

A Unifying Framework for Causal Imitation Learning with Hidden Confounders

Daqian Shao¹, Thomas Kleine Buening², Marta Kwiatkowska¹

¹Department of Computer Science, University of Oxford, UK

²The Alan Turing Institute, UK

Abstract

We propose a general and unifying framework for causal Imitation Learning (IL) with hidden confounders that subsumes several existing confounded IL settings from the literature. Our framework accounts for two types of hidden confounders: (a) those observed by the expert, which thus influence the expert’s policy, and (b) confounding noise hidden to both the expert and the IL algorithm. For additional flexibility, we also introduce a *confounding noise horizon* and time-varying expert-observable hidden variables. We show that causal IL in our framework can be reduced to a set of Conditional Moment Restrictions (CMRs) by leveraging trajectory histories as instruments to learn a history-dependent policy. We propose DML-IL, a novel algorithm that uses instrumental variable regression to solve these CMRs and learn a policy. We provide a bound on the imitation gap for DML-IL, which recovers prior results as special cases. Empirical evaluation on a toy environment with continuous state-action spaces and multiple Mujoco tasks demonstrate that DML-IL outperforms state-of-the-art causal IL algorithms.

1 Introduction

Imitation Learning (IL) has emerged as a prominent paradigm in machine learning, where the objective is to learn a policy that mimics the behaviour of an expert by learning from its demonstrations. While classical IL theory implies that, with infinite data, the IL error should converge to zero and the imitator should be value-equivalent to the expert (Ross et al., 2011, Spencer et al., 2021), it has been observed in practice that IL algorithms often produce incorrect estimates of the expert policy, leading to suboptimal and unsafe behaviours (Lecun et al., 2005, Codevilla et al., 2019, Bansal et al., 2018, Kuefler et al., 2017).

Various attempts have been made to explain the cause of these problems. Previous work has studied spurious correlations within the dataset (de Haan et al., 2019, Codevilla et al., 2019, Pfrommer et al., 2023), temporal noise (Swamy et al., 2022a), the case where the expert has additional information or knowledge (Swamy et al., 2022b, Vuorio et al., 2022, Chen et al., 2019, Choudhury et al., 2017), causal delusions (Ortega and Braun, 2008, Ortega et al., 2021) and covariate shifts (Spencer et al., 2021). However, each of these works considers only a specific aspect of the problem in different settings, and a holistic treatment of the problem and discussion of the connections between these settings is still missing from the literature. In the real world, it is often the case that multiple problems are present simultaneously (e.g., the expert has privileged information and the observed demonstrations are confounded). However, as we demonstrate later, addressing these problems partially or independently is insufficient, and a unified approach is necessary.

Our key observation is that, in fact, all the above settings can be formalised within a unifying framework by considering *hidden confounders*, which are variables present in the environment but not recorded in the demonstrations. Importantly, we distinguish between hidden confounders that can and cannot be observed by the expert.¹ When the expert fully observes the hidden confounders, but the imitator does not, additional information is available to the expert. When the expert also cannot observe the hidden confounders, these hidden confounders act as confounding noise that contaminates the demonstrations, causing spurious correlations and causal delusions. By considering the hidden confounders to include expert-observable and

¹Note that we always assume that the imitator does not observe the hidden confounders.

expert-unobservable parts, we propose a novel and unifying causal IL framework that can interpolate between the two scenarios. This framework not only unifies existing settings, but also enables us to consider a much larger family of problems that are more realistic in practice.

Building on our unifying framework, we aim to develop algorithms that can correctly and efficiently imitate the expert policy, even in the presence of hidden confounders. While interactive IL algorithms, such as DAgger (Ross et al., 2011), have shown promise in addressing specific issues with hidden confounders by allowing direct queries to the expert (Swamy et al., 2022b, Vuorio et al., 2022, Swamy et al., 2022a), this approach relies on access to an interactive expert. However, such an assumption is impractical in many real-world scenarios, where only a fixed set of demonstrations is available. For this reason, here we focus on designing methods that mitigate the effects of hidden confounders while relying solely on (a fixed set of) non-interactive expert demonstrations.

Our key idea is to leverage the trajectory histories as Instrumental Variables (IVs) to break the spurious correlations introduced by expert-unobservable hidden confounders. Moreover, by conditioning on the trajectory history, we can infer information about expert-observable hidden confounders and learn a history-dependent policy that accurately mimics the expert’s behaviour. Importantly, we show that, in our general framework, IL in the presence of hidden confounders can be reformulated as a set of Conditional Moment Restrictions (CMRs)—a well-studied problem in econometrics and causal inference. This reformulation allows us to design practical algorithms with theoretical guarantees on the imitation gap based on efficient algorithms from causal inference for solving CMRs.

Main Contributions.

- We propose a novel unifying framework for causal IL (Section 3). We consider hidden confounders that include expert-observable and expert-unobservable confounding variables to unify and generalise many of the settings in prior work.
- We show that IL in our framework can be reduced to a set of CMRs by leveraging trajectory histories as instruments to learn a history-dependent policy (Section 4).
- We provide a novel algorithm for causal IL (Algorithm 1): based on existing IV regression algorithms, we propose DML-IL to imitate the expert policy in our framework and provide an upper bound on the imitation gap that recovers previous results as special cases (Theorem 4.5).
- Empirically, we show that our algorithms perform well in challenging instances of our general setting, outperforming state-of-the-art methods (Section 5).

1.1 Related Works

Imitation Learning. Imitation learning considers the problem of learning from demonstrations (Pomerleau, 1988, Lecun et al., 2005). Standard IL methods include Behaviour Cloning (Pomerleau, 1988), Inverse RL (Russell, 1998), and adversarial methods (Ho and Ermon, 2016). Interactive IL (Ross et al., 2011) extends standard IL by allowing the imitator to query an interactive expert, facilitating recovery from mistakes. However, in this paper, we do not assume query access to an interactive expert.

Causal Imitation Learning. Recently, it has been shown that IL from offline trajectories can suffer from the existence of latent variables (Ortega et al., 2021), which cause causal delusion. This can be resolved by learning an interventional policy. Following this discovery, various methods (Vuorio et al., 2022, Swamy et al., 2022b) considered IL when the expert has access to the full hidden context that is fixed throughout each episode, but the imitator does not observe the hidden context. They aim to learn an interventional policy through on-policy IL algorithms that require an interactive demonstrator and/or an interactive simulator (e.g., DAgger (Ross et al., 2011)).

Orthogonal to these works, Swamy et al. (2022a) consider latent variables not known to the expert, which act as confounding noise that affects the expert policy, but not the transition dynamics. To address this challenge, the problem is then cast into an IV regression problem. Our work combines and generalises the

above works (Vuorio et al., 2022, Swamy et al., 2022b;a) to allow the latent variables to be only partly known to the expert, evolving through time in each episode and directly affecting both the expert policy and the transition dynamics. Solving this generalisation implies solving the above problems simultaneously.

Causal confusion (de Haan et al., 2019, Pfrommer et al., 2023) considers the situation where the expert’s actions are spuriously correlated with non-causal features of the previous observable states. While it is implicitly assumed that there are no latent variables present in the environment, we can still model this spurious correlation as the existence of hidden confounders that affect both previous states and current expert actions. Slight variations of this setting have been studied in Wen et al. (2020), Spencer et al. (2021), Codevilla et al. (2019). In Appendix A, we explain and discuss how these works can be reduced to special cases of our unifying framework.

From the causal inference perspective (Kumor et al., 2021, Zhang et al., 2020), there have been studies of the theoretical conditions on the causal graph such that the imitator can exactly match the expert performance through backdoor adjustments (*imitability*). Similarly, Ruan et al. (2023) extended imitation conditions and backdoor adjustments to inverse RL. We instead consider a setting where exact imitation is not possible and aim to minimise the imitation gap. Beyond backdoor adjustments, imitability has also been studied theoretically using context-specific independence relations (Jamshidi et al., 2023).

IV Regression and CMRs. In this paper, we transform our causal IL problem into solving a set of CMRs through IVs. Therefore, we briefly introduce IV regression and approaches for solving CMRs. The classic IV regression algorithms mainly consider linear functions (Angrist et al., 1996) and non-linear basis functions (Newey and Powell, 2003, Chen and Christensen, 2018, Singh et al., 2019). More recently, DNNs have been used for function estimation and methods such as DeepIV (Hartford et al., 2017), DeepGMM (Bennett et al., 2019a), AGMM (Dikkala et al., 2020), DFIV (Xu et al., 2020) and DML-IV (Shao et al., 2024) have been proposed.

More generally, IV regression algorithms can be generalised to solve CMRs (Liao et al., 2020, Dikkala et al., 2020, Shao et al., 2024), specifically linear CMRs, where the restrictions are linear functionals of the function of interest. In our paper, the derived CMRs for causal IL are linear, so the above methods can be adopted.

2 Preliminaries: Instrumental Variables and Conditional Moment Restrictions

We first introduce the concept of Instrumental Variables (IVs) and its connection to Conditional Moment Restrictions (CMRs). Consider a structural model that specifies some outcome Y given treatments X :

$$Y = f(X) + \varepsilon(U) \quad \text{with} \quad \mathbb{E}[\varepsilon(U)] = 0, \tag{1}$$

where U is a hidden confounder that affects both X and Y so that $\mathbb{E}[\varepsilon(U) | X] \neq 0$. Due to the presence of this hidden confounder, standard regressions (e.g., ordinary least squares) generally fail to produce consistent estimates of the causal relationship between X on Y , i.e., $\mathbb{E}[Y | do(X)] = f(X)$, where $do(\cdot)$ is the interventional operator (Pearl, 2000). If we only have observational data, a classic technique for learning f is IV regression (Newey and Powell, 2003). An IV Z is an observable variable that satisfies the following conditions:

- *Unconfounded Instrument*: $Z \perp\!\!\!\perp U$;
- *Relevance*: $\mathbb{P}(X|Z)$ is not constant in Z ;
- *Exclusion*: Z does not directly affect Y : $Z \perp\!\!\!\perp Y | (X, U)$.

Using IVs, we are able to formulate the problem of learning f into a set of CMRs (Dikkala et al., 2020), where we aim to solve for f satisfying

$$\mathbb{E}[Y - f(X) | Z] = 0.$$

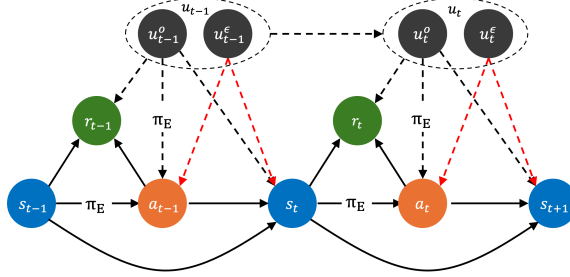


Figure 1: A causal graph of MDPs with hidden confounders, where at each time step the hidden confounder is $u_t = (u_t^o, u_t^e)$. The black dotted lines represent the causal effect of the expert-observable confounder u_t^o , which directly affects a_t because the expert policy can observe u_t^o . It also directly affects s_{t+1} and r_t because otherwise it is irrelevant to the expected return and there is no reason for the expert to consider it. The red dotted lines represent the causal effect of u_t^e that is not observable by the expert, which acts as confounding noise and directly affects the states and actions. u_t^e does not directly affect r_t (following Swamy et al. (2022a)) because the expert policy does not take u_t^e into account, and letting u_t^e directly affect r_t would only add noise to the expected return.

