# Causal Inference out of Control:
# Estimating Performativity without Treatment Randomization

**Gary Cheng** [* 1]   **Moritz Hardt** [2]   **Celestine Mendler-Dünner** [2 3]

## Abstract

Regulators and academics are increasingly interested in the causal effect that algorithmic actions of a digital platform have on user consumption. In pursuit of estimating this effect from observational data, we identify a set of assumptions that permit causal identifiability without assuming randomized platform actions. Our results are applicable to platforms that rely on machine-learning-powered predictions and leverage knowledge from historical data. The key novelty of our approach is to explicitly model the dynamics of consumption over time, exploiting the repeated interaction of digital platforms with their participants to prove our identifiability results. By viewing the platform as a controller acting on a dynamical system, we can show that exogenous variation in consumption and appropriately responsive algorithmic control actions are sufficient for identifying the causal effect of interest. We complement our claims with an analysis of ready-to-use finite sample estimators and empirical investigations. More broadly, our results deriving identifiability conditions tailored to digital platform settings illustrate a fruitful interplay of control theory and causal inference.

## 1. Introduction

How much do advertisements decrease screen time? Do algorithmic recommendations increase consumption of inflammatory content? Does exposure to diverse news sources mitigate political polarization? These are just a few questions that firms, researchers, and regulators ask about digital platforms (Barberá et al., 2015; Brown et al., 2022b). These questions all examine the capacity of a digital platform's algorithmic action to influence users, an effect referred to as performativity (Perdomo et al., 2020).

Estimating performativity is a causal inference problem where the treatment corresponds to the algorithmic action taken by the platform and the outcome variable corresponds to user behavior. We focus on consumption as the user variable of interest and aim to measure how much algorithmic actions of platforms impact consumption. Treatment effect estimation from observational data is challenging in such algorithmic systems because actions are typically driven by observations of past consumption. This feedback loop couples actions and consumption, introducing confounding. The corresponding causal graph is illustrated in Figure 1(a) where $u$ corresponds to the platform action, $x$ to user consumption, and $z$ denotes the confounding set containing all relevant past platform actions and consumption data. The presence of confounding means that we cannot determine whether the correlations between $u$ and $x$ should be attributed to performativity of the actions $u$ or to the common cause $z$.

Injecting independent variation into the platform action $u$—for example, via A/B testing—is sufficient to resolve such confounding. However, randomized experiments may be ethically fraught (Kramer et al., 2014; PNAS, 2014), technically challenging to implement, or prohibitively expensive. Moreover, external investigators often do not have the ability to experimentally intervene in the practices of a platform. For these reasons, we focus this work around observational causal inference and ask: under what conditions can observational designs lead to valid inferences of performativity?

Typically, observational designs build on the premise that there is sufficient variation in the treatment in order to resolve confounding (Rosenbaum & Rubin, 1983; Imbens, 2004). In the context of Figure 1(a), this means that we need to assume overlap between $u$ and $z$—i.e., all strata of the treatment $u$ have positive probability of selection for any possible realization of the confounding set $z$.

At the outset, this assumption seems hard to reconcile with

---

*Work done during internship at the Max Planck Institute for Intelligent Systems, Tübingen. [1]Stanford University Department of Electrical Engineering [2]Max Planck Institute for Intelligent Systems, Tübingen, and Tübingen AI Center [3]ELLIS Institute Tübingen. Correspondence to: Gary Cheng <chenggar@stanford.edu>.

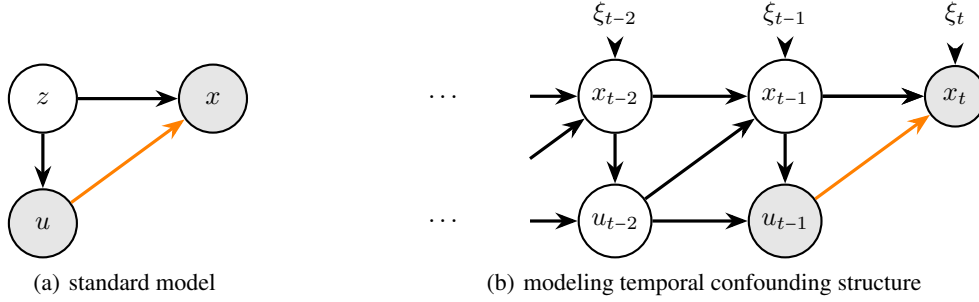(a) standard model      (b) modeling temporal confounding structure

**Figure 1:** The causal inference problem of estimating performativity of algorithmic actions.

practices of digital platforms operating machine learning algorithms. In particular, algorithmic platform actions are typically generated from machine learning models which have been trained on historical observations. As a consequence, past consumption, which is a part of the confounding variable $z$, typically influences future actions—potentially in a deterministic way. Thus, we should expect little independent variation in $u$ for a fixed $z$, meaning that overlap for this choice of $u$ and $z$ is unlikely to hold. Furthermore, because the interactions users have with digital platforms often span multiple time steps, the confounding set $z$ is may be high-dimensional. The high-dimensionality of such a confounding set makes overlap between $u$ and $z$ even harder to satisfy (D'Amour et al., 2017). Together, these properties of algorithmic systems make it challenging to justify valid inferences about performativity in Figure 1(a) from data collected under natural interactions between platforms and participants.

### 1.1. Our work

We propose a more refined causal model—i.e., a different way of representing $u$ and $z$—as a means to bypass the overlap issues we discuss earlier. Using our model, we articulate conditions for *valid* inferences of performativity from observational data, *without* assuming explicit randomization in the platform actions. Our approach is inspired by a control-theoretic view on the problem and explicitly models the repeated interactions between a digital platform and their users across time. We view the platform as a controller repeatedly adapting to changes in user consumption. The structural assumptions on the dynamics translate into structural assumptions on how the variables in the confounding set $z$ relate to each other across time. By exposing this structure, we can trace how variations on consumption propagate through the system, allowing us to design new conditions for causal identifiability tailored to the digital platform setting.

More formally, we work with the model visualized in Figure 1(b). The consumer's features $x_t$ at time $t$ are determined by their previous value $x_{t-1}$, exogenous noise $\xi_t$, as

well as the performative effect of the previous platform action $u_{t-1}$. The platform's action $u_t$ is updated in each step based on the most recent observations of $x_t$ and the action it took previously $u_{t-1}$. We are interested in the performative effect, represented by the orange arrow, and quantified by the following treatment effect function:

$$\mathrm{PE}_t(u, u') \coloneqq \mathbb{E}[x_t \mid \mathrm{do}(u_{t-1} = u)] - \mathbb{E}[x_t \mid \mathrm{do}(u_{t-1} = u')],$$

where $\mathrm{do}(\cdot)$ is the do-operator used to represent interventions in the causal graph (Pearl, 2009).

**Contributions.** Our main contribution is to demonstrate that it is possible to circumvent directly assuming exogenous variations in platform actions $u$ in order to identify $\mathrm{PE}_t$. In particular, our main theoretical claims are the following:

1. In the non-parametric case, we show that a) sufficient exogenous variation on the consumer's consumption at *multiple* consecutive time steps and b) the platform control action being non-degenerate, is necessary and sufficient for causal identifiability of $\mathrm{PE}_t$.

2. In the linear case, we show that the investigator can take advantage of observations of longer roll-outs across time to get identifiability of $\mathrm{PE}_t$ from a *single* consumer-feature perturbation.

To complement our study, we propose a two-stage regression estimator and an adjustment formula estimator for estimating $\mathrm{PE}_t$ from finite samples with theoretical guarantees. We also simulate a recommendation system to empirically test the efficacy of our assumptions at reducing overlap violations. We present more experiments in the Appendix, including ones using real microeconomic data. Taken together our results provide a valuable guidance for the applicability of observational causal inference for investigating performativity on digital platforms.

### 1.2. Background and Related Work

The impact digital platforms have on their users is relevant for diverse applications spanning content recommendation, prediction policy problems, labor markets and social

science research (c.f., Shmueli & Tafti, 2020; Thai et al., 2016; Fleder et al., 2010; Adomavicius et al., 2013; Krauth et al., 2022; Barberá et al., 2015; Brown et al., 2022b; Wagner et al., 2021a). Performative effects influence design choices on the side of the platform (Bottou et al., 2013; Perdomo et al., 2020), produce externalities for platform participants (Wagner et al., 2021b), and offer an important dimension along which to monitor algorithmic systems. A recent work by Hardt et al. (2022) proposed performative power as a formal measure to quantify the extent to which a platform can steer user behavior, relating it to economic power. In this context, our identifiability results provide conditions under which a lower bound on performative power can be assessed from observational data, providing a valuable guardrail for digital market investigations.

**Performativity in Machine Learning.** There is a growing body of work studying the role of performativity in machine learning (c.f., Perdomo et al., 2020; Brown et al., 2022a; Mandal et al., 2023; Eilat & Rosenfeld, 2023; Wang et al., 2023; Hardt & Mendler-Dünner, 2023) and, more specifically, user dynamics in recommender systems (Kalimeris et al., 2021; Chaney et al., 2018; Dean & Morgenstern, 2022). These works typically posit an interaction model between platforms and participants or treat the strength of performativity as a free parameter. Our work complements these investigations by providing an approach to estimate performativity from data. Most related to our work is Mendler-Dünner et al. (2022) who propose 'predicting-from-predictions' as an identification strategy for recovering performative effects under outcome performativity. The approach relies on incongruences in modality to establish causal identifiability, whereas our results exploit interactions across time. Several other works blend performativity and causal inference in the context of strategic classification (e.g., Miller et al., 2020; Shavit et al., 2020; Bechavod et al., 2021; Harris et al., 2022; Horowitz & Rosenfeld, 2023), while they posit performativity, we focus on estimating it.

**Causal Inference.** We build on tools from causal inference (Pearl, 2009) to understand when observational data is sufficient to measure performativity, to deal with potential confounding, and to derive finite-sample estimators. Our work most closely relates to works that handle overlap violations (Chen et al., 2007; Sasaki & Ura, 2017; Yang & Ding, 2018; Petersen et al., 2012a), address adjustment set selection (de Luna et al., 2011) and model time in causal graphs (Blackwell, 2013). However, unlike general causal inference results, our focus lies on measuring the causal effects of algorithm-driven treatments and providing assumptions and results tailored to this setting. This provides many interesting connections but leaves a small intersection with prior work, which we discuss now.

D'Amour et al. (2017) provides an in depth analysis of the restrictive implications of assuming overlap with high-dimensional confounders. Petersen et al. (2012b) discuss how sub-selecting adjustment sets in causal inference can help alleviate challenges of overlap. We apply similar ideas to shrink our adjustment set, motivating our selection using a time-aware interaction model. When studying time-aware causal graphs, most prior work focuses on experimental designs (Klasnja et al., 2015; Dwivedi et al., 2022; Zhang & Bareinboim, 2019). Only a few focus on observational causal inference; some examples include (Blackwell, 2013; Kim et al., 2020; Shah et al., 2022). While these works do not assume control over platform actions, they do presuppose that typical overlap assumptions are satisfied. For example, Shah et al. (2022) focus on a recommender system application and use an exponential family model for the conditional distribution, implicitly assuming overlap. Other time-aware causal inference techniques include difference in difference techniques (Bertrand et al., 2001) and synthetic control techniques (Abadie & Gardeazabal, 2003; Abadie et al., 2010) for longitudinal and panel data. These approaches also implicitly assume overlap—i.e., that treatment and control groups exist or can be synthetically constructed. In contrast, our work focuses on how we can use time-dependent interactions in digital platforms to handle potentially deterministic platform actions and ultimately guarantee identifiabilty from observational data. To the best of our knowledge, the assumptions derived in this work are novel, and there is no prior work exploiting time and the resulting structure in confounding variables to establish identifiability.

**Control Theory.** From a control theory perspective, estimating performativity in our causal model is reminiscent of a system identification problem (Ljung, 2010). However, our problem setup differs from standard system identification results such as Bruder et al. (2018) because we focus on purely observational designs, where we do not choose what platform control actions (i.e., interventions) are taken. Within the system identification literature, we highlight the work of Abbasi-Yadkori & Szepesvari (2011) because of the similarity of their model to the linear model we study in Section 4. Their results hinge on a finite-sample system identification result, similar in spirit to the type of identifiability results found in this paper.

## 2. Model

In our setting, estimating performativity corresponds to quantifying the causal effect of a platform action $u$ on user consumption $x$. The main challenge is that future platform actions are affected by past actions and observations of user behavior, introducing confounding.

Our model—outlined in Figure 1(b)—makes the temporal component of interactions among the confounding variables explicit: We let $x_t \in \mathbb{R}^d$ and $u_t \in \mathbb{R}^p$ denote the consumption and platform action at time step $t$ respectively. We assume for all $t \geq 0$ the dynamics follow

$$x_t = f(x_{t-1}) + g(u_{t-1}) + \xi_t$$
$$u_t = h(x_t) + r(u_{t-1}) \tag{1}$$

with $\xi_t \in \mathbb{R}^d$ modeling potential exogenous variations in $x_t$ and the functions $f : \mathbb{R}^d \to \mathbb{R}^d$, $g : \mathbb{R}^p \to \mathbb{R}^d$, $h : \mathbb{R}^d \to \mathbb{R}^p$, and $r : \mathbb{R}^p \to \mathbb{R}^p$ describe how consumption and platform actions affect one another. We make the following assumption[1] on the exogenous noise:

**Assumption 1** (Mutually Independent Exogenous Variation). *For any $t \geq 1$, the random variable $\xi_t$ is mutually independent of $\xi_k$ for all $k \neq t$ and independent of $(x_0, u_0) \sim P_0$.*

Given these modeling assumptions, we are interested in estimating the treatment effect function $\text{PE}_t$. We note that because our system dynamics (1) are time-invariant and the structural equations for $x$ are assumed to be separable, we have $\text{PE}_t = \text{PE}_{t'}$ for all $t, t'$. Thus, without loss of generality, we will focus on identifying $\text{PE}_T(u, u') = \mathbb{E}[x_T \mid \text{do}(u_{T-1} = u)] - \mathbb{E}[x_T \mid \text{do}(u_{T-1} = u')]$, letting $T$ denote the index we want to estimate performativity for. We use $R_K$ to denote a rollout of the previous $K$ time indices leading up to the chosen time index $T$:

$$R_K \coloneqq (\{x_{T-t}, u_{T-t}\}_{t=1}^K, x_T).$$

We assume access to iid observations of rollouts $R_K$. We will specify $K$ in each result.

