t+1} = p(s_{t+1} | s_t, a_t, \theta)$

▷ True dynamics

Append (a_t, s_{t+1}) to τ .

$t = t + 1$

end for

238 **Naive imitation learning and the conditional policy** To answer the first question, we compare a
239 naive imitation learner to the true conditional policy described in section 2. Even the naive learner
240 needs a memory to be able to adapt to a new instance of the bandit defined by the latent. To provide
241 such a memory, we implement the imitation learner as a recurrent neural network (RNN). From
242 figure 3a, we see that the naive behavioral cloning agent implemented with an RNN learns a policy
243 matching the true conditional policy closely.

244 This results in problems for the policy learned with naive imitation learning when it is deployed
245 in the environment and has to choose the actions itself, as shown in figure 3b. In the figure, it can
246 be seen that the naive imitation learner does not converge to the expert policy during the episode.
247 Furthermore, it closely tracks the action probability of the true conditional policy, suggesting that
248 it has suffered the full impact of the confounding problem. Given that the naive imitator closely
249 matches the true conditional policy at convergence, we can conclude that naive imitation learning is
250 unlikely to be enough to solve confounded imitation learning problems.

251 **Deconfounded imitation learning and the interventional policy** To answer the second question,
252 we implement the deconfounded imitation learner as described in section 4. While in principle the
253 inference model can be learned from the data collected by the imitator, we use a random policy for
254 the data collection. Figures 3a and 3b show that the proposed method (labeled “Deconfounded”)
255 closely matches the true interventional policy both on expert trajectories and online. From this we
256 can conclude that in this simple environment, the proposed method produced an inference model and
257 a conditional policy that together match the interventional policy.

258 Figure 3c shows the imitation loss during training. Notice that the deconfounded algorithm achieves
259 a lower imitation loss than naive BC even though the policy trained with naive BC collapses on the
260 best arm more rapidly on the expert trajectories as shown in figure 3a. This happens because the

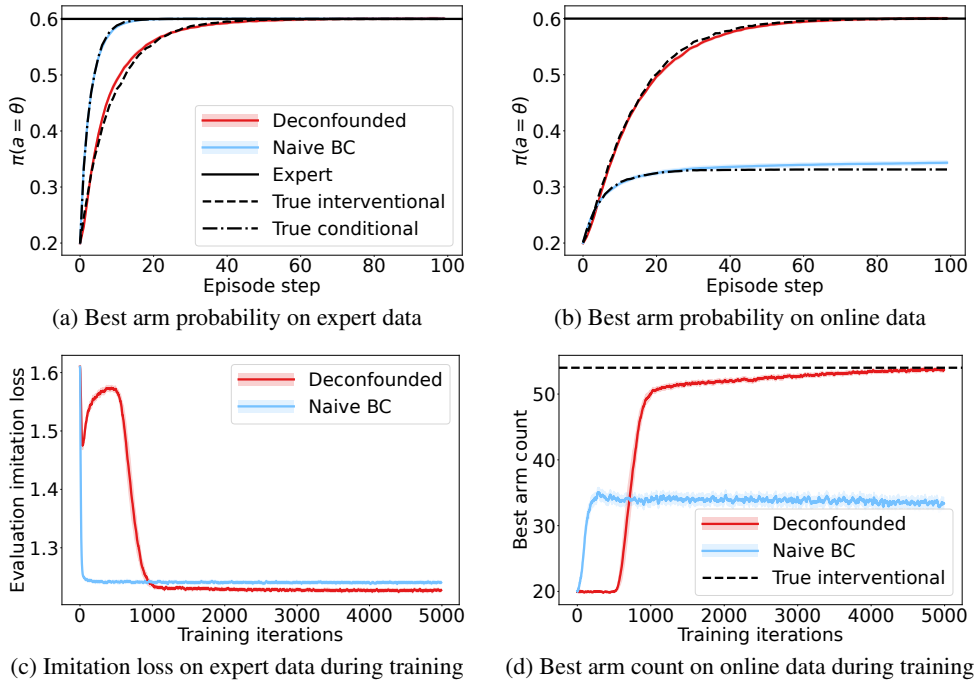


Figure 3: Experiment with imitation learning in a multi-armed bandit problem. For visual clarity, the curves for the learned agents are smoothed averages for sliding window of length 20. The shading shows the standard error of the mean across random seeds. Panel (a) shows the probability of choosing the best arm on evaluation expert trajectories after training. Panel (b) shows the probability of choosing the best arm in the online environment after training. Panel (c) shows the imitation learning loss on evaluation trajectories over the course of training. Panel (d) shows the number of times the policy chose the best arm during an episode over the course of training.

261 deconfounding step in the algorithm conditions the imitator policy on the latent inferred for the whole
 262 episode at every step of the episode.

263 Figure 3d shows the number of times the policies chose the best arm during an episode. This is a
 264 convenient proxy for how closely a policy matches the interventional policy. From the figure, it can
 265 be seen that the deconfounded imitation learning algorithm closely matches the counterfactual policy,
 266 while naive BC converges to lower best arm count.