In our work, we show that trajectory histories can be used as instruments to learn the causal relationship between states and expert actions by transforming the problem of causal IL into a set of CMRs.

3 A Unifying Framework for Causal Imitation Learning

MDPs with Hidden Confounders. In this section, we introduce a novel unifying framework for causal IL in the presence of hidden confounders. We begin by introducing an Markov Decision Process (MDP) with hidden confounders, $(\mathcal{S}, \mathcal{A}, \mathcal{U}, \mathcal{P}, r, \mu_0, T)$, where \mathcal{S} is the state space, \mathcal{A} is the action space and \mathcal{U} is the confounder space. Importantly, parts of the hidden confounders u_t may be available to the expert due to imperfect environment logging and expert knowledge. We model this by segmenting the hidden confounder into two parts $u_t = (u_t^o, u_t^e)$, where u_t^o are observable to the expert and u_t^e are not. Intuitively, u_t^o are additional information that only the expert observes and u_t^e behave as confounding noise in the environment that affects both the state and action.² As a result, the transition function $\mathcal{P}(\cdot \mid s, a, (u^o, u^e))$ depends on both hidden confounders but the reward function $r(s, a, u^o)$ only depends on the state, action and the observable confounder u^o since the confounding noise only directly affects the state and actions. Finally, μ_0 is the initial state distribution and T is the horizon of the problem. We provide a causal graph that illustrates the relationships between variables in Figure 1.

Causal Imitation Learning. We assume that an expert is demonstrating a task following some expert policy π_E (which we will specify later) and we observe a set of $N \geq 1$ expert demonstrations $\{d_1, d_2, \dots, d_N\}$. Each demonstration is a state-action trajectory $(s_1, a_1, \dots, s_T, a_T)$, where, at each time step, we observe the state s_t and the action a_t taken in the environment, and the trajectory follows the transition function $\mathcal{P}(\cdot \mid s_t, a_t, (u_t^o, u_t^e))$. Denote $h_t = (s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t) \in \mathcal{H}$ as the trajectory history at time t , where $\mathcal{H} \subseteq \bigcup_{i=0}^{T-1} (\mathcal{S} \times \mathcal{A})^i \times \mathcal{S}$ is the set of all possible trajectory histories at different time steps. Importantly, we do not observe the reward and the sequence of confounders (u_t^o, u_t^e) .

Given the observed trajectories, our goal is to learn a history-dependent policy $\pi_h : \mathcal{H} \rightarrow \Delta(\mathcal{A})$. We assume that our policy class Π is convex and compact. The Q -function of a policy $\pi_h \in \Pi$ is defined as

²In our framework, we allow the actual actions taken in the environment to be affected by the noise. Noise that only perturbs data records can be considered as a special case of our framework.

$Q_\pi(s_t, a_t, u_t^o) = \mathbb{E}_{\tau \sim \pi_h} [\sum_{t'=t}^T r(s_{t'}, a_{t'}, u_{t'}^o)]$ and the value of a policy is $J(\pi) = \mathbb{E}_{\tau \sim \pi_h} [\sum_{t'=t}^T r(s_{t'}, a_{t'}, u_{t'}^o)]$, where τ is the trajectory following π_h .

This nuanced distinction between u_t^o and u_t^ε is crucial for determining the appropriate method for IL, and we begin with an example to motivate our setting and illustrate the importance of considering $u_t = (u_t^o, u_t^\varepsilon)$.

Example 3.1. Consider a dynamic aeroplane ticket pricing scenario (Wright, 1928, Hartford et al., 2017), where we would like to learn a ticket pricing policy by imitating actual airline prices based on the profit margins set by experts. We have access to information such as destinations, flight time, previous sales, and aeroplane records. However, seasonal demand patterns and events are known to the experts, but are not logged in the dataset. Hence, these hidden confounders are included in u_t^o . In contrast, the experts only determine profit margins as an action because the observed price is confounded (additively) by fluctuations in operating costs such as fuel price and maintenance costs, which are unknown to the expert when they set the profit margin. These hidden confounders act as u_t^ε in our setting. Without an explicit consideration of u_t^o and u_t^ε separately, learning algorithms cannot distinguish between u_t^o and u_t^ε and fail to correctly imitate the expert. We perform experiments on this toy example in Section 5.

In order to learn a policy π_h that can match the performance of π_E , we need to break the spurious correlation between states and expert actions. In fact, this is a causal inference problem of learning the counterfactual expert decision:

$$\begin{aligned} \mathbb{E}[a_t \mid do(s_t, u_t^o)] &= \pi_E(s_t, u_t^o) + \mathbb{E}[u_t^\varepsilon \mid do(s_t, u_t^o)] \\ &= \pi_E(s_t, u_t^o), \end{aligned}$$

which describes what the expert would do if we intervened and placed them in state s_t when observing u_t^o . Here, we assume the confounding noise u_t^ε is additively to the action and zero expectation, which is later formalised and explained in Assumption 3.3. Unfortunately, from the causal inference literature (Shpitser and Pearl, 2008, Pearl, 2000), we know that, without further assumptions, it is impossible to identify π_E .

To determine the minimal assumptions that allow π_E to be identifiable, we first observe that u_t^ε can be correlated for all t , making it impossible to distinguish between the intended actions of the expert and the confounding noise. However, in practice, the effect of confounding noise u_t^ε at t may diminish over time (e.g., wind) or the u_t^ε becomes observable at a future time (e.g, operating costs), which makes the confounding noise at time t and the confounding noise at some future time step t' independent. We formalise this novel notion in terms of a confounding noise horizon k :

Assumption 3.2 (Confounding Horizon). For every t , the confounding noise u_t^ε has a horizon of k where $1 \leq k < T$. More formally, $u_t^\varepsilon \perp\!\!\!\perp u_{t-k}^\varepsilon \forall t > k$.

In addition, we make a standard assumption in causal inference (Pearl, 2000, Hartford et al., 2017, Shao et al., 2024), namely that the confounding noise is additive to the action.

Assumption 3.3 (Additive Noise). The structural equation that generates the actions in the observed trajectories is

$$a_t = \pi_E(s_t, u_t^o) + u_t^\varepsilon, \tag{2}$$

where we assume the confounding noise u_t^ε is additive to a_t and WLOG $\mathbb{E}[u_t^\varepsilon] = 0$ because any non-zero expectation of u_t^ε can be included as a constant in π_E .

Without this additive noise assumption, the causal effect becomes unidentifiable (see, e.g., Balke and Pearl (1994)) and the best we can do is to upper/lower bound it. Next, we show that, with the above two assumptions, it becomes possible to identify the true causal relationship between states and expert actions, and imitate π_E .

4 Causal IL as CMRs

In this section, we demonstrate that performing causal IL in our framework is possible using trajectory histories as instruments. In the next step, we show that the problem can be described as CMRs and propose an effective algorithm to solve it.

The typical target for IL would be the expert policy π_E itself. However, since the expert has access to information, namely u_t^o , which the imitator does not, the best thing an imitator can do is to learn a history-dependent policy π_h that is the closest to the expert. A natural choice is the conditional expectation of $\pi_E(s_t, u_t^o)$ on the history h_t :

$$\pi_h(h_t) := \mathbb{E}_{\mathbb{P}(u_t^o|h_t)}[\pi_E(s_t, u_t^o)] = \mathbb{E}[\pi_E(s_t, u_t^o) | h_t],$$

because the conditional expectation minimizes the least squares criterion (Hastie et al., 2001) and π_h is the best predictor of π_E given h_t . In π_h , the distribution $\mathbb{P}(u_t^o | h_t)$ captures the information about u_t^o that can be inferred from trajectory histories.

Remark 4.1. *Learning π_h is not trivial. Policies learnt naively using behaviour cloning (i.e., $\mathbb{E}[a_t | h_t]$) fail to match π_E . In view of Equation (2), we have that*

$$\begin{aligned} \mathbb{E}[a_t | h_t] &= \mathbb{E}[\pi_E(s_t, u_t^o) | h_t] + \mathbb{E}[u_t^\varepsilon | h_t] \\ &= \pi_h(h_t) + \mathbb{E}[u_t^\varepsilon | h_t], \end{aligned} \tag{3}$$

where $\mathbb{E}[u_t^\varepsilon | h_t] \neq 0$ due to the spurious correlation between u_t^ε and the trajectory history h_t . As a result, $\mathbb{E}[a_t | h_t]$ becomes biased, which can lead to arbitrarily worse performance compared to π_E .

Derivation of CMRs. Leveraging the confounding horizon from Assumption 3.2, it becomes possible to break the spurious correlation using the independence of u_t^ε and u_{t-k}^ε . We propose to use the k -step trajectory history $h_{t-k} = (s_1, a_1, \dots, s_{t-k})$ as an instrument for the current state s_t . Taking the expectation conditional on h_{t-k} in Equation (3) yields

$$\begin{aligned} \mathbb{E}[a_t | h_{t-k}] &= \mathbb{E}[\mathbb{E}[a_t | h_t] | h_{t-k}] \\ &= \mathbb{E}[\pi_h(h_t) | h_{t-k}] + \mathbb{E}[\mathbb{E}[u_t^\varepsilon | h_t] | h_{t-k}] \\ &= \mathbb{E}[\pi_h(h_t) | h_{t-k}] + \mathbb{E}[u_t^\varepsilon | h_{t-k}] \end{aligned}$$

where we use the fact that h_{t-k} is $\sigma(h_t)$ -measurable because $h_{t-k} \subseteq h_t$. Next, recall that $u_t^\varepsilon \perp\!\!\!\perp u_{t-k}^\varepsilon$ by Assumption 3.2, which implies $u_t^\varepsilon \perp\!\!\!\perp h_{t-k}$, so that

$$\begin{aligned} \mathbb{E}[a_t | h_{t-k}] &= \mathbb{E}[\pi_h(h_t) | h_{t-k}] + \mathbb{E}[u_t^\varepsilon] \\ &= \mathbb{E}[\pi_h(h_t) | h_{t-k}]. \end{aligned} \tag{4}$$

As a result, the problem of learning π_h reduces to solving for π_h that satisfies the following identity

$$\mathbb{E}[a_t - \pi_h(h_t) | h_{t-k}] = 0, \tag{5}$$

which is a CMR problem as defined in Section 2. In this case, both a_t and h_t are observed in the confounded expert demonstrations, and h_{t-k} acts as the instrument.