### 2.1. Running example: Video recommendation system

To provide a concrete instantiation of our model, consider an auditor who is interested in estimating the impact of the recommendation algorithm of a video streaming platform—like Twitch or YouTube—on the consumption patterns of their users. Let $x_{1t} \in \mathbb{R}^p$ be some measure of content consumption (e.g., number of hours streamed) for $p$ video categories of interest during week $t$ for a given user. Let $x_{2t} \in \mathbb{R}^{d_z}$ be comprised of measurements about the platform such as revenue per category which could be confounders. We can think of the joint vector $[x_{1t}; x_{2t}] \in \mathbb{R}^d$ as the state variable $x_t$ for $d = p + d_z$. The platform action $u_t \in \mathbb{R}^p$ is a measure of how many videos from the $p$ categories of interest are recommended to a given user during week $t$. The platform interfaces using $u_t$ with the goal of maximizing total profits, which is some deterministic function of $x_t$. The auditor is interested in estimating how the platform action

$u_{t-1}$ impacts the average watch habits $x_{1t}$ of users—i.e., the first $p$ coordinates of $\text{PE}(u, u')$.

Our model postulates that user consumption changes over time based on the recommendations by the algorithm, as well as external factors (e.g., new trends). Formally, taking inspiration from Jambor et al. (2012), we model the dynamics of the system as

$$x_{1t} = f_1(x_{1t-1}, x_{2t-1}) + g_1(u_{t-1}) + \xi_t^{(1)}$$
$$x_{2t} = f_2(x_{1t-1}, x_{2t-1}) + g_2(u_{t-1}) + \xi_t^{(2)}.$$

The function $f_1$ models how much interest users retain in each video category from week to week, as well as the effect of confounders on viewership (e.g., how many hours of viewing time can a competitor poach). The function $f_2$ models how the performance metrics chosen as a target variable by the firm evolve over time, while the function $g_1$ models the platform's ability to control this metric. The auditor wants to estimate the relationship $g_1$ that governs how much consumption increases with the number of recommendations. The noise variables $\xi_t^{(1)}, \xi_t^{(2)}$ allow for natural variation in user preferences, e.g., due to economic developments, independent of past consumption and the platform's recommendations. We can model the platform action as

$$u_t = h(x_{1t}, x_{2t}) + r(u_{t-1}),$$

where $h$ models the platform's algorithm of how viewer statistics and other metrics affect recommendations in the future. The function $r$ models how it regularizes recommendations to avoid overfitting to recent activity.

**Markovian assumption.** Our model implicitly assumes that current platform actions are only affected by the recent past. This assumption seems reasonable for ML-based algorithmic actions, given that digital platforms are constantly retraining ML models on fresh data to improve performance and mitigate distribution shift (Shankar et al., 2022). In addition, the Markovian view of digital platform actions is prevalent in the recommendation system literature. For example, the contextual multi-armed bandit models used to study recommendation systems are Markovian by construction (Langford & Zhang, 2007; Bouneffouf & Rish, 2019). In turn, the Markovian assumption on consumer dynamics is based on the belief that there are few long range causal effects that affect consumption, and the ones that exist—say inherent biases, interests, or habits—can be encoded into all of the states, without blowing up the dimension. Even if the system we are modeling is not Markovian, choosing to use the Markovian assumption as a means to weaken overlap assumptions and select an adjustment set can still be beneficial as long as the errors caused by model misspecification are small, outweighed by the benefits of statistical power. We empirically explore this trade-off in Appendix D.

---

[1] We use Assumption 1 for clarity of exposition; for an alternative and weaker assumption, see Appendix G.

# 3. Identifiability from exogenous variations in user consumption

A quantity is *identifiable* if it can be uniquely determined from the observational data distribution. Conversely, if there exists multiple values of said quantity which are all consistent with the observational data probability distribution, then we say it is *unidentifiable*.

Our goal is to outline necessary and sufficient conditions for identifiability of $\mathrm{PE}_T$ from observations of $u$ and $x$ across time. To provide intuition for the merits of exposing time, let us recall the conditions for identifiability in the general causal graph in Figure 1(a). Classical results from causal inference in the presence of observed confounding (Pearl, 2009) tell us that a sufficient condition for identifiability of the causal effect of $u$ on $x$ is *admissibility* and *overlap*.

**Definition 3.1** (Admissibility). *We say a continuous random variable $Z$ with density $p$ is admissible for adjustment with respect to treatment $U$ and outcome $X$ if the adjustment formula is valid:*

$$\mathbb{E}[X \mid \mathrm{do}(U \coloneqq u)] = \int \mathbb{E}[X \mid U = u, Z = z]\, p(z)\, dz.$$

**Definition 3.2** (Overlap). *For an action $U$ and a confounder $Z$ with well-defined joint density $p$, overlap of $(u, z)$ is satisfied if $p_{U|Z}(u' \mid z') > 0$ for all $u' \in \mathbb{R}^p$ and $z'$ where $p_Z(z') > 0$.*

Intuitively, overlap guarantees that $\mathbb{E}[X \mid U = u, Z = z]$ is well defined, and together with admissibility $\mathbb{E}[X \mid \mathrm{do}(U \coloneqq u)]$ can be uniquely expressed as a function of observational data distributions, and computed via the adjustment formula.

Naturally, $z$ is admissible with respect to $u$ and $x$ in the standard three variable graph of Figure 1(a). The corresponding overlap assumption for identifiability, also known as common support assumption, would require that for any given $z$, the variable $u$ takes on any value with non-zero probability. However, as we argued in the introduction this is not typically the case in digital platform settings where values of $x$ and $u$ from previous time steps confound future actions. In particular, if the new platform action $u$ is a deterministic function of these two variables, overlap of $(u, z)$ is necessarily violated. This holds even if dynamics are Markovian—namely, when $z$ contains only the values of $u$ and $x$ from one previous time step.

Now, let us return to our model that treats $u$ and $x$ separately. The first key observation is that $x_{T-1}$ is admissible with respect to $u_{T-1}$ and $x_T$, see Appendix F.2 for the proof.

**Proposition 1** (Admissibility in our model). *Given the structural equations in (1) and let Assumption 1 hold. Then, $x_{T-1}$ is admissible with respect to $u_{T-1}$ and $x_T$ for any $T \geq 0$.*

Thus, the main challenge for establishing identifiability of $\mathrm{PE}_T$ in our model is to argue about overlap of $(u_{T-1}, x_{T-1})$.

As we will show, this condition can be satisfied, even if $h$ and $r$ in (1) are deterministic functions. Once overlap is given, observations of $R_{K=1}$ are sufficient for identifiability.

We note that although our pathway to showing identifiability in this section is via showing overlap, unlike past work, we will *not* assume at the outset that overlap is satisfied. Instead, we will design digital-platform-specific sufficient conditions for satisfying overlap. We make these sufficient conditions explicit to give us the language to articulate the circumstances when identifiability is possible in settings without a randomized platform action.

## 3.1. Key assumptions

We highlight the two requirements on our dynamical system that will allow us to establish overlap of $(u_{T-1}, x_{T-1})$. The first assumption requires exogenous noise in the system that leads to sufficient variation in consumption $x$ across time.

**Definition 3.3** (Consumption shock). *For a given time step $t \geq 0$ we say there is a consumption shock at time $t$, if the noise $\xi_t$, with density $p_{\xi_t}$, satisfies $p_{\xi_t}(a) > 0$ for all $a \in \mathbb{R}^d$.*

We say the system is exposed to $M$ shocks prior to $T$ if for all $t \in \{T-M, \ldots, T-1\}$, there is a consumption shock. We expect such variations in consumption to naturally occur in the presence of unexpected news events, economic shocks, or new trends. In order to leverage these shocks for the purpose of identifiability, we rely on a second assumption: the platform needs to be sufficiently sensitive to the variations in consumption $x$, so that the consumption shocks propagate into the platform action $u$ at consecutive time steps.

**Definition 3.4** (Responsive platform action). *A platform action is responsive if for all $c \in \mathbb{R}^p$, $r(h(y) + c)$ is a surjective, continuously differentiable map with respect to $y \in \mathbb{R}^d$, and the Jacobian $J \in \mathbb{R}^{p,d}$ of $r(h(y) + c)$ with respect to $y$ satisfies $\mathrm{rank}(J) = \min(p, d)$ for all $y \in \mathbb{R}^d$.*

Intuitively, $r(h(y) + c)$ describes how the current state $y$ affects the next platform action, given that the previous platform action was $c$. This is an assumption which can be verified with enough knowledge of the design of the platform. Consider $r(u) = \alpha u$ as a plausible instantiation, recalling our example from Section 2.1. This corresponds to a model where the platform uses previous platform actions as a regularizer for how they select future actions. This choice of $r$ is surjective. We also expect the number of metrics and confounders which can be affected by platform actions to be large compared to the dimensionality of the platform action—i.e., $d \geq p$. Thus, because $h$ maps to a lower-dimensional space, we can reasonably expect $h$ to be surjective. Finally, because $r(h(y) + c)$ is the composition of $h$ and $r$, surjectivity of $r(h(y) + c)$ follows. The Jacobian rank condition imposes a form of "monotonicity" on $r(h(y) + c)$. In the video recommender system setting

this could correspond to: more views in category $i$ causes more recommendations in category $i$—a plausible assumption in a recommendation system. For another example of a responsive platform action—one based on a paramteric model taking gradient steps—see Appendix A.

Definition 3.4 is inspired by the notion of reachability in control theory. Dean et al. (2019) discusses the connection between reachability and recommender systems; they suggest that recommendation systems should be designed such that users have the ability to "reach" any recommendations they want to see indirectly via the actions they take. This prescription corresponds in spirit to the surjectivity condition of Definition 3.4.

### 3.2. General identifiability result

Building on the definitions from the previous section we are ready to present our main identifiability result. The proof can be found in Appendix F.3.

**Theorem 1.** *Let the dynamical system in* (1) *have a responsive platform action. Let Assumption 1 hold. Fix a $T \geq 2$ and let the auditor observe $R_{K=1}$. Then,*

a) *if the system exhibits $M = 2$ consumption shocks prior to $T$, then the treatment effect function $\mathrm{PE}_T(u, u')$ is identifiable for any $u, u' \in \mathbb{R}^p$.*

b) *There exists a system with $M < 2$ consumption shocks prior $T$, a distribution of $(x_{T-3}, u_{T-3})$, and an invertible $r$ such that for any $f, g, h$, such that for all $u \neq u'$, the function $\mathrm{PE}_T(u, u')$ is unidentifiable.*

This result states that consumption shocks on two preceding state variables are sufficient for the auditor to identify the performative effect from observations. In general, a single consumption shock is not enough for identifiability because $h$ can be a deterministic function. In this case, for any given combination of $x_t, u_{t-1}$, the auditor is only able to see one corresponding value of $u_t$. Thus, the second noise spike is necessary to add another degree of freedom which provides enough variation for overlap. We emphasize that our analysis crucially relies on accounting for how the noise propagates through the system across multiple time steps. In contrast, even multiple consumption shocks do not suffice for identifiability in the time-agnostic, standard causal model of Figure 1(a).

## 4. Identifiability in the linear model

In practice, an auditor may have access to longer rollouts of observations ($K > 1$) or the system may exhibit additional consumption shocks. A natural question is whether this information makes it easier to estimate $\mathrm{PE}_t(u, u')$. We investigate this question in the linear setting, instantiating

our model (1) as follows:

$$
\begin{aligned}
f(x) &\coloneqq Ax & g(u) &\coloneqq Bu \\
h(x) &\coloneqq Cx & r(u) &\coloneqq Du,
\end{aligned}
\tag{2}
$$

where $A \in \mathbb{R}^{d,d}, B \in \mathbb{R}^{d,p}, C \in \mathbb{R}^{p,d}, D \in \mathbb{R}^{p,p}$. The linear dynamics admit a clean characterization of the tradeoff between rollout length and conditions for identifiability. Linear state dynamics is certainly a strong assumption, but it has proven to be a useful approximation in control theory (e.g., Bouabdallah et al., 2004). In this linear setting, identifying the performative effect reduces to identifying the matrix $B$ because $\mathrm{PE}_T(u, u') = B(u' - u)$. We again consider identifiability under consumption shocks. However, for the linear case a weaker definition—one implied by full-support—suffices.

**Definition 4.1** (Fully-spanning consumption shock). *We say there is a fully-spanning consumption shock at time $t$, if $\xi_t$ is such that for all vectors $a \in \mathbb{R}^d$ with $a \neq 0$, $a^\top \xi_t$ is almost surely not a constant.*

We will also replace the responsive platform action assumption (Definition 3.4) with a full rank condition.

**Definition 4.2** (Full-row-rank platform action). *For $M \geq 2$, the platform has a full-row-rank platform action over a span of $M$ steps if $[DC, \ldots, D^{M-1}C]$ has full-row-rank.*

The condition for a full-row rank platform action is related the Kalman rank condition in control theory (Zabczyk, 1992; Gajic). In particular, recall that the Kalman rank condition is satisfied—i.e., the matrix $[B, AB, \ldots, A^{M-1}B]$ is full row rank—if and only if the system is controllable, meaning that given any initial state $x_0$ there exists control actions $u_0, \ldots, u_{M-1}$ which can reach any desired state $x_M$. The main differences in our condition is 1) we view the platform action as the state and the user consumption as the control action–this swaps $A$ with $D$ and $B$ with $C$ —and 2) rank condition matrix starts with $DC$ not $C$. To see the implications of these differences, consider rolling out $u_M$:

$$
u_M = D^M u_0 + C x_M + [DC, \ldots, D^{M-1}C]\begin{bmatrix} x_{M-1} \\ \vdots \\ x_1 \end{bmatrix}.
$$

Interpreting this equation through the lens of controllability, the full-row-rank platform action is satisfied if and only if for any choice of initial platform action $u_0$ and final consumption $x_M$, there exists a consumption profile $x_1, \ldots, x_{M-1}$ that induces any desired final platform action $u_M$.

### 4.1. Benefit of observing longer rollouts

In the linear setting, a full-row rank platform action which spans $M = 2$ time steps is also an expressive platform action (Definition 3.4), and vice versa. That said, Definition 4.2

generalizes Definition 3.4 beyond the $K = 1$ setting. This generalization allows us to characterize the benefits of observing longer rollouts, which we formalize in the following result. We note that unlike the more general result (Theorem 1), this result in the linear setting does not rely on showing overlap conditions are satisfied; instead, we reason about the uniqueness of the data generating process directly. The proof can be found in Appendix F.4.