267 We conclude that our deconfounded imitation learning algorithm is indeed able to learn the inter-
 268 ventional policy, solve the issue of causal delusions faced by naive behavioral cloning, and achieve
 269 near-perfect imitation performance. This comes at the price of requiring exploration data to train the
 270 inference model as well as an increased number of training iterations needed for convergence.

271 6 Related work

272 **Imitation learning** Imitation learning (or learning from demonstration) has a long history with
 273 applications in autonomous driving [Pomerleau, 1988] and robotics [Schaal, 1999]. Standard algo-
 274 rithms include behavioral cloning as well as inverse reinforcement learning Russell [1998], Ng et al.
 275 [2000]. Scaling imitation learning to high-dimensional continuous control problems is more challeng-
 276 ing and has been solved through adversarial methods [Ho and Ermon, 2016, Fu et al., 2017].

277 Imitation learning can suffer from a mismatch between the distributions faced by the expert and
 278 imitator due to the accumulation of errors when rolling out the imitator policy. This issue is commonly
 279 addressed by querying experts during the training [Ross et al., 2011] or by noise insertion in the
 280 expert actions. Note that this issue is qualitatively different from the one we discuss in the paper: it
 281 is a consequence of the limited support of the expert actions and occurs even in the absence of the
 282 latent confounders.

283 **Causality-aware imitation learning** De Haan et al. [2019] diagnose the issue of causal confusion
284 in imitation learning, in which the imitator draws causally wrong conclusions from its inputs because
285 such patterns may be easier to learn. This problem is different from the issue of latent confounders
286 we discuss and may occur even if the expert and imitator have access to the same information.

287 The works most closely related to our paper are Ortega and Braun [2010a,b], Ortega et al. [2021],
288 which point out the issue of latent confounding and causal delusions that we discuss in this paper. In
289 particular, Ortega and Braun [2010b] propose a training algorithm that learns the correct interventional
290 policy. Unlike our algorithm, their approach requires querying experts during the training. However,
291 as we discuss in section 3, their solution has weaker assumptions and in particular also applies to
292 non-Markovian environment dynamics.

293 Kumor et al. [2021] study imitation learning on partially observed structural causal models. They
294 characterize when behavioral cloning conditional on observed adjustment variables can maximize the
295 reward. The environments we consider do not meet these criteria, so behavioral cloning fails to imitate
296 the expert, while our method works under additional assumptions. In an extension, Anonymous
297 [2022] propose that an optimal policy can be recovered from a sub-optimal expert whenever the causal
298 effect from the policy on a hypothetical reward would be identifiable from the observables through the
299 IDENTIFY algorithm [Tian and Pearl, 2002], which is not the case for the environments we consider.

300 Rezende et al. [2020] point out that the same problem appears in partial models, i. e. world models
301 that only use a subset of the (in principle fully observable) state. The part of the state that is not
302 modeled then acts as a confounder for the model predictions. The authors then identify a minimal set
303 of variables that a partial model needs to include in order to not suffer from this confounding issue.
304 While their and our works describe closely related problems, the solutions differ as in our problem
305 we cannot redefine the inputs to include the confounder θ .

306 **Meta-reinforcement learning** Our problem is related to meta-reinforcement learning [Duan et al.,
307 2016, Wang et al., 2016]. Here an agent is trained to perform well in a variety of tasks. Both Rakelly
308 et al. [2019] and Zintgraf et al. [2019] propose meta-RL algorithms that consist of a task encoder
309 and a task-conditional policy, similar to our inference model and latent-conditional imitator policy.
310 One difference to our work is that these tasks can vary through the reward function, which is not
311 compatible with imitation learning because the reward function is not available at test time.

312 Zhou et al. [2019], like our work, consider multiple environments that differ through some latents.
313 They propose agents that learn to probe the environment to determine these latents. The difference to
314 our work is that they work in an RL setting with given reward functions, while our algorithm only
315 requires expert demonstrations.

316 Another closely related research topic is that of meta-imitation learning [Duan et al., 2017], where
317 the aim is to make imitation learning possible from a small number of demonstrations. Despite both
318 problems centering on imitation learning in a distribution of MDPs, the motivations and methods
319 are different. Our work does not consider adapting to new demonstrations, whereas meta-imitation
320 learning does not consider the confounding problem in the demonstrations.

321 7 Conclusion

322 Naive imitation learning algorithms fail in the presence of latent confounders — for instance when
323 the expert has access to more information than the demonstrator. We presented a breakdown of this
324 confounding problem based on when the interventional policy, which solves the confounding issue, is
325 identifiable without query access to the expert. We proposed two algorithms that provably converge
326 to the true interventional policy. In a multi-armed bandit experiment, we demonstrated for one of
327 them that it is able to learn the correct interventional policy, solves the confounding problem that
328 limits naive imitation learning, and converges to the expert behavior.