To make sure the instrument h_{t-k} is valid, we check that it satisfies the conditions of Section 2. Firstly, we have checked that $u_t^\varepsilon \perp\!\!\!\perp h_{t-k}$. Secondly, the environment and the expert policy are non-trivial, which means $\mathbb{P}(h_t | h_{t-k})$ is not constant in h_{t-k} . Finally, h_{t-k} indeed only affects a_t through s_t by the Markovian property. However, the strength of the instrument, which informally represents the correlation between the instrument h_{t-k} and h_t , plays an important role in how well we can identify $\pi_h(h_t)$ by solving the CMRs in Equation (5). In particular, we see that, as the confounding horizon k increases, the correlation between h_{t-k} and h_t weakens and h_{t-k} becomes a weaker instrument. This means that it is less able to identify π_h via the CMR in Equation (5) and the final learnt imitator will have poorer performance. This is confirmed theoretically in Proposition 4.3 and experimentally in Section 5, and we will formalise this notion of instrument strength in Section 4.2.

Algorithm 1 DML-IL

- 1: **input** Dataset \mathcal{D}_E of expert demonstrations, Confounding noise horizon k
 - 2: Initialize the roll-out model \hat{M} as a Gaussian mixture model
 - 3: **repeat**
 - 4: Sample (h_t, a_t) from data \mathcal{D}_E
 - 5: Fit the roll-out model $(h_t, a_t) \sim \hat{M}(h_{t-k})$ to maximize the log likelihood
 - 6: **until** convergence
 - 7: Initialize the expert model $\hat{\pi}_h$ as a neural network
 - 8: **repeat**
 - 9: Sample h_{t-k} from \mathcal{D}_E
 - 10: Generate \hat{h}_t and \hat{a}_t using the roll-out model \hat{M}
 - 11: Update $\hat{\pi}_h$ to minimise the loss $\ell := \|\hat{a}_t - \hat{\pi}_h(\hat{h}_t)\|_2$
 - 12: **until** convergence
 - 13: **return** A history-dependent imitator policy $\hat{\pi}_h$
-

4.1 Practical Algorithms for Solving the CMRs

There are various techniques (Shao et al., 2024, Bennett et al., 2019b, Xu et al., 2020, Dikkala et al., 2020) for solving the CMRs $\mathbb{E}[a_t|h_{t-k}] = \mathbb{E}[\pi_h(h_t)|h_{t-k}]$. Here, the *CMR error* that we aim to minimise is given by

$$\sqrt{\mathbb{E}[\mathbb{E}[a_t - \hat{\pi}_h(h_t)|h_{t-k}]^2]} = \|\mathbb{E}[a_t - \hat{\pi}_h(h_t)|h_{t-k}]\|_2.$$

In Algorithm 1, we introduce DML-IL, an algorithm adapted from the IV regression algorithm DML-IV (Shao et al., 2024)³, which solves our CMRs by minimising the CMR error. The first part of the algorithm (line 3-7) learns a roll-out model \hat{M} that generates a trajectory k steps ahead given h_{t-k} . Then, the roll-out model \hat{M} is used to train the policy model $\hat{\pi}_h$ (line 8-13). $\hat{\pi}_h$ takes the generated trajectory \hat{h}_t from $\hat{M}(h_{t-k})$ as inputs, and minimises the mean squared error to the next action. Using generated trajectories is crucial in breaking the spurious correlation caused by u_t^ε between past states and actions, and using the trajectory history before h_{t-k} allows the imitator to infer information about u_t^o .

DML-IL can also be implemented with K -fold cross-fitting, where the dataset is partitioned into K folds, with each fold alternately used to train $\hat{\pi}_h$ and the remaining folds to train \hat{M} . This ensures unbiased estimation and improves the stability of training. The base IV algorithm DML-IV with K -fold cross-fitting is theoretically shown to converge at the rate of $O(N^{-1/2})$ (Shao et al., 2024), where N is the sample size, under regularity conditions. DML-IL with K -fold cross-fitting (see Appendix D.1 for details) will thus inherit this convergence rate guarantee.

Note that Algorithm 1 requires the confounding noise horizon k as input. While the exact value of k can be difficult to obtain in reality, any upper bound \bar{k} of k is sufficient to guarantee the correctness of Algorithm 1, since $h_{t-\bar{k}}$ is also a valid instrument. Ideally, we would like a data-driven approach to determine k . Unfortunately, it is generally intractable to empirically verify whether h_{t-k} is a valid instrument from a static dataset, especially the unconfounded instrument condition (i.e., $h_{t-k} \perp\!\!\!\perp u_t^\varepsilon$). Therefore, we rely on the user to provide a sensible choice of \bar{k} based on the environment that does not substantially overestimate k .

4.2 Theoretical Analysis

In this section, we derive theoretical guarantees for our algorithm, focusing on the imitation gap and its relationship with existing work.

On a high level, in order to bound the imitation gap of the learnt policy $\hat{\pi}_h$, i.e., $J(\pi_E) - J(\hat{\pi}_h)$, we need to control:

³DML stands for double machine learning (Chernozhukov et al., 2018), which is a statistical technique to ensure fast convergence rate for two-step regression, as is the case in Algorithm 1.

- (i) The amount of information about the hidden confounders that can be inferred from trajectory histories;
- (ii) The ill-posedness (or identifiability) of the set of CMRs, which intuitively measures the strength of the instrument h_{t-k} ;
- (iii) The disturbance of the confounding noise to the states and actions at test time.

These factors are all determined by the environment and the expert policy. To control (i), we measure how much information about u_t^o is captured by the trajectory history h_t by analysing the Total Variation (TV) distance between the distribution of u_t^o and $\mathbb{E}[u_t^o|h_t]$ along the trajectories of π_E . To control (ii) and (iii), we need to introduce the following two key concepts.

Definition 4.2 (The ill-posedness of CMRs (Dikkala et al., 2020, Chen and Pouzo, 2012)). Given the derived CMRs in Equation (5), for a policy $\pi \in \Pi$, $\|\pi_E - \pi\|_2$ is the root mean squared error to the expert and $\|\mathbb{E}[a_t - \pi(s_t)|s_{t-k}]\|_2$ is the CMR error we aim to minimise. Then, the *ill-posedness* $\nu(\Pi, k)$ of the policy space with confounding noise horizon k is given by

$$\nu(\Pi, k) = \sup_{\pi \in \Pi} \frac{\|\pi_E - \pi\|_2}{\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2}.$$

The ill-posedness $\nu(\Pi, k)$ measures the strength of the instrument where a higher $\nu(\Pi, k)$ indicates a weaker instrument. It bounds the ratio between the learning error of the imitator following our CMR objective and its L_2 error to the expert policy.

As discussed previously, intuitively, the strength of the instrument would decrease as the confounding horizon k increases. This is in fact true and is confirmed by the following proposition. The proof is deferred to Appendix B.1.

Proposition 4.3. *The ill-posedness $\nu(\Pi, k)$ is monotonically increasing as the confounded horizon k increases.*

Next, we introduce the notion of c-TV stability.

Definition 4.4 (c-total variation stability (Bassily et al., 2021, Swamy et al., 2022a)). Let $P(X)$ be the distribution of a random variable $X : \Omega \rightarrow \mathcal{X}$. $P(X)$ is c-TV stable if for $a_1, a_2 \in \mathcal{X}$ and $\Delta > 0$,

$$\|a_1 - a_2\| \leq \Delta \implies \delta_{TV}(a_1 + X, a_2 + X) \leq c\Delta.$$

where $\|\cdot\|$ is some norm defined on \mathcal{X} and δ_{TV} is the total variation distance.

A wide range of distributions are c-TV stable. For example, standard normal distributions are $\frac{1}{2}$ -TV stable. We apply this notion to the distribution over u_t^ε to bound the disturbance it induces in the trajectory and the expected return.

With the notion of ill-posedness and c-TV stability, we can now analyse and upper bound the imitation gap $J(\pi_E) - J(\hat{\pi}_h)$ by controlling the three components (i) – (iii) discussed above. The full proof is deferred to Appendix B.2.

Theorem 4.5 (Imitation Gap Bound). *Let $\hat{\pi}_h$ be the learnt policy with CMR error ε and let $\nu(\Pi, k)$ be the ill-posedness of the problem. Assume that $\delta_{TV}(u_t^o, \mathbb{E}_{\pi_E}[u_t^o|h_t]) \leq \delta$ for $\delta \in \mathbb{R}^+$, $P(u_t^\varepsilon)$ is c-TV stable and π_E is deterministic. Then, the imitation gap is upper bounded by*

$$J(\pi_E) - J(\hat{\pi}_h) \leq T^2(c\varepsilon\nu(\Pi, k) + 2\delta) = \mathcal{O}(T^2(\delta + \varepsilon)).$$

This upper bound scales at the rate of T^2 , which aligns with the expected behaviour of imitation learning without an interactive expert (Ross and Bagnell, 2010). Next, we show that the upper bounds on the imitation gap from prior work (Swamy et al., 2022a;b) are special cases of Theorem 4.5. The proofs are deferred to Appendix B.3.

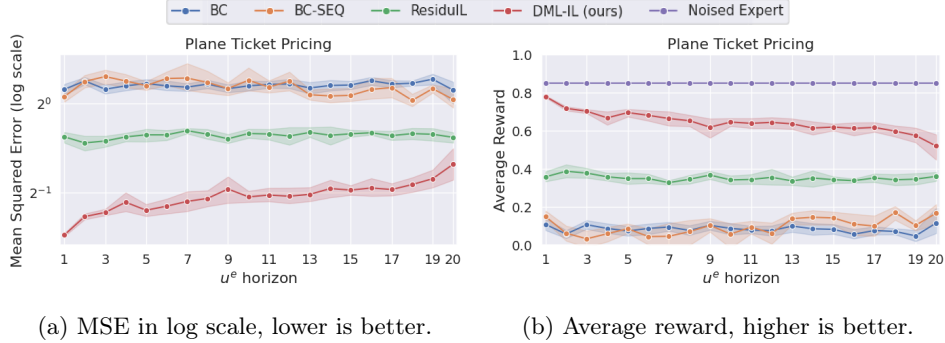


Figure 2: The MSE between the learnt policy and the expert, and the average reward, in the plane ticket environment (Example 3.1).

Corollary 4.6. *In the special case that $u_t^o = 0$, i.e., there are no expert-observable confounders, or $u_t^o = \mathbb{E}_{\pi_E}[u_t^o|h_t]$, i.e., u_t^o is $\sigma(h_t)$ measurable (all information about u_t^o is contained in the history), the imitation gap is upper bounded by*

$$J(\pi_E) - J(\hat{\pi}_h) \leq T^2(c\varepsilon\nu(\Pi, k)) = \mathcal{O}(T^2\varepsilon),$$

which coincides with Theorem 5.1 of Swamy et al. (2022a).