**Theorem 2.** *Consider the dynamical system in* (1) *with linear functions $f, g, h, r$ defined in* (2). *Let Assumption 1 hold. Fix a time step $T \geq K + 1$, let the auditor observe iid samples of $R_K$. Let there be a fully-spanning consumption shock at time step $T - K$. Then,*

  a) *if $K = 1$, then for any $A, B, C, D$, there exists a distribution over $(x_{T-2}, u_{T-2})$ such that $\mathrm{PE}_T(u, u')$ is unidentifiable.*

  b) *if $K \geq 2$, then full-row-rank platform action over the span of $K$ steps is sufficient for identifiability of $\mathrm{PE}_T(u, u')$ for any $u, u'$.*

  c) *if $K \geq 2$, $x_{T-K-1} = u_{T-K-1} = 0$, and $\xi_t = 0$ for $t \geq T - K + 1$, then full-row-rank platform action over the span of $K$ steps is necessary for identifiability of $\mathrm{PE}_T(u, u')$ for any $u, u'$.*

Theorem 2 characterizes the tradeoff between identifiability, length of the observed rollout, and rank conditions on the platform dynamics matrices in the linear setting. We see that one consumption shock is not enough to identify the performative effect from only observations of $R_{K=1}$—just like in the general setting—but given additional observations $R_{K \geq 2}$ one consumption shock suffices in the linear case. Moreover, as $K$ gets larger, the rank assumptions required become easier to satisfy, allowing for more poorly conditioned dynamical systems to be identifiable.

More broadly, this result implies that an auditor can leverage longer interaction sequences between consumer and platform to make it easier to identify performativity. We note that the proof technique we use to show Theorem 2 can be extended to analyze linear systems where there is more than one consumption shock; however, unifying this result with settings where only a subset of a long rollout is observed is an open question we defer to future work. Preliminary empirical investigations Appendix E suggest that—similar to longer rollouts—more consumption shocks also weaken the required conditions for identifiability. We also defer the extension of the proof techniques used in proving Theorem 2—which do not rely on showing overlap conditions are satisfied—to non-linear settings for future work.

## 5. Estimation from finite samples

In practice, an auditor will only have access to a finite number of observations. We discuss two finite-sample estimators

for measuring $\mathrm{PE}_T$.

### 5.1. Adjustment formula estimator

The adjustment formula offers a direct way to estimate $\mathrm{PE}_T$ from observations of $R_{K=1}$ in a non-parametric way. In particular, we replace the population conditional expectations and probabilities with empirical estimates. Letting $\mathcal{S}_u$ denote the observations of $(u_{t-1}, x_{t-1})$ where $u_{t-1} = u$ and $\widehat{E}, \widehat{P}$ denote empirical estimates of expectations and probabilities respectively, the adjustment formula estimator is defined as

$$\sum_{x \in \mathcal{S}_u} \widehat{E}[x_t | u_{t-1} = u, x_{t-1} = x] \widehat{P}(x_{t-1} = x) \qquad (3)$$

Finite-sample analysis of the adjustment formula estimator can be found in Appendix B. This estimator is also applicable to the standard causal model in Figure 1(a), though at the cost of being ill-defined if overlap conditions are violated—i.e., when some pairs of $(u, x)$ have no observations.

### 5.2. Two-stage regression estimator

We also introduce a two-stage regression estimator tailored to the time-aware structure of our data generation model. This estimator is applicable if observations of $R_{K=2}$ are available, and can always be computed, even if the overlap conditions needed for theoretical guarantees do not hold. We will analyze the two-stage regression estimator in the linear setting from Section 4 and without loss of generality, we set $T = 3$. Our results can be generalized to settings where $f, g, h, r$ are from a non-linear function class (e.g., via Rademacher complexity arguments), but we focus on the simple linear setting for the sake of clarity. Recall that recovering $B$ is sufficient to estimate $\mathrm{PE}_T(u, u')$ in the linear case, as

$$\mathrm{PE}_T(u, u') = B(u' - u).$$

We let $x_t^{(k)}, u_t^{(k)}, \xi_t^{(k)}$ denote the $k$th observations of $x_t$, $u_t$, and $\xi_t$ respectively. Let $X_t \in \mathbb{R}^{d,n}$, $U_t \in \mathbb{R}^{p,n}$, and $E_t \in \mathbb{R}^{d,n}$ be matrices that comprise the $n$ samples of $x_t$, $u_t$, and $\xi_t$ respectively. The two-stage regression estimator is defined as $\widehat{B}$, where

$$\widehat{C} := \underset{C \in \mathbb{R}^{p,d}}{\mathrm{argmin}} \, \frac{1}{2n} \|U_1 - CX_1\|_{\mathrm{Fr}}^2$$

$$\widehat{H} := \underset{H \in \mathbb{R}^{d,d}}{\mathrm{argmin}} \, \frac{1}{2n} \|X_2 - HX_1\|_{\mathrm{Fr}}^2$$

$$\widehat{B} := \underset{B \in \mathbb{R}^{d,p}}{\mathrm{argmin}} \, \frac{1}{2n} \left\|X_3 - \widehat{H}X_2 - B(U_2 - \widehat{C}X_2)\right\|_{\mathrm{Fr}}^2.$$

We need the following assumption in order to provide our convergence guarantees for this estimator.

**Assumption 2** ($\rho$-Bounded System Dynamics). *The linear dynamical system specified by* (1) *and* (2) *has $\rho$-Bounded System Dynamics if* $\|A + BC\|_{\text{op}} \leq \rho\sigma_{\min}(DC)$.

Intuitively, Assumption 2 ensures that the magnitude of state and platform actions are of the same scale. We will use the notation $\kappa_A \coloneqq \sigma_{\max}(A)/\sigma_{\min}(A)$ to denote the condition number of a matrix $A$ and $\hat{\Sigma}_1 \coloneqq \frac{1}{n}\sum_{k=1}^n \xi_1^{(k)}(\xi_1^{(k)})^\top$ to denote the sample covariance of $\xi_1$. Theorem 3 provides a convergence result for the two-stage regression estimator of $B$, assuming $\xi_3 = 0$; the proof can be found in Appendix F.5.

**Theorem 3.** *Consider the dynamical system in* (1) *with $x_0 = u_0 = \xi_3 = 0$, with $f, g, h, r$ defined in* (2)*, and with full-row-rank platform action over the span of $K = 2$ steps. Let the auditor observe $n$ iid samples of $R_{K=2}$. Let $\mathbb{E}\|\xi_2\|_2^2 = \sigma_2^2 d$, and Assumption 2 hold. Let $\mathcal{G}$ be the event where $X_1 X_1^\top$ is invertible. If $\mathbb{E}\left[\kappa_{\hat{\Sigma}_1}^2 \lambda_{\min}(\hat{\Sigma}_1)^{-1}\right] \leq \tau_1$, then*

$$\frac{1}{pd}\mathbb{E}\left[\|\widehat{B} - B\|_{\text{Fr}}^2 \mid \mathcal{G}\right] \leq \frac{\sigma_2^2 \rho^2 \kappa_{DC}^2 \tau_1}{n}.$$

To illustrate this result, consider a simple Gaussian noise example. Suppose $\xi_1$ and $\xi_2$ are drawn iid from $\mathsf{N}(0, \sigma_1^2 I_d)$ and $p = d$. We have $\mathbb{E}\|\xi_2\|_2^2 = \sigma_2^2 d$. For $n \geq d$, $\hat{\Sigma}_1$ is almost surely invertible. $(X_1 X_1^\top/\sigma_1^2)^{-1}$ has an inverse Wishart distribution and thus, $\mathbb{E}[\hat{\Sigma}_1^{-1}] = \frac{n}{(n-d-1)\sigma_1^2} I_d$ for $n > d + 1$. Theorem 3 gives $\mathbb{E}\|\widehat{B} - B\|_{\text{Fr}}^2 \leq \frac{d^2 \sigma_2^2 \rho^2}{(n-d-1)\sigma_1^2}$, scaling like the linear regression error rate.

# 6. Empirical investigations

Now we use semi-synthetic recommender system data to empirically investigate the effectiveness of our modeling assumptions in alleviating overlap violations. We focus on the adjustment formula estimator here, deferring experiments about the two-stage regression estimator to the Appendix. We simulate an online platform with items, recommendations, users, and ratings using RecLab (Krauth et al., 2020) with the goal of estimating the effect the recommended topic has on user ratings.

**RecLab Simulator.** In the RecLab Topic environment, we create 1000 items that the platform can recommend to 2000 users. Each item has a topic attribute $w \in \{0, 1, 2, 3\}$. The simulator is such that for any given topic $w$, user $s$ has a preference $b_{s,w} \sim \text{Unif}(0.5, 5.5)$. When user $s$ is recommended an item of topic $w$, they will rate it as $b_{s,w} + \mathcal{N}(0, 0.01)$ clipped to be between 1 and 5. If the user has been recommended a item of topic $w$ within its last 3 recommendations, they are in an "echo chamber" and will rate the item as $b_{s,w} + \mathcal{N}(0, 0.01) + 1$. We simulate $N$ time steps, where a time step corresponds to recommending each of the 2000 users one item, having them rate the item, and updating the

| | Recommender | Estimated $TE_1$ | Fraction of undefined terms | Probability mass of undefined terms |
|---|---|---|---|---|
| K=1 | Random | 4.58 | 0 / 3 | 0.0% |
| | EASE (no update) | 4.65 | 0 / 3 | 0.0% |
| | EASE | 4.68 | 0 / 3 | 0.0% |
| K=2 | Random | 4.56 | 0 / 9 | 0.0% |
| | EASE (no update) | 4.66 | 0 / 9 | 0.0% |
| | EASE | 4.69 | 0 / 9 | 0.0% |
| K=3 | Random | N/A | 4 / 27 | 4.1% |
| | EASE (no update) | N/A | 5 / 26 | 5.2% |
| | EASE | N/A | 4 / 27 | 6.6% |

**Table 1.** $TE_1$ estimated using the adjustment formula estimator with $N = 751$

model. We will use two different recommender systems provided by RecLab for choosing the items recommended to users: the EASE recommender (Steck, 2019) and a random recommender. We initiate the EASE recommender by simulating 100 cycles. In the 'no-update' variation, we fix the EASE recommender weights to their initial value, and in the other variation, we continue to update the weights during the following $N$ cycles based on user ratings. We will vary $N$ in this experiment to explore how the number of samples affects overlap.

**Causal effect of interest.** We will focus on one particular user's ratings[2]. Let $x_t \in [1, 5]$ will denote this user's rating of the item recommended to them at time step $t$. Similarly, $u_t \in \{0, 1, 2, 3\}$ will denote the topic of the item recommended to our user of interest at time step $t$. We aim to estimate $TE_1 \coloneqq \mathbb{E}[x_t \mid \text{do}(u_{t-1} = 1)]$ from which we can construct the performative effect the recommendation of topic 1 has on the user's rating. We look at a sliding window over the data $\{(x_{t-K}, \ldots, x_t, u_{t-K}, \ldots, u_{t-1})\}_{t=K}^N$. We treat these samples as the iid observations of $R_K$. We will let $z_{t-1} \coloneqq (x_{t-K}, \ldots, x_{t-1})$ be the confounders. We will vary $K$ to explore how the size of the confounding set affects estimation. Note that $K = 3$ is the well-specified setting as users have a memory length of 3. We will use the adjustment formula estimator, as the simulated dynamics are non-separable.

## 6.1. Coverage

We compute the adjustment formula estimator, described in (3), for each recommender system simulated for $N = 751$ time steps. We refer to the pairs of $(u, z)$ for which there are no observations to calculate the conditional expectation in (3) as the undefined terms. In Table 1, we report the number of undefined terms observed for various choices of $K$ and recommender system. We also report the probability mass of $z_{t-1}$ corresponding to these undefined terms. When $K = 3$, the dimensionality becomes too large, and there are not enough samples to get coverage. The esti-

---

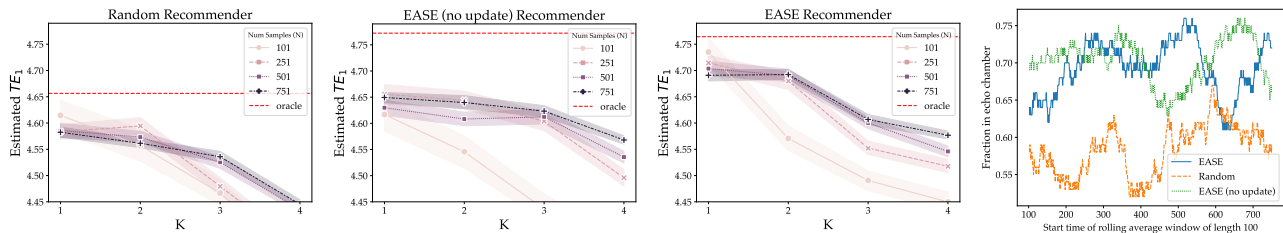[2]the user considered has $b_{s,w} = [3.244, 4.076, 3.514, 3.224]$

**Figure 2.** The 3 plots on the left are bootstrapped estimates of $\text{TE}_1$ with 95% confidence intervals using the adjustment formula estimator for different recommenders. Right-most plot is a rolling average of how often the user of interest is in an echo chamber.

mates of $\text{TE}_1$ between $K = 1$ and $K = 2$ are comparable, suggesting that the Markovian assumption does not hinder accuracy. We note that for $K = 2$, there are often not enough samples to estimate the treatment effects of the other topics $\text{TE}_0, \text{TE}_2, \text{TE}_3$. Our findings are consistent with the combinatorial nature of calculating terms within the adjustment formula estimator: if $n$ samples are needed for overlap when $K = 1$, then $\approx n^K$ samples are needed for overlap when $K > 1$.