329 While our work makes progress theoretically, the empirical demonstration is limited to a simple toy
330 problem. In future work we aim to scale up the experiments to simulated robotics environments. We
331 also plan to adapt the algorithm to other imitation learning algorithms and to study learning objectives
332 that incentivize efficient exploration. A scaled-up algorithm for deconfounded imitation learning may
333 be an important stepping stone on the way to general learning algorithms for control problems, with
334 potentially large impact in the context of robotics and autonomous driving.

335 References

- 336 David Abel, Will Dabney, Anna Harutyunyan, Mark K Ho, Michael Littman, Doina Precup, and
337 Satinder Singh. On the expressivity of markov reward. *Advances in Neural Information Processing*
338 *Systems*, 34:7799–7812, 2021.
- 339 Anonymous. Causal imitation learning via inverse reinforcement learning. https://openreview.net/pdf?id=B-z41MBL_tH, 2022. URL https://openreview.net/pdf?id=B-z41MBL_tH.
340 Accessed: 2022-9-30.
- 342 Pim de Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. *Adv.*
343 *Neural Inf. Process. Syst.*, 32, 2019.
- 344 Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RI^2 : Fast
345 reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- 346 Yan Duan, Marcin Andrychowicz, Bradly C Stadie, Jonathan Ho, Jonas Schneider, Ilya Sutskever,
347 Pieter Abbeel, and Wojciech Zaremba. One-Shot imitation learning. March 2017.
- 348 Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforce-
349 ment learning. October 2017.
- 350 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
351 *information processing systems*, 29, 2016.
- 352 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
353 *arXiv:1312.6114*, 2013.
- 354 Daniel Kumor, Junzhe Zhang, and Elias Bareinboim. Sequential causal imitation learning
355 with unobserved confounders. October 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=o6-k168bBD8)
356 [o6-k168bBD8](https://openreview.net/forum?id=o6-k168bBD8).
- 357 Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*,
358 volume 1, page 2, 2000.
- 359 PA Ortega and DA Braun. A bayesian rule for adaptive control based on causal interventions. In *Third*
360 *Conference on Artificial General Intelligence (AGI 2010)*, pages 121–126. Atlantis Press, 2010a.
- 361 Pedro A Ortega and Daniel A Braun. A minimum relative entropy principle for learning and acting.
362 *Journal of Artificial Intelligence Research*, 38:475–511, 2010b.
- 363 Pedro A Ortega, Markus Kunesch, Grégoire Delétang, Tim Genewein, Jordi Grau-Moya, Joel Veness,
364 Jonas Buchli, Jonas Degraeve, Bilal Piot, Julien Perolat, Tom Everitt, Corentin Tallec, Emilio
365 Parisotto, Tom Erez, Yutian Chen, Scott Reed, Marcus Hutter, Nando de Freitas, and Shane Legg.
366 Shaking the foundations: delusions in sequence models for interaction and control. October 2021.
- 367 Judea Pearl. *Causality*. Cambridge University Press, September 2009. URL [https://www.](https://www.cambridge.org/core/books/causality/B0046844FAE10CBF274D4ACBDAEB5F5B)
368 [cambridge.org/core/books/causality/B0046844FAE10CBF274D4ACBDAEB5F5B](https://www.cambridge.org/core/books/causality/B0046844FAE10CBF274D4ACBDAEB5F5B).
- 369 Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Adv. Neural Inf. Process. Syst.*,
370 1988.
- 371 Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy
372 meta-reinforcement learning via probabilistic context variables. In *International conference on*
373 *machine learning*, pages 5331–5340. PMLR, 2019.
- 374 Danilo J Rezende, Ivo Danihelka, George Papamakarios, Nan Rosemary Ke, Ray Jiang, Theophane
375 Weber, Karol Gregor, Hamza Merzic, Fabio Viola, Jane Wang, et al. Causally correct partial
376 models for reinforcement learning. *arXiv preprint arXiv:2002.02836*, 2020.
- 377 Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
378 prediction to No-Regret online learning. In Geoffrey Gordon, David Dunson, and Miroslav
379 Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence*
380 *and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 627–635, Fort
381 Lauderdale, FL, USA, 2011. PMLR.

- 382 Stuart Russell. Learning agents for uncertain environments (extended abstract). In *Proceedings of*
383 *the eleventh annual conference on Computational learning theory - COLT' 98*, New York, New
384 York, USA, 1998. ACM Press.
- 385 S Schaal. Is imitation learning the route to humanoid robots? *Trends Cogn. Sci.*, 3(6):233–242, June
386 1999.
- 387 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
388 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 389 Jin Tian and Judea Pearl. *A general identification condition for causal effects*. eScholarship, University
390 of California, 2002. URL <https://www.aaai.org/Papers/AAAI/2002/AAAI02-085.pdf>.
- 391 Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos,
392 Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv*
393 *preprint arXiv:1611.05763*, 2016.
- 394 Wenxuan Zhou, Lerrel Pinto, and Abhinav Gupta. Environment probing interaction policies. *arXiv*
395 *preprint arXiv:1907.11740*, 2019.
- 396 Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and
397 Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning.
398 *arXiv preprint arXiv:1910.08348*, 2019.