When there are no hidden confounders, i.e., $u_t^\varepsilon = 0$, our framework is reduced to that of Swamy et al. (2022b). However, Swamy et al. (2022b) provided an abstract bound that directly uses the supremum of key components in the imitation gap over all possible Q functions to bound the imitation gap. We further extend and concretise the bound using the learning error ε and the TV distance bound δ instead of relying on the suprema.

Corollary 4.7. *In the special case that $u_t^\varepsilon = 0$, if the learnt policy has optimisation error ε , the imitation gap is upper bounded by*

$$J(\pi_E) - J(\hat{\pi}_h) \leq T^2 \left(\frac{2}{\sqrt{\dim(A)}} \varepsilon + 2\delta \right),$$

which is a concrete bound that extends the abstract bound in Theorem 5.4 of Swamy et al. (2022b).

Remark 4.8. *If both u_t^ε and u_t^o are zero, we then recover the classic setting of IL without confounders (Ross and Bagnell, 2010), and the imitation gap bound is $T^2\varepsilon$, where ε is the optimisation error of the algorithm.*

5 Experiments

In this section, we empirically evaluate the performance of Algorithm 1 (DML-IL) on the toy environment with continuous state and action spaces introduced in Example 3.1 and Mujoco environments: Ant, Half Cheetah and Hopper. We compare with the following existing methods: Behavioural Cloning (BC), which naively minimises $\mathbb{E}[-\log \pi(a_t|s_t)]$; BC-SEQ (Swamy et al., 2022b), which learns a history-dependent policy to handle hidden contexts observable to the expert; ResiduIL (Swamy et al., 2022a), which we adapt to our setting by providing h_{t-k} as instruments to learn a history-independent policy; and the noised expert, which is the performance of the expert when put in the confounded environment, and corresponds to the maximally achievable performance. In Appendix E, we also include additional evaluations of using other IV regression algorithms, including DFIV (Xu et al., 2020) and DeepGMM (Bennett et al., 2019a), as the core CMR solver but found inconsistent and subpar performance.

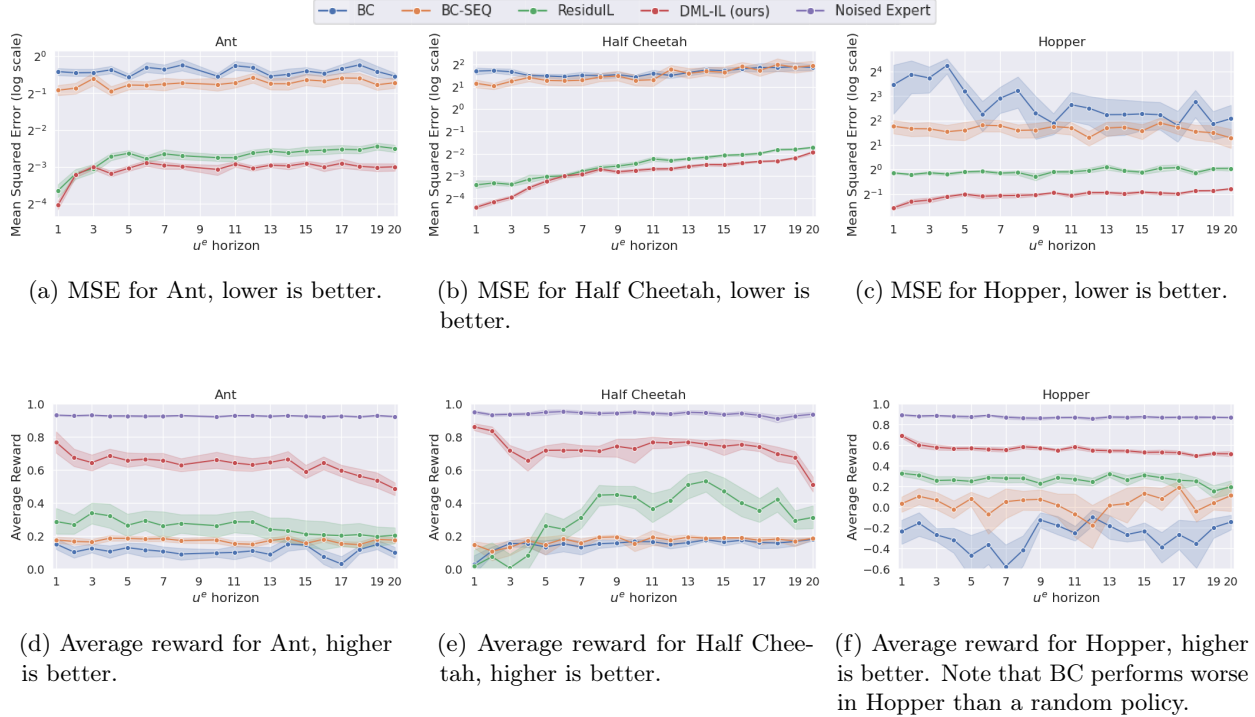


Figure 3: The MSE between the learnt policy and expert, and the average reward, in Mujoco environments.

We train imitators with 20000 samples (40 trajectories of 500 steps each) of the expert trajectory using each algorithm and report the average reward when tested online in their respective environments. The reward is scaled such that 1 is the performance of the un-noised expert, and 0 is a random policy. We also report the Mean Squared Error (MSE) between imitator’s and expert’s actions. The purpose of evaluating the MSE is to assess how well the imitator learnt from the expert, and importantly whether the confounding noise problem is mitigated. When the confounding noise u^ε is not handled, we should expect to observe a much higher MSE. All results are plotted with one standard deviation as a shaded area. In addition, we vary the confounding noise horizon k from 1 to 20 in order to increase the difficulty of the problem with weaker instruments h_{t-k} .

5.1 Plane Ticket Pricing Environment

Experimental Setup. We first consider the plane ticket pricing environment described in Example 3.1. The confounding noise u^ε are operation costs and u_t^o are seasonal demand patterns and events. We set u_t^o to continuously vary with a rate of change of approximately every 30 steps. Details on this environment are provided in Appendix C.1.

Results. The results are presented in Figure 2. DML-IL performed the best with the lowest MSE and the highest average reward that is closest to the expert, especially when the u_t^ε horizon is 1. This implies that DML-IL is successful in handling both u_t^ε and u_t^o . ResiduIL is able to reduce the confounding effect of u_t^ε evident by the lower MSE compared to the two other methods that do not deal with u_t^ε . However, since it does not explicitly consider u_t^o , the imitator has no information on u_t^o and the best it can do is to assume some average value (or expectation) of u_t^o . Therefore, while ResiduIL still achieves some reward, its performance gap to DML-IL is due to the lack of consideration of u_t^o . Both BC and BC-SEQ fail completely in the presence of confounding noise u_t^ε , with orders of magnitude higher MSE and average reward close to a random policy. From the similar performance of BC-SEQ and BC, we see that the use of trajectory

history to infer u_t^o is not helpful when the confounding noise is not handled explicitly. This demonstrates that considering the effect of u_t^ε and u_t^o partially is insufficient to learn a well performing imitator under the general setting.

In addition, as the confounding noise horizon k increases (x-axis), the performance of DML-IL drops. This coincides with the fact that the instrument is weaker and less information regarding u_t^o can be captured from h_{t-k} as k increases. When $k = 20$, it can be seen that the performance of DML-IL is close to that of ResiduIL, which does not consider the effect of u_t^o , because very limited information about the current expert-observable confounder u_t^o can be inferred from history 20 steps ago.

5.2 Mujoco Environments

Experimental Setup. In Figure 3, we consider the Mujoco tasks. While the original tasks do not have hidden variables, we modify the environment to introduce u_t^ε and u_t^o . Specifically, instead of controlling the ant, half cheetah and hopper, respectively, to travel as fast as possible, the goal is to control the agent to travel at a target speed that is varying throughout an episode. This target speed is u_t^o , which is observed by the expert but not recorded in the dataset. In addition, we add confounding noise u_t^ε to s_t and a_t to mimic the confounding noise such as wind. Additional details about the modification made to the environments are provided in Appendix C.2.

Results. DML-IL outperforms other methods in all three Mujoco environments as shown in Figure 3. Similarly to the plane ticket environment, ResiduIL is effective in removing the confounding noise but fails to match the average reward of DML-IL due to the lack of knowledge of u_t^o . BC and BC-SEQ have much higher MSE and fail to learn meaningful policies. As u_t^ε horizon increases, the performance of DML-IL drops as expected due to weaker instruments and limited inferrable information regarding u_t^o , especially for the Ant and Half Cheetah tasks.

6 Conclusion

In this paper, we proposed an unifying framework for confounded IL with hidden confounders that unifies and extends previous confounded IL settings. Specifically, we considered hidden confounders to be partially observable to the expert, and demonstrated that causal IL under this framework can be reduced to a set of CMRs with the trajectory histories as instruments. We proposed DML-IL, a novel algorithm to solve these CMRs and learn an imitator. We provided bounds on the imitation gap for the learnt imitator. Finally, we empirically evaluated DML-IL on multiple tasks including Mujoco environments and demonstrated state-of-the-art performance against other causal IL algorithms.

Acknowledgments

This work was supported by the EPSRC Prosperity Partnership FAIR (grant number EP/V056883/1). DS acknowledges funding from the Turing Institute and Accenture collaboration. MK receives funding from the ERC under the European Union’s Horizon 2020 research and innovation programme (FUN2MODEL, grant agreement No. 834115).

References

- Angrist, J. D., Imbens, G. W., and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91:444–455.
- Balke, A. and Pearl, J. (1994). Counterfactual probabilities: Computational methods, bounds and applications. *Uncertainty Proceedings 1994*, pages 46–54.