### 6.2. Identification

We now modify the adjustment formula estimator to silently fail when overlap does not hold. In particular, we use the empirical average rating $\widehat{E}[x_t]$ in place of $\widehat{E}[x_t | u_{t-1} = u, x_{t-1} = x]$ for $(u, z)$ with no observed samples in (3). We bootstrap this estimator with 40 bootstrap samples and plot our results on estimating $\text{TE}_1$ for various choices of $K$ and $N$ in Figure 2. The random recommender system serves as a sanity check. Because there are no confounders by construction, the treatment effect is just the expected rating $4.076 + (1 - (3/4)^3) = 4.65$, which is close to the $4.61$ we estimate when $K = 1$. For both variations of the EASE recommender system, we estimate the ground truth treatment effect by exploiting our knowledge of user rating dynamics. In particular, we compute the fraction of the total 751 time steps the user is in an echo chamber $n/751$ and then compute an estimate of the expected rating $4.076 + n/751$. This baseline computation assume time-invariant dynamics, which is amenable to the EASE (no update) recommender. We also plot a sliding window average of the time our user of interest is stuck in an echo chamber in Figure 2 (right). We see that this fraction stays fairly constant across time, providing evidence that a time-invariant assumption is reasonable. The estimated treatment effect of $4.74$ reported with $K = 1$ for the EASE is close to the ground truth treatment effect of $4.76$. There is a gap between the estimated and ground truth treatment effects for the EASE (no update) recommender, but the estimate is still relatively higher than the estimate of the treatment effect of the random recommender. Our results suggest that our method can effectively capture EASE's superiority in achieving higher user engagement relative to the random model and also reveals that

EASE learns to take advantage of echo chamber effects to boost user ratings. The treatment effect dips down for larger $K$ in all of the plots because coverage is poor, the estimator silently fails, causing it to underestimate the treatment effect, corroborating our findings from Table 1. This being said, we also see that as the number of observations $N$ increases, this deflation due to coverage mitigates as expected, with the treatment effect gradually moving up towards the treatment effect estimated with $K = 1$.

### 6.3. Ablations

We report several additional experiments in the Appendix. In Appendix C we experiment on real time-series data from a related microeconomic use-case. In Appendix D we investigate how the Markovian assumption we use performs in high-dimensional non-Markovian models using synthetic data. In both of these experiments, we again observe that the benefits of reducing overlap violations outweigh the costs arising from potential model misspecification in terms of minimizing estimation error.

## 7. Conclusion

In this work, we discuss sufficient conditions for the identifiability of causal effects from natural interactions of users with digital platforms. We explain why the standard causal model does not satisfy the typical overlap assumptions required for causal inference of performative effects from observational data. As a solution, we propose a time-aware causal model, and we demonstrate on this model that valid inferences of performativity are possible even when algorithmic actions are chosen deterministically based on historical data. We exploit repeated interactions between participants and the platform as well as natural variations in participant behavior without assuming randomization in the platform action. We provide a finite-sample estimator tailored to this digital platform setting, and we show how observations of longer interaction sequences can be beneficial for estimation. More broadly, our work demonstrates how connecting causal inference and control theory helps justify the use of observational causal inference in digital platforms, providing an important guardrail for digital market investigations.

## Impact Statement

## Acknowledgements

## References

Abadie, A. and Gardeazabal, J. The economic costs of conflict: A case study of the basque country. 2003.

Abadie, A., Diamond, A., and Hainmueller, J. Synthetic control methods for comparative case studies: Estimating the effect of california's tobacco control program. *Journal of the American Statistical Association*, 105(490):493–505, 2010.

Abbasi-Yadkori, Y. and Szepesvari, C. Regret bounds for the adaptive control of linear quadratic systems. In *COLT*, 2011.

Adomavicius, G., Bockstedt, J. C., Curley, S. P., and Zhang, J. Do recommender systems manipulate consumer preferences? a study of anchoring effects. *Information Systems Research*, 24(4):956–975, 2013.

Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., and Bonneau, R. Tweeting from left to right. *Psychological Science*, 26:1531 – 1542, 2015.

Bechavod, Y., Ligett, K., Wu, S., and Ziani, J. Gaming helps! learning from strategic interactions in natural dynamics. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130, pp. 1234–1242, 2021.

Bertrand, M., Duflo, E., and Mullainathan, S. How much should we trust differences-in-differences estimates? *Experimental & Empirical Studies eJournal*, 2001.

Blackwell, M. A framework for dynamic causal inference in political science. *American Journal of Political Science*, 57:504–520, 2013.

Bottou, L., Peters, J., Quiñonero Candela, J., Charles, D. X., Chickering, D. M., Portugaly, E., Ray, D., Simard, P., and Snelson, E. Counterfactual reasoning and learning systems: the example of computational advertising. 14 (1):3207–3260, 2013. ISSN 1532-4435.

Bouabdallah, S., Noth, A., and Siegwart, R. PID vs LQ control techniques applied to an indoor micro quadrotor. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, volume 3, 2004.

Bouneffouf, D. and Rish, I. A survey on practical applications of multi-armed and contextual bandits. *ArXiv*, abs/1904.10040, 2019.

Brown, G., Hod, S., and Kalemaj, I. Performative prediction in a stateful world. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151, pp. 6045–6061, 2022a.

Brown, M., Bisbee, J. H., Lai, A., Bonneau, R., Nagler, J., and Tucker, J. A. Echo chambers, rabbit holes, and algorithmic bias: How youtube recommends content to real users. *SSRN Electronic Journal*, 2022b.

Bruder, D., Remy, C. D., and Vasudevan, R. Nonlinear system identification of soft robot dynamics using koopman operator theory. *2019 International Conference on Robotics and Automation (ICRA)*, pp. 6244–6250, 2018.

Chaney, A. J.-B., Stewart, B. M., and Engelhardt, B. E. How algorithmic confounding in recommendation systems increases homogeneity and decreases utility. *Proceedings of the 12th ACM Conference on Recommender Systems*, 2018.

Chen, X., Hong, H., and Tarozzi, A. Semiparametric efficiency in GMM models with auxiliary data. *Annals of Statistics*, 36:808–843, 2007.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. M. Double/debiased machine learning for treatment and structural parameters. *Econometrics: Econometric & Statistical Methods - Special Topics eJournal*, 2017.

D'Amour, A., Ding, P., Feller, A., Lei, L., and Sekhon, J. S. Overlap in observational studies with high-dimensional covariates. *Journal of Econometrics*, 2017.

de Luna, X., Warenbaum, I., and Richardson, T. S. Covariate selection for the nonparametric estimation of an average treatment effect. *Biometrika*, 98:861–875, 2011.

Dean, S. and Morgenstern, J. H. Preference dynamics under personalized recommendations. *Proceedings of the 23rd ACM Conference on Economics and Computation*, 2022.

Dean, S., Rich, S., and Recht, B. Recommendations and user agency: the reachability of collaboratively-filtered information. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 2019.

Dwivedi, R., Murphy, S. A., and Shah, D. Counterfactual inference for sequential experimental design. *ArXiv*, abs/2202.06891, 2022.

Eilat, I. and Rosenfeld, N. Performative recommendation: diversifying content via strategic incentives. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.

Fleder, D., Hosanagar, K., and buja, a. Recommender systems and their effects on consumers: the fragmentation debate. 06 2010.

Gajic, Z. URL https://www.ece.rutgers.edu/~gajic/psfiles/chap5traCO.pdf.

Hardt, M. and Mendler-Dünner, C. Performative prediction: Past and future, 2023.

Hardt, M., Jagadeesan, M., and Mendler-Dünner, C. Performative power. In *Advances in Neural Information Processing Systems*, 2022.

Harris, K., Ngo, D. D. T., Stapleton, L., Heidari, H., and Wu, S. Strategic instrumental variable regression: Recovering causal relationships from strategic responses. In *International Conference on Machine Learning*, pp. 8502–8522. PMLR, 2022.

Horowitz, G. and Rosenfeld, N. Causal strategic classification: A tale of two shifts. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 13233–13253. PMLR, 23–29 Jul 2023.

Imbens, G. W. Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review. *The Review of Economics and Statistics*, 86(1):4–29, 02 2004.

Jambor, T., Wang, J., and Lathia, N. Using control theory for stable and efficient recommender systems. *Proceedings of the 21st international conference on World Wide Web*, 2012.

Jin, C., Netrapalli, P., Ge, R., Kakade, S. M., and Jordan, M. I. A short note on concentration inequalities for random vectors with subgaussian norm. *ArXiv*, abs/1902.03736, 2019.

Kalimeris, D., Bhagat, S., Kalyanaraman, S., and Weinsberg, U. Preference amplification in recommender systems. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021.

Kiggins, J. Avocado prices. https://www.kaggle.com/datasets/neuromusic/avocado-prices, 2018.

Kim, I. S., Rauh, A., Wang, E. H., and Imai, K. Matching methods for causal inference with time-series cross-sectional data. *American Journal of Political Science*, 2020.

Klasnja, P. V., Hekler, E. B., Shiffman, S., Boruvka, A., Almirall, D., Tewari, A., and Murphy, S. A. Micro-randomized trials : An experimental design for developing just-intime adaptive interventions. 2015.

Kramer, A. D. I., Guillory, J., and Hancock, J. T. Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences of the United States of America*, 2014.

Krauth, K., Dean, S., Zhao, A., Guo, W., Curmei, M., Recht, B., and Jordan, M. I. Do offline metrics predict online performance in recommender systems? *ArXiv*, abs/2011.07931, 2020.

Krauth, K., Wang, Y., and Jordan, M. Breaking feedback loops in recommender systems with causal inference. *ArXiv*, abs/2207.01616, 2022.

Langford, J. and Zhang, T. The epoch-greedy algorithm for contextual multi-armed bandits. In *Advances in Neural Information Processing Systems*, 2007.

Ljung, L. Perspectives on system identification. *Annual Reviews in Control*, 34(1):1–12, 2010.

Lowther, G. Is the image of a null set under a differentiable map always null? Mathematics Stack Exchange. URL https://math.stackexchange.com/q/59115. (version: 2011-08-25).

Mandal, D., Triantafyllou, S., and Radanovic, G. Performative reinforcement learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 23642–23680. PMLR, 23–29 Jul 2023.

Mendler-Dünner, C., Ding, F., and Wang, Y. Anticipating performativity by predicting from predictions. In *Advances in Neural Information Processing Systems*, 2022.

Miller, J., Milli, S., and Hardt, M. Strategic classification is causal modeling in disguise. In *Proceedings of the 37th International Conference on Machine Learning*, volume

119 of *Proceedings of Machine Learning Research*, pp. 6917–6926. PMLR, 13–18 Jul 2020.

Pearl, J. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd edition, 2009. ISBN 052189560X.

Perdomo, J., Zrnic, T., Mendler-Dünner, C., and Hardt, M. Performative prediction. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.

Petersen, M. L., Porter, K. E., Gruber, S., Wang, Y., and van der Laan, M. J. Diagnosing and responding to violations in the positivity assumption. *Statistical Methods in Medical Research*, 2012a.

Petersen, M. L., Porter, K. E., Gruber, S., Wang, Y., and van der Laan, M. J. Diagnosing and responding to violations in the positivity assumption. *Statistical Methods in Medical Research*, 21(1):31–54, 2012b.

PNAS. Editorial expression of concern: Experimental evidence of massivescale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, 2014.

Rosenbaum, P. R. and Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70:41–55, 1983.

Santos, A. If $f \in \mathcal{C}^1$ and $\{\nabla f = 0\}$ has lebesgue measure $0$, then $\{f \in b\}$ has lebesgue measure $0$ for all borel measurable $b \subset \mathbb{R}$ with lebesgue measure $0$. Mathematics Stack Exchange. URL https://math.stackexchange.com/q/3216190. (version: 2019-07-05).

Sasaki, Y. and Ura, T. Inference for moments of ratios with robustness against large trimming bias and unknown convergence rate. *arXiv: Methodology*, 2017.

Shah, A., Dwivedi, R., Shah, D., and Wornell, G. W. On counterfactual inference with unobserved confounding. *ArXiv*, abs/2211.08209, 2022.

Shankar, S., Garcia, R., Hellerstein, J. M., and Parameswaran, A. G. Operationalizing machine learning: An interview study. *ArXiv*, abs/2209.09125, 2022.

Shavit, Y., Edelman, B. L., and Axelrod, B. Causal strategic linear regression. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20, 2020.

Shmueli, G. and Tafti, A. "Improving" prediction of human behavior using behavior modification. *Arxiv:2008.12138*, 2020.

Steck, H. Embarrassingly shallow autoencoders for sparse data. *The World Wide Web Conference*, 2019.

Thai, J., Laurent-Brouty, N., and Bayen, A. M. Negative externalities of gps-enabled routing applications: A game theoretical approach. In *IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2016.

Wagner, C., Strohmaier, M., Olteanu, A., Kıcıman, E., Contractor, N., and Eliassi-Rad, T. Measuring algorithmically infused societies. *Nature*, 595(7866):197–204, 2021a.

Wagner, C., Strohmaier, M., Olteanu, A., Kıcıman, E., Contractor, N., and Eliassi-Rad, T. Measuring algorithmically infused societies. *Nature*, 595(7866):197–204, 2021b.

Wang, X., Yau, C.-Y., and Wai, H. T. Network effects in performative prediction games. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 36514–36540, 2023.

Yang, S. and Ding, P. Asymptotic inference of causal effects with observational studies trimmed by the estimated propensity scores. *Biometrika*, 105:487–493, 2018.

Zabczyk, J. Mathematical control theory - an introduction. In *Systems & Control: Foundations & Applications*, 1992.

Zhang, J. and Bareinboim, E. Near-optimal reinforcement learning in dynamic treatment regimes. In *Neural Information Processing Systems*, 2019.

## A. Parametric recommender

To illustrate an example of a responsive platform action, let's consider a setting where a user is iteratively interacting with a recommender system parameterized by weights $\theta \in \mathbb{R}^d$. Let $z \in \mathbb{R}^d$ be some fixed, time-invariant feature about the user. Let the state at time $t$ denoted as $x_t \in \mathbb{R}$ be the number of ads the user clicked on at time $t$. At time $t$, the recommender system outputs $u_t = \theta_t^T z$, a scalar quantifying the number of ads to serve to the user. The recommender system's goal is select $\theta$ to minimize $\ell(\theta; x) = (\gamma z^T \theta - x)^2/2$ for some $\gamma > 0$ which is an estimate of the click-through-rate for an ad shown to a given user. In other words, the recommender system wants to serve ads in proportion to the number ads they click on. Intuitively, if they serve more, this hurts the user experience; if they serve less, then they are losing ad revenue. The recommender system will update its weights at each time step with gradient descent: $\theta_t = \theta_{t-1} - \alpha \nabla_\theta \ell(\theta; x_t) = \theta_{t-1} - \alpha \gamma z (\gamma z^T \theta_{t-1} - x_t)$. Unrolling $u_t$, we have

$$u_t = \gamma z^T \theta_{t-1} - \alpha \gamma^2 \|z\|_2^2 (\gamma z^T \theta_{t-1} - x_t)$$
$$= (1 - \alpha \gamma^2 \|z\|_2^2) u_{t-1} + \alpha \gamma^2 \|z\|_2^2 x_t$$

We have that $h(x) = \alpha \gamma^2 \|z\|_2^2 x$ and $r(u) = (1 - \alpha \gamma^2 \|z\|_2^2)u$. As long as $\alpha \gamma^2 \|z\|_2^2 \neq 0$ and $\alpha \gamma^2 \|z\|_2^2 \neq 1$, then $h$ and $r$ are surjective functions. This implies that $r(h(x) + c)$ is surjective for any choice of $x$ and $c$, which means that this recommender system is a responsive platform action.