- Bansal, M., Krizhevsky, A., and Ogale, A. (2018). Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. *Robotics: Science and Systems*, 5:5986.
- Bassily, R., Steinke, T., Nissim, K., Stemmer, U., Smith, A., and Ullman, J. (2021). Algorithmic stability for adaptive data analysis. *SIAM Journal on Computing*, pages 1046–1059.
- Bennett, A., Kallus, N., and Schnabel, T. (2019a). Deep generalized method of moments for instrumental variable analysis. *Advances in Neural Information Processing Systems*, 32.
- Bennett, A., Kallus, N., and Schnabel, T. (2019b). Deep generalized method of moments for instrumental variable analysis. *Advances in Neural Information Processing Systems*, 32.
- Chen, D., Zhou, B., Koltun, V., and Krähenbühl, P. (2019). Learning by cheating. *Proceedings of Machine Learning Research*, 100:66–75.
- Chen, X. and Christensen, T. M. (2018). Optimal sup-norm rates and uniform inference on nonlinear functionals of nonparametric iv regression. *Quantitative Economics*, 9:39–84.
- Chen, X. and Pouzo, D. (2012). Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals. *Econometrica*, 80:277–321.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68.
- Choudhury, S., Bhardwaj, M., Arora, S., Kapoor, A., Ranade, G., Scherer, S., and Dey, D. (2017). Data-driven planning via imitation learning. *International Journal of Robotics Research*, 37:1632–1672.
- Codevilla, F., Santana, E., Lopez, A., and Gaidon, A. (2019). Exploring the limitations of behavior cloning for autonomous driving. *Proceedings of the IEEE International Conference on Computer Vision*, 2019-October:9328–9337.
- de Haan, P., Jayaraman, D., and Levine, S. (2019). Causal confusion in imitation learning. *Advances in Neural Information Processing Systems*, 32.
- Dikkala, N., Lewis, G., Mackey, L., and Syrgkanis, V. (2020). Minimax estimation of conditional moment models. *Advances in Neural Information Processing Systems*, 2020-December.
- Hartford, J., Lewis, G., Leyton-Brown, K., and Taddy, M. (2017). Deep IV: A flexible approach for counterfactual prediction. *Proceedings of the 34th International Conference on Machine Learning*.
- Hastie, T., Tibshirani, R., and Friedman, J. (2001). *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA.
- Ho, J. and Ermon, S. (2016). Generative adversarial imitation learning. *Advances in Neural Information Processing Systems*, pages 4572–4580.
- Jamshidi, F., Akbari, S., and Kiyavash, N. (2023). Causal imitability under context-specific independence relations.
- Kakade, S. and Langford, J. (2002). Approximately optimal approximate reinforcement learning. *International Conference on Machine Learning*.
- Kuefler, A., Morton, J., Wheeler, T., and Kochenderfer, M. (2017). Imitating driver behavior with generative adversarial networks. *IEEE Intelligent Vehicles Symposium, Proceedings*, 5:204–211.
- Kumor, D., Zhang, J., and Bareinboim, E. (2021). Sequential causal imitation learning with unobserved confounders. *Proceedings of the 35th Conference on Neural Information Processing Systems*.

- Lecun, Y., Muller, U., Ben, J., Cosatto, E., and Flepp, B. (2005). Off-road obstacle avoidance through end-to-end learning. *Advances in Neural Information Processing Systems*, 18.
- Liao, L., Chen, Y. L., Yang, Z., Dai, B., Wang, Z., and Kolar, M. (2020). Provably efficient neural estimation of structural equation model: An adversarial approach. *Advances in Neural Information Processing Systems*, 2020-December.
- Newey, W. K. and Powell, J. L. (2003). Instrumental variable estimation of nonparametric models. *Econometrica*, 71:1565–1578.
- Ortega, P. A. and Braun, D. A. (2008). A minimum relative entropy principle for learning and acting. *Journal of Artificial Intelligence Research*, 38:475–511.
- Ortega, P. A., Kunesch, M., Delétang, G., Genewein, T., Grau-Moya, J., Veness, J., Buchli, J., Degraeve, J., Piot, B., Perolat, J., Everitt, T., Tallec, C., Parisotto, E., Erez, T., Chen, Y., Reed, S., Hutter, M., de Freitas, N., and Legg, S. (2021). Shaking the foundations: delusions in sequence models for interaction and control.
- Pearl, J. (2000). Causality: Models, reasoning, and inference. *Econometric Theory*.
- Pfrommer, S., Bai, Y., Lee, H., and Sojoudi, S. (2023). Initial state interventions for deconfounded imitation learning.
- Pomerleau, D. A. (1988). Alvin: An autonomous land vehicle in a neural network. *Advances in Neural Information Processing Systems*, 1.
- Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N. (2021). Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8.
- Ross, S. and Bagnell, J. A. (2010). Efficient reductions for imitation learning. *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*.
- Ross, S., Gordon, G. J., and Bagnell, J. A. (2011). A reduction of imitation learning and structured prediction to no-regret online learning. *Journal of Machine Learning Research*, 15:627–635.
- Ruan, K., Zhang, J., Di, X., and Bareinboim, E. (2023). Causal imitation learning via inverse reinforcement learning. *Proceedings at the International Conference on Learning Representations*.
- Russell, S. (1998). Learning agents for uncertain environments (extended abstract). *In The Eleventh Annual Conference on Computational Learning Theory*.
- Shao, D., Soleymani, A., Quinzan, F., and Kwiatkowska, M. (2024). Learning decision policies with instrumental variables through double machine learning. *Proceedings of the International Conference on Machine Learning*.
- Shpitser, I. and Pearl, J. (2008). Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research*, 9:1941–1979.
- Singh, R., Sahani, M., and Gretton, A. (2019). Kernel instrumental variable regression. *Advances in Neural Information Processing Systems*, 32.
- Spencer, J., Choudhury, S., Venkatraman, A., Ziebart, B., and Bagnell, J. A. (2021). Feedback in imitation learning: The three regimes of covariate shift.
- Swamy, G., Choudhury, S., Bagnell, J. A., and Wu, Z. S. (2022a). Causal imitation learning under temporally correlated noise. *Proceedings of Machine Learning Research*, 162:20877–20890.

- Swamy, G., Choudhury, S., Bagnell, J. A., and Wu, Z. S. (2022b). Sequence model imitation learning with unobserved contexts. *Advances in Neural Information Processing Systems*, 35.
- Vuorio, R., Brehmer, J., Ackermann, H., Dijkman, D., Cohen, T., and de Haan, P. (2022). Deconfounded imitation learning.
- Wen, C., Lin, J., Darrell, T., Jayaraman, D., and Gao, Y. (2020). Fighting copycat agents in behavioral cloning from observation histories. *Advances in Neural Information Processing Systems*, 2020-December.
- Wright, P. G. (1928). The tariff on animal and vegetable oils. <https://doi.org/10.1086/254144>, 38:619–620.
- Xu, L., Chen, Y., Srinivasan, S., de Freitas, N., Doucet, A., and Gretton, A. (2020). Learning deep features in instrumental variable regression. *ICLR 2021 - 9th International Conference on Learning Representations*.
- Zhang, J., Kumor, D., and Bareinboim, E. (2020). Causal imitation learning with unobserved confounders. *Proceedings of the 34th Conference on Neural Information Processing Systems*.

A Reducing Our Unifying Framework to Related Literature

In this section, we discuss how the various previous works can be obtained as special cases of our unifying framework.

A.1 Temporally Correlated Noise (Swamy et al., 2022a)

The Temporally Correlated Noise (TCN) proposed in Swamy et al. (2022a) is a special case of our setting where $u^o = 0$ and only the confounding noise u^ε is present. Following Equation 14-17 of Swamy et al. (2022a), their setting can be summarised as

$$\begin{aligned} s_t &= \mathcal{T}(s_{t-1}, a_{t-1}) \\ &= \mathcal{T}(s_{t-1}, \pi_E(s_{t-1}) + u_{t-1} + u_{t-2}) \\ a_t &= \pi_E(s_t) + u_t + u_{t-1}, \end{aligned}$$

where \mathcal{T} is the transition function and u_t are the TCN. It can be seen that TCN is the confounding noise u^ε since the expert policy doesn't take it into account and it affects (or confounds) both the state and action.

It can be seen that this is a special case of our framework when $u_t^o = 0$, where $a_t = \pi_E(s_t) + \varepsilon(u_t^\varepsilon)$ from Equation (2), and more specifically when the confounding noise horizon in Theorem 3.2 is 2. In addition, the theoretical results in Swamy et al. (2022a) can be deduced from our main results as shown in Corollary 4.7.

A.2 Unobserved Contexts (Swamy et al., 2022b)

The setting considered by Swamy et al. (2022b) is a special case of our setting when $u^\varepsilon = 0$ and only u^o are present. Following Section 3 of Swamy et al. (2022b), their setting can be summarised as

$$\begin{aligned} \mathcal{T} &: \mathcal{S} \times \mathcal{A} \times C \rightarrow D(\mathcal{S}) \\ \nabla &: \mathcal{S} \times \mathcal{A} \times C \rightarrow [-1, 1] \\ a_t &= \pi_E(s_t, c) \end{aligned}$$

where $c \in C$ is the context, which is assumed to be fixed throughout an episode. There are no hidden confounders in this setting and the context c is included in u^o under our framework. Note that in our setting we also allow u^o to be varying throughout an episode. In addition, the theoretical results in Swamy et al. (2022b) can be deduced from our main results as shown in Corollary 4.6.

A.3 Imitation Learning with Latent Confounders (Vuorio et al., 2022)

The setting considered by Vuorio et al. (2022) is also a special case of our setting when $u^\varepsilon = 0$ and only u^o are present, which is very similar to Swamy et al. (2022b). In Section 2.2 of Vuorio et al. (2022), they introduced a latent variable $\theta \in \Theta$ that is fixed throughout an episode and $a_t = \pi_E(s_t, \theta)$. There are no hidden confounders in this setting and the latent variable θ is included in u^o in our framework. No theoretical imitation gap bounds are provided in Vuorio et al. (2022). However, Corollary 4.6 can be directly applied to their setting and bound the imitation gap.

A.4 Causal Delusion and Confusion (Ortega et al., 2021, de Haan et al., 2019, Pfrommer et al., 2023, Spencer et al., 2021, Wen et al., 2020)

The concept of causal delusion (Ortega et al., 2021) and confusion is widely studied in the literature (de Haan et al., 2019, Pfrommer et al., 2023, Spencer et al., 2021, Wen et al., 2020) from different perspectives. A classic example of causal confusion is learning to break in an autonomous driving scenario. The states are images with full view of the dashboard and the road conditions. The break indicator in this scenario is

the confounding variable that correlates with the action of breaking in subsequent steps, which causes the imitator to learn to break if the break indicator light is already on. Therefore, another name for this problem is the latching problem, where the imitator latches to spurious correlations between current action and the trajectory history. In the setting of Ortega et al. (2021), this is explicitly modelled as latent variables that affect both the action and state, causing spurious correlation between them and confusing the imitator. In other settings (de Haan et al., 2019, Pfrommer et al., 2023, Spencer et al., 2021, Wen et al., 2020), there are no explicit unobserved confounders, but the nuisance correlation between the previous states and actions can be modelled as the existence of hidden confounders u^ε in our framework. Specifically, in de Haan et al. (2019), x_{t-1} and a_{t-1} are considered confounders that affect the state variable x_t , which causes a spurious correlation between previous state action pairs and a_t . The spurious correlation between variables is typically modelled as the existence of a hidden confounder u^ε that affects both variables in causal modelling. For example, the actual hazard or event that causes the expert to break will be the hidden confounder u^ε that affects both the break and the break indicator.

However, despite the fact that this setting can be considered a special case of our general framework, we stress that the concrete and practical problems considered in de Haan et al. (2019), Pfrommer et al. (2023), Spencer et al. (2021), Wen et al. (2020) are different from ours, where they assumed implicitly that the hidden confounders u^ε are embedded in the observations or outright observed.

B Proofs of Main Results

In this section, we provide the proofs for the main results and corollaries in this paper.

B.1 Proof of Propositions

Proposition 4.3: The ill-posedness $\nu(\Pi, k)$ is monotonically increasing as the confounded horizon k increases.

Proof. From definition, we have that

$$\nu(\Pi, k) = \sup_{\pi \in \Pi} \frac{\|\pi_E - \pi\|_2}{\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2}.$$

We would like to show for each $\pi \in \Pi$, $\frac{\|\pi_E - \pi\|_2}{\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2}$ is increasing as k increases, which would imply that $\nu(\Pi, k)$ is increasing. For each $\pi \in \Pi$, we see that the numerator is constant w.r.t the horizon k . Therefore, it is enough to check that for each $\pi \in \Pi$, the denominator $\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2$ decreases as k increases. For any two integer horizon $k_1 > k_2$,

$$\mathbb{E}[a_t - \pi(h_t)|h_{t-k_1}]^2 = \mathbb{E}[\mathbb{E}[a_t - \pi(h_t)|h_{t-k_2}]|h_{t-k_1}]^2 \tag{6}$$

$$\leq \mathbb{E}[\mathbb{E}[a_t - \pi(h_t)|h_{t-k_2}]^2|h_{t-k_1}] \tag{7}$$

$$= \mathbb{E}[a_t - \pi(h_t)|h_{t-k_2}]^2 \tag{8}$$

by the tower property of conditional expectation as $\sigma(h_{t-k_1}) \subseteq \sigma(h_{t-k_2})$, Jensen's inequality for conditional expectations, and the fact that $\mathbb{E}[a_t - \pi(h_t)|h_{t-k_2}]^2$ is h_{t-k_1} measurable, respectively for each line. Therefore, we have that $\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]$ is decreasing, which implies $\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2$ is decreasing and $\nu(\Pi, k)$ is increasing as k increases, which completes the proof. \square

B.2 Main results for guarantees on the imitation gap

Theorem 4.5: Let $\hat{\pi}_h$ be the learnt policy with CMR error ε and let $\nu(\Pi, k)$ be the ill-posedness of the problem. Assume that $\delta_{TV}(u_t^\varepsilon, \mathbb{E}_{\pi_E}[u_t^\varepsilon|h_t]) \leq \delta$ for $\delta \in \mathbb{R}^+$, $P(u_t^\varepsilon)$ is c-TV stable and π_E is deterministic. Then, the imitation gap is upper bounded by

$$J(\pi_E) - J(\hat{\pi}_h) \leq T^2(c\varepsilon\nu(\Pi, k) + 2\delta) = \mathcal{O}(T^2(\delta + \varepsilon)).$$

Proof of Theorem 4.5. Recall that $J(\pi)$ is the expected reward following π , and we would like to bound the performance gap $J(\pi_E) - J(\hat{\pi}_h)$ between the expert policy π_E and the learned history-dependent policy $\hat{\pi}_h$. Let $Q_{\hat{\pi}_h}(s_t, a_t, u_t^o)$ be the Q-function of $\hat{\pi}_h$. Using the Performance Difference Lemma (Kakade and Langford, 2002), we have that for any Q-function $\tilde{Q}(h_t, a_t)$ that takes in the trajectory history h_t and action a_t ,

$$\begin{aligned}
J(\pi_E) - J(\hat{\pi}_h) &= \mathbb{E}_{\tau \sim \pi_E} \left[\sum_{t=1}^T Q_{\hat{\pi}_h}(s_t, a_t, u_t^o) - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h}(s_t, a, u_t^o)] \right] \\
&= \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h}(s_t, a_t, u_t^o) - \tilde{Q}(h_t, a_t) + \tilde{Q}(h_t, a_t) - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q} + \tilde{Q}]] \\
&= \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [\tilde{Q} - \mathbb{E}_{a \sim \hat{\pi}_h} [\tilde{Q}]] + \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q} - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}]] \tag{9}
\end{aligned}$$

We first bound the second part of Equation (9). Denote by δ_{TV} the total variation distance. For two distributions P, Q , recall the property of total variation distance for bounding the difference in expectations:

$$|\mathbb{E}_P[f(x)] - \mathbb{E}_Q[f(x)]| \leq \|f\|_\infty \delta_{TV}(P, Q).$$

In order to bound the second part of Equation (9), for any Q function, consider inferred \tilde{Q} using the conditional expectation of u^o on the history h ,

$$\tilde{Q}(h_t, a_t) := Q(s_t, a_t, \mathbb{E}_{\tau \sim \pi_E} [u_t^o | h_t]),$$

where note that $s_t \in h_t$. We have that, when the transition trajectory $(s_t, u_t^o, u_t^\varepsilon, r_t) \sim \pi_E$ follows the expert policy, for any action $\hat{a} \sim \pi$ following some policy π (in our case, it can be π_E or $\hat{\pi}_h$),

$$\begin{aligned}
|\mathbb{E}_{\tau \sim \pi_E, \hat{a} \sim \pi} [Q(s_t, \hat{a}, u_t) - \tilde{Q}(h_t, \hat{a})]| &= |\mathbb{E}_{\tau \sim \pi_E, \hat{a} \sim \pi} [Q(s_t, \hat{a}, u_t^o) - Q(s_t, \hat{a}, \mathbb{E}_{\tau \sim \pi_E} [u_t^o | h_t])]| \\
&= |\mathbb{E}_{u_t^o \sim \pi_E} [\mathbb{E}_{\pi_E, \pi} [Q(s_t, \hat{a}, u_t^o) | u_t^o] - \mathbb{E}_{u_t^o | h_t \sim \pi_E} [\mathbb{E}_{\pi_E, \pi} [Q(s_t, \hat{a}, u_t^o) | u_t^o]]]| \tag{10} \\
&\leq \|\mathbb{E}_{\pi_E, \pi} [Q(s_t, \hat{a}, u_t^o) | u_t^o]\|_\infty \delta_{TV}(u_t^o, \mathbb{E}_{\pi_E} [u_t^o | h_t]) \tag{11} \\
&\leq T \cdot \delta_{TV}(u_t^o, \mathbb{E}_{\pi_E} [u_t^o | h_t]) \tag{12} \\
&\leq T\delta \tag{13}
\end{aligned}$$

where Equation (10) uses the tower property of expectations, Equation (11) uses the total variation distance bound for bounded functions, Equation (12) uses the fact that the Q function is bounded by T and Equation (13) uses the condition that $\delta_{TV}(u_t^o, \mathbb{E}_{\pi_E} [u_t^o | h_t]) \leq \delta$ in the theorem statement. Since Equation (9) holds for any choice of \tilde{Q} , we choose $\tilde{Q}_{\hat{\pi}_h}(h_t, a_t) := Q_{\hat{\pi}_h}(s_t, a_t, \mathbb{E}_{\tau \sim \pi_E} [u_t^o | h_t])$ such that we can apply Equation (13) twice to bound the second part of Equation (9):

$$\begin{aligned}
\mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]] &\leq \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h} + |\mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]]| \\
&= \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}] + |\mathbb{E}_{s_t, u_t \sim \pi_E, a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]| \\
&\leq |\mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]| + T\delta \tag{14} \\
&\leq 2T\delta
\end{aligned}$$

where Equation (14) holds by applying Equation (13) because the expectation of the trajectories (and their transitions) are over π_E , and the actions which are used only as arguments into the Q function are sampled from $\hat{\pi}_h$.

Next, we bound the first part of Equation (9). Recall that the ill-posedness of the problem for a policy class Π is

$$\nu(\Pi, k) = \sup_{\pi \in \Pi} \frac{\|\pi_E - \pi\|_2}{\|\mathbb{E}[a_t - \pi(h_t)|h_{t-k}]\|_2}$$

where $\|\pi_E - \pi\|_2$ is the RMSE and $\|\mathbb{E}[a_t - \pi(s_t)|s_{t-k}]\|_2$ is the CMR error from our algorithm. Since the learned policy $\hat{\pi}_h$ have CMR error of ε , we have that

$$\|\pi_E - \hat{\pi}_h\|_2 \leq \nu(\Pi, k) \|\mathbb{E}[a_t - \hat{\pi}_h(h_t)|h_{t-k}]\|_2 \leq \nu(\Pi, k)\varepsilon$$

Next, recall that c-total variation stability of a distribution $P(u^\varepsilon)$ where $u^\varepsilon \in A$ for some space A implies for two elements $a_1, a_2 \in A$,

$$\|a_1 - a_2\|_2 \leq \Delta \implies \delta_{TV}(a_1 + u^\varepsilon, a_2 + u^\varepsilon) \leq c\Delta.$$

Since $P(u_t^\varepsilon)$ is c-TV stable w.r.t the action space A , we have that for all history trajectories $h_t \in H$ (note that $s_t \in h_t$)

$$\delta_{TV}(\pi_E(s_t) + u_t^\varepsilon, \hat{\pi}_h(h_t) + u_t^\varepsilon) \leq c\|\pi_E(s_t) - \hat{\pi}_h(h_t)\|_2.$$

Then, we have that by Jensen's inequality,

$$\begin{aligned} \mathbb{E}_{h_t \sim \pi_E} [\delta_{TV}(\pi_E(s_t) + u_t^\varepsilon, \hat{\pi}_h(h_t) + u_t^\varepsilon)]^2 &\leq \mathbb{E}_{h_t \sim \pi_E} [\delta_{TV}(\pi_E(s_t) + u_t^\varepsilon, \hat{\pi}_h(h_t) + u_t^\varepsilon)^2] \\ \implies \mathbb{E}_{h_t \sim \pi_E} [\delta_{TV}(\pi_E(s_t) + u_t^\varepsilon, \hat{\pi}_h(h_t) + u_t^\varepsilon)] &\leq \sqrt{\mathbb{E}_{h_t \sim \pi_E} [\delta_{TV}(\pi_E(s_t) + u_t^\varepsilon, \hat{\pi}_h(h_t) + u_t^\varepsilon)^2]} \\ &\leq \sqrt{c^2 \mathbb{E}_{h_t \sim \pi_E} [\|\pi_E(s_t) - \hat{\pi}_h(h_t)\|_2^2]} \\ &= c\|\pi_E - \hat{\pi}_h\|_2 \leq c\varepsilon\nu(\Pi, k) \end{aligned}$$