## B. Adjustment formula estimator

Admissibility of the dynamical system we are studying (Proposition 1) makes estimating the adjustment formula (Definition 3.1) sufficient for estimating the performative effect. Since $x$ and $u$ can take on continuous values we start with discretizations of $\mathbb{R}^d$ and $\mathbb{R}^p$ denoted as finite collections of bounded, non-intersecting sets $\mathcal{N} \coloneqq \{\mathcal{X}_\alpha\}$ and $\mathcal{M} \coloneqq \{\mathcal{U}_\beta\}$ indexed by $\alpha$ and $\beta$ respectively. Suppose that every element of $\mathcal{N}$ and $\mathcal{M}$ has diameter at most $\varepsilon/2$ and has Lebesgue measure greater than 0. For a point $x \in \cup \mathcal{N}$, define $\alpha(x)$ such that $x \in \mathcal{X}_{\alpha(x)}$. Define $\beta(u)$ respectively. We will assume we have $n$ samples of the form $\mathcal{D}^n = \{(x_1^{(k)}, u_1^{(k)}, x_2^{(k)})\}_{k=1}^n$, where every sample is drawn iid from (1). With these quantities, we form estimates of the components of the adjustment formula; here without loss of generality, we set $T = 2$.

$$\widehat{\mathbb{E}}[x_2 \mid u_1 \in \mathcal{U}, x_1 \in \mathcal{X}] \coloneqq \frac{\sum_{k \in [n]} x_2^{(k)} \mathbf{1}\left\{u_1^{(k)} \in \mathcal{U}, x_1^{(k)} \in \mathcal{X}\right\}}{\sum_{k \in [n]} \mathbf{1}\left\{u_1^{(k)} \in \mathcal{U}, x_1^{(k)} \in \mathcal{X}\right\}}$$

$$\widehat{P}(x_1 \in \mathcal{X}) \coloneqq \frac{1}{n} \sum_{k=1}^n \mathbf{1}\left\{x_1^{(k)} \in \mathcal{X}\right\}.$$

After combining, we have an estimate of the performative effect:

$$\hat{x}_2(u) \coloneqq \sum_\alpha \widehat{\mathbb{E}}[x_2 \mid u_1 \in \mathcal{U}_{\beta(u)}, x_1 \in \mathcal{X}_\alpha] \widehat{P}(x_1 \in \mathcal{X}_\alpha).$$

We will need some mild assumptions to prove a guarantee on the estimator. Our first assumption controls how much previous user state and platform actions affect future state actions. The magnitude of the effect must be bounded in proportion to the inputs.

**Assumption 3.** *The relationship between $x_2$ and $x_1, u_1$ is L-Lipschitz continuous in the sense that for any $w, w' \in \mathbb{R}^d$ and $u, u' \in \mathbb{R}^p$, and with $v \coloneqq [u^\top, w^\top]^\top$, it holds that*

$$\|\mathbb{E}[x_2 | u_1 = u, x_1 = w] - \mathbb{E}[x_2 | u_1 = u', x_1 = w']\| \le L \|v - v'\|.$$

We also need to control how far the discretized conditional expectation $\mathbb{E}[x_2 \mid u_1 \in \mathcal{U}, x_1 \in \mathcal{X}]$ deviates from $\mathbb{E}[x_2 | u_1 = u, x_1 = x]$. To do this, we impose a regularity condition on the conditional distribution.

**Assumption 4.** *Let $w, w' \in \mathbb{R}^d$ and $u, u' \in \mathbb{R}^p$, and with $v \coloneqq [u^\top, w^\top]^\top$ be such that $\|v - v'\| \le \varepsilon$. Then, for any $x \in \cup \mathcal{N}$, the following condition on the density $p$ holds for some $\eta(\varepsilon) \in (0, 1)$ such that $\lim_{\varepsilon \to 0} \eta(\varepsilon) = 0$:*

$$1 - \eta(\varepsilon) \le \frac{p(u_1 = u, x_1 = w | x_2 = x)}{p(u_1 = u', x_1 = w' | x_2 = x)} \le 1 + \eta(\varepsilon).$$

This assumption ensures that the conditional distribution is "stable" in any $\varepsilon$-neighborhood. Finally, we need one more assumption which guarantees we obtain enough samples for every slice of data. Assumption 5 is defined with respect to the variables: cover granularity $\varepsilon > 0$, error tolerances $\delta \in (0,1)$ and $\gamma > 0$, and failure probability tolerance $\rho \in (0,1)$.

**Assumption 5.** *Let $n_{\mathcal{U},\mathcal{X}} \coloneqq \sum_{k\in[n]} \mathbf{1}\left\{u_1^{(k)} \in \mathcal{U}, x_1^{(k)} \in \mathcal{X}\right\}$. Let $n_{\mathcal{U},\mathcal{X}} \geq \frac{2d\sigma^2}{\gamma^2} \log(4|\mathcal{N}|/\rho)$ for all $\mathcal{X} \in \mathcal{N}$ and $\mathcal{U} \in \mathcal{M}$. Further let $n \geq \max_{\mathcal{X}\in\mathcal{N}} \frac{1}{2\delta^2 P(x_1\in\mathcal{X})^2} \log(4|\mathcal{N}|/\rho)$.*

We present our convergence result now in Theorem 4.

**Theorem 4.** *Consider the dynamical system in (1) with any arbitrary $P_{-1}$. Let the auditor observe $n$ iid samples of $(x_1, u_1, x_2)$. Suppose $x_2$ is $\sigma^2$-subgaussian conditioned on $u_1$ and $x_1$. Let $\mathbb{E}[\xi_2|x_1 = x, u_1 = w] = 0$, $\mathbb{E}[\|\xi_2\| \mid x_1 = x, u_1 = w] \leq c_1$ for all $x \in \mathbb{R}^d$ and $w \in \mathbb{R}^p$. Let $f$ and $g$ be continuous functions, and define $R$ such that $\sup_{x\in\cup\mathcal{N}, w\in\cup\mathcal{M}} \max(\|f(x)\|, \|g(w)\|) \leq R$. Let the conditions of Theorem 1 hold, Assumption 3 hold with $L$, Assumption 4 hold with $\eta$, and Assumption 5 hold. For any specified $u \in \cup\mathcal{M}$ with probability at least $1 - \rho$, the following holds*

$$\|\hat{x}_2(u) - \mathbb{E}[x_2|do(u_1 \coloneqq u)]\| \leq \delta\gamma + 2\delta R + \gamma + \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)}(2R + c_1)$$
$$+ L\varepsilon + \mathbb{E}[\|f(x_1)\| \mathbf{1}\{x_1 \notin \cup\mathcal{N}\}] + (1 - P_{x_1}(\cup\mathcal{N}))R.$$

The proof of Theorem 4 can be found in Appendix F.6. Let us go through all the terms in the bound, to verify that they can all be made arbitrarily small (with sufficient samples). $\delta$ and $\gamma$ can be made smaller, so long as the auditor receives proportionally enough samples. The auditor can create a finer discretization to make $\varepsilon$ smaller and therefore $\eta$ smaller as well. If we assume that $\mathbb{E}[\|f(x_1)\|] \leq \infty$, then the last two terms tend to zero as the auditor's approximation of $\mathbb{R}^d$—i.e., $\cup\mathcal{N}$—covers more of the space.

## C. Estimating price elasticity of demand from real data

Besides digital platforms, our model also applies to some economic settings. Micro-economists are often interested in estimating the effect product prices have on demand, termed the *price elasticity of demand*. If we model product demand using $x_t$ and model product prices using $u_t$, then the price elasticity of demand is precisely the performative effect. Confounders like product quality can be accounted for in the state variable $x_t$. We use this setting to perform additional experiments.

### C.1. Setup

We apply our model to the task of estimating the price elasticity of demand (PED) from time series data. here, we are interested in how the price (platform action) affects the demand (consumption). We use an avocado time series dataset (Kiggins, 2018) that consists of biweekly measurements of the prices of avocados and the amount of avocados purchased by region in the US from 2015 to 2018. The avocado time series dataset is comprised of several time series spanning different regions of the United States. To construct the dataset we are operating on, we combine data from two regions—Southeast and Great Lakes—chosen by pricing and demand similarity. For a week $t \in [N]$, $u_t$ corresponds to the logged average avocado price, and $x_t$ corresponds the logged number of avocados purchased. We posit the following model:

$$x_t = \tilde{f}(z_{t-1}) + g(u_{t-1}).$$

where $z_{t-1}$ denotes the set of confounding variable that we adjust for, which we will specify shortly. In this model, the PED is defined as $\nabla g$. This quantity is a curve if the function $g$ is non-linear; however, in this section, we will assume that $g$ is linear, which reduces the problem of estimating the PED into one of estimating a scalar. We will analyze three estimators: adjustment formula estimator, random forest double ML (RF-DML), and linear regression double ML (LR-DML).

### C.2. Implementation details of estimators

We outline how our three estimators—the adjustment formula estimator, random forest double ML (RF-DML), and linear regression double ML (LR-DML)—are implemented. The adjustment formula estimator relies on computing (3) on a discretized platform action and consumption variables. The discretization is important to ensure overlap over confounder and treatment variables, as the adjustment formula estimator is not well defined without overlap. In particular, let $\{\mathcal{Z}_\gamma\}_\gamma, \{\mathcal{U}_\beta\}_\beta$
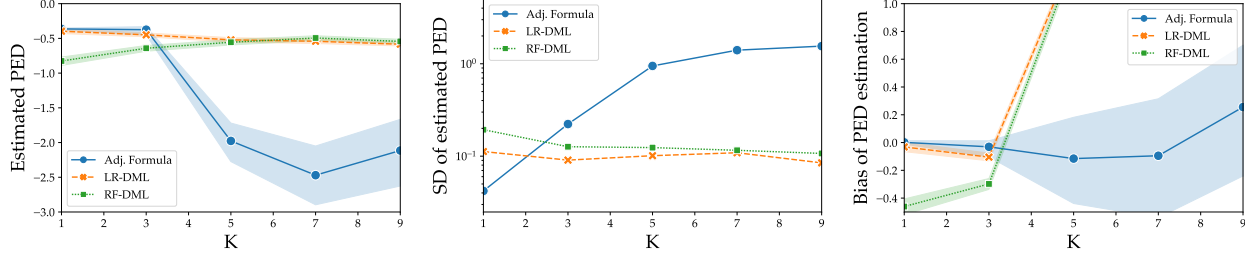
**Figure 3.** (left) Bootstrapped estimates of PED with 95% confidence intervals. (middle) Standard deviation of each estimator. (right) Bias defined as $\mathbb{E}\big[\Psi(Y^K)\big] - \Psi_a(y^K)$.

denote discretizations of the confounders $z_{t-1}$ and platform action $u_{t-1}$, and let $\beta(u)$ be such that $u \in \mathcal{U}_{\beta(u)}$. We define the adjustment formula estimator as:

$$\hat{x}(u) \coloneqq \sum_\gamma \left[ \frac{\sum_t x_{t+1} \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)} \right\}}{\sum_t \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)} \right\}} \right] \frac{\sum_t \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma \right\}}{n}.$$

Detailed discussion and theoretical guarantees regarding the adjustment formula can be found in Appendix B. We discretize the logged price into two buckets: $(-0.479, 0.131], (0.131, 0.683]$ and the logged demand into two buckets $(14.539, 15.014], (15.014, 15.837]$. After using the adjustment formula estimator to estimate the effect price has on demand, we then use this estimator to assign predicted demands to all of the prices observed in the dataset. We then use linear regression to estimate the slope of the relationship between predicted demand and price—this is what we refer to as the adjustment formula estimate of the PED. This approach is motivated by methods suggested by Petersen et al. (2012a). In our experiments, we let the adjustment formula estimator silently fail when overlap does not hold. This means for terms in the adjustment formula $\hat{x}$ defined above where $\sum_t \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma \right\} > 0$ and $\sum_t \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)} \right\} = 0$, we set $\frac{\sum_t x_{t+1} \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)} \right\}}{\sum_t \mathbf{1}\left\{ z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)} \right\}}$ equal to $0$. This modification could cause the adjustment formula estimator to underestimate the PED for large $K$, potentially causing the bias to spike for larger $K$ for LR-DML and RF-DML in Figure 3. We explore this issue further later.

The double machine learning approach (Chernozhukov et al., 2017) first uses half of the training data to residualize the confounders out of the treatment and effect. More specifically, let's split our samples of $x_t$ stored in a vector $X$ into two parts $X_a, X_b$. Let's do the same for $z_{t-1}$ and $u_{t-1}$ to form $Z_a, Z_b$ and $U_a$ and $U_b$ respectively. DML proceeds by solving

$$h = \underset{h' \in \mathcal{F}}{\operatorname{argmin}} L(h'(Z_a), U_a)$$
$$f = \underset{f' \in \mathcal{F}}{\operatorname{argmin}} L(f'(Z_a), X_a)$$
$$g = (U_b - h(Z_b))^\top U_b (U_b - h(Z_b))^\top (X_b - f(Z_b)),$$

where the first two steps are a residualizing procedure; here, for a function $h'$ the notation $h'(Z_a)$ denotes applying $h'$ on each row of $Z_a$. For LR-DML, the residualizing procedure uses linear regression; i.e., $\mathcal{F}$ corresponds to the class of linear functions and loss function $L$ is squared loss. For RF-DML, the residualizing procedure uses a random forest model implemented using sci-kit-learn 1.2.0. Then, in the second step, both RF-DML and LR-DML use the other half of the training data to perform a slightly modified version linear regression—see Equation (1.5) in Chernozhukov et al. (2017)—on the residualized treatment and residualized effect. The slope of this estimated line is the estimated PED.

### C.3. Empirical findings

**Varying the adjustment set to characterize overlap violations.** Our focus in this section is to investigate whether the Markovian assumption on the system dynamics our model posits actually mitigates overlap violations. We will vary the size of the confounding set to measure the variance and bias of different estimators as a proxy for overlap violations. We look at a sliding window over the data $\{(x_{t-K}, \ldots, x_t, u_{t-K}, \ldots, u_{t-1})\}_{t=K}^N$. We treat these samples as the iid observations of $R_K$ that the auditor observes. We will use $u_{t-1}$ as the treatment variable, $z_{t-1} \coloneqq (x_{t-K}, \ldots, x_{t-1})$ as the confounders, $x_t$ as the outcome. We will vary $K$ to explore how the size of the confounding set affects estimation.
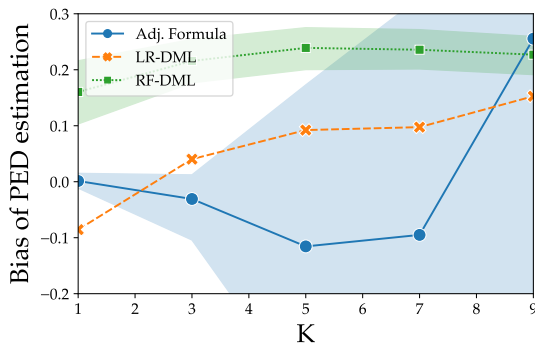
15

**Figure 4:** Bias defined as $\mathbb{E}[\Psi(Y^K)] - \Psi(y^K)$.

**Effect of shrinking adjustment set on estimator variance.** We bootstrap the adjustment formula estimator and the two double ML estimators. We find that the number of confounders heavily affects the bootstrapped variance of the PED estimators, suggesting that the Markovian modeling assumption (i.e., setting $K = 1$) used by our theory is also useful in practice. We report the predicted PED for all of the estimators in Figure 3 (left). For each estimator, we bootstrap the dataset 40 times to form confidence intervals. We report the standard deviation of the bootstrapped estimates in Figure 3 (middle). We see that the variance of the adjustment formula estimator increases as the number of confounders increases. The RF-DML and LR-DML variance curves are fairly stable with respect to $K$, suggesting that our Markovian assumption does not affect the variance of those estimators by much.

**Effect of shrinking adjustment set on estimator bias.** Stronger assumptions enable identifiability, but they come at a price of potential modeling errors. We have motivated our Markovian assumption theoretically, and now we want to understand how well they reflect reality. We use the bootstrapping technique proposed by Petersen et al. (2012a) for testing the bias of our estimators, which we describe now. Let $\Psi$ be the estimator of the PED we are testing, and $\Psi_a$ be the adjustment formula estimator of the PED. Further, let $y^K$ denote the avocado dataset for sequences of length $K$, and let $Y^K$ denote a bootstrapped sample constructed from $y^K$. We plot an empirical estimate of $\mathbb{E}[\Psi(Y^K)] - \Psi_a(y^K)$ using 40 bootstrap samples with confidence intervals in Figure 3 (right). We see that the adjustment formula and LR-DML estimators have small bias for small values of $K$, and all estimators have larger bias for large values of $K$. Recall that the modification to the adjustment formula could cause the adjustment formula estimator to underestimate the PED for large $K$, potentially causing the bias to spike for larger $K$ for LR-DML and RF-DML. In order to account for this, we also plot the estimated bias defined as $\mathbb{E}[\Psi(Y^K)] - \Psi(y^K)$—where we replace $\Psi_a$ with $\Psi$—in Figure 4. We see that the bias still increases as $K$ gets larger, suggesting that more confounders does in fact increases the bias of the estimator.

**Importance of shrinking adjustment set for overlap.** We report what the adjustment formula estimator estimates for a discretized treatment $u$ in Table 2. A "Low" price in the treatment column corresponds to the logged price bucket $(-0.479, 0.131]$. A "High" price corresponds to $(0.131, 0.683]$. In this experiment, we do not use the modified adjustment formula estimator that silently fails when overlap is not satisfied, which we introduced earlier in this section. Instead, we will report an explicit "N/A" when overlap fails. The "Fraction of undefined terms" column corresponds to the number of $\gamma$ values where $\sum_t \mathbf{1}\{z_t \in \mathcal{Z}_\gamma\} > 0$ and $\sum_t \mathbf{1}\{z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)}\} = 0$ over the total number of values of $\gamma$ where $\sum_t \mathbf{1}\{z_t \in \mathcal{Z}_\gamma\} > 0$. If "Fraction of undefined terms" is non-zero, then $\hat{x}(u)$ is not well defined. $N/A$ denotes when this occurs. The entries of "Probability mass of undefined terms" column is equal to $\sum_\gamma \frac{\sum_t \mathbf{1}\{z_t \in \mathcal{Z}_\gamma\}}{n} \mathbf{1}\{\sum_t \mathbf{1}\{z_t \in \mathcal{Z}_\gamma, u_t \in \mathcal{U}_{\beta(u)}\} = 0\}$. We can see that as $K$ gets larger, the number of undefined estimates, the relative fraction of undefined values, and the mass of said values gets larger. This preliminary analysis already suggests that there are overlap issues as $K$ gets larger.

**Conclusion.** Our experiments suggest that our Markovian assumption (i.e., $K = 1$) mitigates overlap issues while still accurately modeling reality. We believe the increase (with $K$) in bias and variance of the estimators is caused by overlap issues; as $K$ gets larger, the dimension of the confounders gets larger, making overlap harder to satisfy.

16

| | Price (Intervention $u$) | Estimated Effect on demand | Fraction of undefined terms | Probability mass of undefined terms |
|---|---|---|---|---|
| K=1 | High | 14.95 | 0 / 2 | 0.0% |
| | Low | 15.11 | 0 / 2 | 0.0% |
| K=3 | High | 14.96 | 0 / 8 | 0.0% |
| | Low | 15.11 | 0 / 8 | 0.0% |
| K=5 | High | N/A | 5 / 31 | 4.6% |
| | Low | 15.10 | 0 / 31 | 0.0% |
| K=7 | High | N/A | 36 / 89 | 15.1% |
| | Low | N/A | 16 / 89 | 9.0% |
| K=9 | High | N/A | 73 / 145 | 25.9% |
| | Low | N/A | 40 / 145 | 19.7% |

**Table 2:** Adjustment formula estimated effects on avocado demand for price interventions.

## D. Markovian assumption in Non-Markovian models

Although the system we are studying may not be Markovian—even in spite of the qualitative arguments we made in Section 2—the Markovian assumption may still be preferred over the model-free overlap assumption normally used, if the misspecification errors are small. In this section, we empirically quantify the trade-off between model-free overlap assumptions and Markovian assumption in a dynamical system which progressively become less Markovian. In particular, we compare the two-stage regression estimator from Section 5 against the double ML estimator (Chernozhukov et al., 2017) which we discuss in Appendix C.2. Because the former estimator makes a Markovian assumption and the latter estimator makes a blanket overlap assumption, we will use whichever estimator does better as a proxy of which assumption is preferable.

We let $\xi_t \in \mathbb{R}^d$ are drawn iid from $N(0, I)$. In this section, we have 256 samples of $(x_0, u_0, x_1, u_1, x_2)$ drawn iid from a dynamical system defined as

$$
\begin{aligned}
x_0 &\sim \xi_0, & u_0 &= Ax_0 \\
x_1 &= Ax_0 + Bu_0 + \xi_1, & u_1 &= Cx_1 + Du_0 + Rx_0 \\
x_2 &= Ax_1 + Bu_1 + Px_0 + Qu_0 + \xi_2.
\end{aligned}
$$

$A, B, C, D, P, Q, R$ are all in $\mathbb{R}^{d,d}$ and constructed such that the non-zero eigenvalues of $A, B, C, D$ are all roughly around 1 and the non-zero eigenvalues of $P, Q, R$ are all roughly around $\lambda$. To generate $P$, we first generate a random matrix $W$ in $\mathbb{R}^{d,m}$ with $m = 100$ and with independent standard Gaussians as its entries, and then we set $P = \lambda WW^T/m$. We generate $Q, R$ in the same way, and we generate $A, B, C, D$ in the same way, except with $\lambda = 1$. Our results look similar for $m = 1000$. In this experiment we will vary $\lambda$ and $d$. The larger $\lambda$ is, the further away this system is from the assumption of a Markovian system we propose in our paper. The larger $d$ is, the more susceptible to overlap violations this system becomes. As we vary $\lambda$ and $d$, we compare the error in estimating the performative effect ($B$) between the two stage regression estimator (LR-2SR and RF-2SR) we propose in Section 5—which uses a Markovian assumption—and the double ML estimator (LR-DML and RF-DML) discussed in Appendix C.2, which instead treats $(x_0, u_0, x_0)$ all as confounders. LR-2SR stands for Linear Regression-two-stage regression and corresponds exactly to the estimator we develop in Section 5. RF-2SR stands for Random forest-two-stage regression and also is a version of the two-stage regression we develop in Section 5, except that it uses the random forest model from sci-kit-learn version 1.2.0 to do the first stage of regression.

We plot our results in Figure 5. For all $\lambda \in [0, 0.1, 1]$, the Markovian approach (2SR) outperforms the model-free (DML) approach for $d \geq n/32$. For $\lambda \in [10, 100, 1000]$, the Markovian approach (2SR) outperforms the model-free (DML) approach for $d \geq 3n/10$. In other words, unless there are a lot of samples ($n \geq 10d/3$) and the misspecified matrices $P, Q, R$ have eigenvalues at least 10 times larger than that of $A, B, C, D$, our Markovian framework outperforms the model-free approach.

## E. Benefit of additional consumption shocks

Instead of observing a longer rollout, more consumption shocks is also another way to weaken the conditions required for identifiability. To explore the benefit of additional consumption shocks, we use an illustrative simulation.
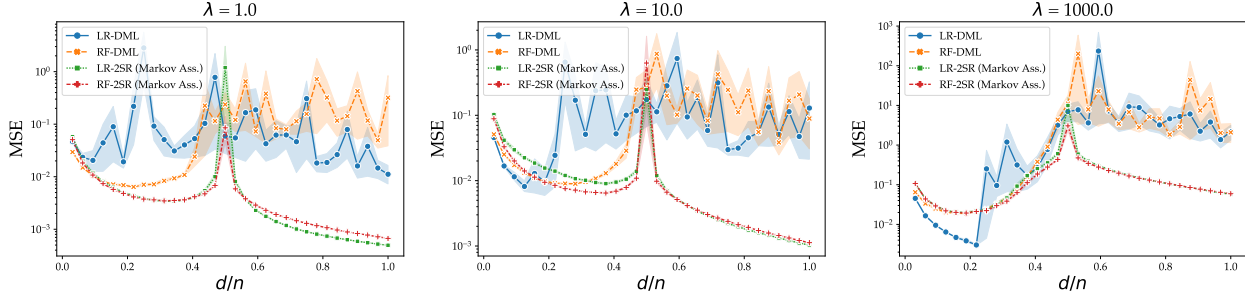
**Figure 5.** Mean-squared error plotted against the $d/n$ for various choices of eigenvalue scaling denoted by $\lambda$. LR-2SR and RF-2SR make Markovian assumptions on the dynamical system.
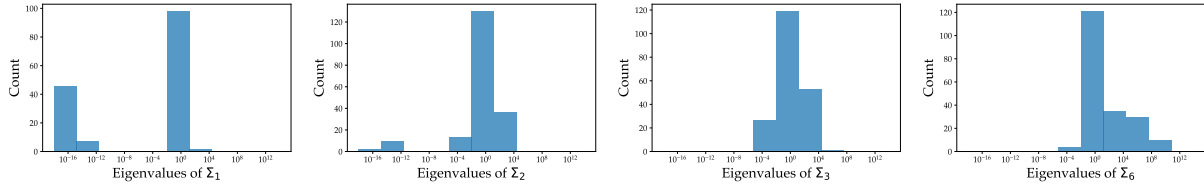


**Figure 6:** Histograms of eigenvalues of $\Sigma_t$ as defined in Section E.

We instantiate the linear dynamics (2) as follows. We consider the symmetric case where $d = p$. To generate $B$, we sample a random matrix $W$ in $\mathbb{R}^{d,n}$ for $n \gg d$ with independent standard Gaussians as its entries, and we set $B = WW^\top/n$. We repeat this process to generated $A$ and $D$. This way of generating our dynamics matrices ensures the matrices are well conditioned. We generate $C$ the same except by instead setting $W \in \mathbb{R}^{d,r}$ for $r < d$, making $C$ rank $r$ instead of rank $d$. We set $d = 100$, $n = 2000$, and $r = 80$. We note that in this system, $\operatorname{rank} DC = 80 = r < d$ and $\operatorname{rank} C = 80 = r < d$. We let $\xi_t \overset{d}{=} \mathsf{N}(0, I)$ for all $t$ and starts simulating the dynamics from $x_0 = u_0 = 0$.

This system, $B$ is identifiable from observations of triplets $R_{K=1}^{t+1} \coloneqq (x_t, u_t, x_{t+1})$ whenever the covariance matrix of $(x_t, u_t)$ has full rank. We can explicitly write down the covariance matrix of $(x_t, u_t)$ for our linear system as

$$\Sigma_t = J\Sigma_{t-1}J^\top + MM^\top = \sum_{k=0}^{t-1} (J^k)MM^\top(J^k)^\top$$

$$J \coloneqq \begin{bmatrix} A & B \\ CA & CB + D \end{bmatrix} \qquad M \coloneqq \begin{bmatrix} I \\ C \end{bmatrix}.$$

(4)

In the following we simulate $\Sigma_t$ for different $t$ and inspect the eigenvalues to determine whether identifiability from triplets $R_{K=1}^{t+1}$ is achieved. The larger $t$ the more consumption shocks preceed the observed triplet. We illustrate the histogram of the eigenvalues of $\Sigma_t$ for different $t$ in Figure 6. We see that the eigenvalues of get larger as more time passes: e.g., the eigenvalue mass of $\Sigma_6$ is further to the right of the eigenvalue mass of $\Sigma_3$ in Figure 6. This suggests that more noise spikes over more time steps make the observations better conditioned. As we will see in the next section, this makes estimating the performative effect provably easier for the auditor in practice.

Theorem 2 says that observing $R_{K=1}^2$—with only one consumption shock—is not sufficient for identifiability, as $\operatorname{rank} DC$ is not full row rank. This is consistent with the eigenvalue histogram of $\Sigma_2$ in Figure 6 as there are still 0 eigenvalues. However, Theorem 2 is not prescriptive for $R_{K=1}^t$ for settings where there are several consumption shocks and $t > 2$. However, when we inspect the eigenvalues for two consumption shocks, corresponding to observations of $R_{K=1}^3$. By Theorem 1 two consumption shocks are sufficient for identifiability in our system. This is also consistent with the eigenvalue histogram of $\Sigma_3$ in Figure 6, as all eigenvalues are bounded away from 0 at that time step.

An interesting question for future work is to formally unify the setting where we observe a longer rollout with one consumption shock and this setting where we observe a short rollout $R_{K=1}$ with several consumption shocks.