Therefore, by applying the total variation distance bound for expectations of $\tilde{Q}_{\hat{\pi}_h}$ over different distributions of action a_t , we have that

$$\mathbb{E}_{\tau \sim \pi_E} [\tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [\tilde{Q}_{\hat{\pi}_h}]] = \mathbb{E}_{\tau \sim \pi_E} [\tilde{Q}_{\hat{\pi}_h}(h_t, a_t) - \mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \hat{\pi}_h(h_t))]] \quad (15)$$

$$= \mathbb{E}_{h_t \sim \pi_E} [\mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \pi_E(s_t) + u_t^\varepsilon)] - \mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \hat{\pi}_h(h_t) + u_t^\varepsilon)]] \quad (16)$$

$$\leq \|\tilde{Q}_{\hat{\pi}_h}\|_\infty \mathbb{E}_{h_t \sim \pi_E} [\delta_{TV}(F(\pi_E(s_t) + u_t^\varepsilon), F(\hat{\pi}_h(h_t) + u_t^\varepsilon))] \quad (17)$$

$$\leq Tc\varepsilon\nu(\Pi, k) \quad (18)$$

Combining all of above, we see that from Equation (9), by selecting $\tilde{Q}_{\hat{\pi}_h}(h_t, a_t) := Q_{\hat{\pi}_h}(s_t, a_t, \mathbb{E}_{\tau \sim \pi_E}[u_t^\varepsilon | h_t])$, the imitation gap can be bounded by

$$J(\pi_E) - J(\hat{\pi}_h) = \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [\tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [\tilde{Q}_{\hat{\pi}_h}]] + \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]] \quad (19)$$

$$\leq \sum_{t=1}^T Tc\varepsilon\nu(\Pi, k) + \sum_{t=1}^T 2T\delta \quad (20)$$

$$\leq T \cdot (Tc\varepsilon\nu(\Pi, k) + 2T\delta) \quad (21)$$

$$= T^2(c\varepsilon\nu(\Pi, k) + 2\delta) = \mathcal{O}(T^2(\varepsilon + \delta)), \quad (22)$$

which concludes the proof. \square

B.3 Proofs of Corollaries

Corollary 4.6: In the special case that $u_t^o = 0$, meaning that there is no confounder observable to the expert, or $u_t^o = \mathbb{E}_{\pi_E}[u_t^o|h_t]$, meaning that u_t^o is $\sigma(h_t)$ measurable (all information regarding u_t^o is represented in the history), the imitation gap bound is $T^2(c\varepsilon\nu(\Pi, k))$, which coincides with Theorem 5.1 of Swamy et al. (2022a).

Proof. If $u_t^o = 0$, then we have $u_t^o = \mathbb{E}_{\pi_E}[u_t^o|h_t]$ since u_t^o is a constant. If $u_t^o = \mathbb{E}_{\pi_E}[u_t^o|h_t]$, we have that

$$\delta_{TV}(u_t^o, \mathbb{E}_{\pi_E}[u_t^o|h_t]) = \delta_{TV}(u_t^o, u_t^o) \leq 0$$

By plugging $\delta = 0$ into Theorem 4.5, we have that $J(\pi_E) - J(\hat{\pi}_h) \leq T^2(c\varepsilon\nu(\Pi, k))$, which is the same as the imitation gap derived in Swamy et al. (2022a) and completes the proof. \square

Corollary 4.7: In the special case that $u_t^\varepsilon = 0$, if the learned policy via supervised BC have error ε , then the imitation gap bound is $T^2(\frac{2}{\sqrt{\dim(A)}}\varepsilon + 2\delta)$, which is a concrete bound that extends the abstract bound in Theorem 5.4 of Swamy et al. (2022b).

Proof. In Theorem 5.4 of Swamy et al. (2022b), for the offline case, which is the setting we are considering (as opposed to the online settings), they defined the following quantities for bounding the imitation gap in a very general fashion,

$$\begin{aligned} \varepsilon_{\text{off}} &:= \sup_{\tilde{Q}} \mathbb{E}_{\tau \sim \pi_E}[\tilde{Q} - \mathbb{E}_{a \sim \hat{\pi}_h}[\tilde{Q}]] \\ \delta_{\text{off}} &:= \sup_{Q \times \tilde{Q}} \mathbb{E}_{\tau \sim \pi_E}[Q_{\hat{\pi}_h} - \tilde{Q} - \mathbb{E}_{a \sim \hat{\pi}_h}[Q_{\hat{\pi}_h} - \tilde{Q}]]. \end{aligned}$$

The imitation gap by Theorem 5.4 in Swamy et al. (2022b) under the assumption that $u_t^\varepsilon = 0$ is $T^2(\varepsilon_{\text{off}} + \delta_{\text{off}})$, which can also be deduced from Equation (9) by naively applying the above supremum. To obtain a concrete bound, we can provide a tighter bound for $\mathbb{E}_{\tau \sim \pi_E}[Q_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h}[Q_{\hat{\pi}_h}]]$, which is the first part of Equation (9), given that $u_t^\varepsilon = 0$.

For two elements $a_1, a_2 \in A$, we have that by Cauchy–Schwarz,

$$\delta_{TV}(a_1 + 0, a_2 + 0) = \frac{1}{2} \|a_1 - a_2\|_1 \leq \frac{\sqrt{\dim(A)}}{2} \|a_1 - a_2\|_2.$$

Then, we have that

$$\|a_1 - a_2\|_2 \leq \Delta \implies \delta_{TV}(a_1, a_2) \leq \frac{2}{\sqrt{\dim(A)}} \Delta$$

so that by Theorem 4.5,

$$\mathbb{E}_{\tau \sim \pi_E}[\tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h}[\tilde{Q}_{\hat{\pi}_h}]] = \mathbb{E}_{\tau \sim \pi_E}[\tilde{Q}_{\hat{\pi}_h}(h_t, a_t) - \mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \hat{\pi}_h(h_t))]] \quad (23)$$

$$= \mathbb{E}_{h_t \sim \pi_E}[\mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \pi_E(s_t))] - \mathbb{E}[\tilde{Q}_{\hat{\pi}_h}(h_t, \hat{\pi}_h(h_t))]] \quad (24)$$

$$\leq \|\tilde{Q}_{\hat{\pi}_h}\|_\infty \frac{2}{\sqrt{\dim(A)}} \|\pi_E - \pi\|_2 \quad (25)$$

$$\leq T \frac{2}{\sqrt{\dim(A)}} \varepsilon, \quad (26)$$

since when $u_t^\varepsilon = 0$ the learning error via supervised learning is $\varepsilon := \|\pi_E - \pi\|_2$. Therefore, the final imitation bound following Theorem 4.5 is

$$J(\pi_E) - J(\hat{\pi}_h) = \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [\tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [\tilde{Q}_{\hat{\pi}_h}]] + \sum_{t=1}^T \mathbb{E}_{\tau \sim \pi_E} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h} - \mathbb{E}_{a \sim \hat{\pi}_h} [Q_{\hat{\pi}_h} - \tilde{Q}_{\hat{\pi}_h}]] \quad (27)$$

$$\leq \sum_{t=1}^T T \frac{2}{\sqrt{\dim(A)}} \varepsilon + \sum_{t=1}^T 2T\delta \quad (28)$$

$$= T^2 \left(\frac{2}{\sqrt{\dim(A)}} \varepsilon + 2\delta \right). \quad (29)$$

This bound is a concrete bound, obtained through detailed analysis of the problem at hand, that coincides with the abstract bound $T^2(\varepsilon_{\text{off}} + \delta_{\text{off}})$ provided in Theorem 5.4 of Swamy et al. (2022a). Note that this bound is independent of the ill-posedness $\nu(\Pi, k)$ and the c -TV stability of u_t^ε , which are present in the bound of Theorem 4.5, because of the lack of hidden confounders u_t^ε . \square

C Environments and Tasks

C.1 Dynamic Aeroplane Ticket Pricing

Here, we provide details regarding the dynamic aeroplane ticket pricing environment introduced in Example 3.1. The environment and the expert policy are defined as follows:

$$\mathcal{S} := \mathbb{R} \quad (30)$$

$$\mathcal{A} := [-1, 1] \quad (31)$$

$$s_t = \text{sign}(s) \cdot u_t^o - u_t^\varepsilon \quad (32)$$

$$\pi_E = \text{clip}(-s/u_t^o, -1, 1) \quad (33)$$

$$a_t = \pi_E + 10 \cdot u_t^\varepsilon \quad (34)$$

$$u_t^o = \text{mean}(p_t \sim \text{Unif}[-1, 1], p_{t-1}, \dots, p_{t-M}) \quad (35)$$

$$u_t^\varepsilon = \text{mean}(q_t \sim \text{Normal}(0, 0.1 \cdot \sqrt{k}), q_{t-1}, \dots, q_{t-k+1}) \quad (36)$$

where M is the influence horizon of the expert-observable u^o , which we set to 30. The states s_t are the profits at each time step, and the actions a_t are the final ticket price. u_t^o represent the seasonal patterns, where the expert π_E will try to adjust the price accordingly. u_t^ε represent the operating costs, which are additive both to the profit and price. Both u_t^o and u_t^ε are the mean over a set of i.i.d samples, q_t and p_t , and vary across the time steps by updating the elements in the set at each time step. This constructions allows u_t^ε and u_{t-k}^ε to be independent since all set elements q_t will be re-sampled from time step $t - k$ to t . We multiply the standard deviation of q_t by \sqrt{k} to make sure u_t^ε , which is the average over k i.i.d variables, have the same standard deviation for all choices of k .

C.2 Mujoco Environments

We evaluate DML-IL on three Mujoco environments: Ant, Half Cheetah and Hopper. The original tasks do not contain hidden variables, so we modify the environment to introduce u^ε and u^o . We use the default transition, state and action space defined in the Mujoco environment. However, we changed the task objectives by altering the reward function and added confounding noise to both the state and action. Specifically, instead of controlling the ant, half cheetah and hopper, respectively, to travel as fast as possible, the goal is to control the agent to travel at a target speed that is varying throughout an episode. This target speed is u^o , which is observed by the expert but not recorded in the dataset. In addition, we add confounding noise u_t^ε to

s_t and a_t to mimic the environment noise such as wind noise. In all cases, the target speed u_t^o , confounding noise u_t^ε and the action a_t are generated as follows:

$$a_t = \pi_E + 20 \cdot u_t^\varepsilon \tag{37}$$

$$u_t^o = \text{mean}(p_t \sim \text{Unif}[-2, 4], p_{t-1}, \dots, p_{t-M}) \tag{38}$$

$$u_t^\varepsilon = \text{mean}(q_t \sim \text{Normal}(0, 0.01 \cdot \sqrt{k}), q_{t-1}, \dots, q_{t-k+1}) \tag{39}$$

where $M = 30$, the state transitions follow the default Mujoco environment and the expert policy π_E is learned online in the environment. u_t^o and u_t^ε follow the aeroplane ticket pricing environment to be the average over a queue of i.i.d random variables. The reward is defined to be the $1_{\text{healthy}} - (\text{current velocity} - u_t^o)^2 - \text{control loss}$, where 1_{healthy} gives reward 1 as long as the agent is in a healthy state as defined in the Mujoco documentation. The second penalty term penalises deviation between the current agent’s velocity and the target velocity u_t^o . The control loss term is also as defined in default Mujoco, which is $0.1 * \sum(a_t^2)$ at each step to regularize the size of actions.