# F. Proofs

## F.1. Auxiliary results

**Lemma F.1** (Multivariate change of variables). *Let $X$ be a random variable with density $p_X$ and let $Y = g(X)$ where $g$ is an invertible mapping with Jacobian $J_g$, then $p_Y(a) = p_X(g^{-1}(a))|J_g(a)|^{-1}$.*

**Proof**

$$P(Y \in A) = P(X \in g^{-1}(A)) = \int_{g^{-1}(A)} p_X(x)dx = \int_A p_X(g^{-1}(x))|J_{g^{-1}}(x)|dx$$
$$= \int_A p_X(g^{-1}(x))|J_g(x)|^{-1}dx.$$

The definition of density gives the result. $\square$

**Definition F.1** (Lusin's (N) condition). *A function $f : \mathbb{R}^d \to \mathbb{R}^p$ satisfies Lusin's (N) condition if for every Lebesgue-measure 0 set $A \subset \mathbb{R}^d$, $f(A)$ has Lebesgue-measure 0.*

**Definition F.2** (Non-singular measurable transformation). *A function $f : \mathbb{R}^d \to \mathbb{R}^p$ is a non-singular measurable transformation if for every Lebesgue-measure 0 set $A \subset \mathbb{R}^p$, the preimage of $A$, $f^{-1}(A)$ has Lebesgue-measure 0.*

**Lemma F.2.** *For a measurable function $h : \mathbb{R}^d \to \mathbb{R}^p$, let $h^{-1}$ denote the preimage. Let $h$ be a non-singular measurable transformation which satisfies Lusin's (N) condition. Let $X$ be a $\mathbb{R}^d$-valued random variable with measure $P_X$ and density $p_X$, and let $Y := h(X)$ be a $\mathbb{R}^p$-valued random variable. Then the following is true:*

1. *$P_Y$ has a density $p_Y$ with respect to the Lebesgue measure.*

2. *if $p_X(a) > 0$ for almost all $a \in \mathbb{R}^d$ with respect to the Lebesgue measure, then $p_Y(b) > 0$ for almost all $b \in \mathbb{R}^p$ with respect to the Lebesgue measure.*

**Proof**   Recall that a $\sigma$-finite measure $\nu$ has a density with respect to $\sigma$-finite measure $\mu$ if and only if $\nu$ is absolutely continuous with respect to $\mu$ (denoted as $\nu \ll \mu$).

We prove the first point first. We will show that the measure of $Y$, $P_Y$, is absolutely continuous with respect to the Lebesgue measure $\lambda$. Let $A \subset \mathbb{R}^p$ be such that $\lambda(A) = 0$, then

$$\lambda(A) = 0 \implies \lambda(h^{-1}(A)) = 0 \implies P_X(h^{-1}(A)) = 0 \implies P_Y(A) = 0.$$

The first implication is because $h$ is a non-singular measurable transformation. The second implication is because $P_x \ll \lambda$ as $X$ has a density with respect to $\lambda$.

To prove the second point, we first show that $p_X(a) > 0$ for all $a \in \mathbb{R}^d$ implies $P_X \gg \lambda$. To see this, observe that for any $A$, $\lambda(A) = \int_A \frac{1}{p_X(y)} p_X(y)dy\lambda = \int_A \frac{1}{p_X(y)} P_X(dy)$. With this we show that $P_Y \gg \lambda$. Let $A \subset \mathbb{R}^p$ be such that $P_Y(A) = 0$, then

$$P_Y(A) = 0 \implies P_X(h^{-1}(B)) = 0 \implies \lambda(h^{-1}(B)) = 0 \implies \lambda(B) = 0.$$

The second implication is because $P_X \gg \lambda$ and the third implication is because $h$ satisfies Lucin's condition. We prove that $p_Y > 0$ almost everywhere by contradiction. Because $P_Y$ and $\lambda$ are mutually absolutely continuous, there exists $q$ such that $\lambda = qP_Y$. Then because $P_Y = p_Y\lambda$, $\lambda = qp_Y\lambda$. Thus, $qp_Y$ must equal 1 almost everywhere with resepct to the Lebesgue measure, $p_Y$ must be non-zero almost everywhere. $\square$

## F.2. Proof of Proposition 1

Without loss of generality we consider $T = 2$. Recall that the do action alters the data generation model by deleting incoming edges into $u_1$.

$$\mathbb{E}[x_2|do(u_1 := u)] = \mathbb{E}[f(x_1) + g(u) + \xi_2]$$
$$= \int \mathbb{E}[f(z) + g(u) + \xi_2 \mid x_1 = z]p_{x_1}(z)dz$$
$$= \int \mathbb{E}[f(x_1) + g(u_1) + \xi_2 \mid u_1 = u, x_1 = z]p_{x_1}(z)dz$$
$$= \int \mathbb{E}[x_2 \mid u_1 = u, x_1 = z]p_{x_1}(z)dz.$$

The second and third equalities use the fact that $\xi_2$ is independent of $x_1, u_1$.

## F.3. Proof of Theorem 1

Without loss of generality we will set $T = 2$ in this proof.

### F.3.1. PART 1: IDENTIFIABILITY

**Showing overlap**   We will first show that $(x_1, u_1)$ has full support, which automatically implies overlap. Let $z_{-1} :=(u_{-1}, x_{-1})$. Because

$$p_{u_1,x_1}(u,x) = \int p_{u_1,x_1|z_{-1}=z}(u,x)p_{z_{-1}}(z)dz,$$

it suffices to show that $p_{u_1,x_1|z_{-1}=z}$ has full support for any $z \in \mathbb{R}^{p+d}$. For this reason, in this proof, we fix $z_{-1}$—i.e., $u_{-1}$ and $x_{-1}$ will be treated like constants—and for notional simplicity, we omit explicitly conditioning on the event $z_{-1} = z$. Let $c := r(u_{-1})$, $d := g(u_{-1}) + f(x_{-1})$, and $\xi_0' := \xi_0 + d$. Observe that $(\xi_0', \xi_1)$ still has full support. Using this modified notation, we have

$$x_0 = \xi_0'$$
$$u_0 = h(\xi_0') + c$$
$$x_1 = f(\xi_0') + g(h(\xi_0') + c) + \xi_1$$
$$u_1 = h(x_1) + r(h(\xi_0') + c)$$

We first show that $x_1$ has full support. Recall $\xi_1$ has positive density over $\mathbb{R}^d$. Because addition by a constant is an invertible, differentiable function, Lemma F.1 implies that $f(\xi_0') + g(h(\xi_0') + c) + \xi_1|\xi_0'$ has positive density over $\mathbb{R}^d$. Since $\xi_0'$ also has positive density over $\mathbb{R}^d$, integration tells us that $x_1 = f(\xi_0') + g(h(\xi_0') + c) + \xi_1$ has positive density over $\mathbb{R}^d$.

Because $p_{x_1,u_1} = p_{u_1|x_1}p_{x_1}$ and $x_1$ has full support, it suffices to show that $u_1|x_1$ has full support over $\mathbb{R}^p \times \mathbb{R}^d$. Let $q_c(y) := r(h(y) + c)$. It is sufficient to show that $p_{q_c(\xi_0')|x_1}$ is positive everywhere. To see this, observe that $u_1|x_1 = h(x_1) + q_c(\xi_0')|x_1$. Because addition by a constant is an invertible, differentiable function, if $q_c(\xi_0')|x_1$ had positive density everywhere, then Lemma F.1 tells us that $u_1|x_1$ would have positive density everywhere. One can show that the class of continuously differentiable, surjective functions with either full row-rank or full column rank Jacobian satisfy Definitions F.1 and F.2 (Santos; Lowther). Thus, Definition 3.4 holds, the conditions of Lemma F.2 hold, and thus, it suffices to show $\xi_0'|x_1$ has positive density everywhere. We observe that

$$p_{\xi_0'|x_1}(a,b) = \frac{p_{x_1|\xi_0'}(b,a)p_{\xi_0'}(a)}{p_{x_1}(b)}.$$

Since $x_1$ has full support, the denominator is positive. Since $\xi_0'$ has full support, $p_{\xi_0'}(a) > 0$ as well. Finally, we had already shown earlier in the proof that $x_1|\xi_0'$ (i.e., $f(\xi_0') + g(h(\xi_0') + c) + \xi_1|\xi_0'$) has positive density everywhere as well.

**Concluding argument**   Because $p_{u_1,x_1}$ is positive everywhere, $\mathbb{E}[x_2 \mid u_1 = u, x_1 = z]$ is well defined. Additionally, because $x_1$ has density, $\int_z \mathbb{E}[x_2 \mid u_1 = u, x_1 = z]p_{x_1}(z)dz$ is well defined as well. Finally because our model is admissible as stated in Proposition 1, $\mathbb{E}[x_2 \mid do(u_1 := u)] = \int_z \mathbb{E}[x_2 \mid u_1 = u, x_1 = z]p_{x_1}(z)dz$. The right hand side of this relationship is well defined and can be computed from knowledge of the distribution of $(x_1, u_1, x_2)$; thus, $\mathbb{E}[x_2 \mid do(u_1 := u)]$ can be computed from the distribution of observations $(x_1, u_1, x_2)$. Because this quantity identifiable, the performative effect $\mathrm{PE}(u, u')$ is also identifiable for any $u, u' \in \mathbb{R}^d$.

20

F.3.2. PART 2: UNIDENTIFIABILITY

**Case 1** $\xi_0 \stackrel{a.s.}{=} 0$**:** We consider the case where we observe $x_2, x_1, u_1$ and $\xi_0 \stackrel{a.s.}{=} 0$. Let $P_0$ be the point mass over the 0 vector; i.e., $x_0 = u_0 = 0$. Define a measurable function $\Delta : \mathbb{R}^p \to \mathbb{R}^d$ such that $\Delta \neq 0$. Let $r$ be the identity function. For any functions $f, g, h$, define, $\hat{f}(a) := f(a) + \Delta(h(a))$, $\hat{g}(b) := g(b) - \Delta(b)$, and $\hat{h}(c) = h(c)$. For noise variables $(\xi_1, \xi_2)$, let $(\hat{\xi}_1, \hat{\xi}_2)$ an identically distributed copy. Let $R_2 = (x_1, u_1, x_2)$ be sampled according to the dynamics specified by (1) using the functions $f, g, h$, noise variables $(\xi_1, \xi_2)$, and with initial conditions $x_0 = u_0 = 0$. Let $\hat{R}_2 = (\hat{x}_1, \hat{u}_1, \hat{x}_2)$ be sampled according to the dynamics specified by (1) using the functions $\hat{f}, \hat{g}, \hat{h}$ in place of $f, g, h$, noise variables $(\hat{\xi}_1, \hat{\xi}_2)$ in place of $(\xi_1, \xi_2)$, and with initial conditions $\hat{x}_0 = \hat{u}_0 = 0$. We see that

$$x_1 \stackrel{d}{=} \xi_1 \stackrel{d}{=} \hat{x}_1$$
$$u_1 \stackrel{d}{=} h(\xi_1) \stackrel{d}{=} \hat{h}(\xi_1) \stackrel{d}{=} \hat{u}_1$$
$$x_2 \stackrel{d}{=} f(\xi_1) + g(h(\xi_1)) + \xi_2$$
$$\stackrel{d}{=} f(\xi_1) + \Delta(h(\xi_1)) + g(h(\xi_1)) - \Delta(h(\xi_1)) + \xi_2$$
$$\stackrel{d}{=} \hat{f}(\xi_1) + \hat{g}(\hat{h}(\xi_1)) + \xi_2 \stackrel{d}{=} \hat{x}_2.$$

**Case 2** $\xi_1 \stackrel{a.s.}{=} 0$**:** We consider the case where we observe $x_2, x_1, u_1$ and $\xi_1 \stackrel{a.s.}{=} 0$. Let $P_{-1}$ be the point mass over the 0 vector; i.e., $x_{-1} = u_{-1} = 0$. Let $r$ be the identity function. For any functions $f, g, h$, define, $\hat{f}(a) := f(a) + \Delta(h(a))$, $\hat{g}(b) := g(b) - \Delta(b)$, and $\hat{h}(c) = h(c)$. For noise variables $(\xi_0, \xi_2)$, let $(\hat{\xi}_0, \hat{\xi}_2)$ an identically distributed copy.

Let $R_2 = (x_1, u_1, x_2)$ be sampled according to the dynamics specified by (1) using the functions $f, g, h, r$, noise variables $(\xi_0, \xi_2)$, and with initial conditions $x_{-1} = u_{-1} = 0$. Let $\hat{R}_2 = (\hat{x}_1, \hat{u}_1, \hat{x}_2)$ be sampled according to the dynamics specified by (1) using the functions $\hat{f}, \hat{g}, \hat{h}, r$ in place of $f, g, h, r$, noise variables $(\hat{\xi}_0, \hat{\xi}_2)$ in place of $(\xi_0, \xi_2)$, and with initial conditions $\hat{x}_{-1} = \hat{u}_{-1} = 0$.

Following the steps from the first case, we know that $(x_0, u_0, x_1) \stackrel{d}{=} \hat{x}_0, \hat{u}_0, \hat{x}_1$. We have the following equalities

$$u_1 \stackrel{d}{=} h(x_1) + r(u_0) \stackrel{d}{=} \hat{h}(\hat{x}_1) + r(\hat{u}_0) \stackrel{d}{=} \hat{u}_1$$
$$x_2 \stackrel{d}{=} f(x_1) + g(h(x_1)) + \xi_2$$
$$\stackrel{d}{=} f(x_1) + \Delta(h(x_1)) + g(h(x_1)) - \Delta(h(x_1)) + \xi_2$$
$$\stackrel{d}{=} \hat{f}(x_1) + \hat{g}(\hat{h}(x_1)) + \xi_2 \stackrel{d}{=} \hat{x}_2.$$

## F.4. Proof of Theorem 2

F.4.1. SUPPORTING LEMMAS

We first outline a series of helpful supporting lemmas. This first lemma draws an equivalence between matrices and the probability distributions induced by these matrices, allowing us to reason about one by reasoning about the other.