C.2.1 Ant

In the Ant environment, we follow the gym implementation ⁴ with 8-dimensional action space and 28-dimensional observable state space, where the agent’s position is also included in the state space. Since the target speed u_t^o is not recorded in the trajectory dataset, we scale the current position of the agent with respect to the target speed, $pos_t' = pos_{t-1} + \frac{pos_t - pos_{t-1}}{u_t^o}$, and use the new agent position pos_t' in the observed states. This allows the imitator to infer information regarding u_t^o from trajectory history, namely from the rate of change in the past positions.

C.2.2 Half Cheetah

In the Half Cheetah environment, we follow the gym implementation ⁵ with 6-dimensional action space and 18-dimensional observable state space, where the agent’s position is also included in the state space. Similarly to the Ant environment, we scale the current position of the agent to $pos_t' = pos_{t-1} + \frac{pos_t - pos_{t-1}}{u_t^o}$ such that the imitator can infer information regarding u_t^o from trajectory history.

C.2.3 Hopper

In the Hopper environment, we follow the gym implementation ⁶ with 3-dimensional action space and 12-dimensional observable state space, where the agent’s position is also included in the state space. Similarly to the Ant environment, we scale the current position of the agent to $pos_t' = pos_{t-1} + \frac{pos_t - pos_{t-1}}{u_t^o}$ such that the imitator can infer information regarding u_t^o from trajectory history.

D Implementation Details

D.1 DML-IL with K -Fold Cross-Fitting

Here, we outline DML-IL with K -fold cross-fitting, which ensures unbiased estimation and improves training stability. The algorithm is shown in Algorithm 2. The dataset is partitioned into K folds based on the trajectory index. For each fold, we use the leave-out data, that is, indices $I_k^c := [N] \setminus I_k$, to train separate roll-out models \hat{M}_i for $i \in [1..K]$. Then, to train a single expert model $\hat{\pi}_h$, we sample the trajectory history h_{t-k} from each fold and use the roll-out model trained with the leave-out data to complete the trajectory

⁴Ant environment: <https://www.gymnasium.dev/environments/mujoco/ant/>

⁵Half Cheetah environment: https://www.gymnasium.dev/environments/mujoco/half_cheetah/

⁶Hopper environment: <https://www.gymnasium.dev/environments/mujoco/hopper/>

Algorithm 2 DML-IL with K -fold cross-fitting

Input: Dataset \mathcal{D}_E of expert demonstrations, Confounding noise horizon k , number of folds K for cross-fitting
Output: A history-dependent imitator policy $\hat{\pi}_h$
Get a partition $(I_k)_{k=1}^K$ of dataset indices $[N]$ of trajectories
for $k = 1$ **to** K **do**
 $I_k^c := [N] \setminus I_k$
 Initialize the roll-out model \hat{M}_i as a mixture of Gaussians model
 repeat
 Sample (h_t, a_t) from data $\{(\mathcal{D}_{E,i}) : i \in I_k^c\}$
 Fit the roll-out model $(h_t, a_t) \sim \hat{M}_i(h_{t-k})$ to maximize log likelihood
 until convergence
end for
Initialize the expert model $\hat{\pi}_h$ as a neural network
repeat
 for $k = 1$ **to** K **do**
 Sample h_{t-k} from $\{(\mathcal{D}_{E,i}) : i \in I_k\}$
 Generate \hat{h}_t and \hat{a}_t using the roll-out model \hat{M}_i
 Update $\hat{\pi}_h$ to minimise the loss $\ell := \|\hat{a}_t - \hat{\pi}_h(\hat{h}_t)\|_2$
 end for
until convergence

and train $\hat{\pi}_h$. This technique is very important in Double Machine Learning (DML) literature (Shao et al., 2024, Chernozhukov et al., 2018) for it provides both empirical stability and theoretical guarantees. The base IV regression algorithm DML-IV with K -fold cross-fitting is theoretically shown to converge at the rate of $O(N^{-1/2})$ (Shao et al., 2024), where N is the sample size, under technical regularity and identifiability conditions (see Shao et al. (2024) for the technical conditions). These conditions are typically assumed for similar theoretical analyses, and DML-IL with K -fold cross-fitting will thus inherit this convergence rate guarantee if the regularity conditions are satisfied.

D.2 Expert Training

The expert in the aeroplane ticket pricing environment is explicitly hand crafted. For the Mujoco environments, we used the Stable-Baselines3 (Raffin et al., 2021) implementation of soft actor-critic (SAC) and the default hyperparameters for each task outlined by Stable-Baseline3. The expert policy is an MLP with two hidden layers of size 256 and ReLU activations, and we train the expert for 10^7 steps.

D.3 Imitator Training

With the expert policy π_E , we generate 40 expert trajectories, each of 500 steps, following our previously defined environments. Specifically, the confounding noise is added to the state and actions and crucially u_t^o is not recorded in the trajectories. The naive BC directly learns $\mathbb{E}[a_t | s_t]$ via supervised learning. ResiduIL mainly follows the implementation of Swamy et al. (2022a), where we adopt it to allow longer confounding horizon $k > 1$. For DML-IL and BC-SEQ, a history-dependent policy is used, where we fixed the look-back length to be $k + 3$, where k is the confounding horizon. BC-SEQ then just learns $\mathbb{E}[a_t | h_t]$ via supervised learning, and DML-IL is implemented with K -fold following Algorithm 2. The policy network architecture for BC, BC-SEQ and ResiduIL are 2 layer MLPs with 256 hidden size. The policy network $\hat{\pi}_h$ and the mixture of Gaussians roll-out model \hat{M} for DML-IL have similar architecture, with details provided in Table 1. We use AdamW optimizer with weight decay of 10^{-4} and learning rate of 10^{-4} . The batch size is 64 and each model is trained for 150 epochs, which is sufficient for their convergence.

Table 1: Network architecture for DML-IL. For mixture of Gaussians output, we report the number of components. No dropout is used.

(a) Roll-out model \hat{M}		(b) Policy model $\hat{\pi}_h$	
Layer Type	Configuration	Layer Type	Configuration
Input	state dim \times 3	Input	state dim \times (k+3)
FC + ReLU	Out: 256	FC + ReLU	Out: 256
FC + ReLU	Out: 256	FC + ReLU	Out: 256
MixtureGaussian	5 components; Out: state dim \times k	FC	Out: action dim

D.4 Imitator Evaluation

The trained imitator is then evaluated for 50 episodes, each 500 steps in the respective confounded environments. The average reward and the mean squared error between the imitator’s action and the expert’s action are recorded.

E Adopting other IV regression algorithms

In this paper, we have transformed causal IL with hidden confounders into a set of CMRs as defined in Equation (5). Therefore, in principle many IV regression algorithms can be adopted to solve our CMRs. We also experimented with other IV regression algorithms that have been previously shown to be practical (Shao et al., 2024) for different tasks and high-dimensional input. Specifically, we experimented with DFIV (Xu et al., 2020), which is an iterative algorithm that integrates the training of two models that depend on each other, and DeepGMM (Bennett et al., 2019a), which solves a minimax game by optimising two models adversarially. Note that DeepIV (Hartford et al., 2017) can be considered a special case of DML-IV (Shao et al., 2024), so we did not reimplement it. The additional results for using DFIV and DeepGMM as the CMRs solver are provided in Figure 4 and Figure 5. It can be seen from Figure 4 that only DFIV achieves good performance in the aeroplane ticket pricing environment, surpassing the performance of ResiduIL. For the Ant Mujoco task in Figure 5, both DFIV and DeepGMM fail to learn good policies, with only slightly lower MSE than BC and BC-SEQ. We think this is due to the high-dimensional state and action spaces and the inherent instability in the DFIV and DeepGMM algorithms. For DFIV, the interleaving of training of two models causes highly non-stationary training targets for both models, and, for DeepGMM, the adversarial training procedure of two models is similar to that of generative adversarial Networks (GANs), which is known to be unstable and difficult to train.

We conclude that solving the CMRs for an imitator policy can be sensitive to the choice of solver as well as the choice of hyperparameters. In addition, some IV regression algorithms do not work well with high dimension inputs. Our IV algorithm of choice, DML-IV, provides a robust base for the DML-IL algorithm that demonstrated good performance across all tasks and environments.

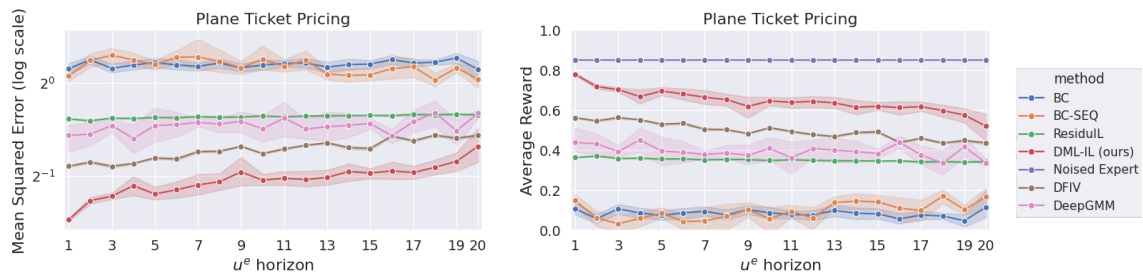


Figure 4: Additional results for the MSE between learnt policy and expert, and the average reward, in the plane ticket environment (Example 3.1), with DFIV and DeepGMM as the CMRs solver.

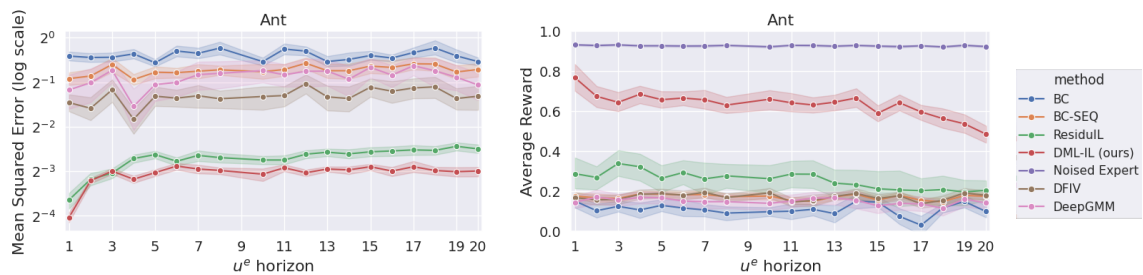


Figure 5: Additional results for the MSE between learnt policy and expert, and the average reward, Ant Mujoco environment, with DFIV and DeepGMM as the CMRs solver.