**Lemma F.3.** *Let $\{\xi_i\}_{i=1}^n$ be a set of mutually independent random vectors in $\mathbb{R}^d$ with full span. Let $\{A_i\}_{i=1}^n$ be a set of deterministic matrices in $\mathbb{R}^{d,p}$. Let $v \in \mathbb{R}^d$ be a random vector in $\mathbb{R}^d$ mutually independent of $\{\xi_i\}_{i=1}^n$. $A_i = 0$ for all $i \in [n]$ and $v \stackrel{a.s.}{=} 0$ if and only if $v + \sum_{i=1}^n A_i \xi_i \stackrel{a.s.}{=} 0$.*

**Proof** The left to right direction is obvious. We now prove the right to left direction by cases. Suppose $v$ is almost surely a constant vector. Suppose that only one $j \in [n]$ such that $A_j \neq 0$, then its not possible that $A_j \xi_j \stackrel{a.s.}{=} -v$ by definition of full span. Suppose there exists $j, k \in [n]$ such that $A_j \neq 0$ and $A_k \neq 0$. This means that $A_j \xi_j$ is almost surely not a constant. We also know that conditioned on $\{A_i \xi_i\}_{i \neq j}$, $A_j \xi_j$ is almost surely a constant. This implies that $P_{A_j \xi_j} \neq P_{A_j \xi_j | \{A_i \xi_i\}_{i \neq j}}$ which contradicts the assumption of mutual independence. Suppose $v$ is almost surely not a constant vector. Then $P_v \neq P_{v | \{A_i \xi_i\}_{i \in [n]}}$ as $v$ is almost surely a constant vector conditioned on $\{A_i \xi_i\}_{i \in [n]}$. This contradicts mutual independence. $\square$

For our next lemma and for the rest of the proof, we need to define some notation. Consider the following variables:

$$\Theta_x := \begin{bmatrix} A^\top \\ B^\top \end{bmatrix} \qquad \Theta_u := \begin{bmatrix} C^\top \\ D^\top \end{bmatrix} \qquad \xi_x := \begin{bmatrix} \xi_1^\top \\ \vdots \\ 0 \end{bmatrix} \qquad \xi_u := 0.$$

Let $\hat{\Theta}_x, \hat{\Theta}_u$ be defined with respect to $\hat{A}, \hat{B}, \hat{C}, \hat{D}$. Let $\hat{\xi}_x \overset{d}{=} \xi_x$ and $\hat{\xi}_u \overset{d}{=} \xi_u$. Let $(A, B, C, D)$ and $\xi_x, \xi_u$ induce $P_T$ and let $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$ and $\hat{\xi}_x, \hat{\xi}_u$ induce $\hat{P}_T$. Let $x := [x_1, \ldots, x_T]^\top$ and $u := [u_1, \ldots, u_{T-1}]^\top$ be observations from $P_T$ and let $\hat{x}$ and $\hat{u}$ defined with hat variables be observations from $\hat{P}_T$. Finally let $z := (x, u)$ and $\hat{z} := (\hat{x}, \hat{u})$. Finally, we define matrices $Q_x, Q_u, \hat{Q}_x$, and $\hat{Q}_u$ such that the following relationships hold

$$x - \xi_x \overset{d}{=} Q_x \Theta_x \qquad u - \xi_u \overset{d}{=} Q_u \Theta_u$$
$$\hat{x} - \hat{\xi}_x \overset{d}{=} \hat{Q}_x \hat{\Theta}_x \qquad \hat{u} - \hat{\xi}_u \overset{d}{=} \hat{Q}_u \hat{\Theta}_u.$$

Our next lemma translates relationships about one set of dynamics matrices into relationships about the other set of dynamics relationships.

**Lemma F.4.** *If* $x - \xi_x \overset{d}{=} Q_x \hat{\Theta}_x$, *then* $x - \xi_x \overset{d}{=} Q_x \hat{\Theta}_x \overset{d}{=} \hat{Q}_x \hat{\Theta}_x \overset{d}{=} \hat{x} - \hat{\xi}_x$. *Similarly, if* $u - \xi_u \overset{d}{=} Q_u \hat{\Theta}_u$, *then* $u - \xi_u \overset{d}{=} Q_u \hat{\Theta}_u \overset{d}{=} \hat{Q}_u \hat{\Theta}_u \overset{d}{=} \hat{u} - \hat{\xi}_u$.

**Proof**    Recall that the random variables in the vector $z$ corresponds to nodes in the causal directed acyclic graph shown in Figure 1(b). Define $\sigma : \mathbb{Z} \to \mathbb{Z}$ such that $z_{\sigma(i)}$ is in sorted DAG order with respect to the DAG in Figure 1(b) (i.e, the parents of $z_{\sigma(i)}$ have $\sigma$ indices smaller than $\sigma(i)$ and its children have $\sigma$ indices larger than $\sigma(i)$). We proceed inductively to show that $z_{\sigma(i)} \overset{d}{=} \hat{z}_{\sigma(i)}$.

*Base case:* $z_{\sigma(1)} \overset{d}{=} L(\xi, \hat{\Theta}_x, \hat{\Theta}_u)$, where $L$ is some function, linear in each of its inputs. Since $\xi \overset{d}{=} \hat{\xi}$, we have that $z_{\sigma(1)} \overset{d}{=} L(\hat{\xi}, \hat{\Theta}_x, \hat{\Theta}_u) = \hat{z}_{\sigma(1)}$; the last equality follows from definition.

*Inductive step:* suppose $z_{\sigma(j)} \overset{d}{=} \hat{z}_{\sigma(j)}$ jointly over all $j$. We know that $z_{\sigma(j+1)} \overset{d}{=} L(\{z_{\sigma(i)}\}_{i<j}, \xi, \hat{\Theta}_x, \hat{\Theta}_u)$ where $L$ is linear in $\{z_{\sigma(i)}\}_{i<j}$, linear in $\xi$, linear with respect to $\hat{\Theta}_x$, and linear in $\hat{\Theta}_u$. By the inductive hypothesis we know that $z_{\sigma(j+1)} \overset{d}{=} L(\{z_{\sigma(i)}\}_{i<j}, \xi, \hat{\Theta}_x, \hat{\Theta}_u)$ which in turn is equal in distribution to $L(\{\hat{z}_{\sigma(i)}\}_{i<j}, \hat{\xi}, \hat{\Theta}_x, \hat{\Theta}_u) \overset{d}{=} \hat{z}_{\sigma(j+1)}$, as all the inputs to the function are equal in distribution.

Because the entries of $Q$ ($\hat{Q}$ respectively) are comprised of entries of $z$ ($\hat{z}$ respectively), we have that $\hat{Q} \overset{d}{=} Q$. This proves the desired result. $\qquad \square$

### F.4.2. PART 1: UNIDENTIFIABILITY WHEN $K = 1$

Without loss of generality, let $T = 2$. The proof of this result proceeds exactly as the proof of the unidentifiability result in Theorem 1 in Appendix F.3.2 except with $f, g, h, r$ defined as in Equation (2) and with $\Delta : \mathbb{R}^p \to \mathbb{R}^d$ set to any linear function $\Delta(x) = Wx$ where $W \in \mathbb{R}^{d,p}$ is such that $W \neq 0$.

### F.4.3. PARTS 2 AND 3: IDENTIFIABILITY WHEN $K \geq 2$

Now that we have established our supporting lemmas, we can now prove our desired result. Without loss of generality, we will set $T = K + 1$.

**Necessity and sufficiency when $x_0 = u_0 = \xi_t = 0$ for $t \geq 2$.**    Let $X_t \in \mathbb{R}^{d,d}$ and $U_t \in \mathbb{R}^{p,d}$ be defined such that $x_t = X_t \xi_1$ and $u_t = U_t \xi_1$. Further define the following random matrix:

$$Q_x := \begin{bmatrix} x_0^\top & u_0^\top \\ x_1^\top & u_1^\top \\ \vdots & \vdots \\ x_{T-1}^\top & u_{T-1}^\top \end{bmatrix}.$$

Define hat versions of all variables accordingly. We have that $x - \xi_x \overset{d}{=} \hat{x} - \hat{\xi}_x$ and $u - \xi_u \overset{d}{=} \hat{u} - \hat{\xi}_u$. Moreover, $Q_x$ is comprised of entries of $x$ and $u$, $Q_x \overset{d}{=} \hat{Q}_x$ (jointly). Thus,

$$x - \xi_x \overset{d}{=} \hat{x} - \xi_x \overset{d}{=} \hat{Q}_x \hat{\Theta}_x \overset{d}{=} Q_x \hat{\Theta}_x$$
$$u - \xi_u \overset{d}{=} \hat{u} - \xi_u \overset{d}{=} \hat{Q}_u \hat{\Theta}_u \overset{d}{=} Q_u \hat{\Theta}_u. \tag{5}$$

Finally, defining the fixed matrices $X := [X_2, \ldots, X_T]^\top$, $U := [U_2, \ldots, U_{T-1}]^\top$, and

$$Q_X := \begin{bmatrix} X_1^\top & U_1^\top \\ \vdots & \vdots \\ X_{T-1}^\top & U_{T-1}^\top \end{bmatrix},$$

we can rewrite (5) as

$$\begin{bmatrix} \xi_1^\top \\ \vdots \\ \xi_1^\top \end{bmatrix} \odot X = \begin{bmatrix} \xi_1^\top \\ \vdots \\ \xi_1^\top \end{bmatrix} \odot Q_X \hat{\Theta}_x. \tag{6}$$

Note, that in this reparameterization, we omit the $x_1 - \xi_1 = Ax_0 + Bu_0$, as these terms are equal to 0. Using Lemma F.3 we know the above equality holds if and only if the following holds

$$X = Q_X \hat{\Theta}_x. \tag{7}$$

Lemma F.4 tells us $B$ is identifiable if and only if the entries of $\Theta_x$ corresponding to $B$ is unique (7). Indeed, if there exists two solutions $(\hat{\Theta}_x, \hat{\Theta}_u) \neq (\Theta_x, \Theta_u)$ such that $B \neq \hat{B}$, we can use Lemma F.4 to show that $\hat{P}_T = P_T$; i.e., the system is not identifiable. The other direction is trivial, as $B$ being identifiable implies that $B$ is unique.

We now give equivalent conditions for when $B$ is unique. Let $\mathcal{S} := \{e_j\}_{j=d+1}^{d+p}$ where $e_j$ is the $j$th standard basis vector in $\mathbb{R}^{d+p}$. $B$ is unique (i.e., $\hat{B} = B$) if an only if $\text{null}(Q_X) \perp \text{span}(\mathcal{S})$. Indeed suppose $v \in (Q_X)$ is such that $v$ is not orthogonal to $\text{span}(\mathcal{S})$, then $\hat{\Theta}_x = \Theta_x + v\mathbf{1}^\top$ is also a solution to (7); moreover, $\hat{B} \neq B$ because $v$ is not orthogonal to $\text{span}(\mathcal{S})$. Conversely suppose for all $v \in (Q_X)$, $v$ is orthogonal to $\text{span}(\mathcal{S})$. Then, any alternative solution $\hat{\Theta}_x \neq \Theta_x$ must satisfy $\mathcal{C}(\hat{\Theta}_x - \Theta_x) \perp \text{span}(\mathcal{S})$, where $\mathcal{C}$ denotes the column span, which implies that $\hat{B} = B$.

Note that if $M$ is a full rank matrix, $MQ$ has the same null space as $Q$. Further observe that by using elementary row operations, we know that there exists full rank square matrices $M_1$ and $M_2$ such that

$$Q_X = \begin{bmatrix} I & C^\top \\ X_2^\top & (CX_2 + DU_1)^\top \\ \vdots & \vdots \\ X_{T-1}^\top & (CX_{T-1} + DU_{T-2})^\top \end{bmatrix} = M_1 \begin{bmatrix} I & C^\top \\ 0 & (DU_1)^\top \\ \vdots & \vdots \\ 0 & (DU_{T-2})^\top \end{bmatrix} = M_2 \begin{bmatrix} I & C^\top \\ 0 & (DC)^\top \\ \vdots & \vdots \\ 0 & (D^{T-2}C)^\top \end{bmatrix}.$$

$M_1$ and $M_2$ are products of full rank matrices corresponding to elementary row operations. $M_2$ is constructed by repeatedly applying the fact $U_t = CX_t + DU_{t-1}$. Thus, $B$ is unique if and only if $\text{null}(\tilde{Q}_X) \perp \text{span}(\mathcal{S})$ where

$$\tilde{Q}_X := \begin{bmatrix} I & C^\top \\ 0 & (DC)^\top \\ \vdots & \vdots \\ 0 & (D^{T-2}C)^\top \end{bmatrix}.$$

This is equivalent to $\text{span}(\mathcal{S}) \subset \mathcal{R}(\tilde{Q}_X)$, where $\mathcal{R}$ denotes row span, which is then equivalent to $[DC, \ldots, D^{T-2}C]$ being full row rank (recall $T = K + 1$). Tracing back all the if and only if statements gives the result.

**Sufficiency even when $x_0 \neq 0$ and $u_0 \neq 0$.** In this setting, the proof for Claim 1 holds up to Equation (6). Equation (6) changes to the following

$$\begin{bmatrix} \xi_1^\top \\ \vdots \\ \xi_1^\top \end{bmatrix} \odot X + w_1(x_0, u_0, \xi_{>1}) = \begin{bmatrix} \xi_1^\top \\ \vdots \\ \xi_1^\top \end{bmatrix} \odot Q_X \hat{\Theta}_x + w_2(x_0, u_0, \xi_{>1}).$$

By Lemma F.3, we know that these equalities hold if and only if Equation (7) holds, $w_1(x_0, u_0, \xi_{>1}) = w_2(x_0, u_0, \xi_{>1})$ holds. $\mathrm{null}(Q_X) \perp \mathrm{span}(\mathcal{S})$ *suffices* (but is no longer necessary as there is one other relationships we are not accounting for) in showing there is a unique B in any solution of the linear system in Equation (7). The rest of the argument in Claim 1 follows identically.

## F.5. Proof of Theorem 3

We first introduce a helpful supporting lemma.

**Lemma F.5.** *Suppose $n$ samples are drawn iid from $P_2$. If $X_1 X_1^\top$ is invertable, then $\widehat{C} = C$ and $\widehat{H} = A + BC + E_2 X_1^\top (X_1 X_1^\top)^{-1}$. If $X_1 X_1^\top$ is invertable and $DCX_1 X_1^\top C^\top D^\top$ is invertable, then $\widehat{B} = B - E_2 X_1^\top (X_1 X_1^\top)^{-1} X_2 (DCX_1)^\top (DCX_1 X_1^\top C^\top D^\top)^{-1}$.*

**Proof** Substituting $U_1 = CX_1$ and $X_2 = (A + BC)X_1 + E_2$ into the closed form solutions of $\widehat{C}$ and $\widehat{H}$ respectively gives the first result.

To get the second result, we use the fact that $\widehat{C} = C$ and $\widehat{H} = A + BC + E_2 X_1^\top (X_1 X_1^\top)^{-1}$ by the first result. We observe that $U_2 = CX_2 + DU_1 = CX_2 + DCX_1$ to get that $\widehat{B} = (X_3 - \widehat{H}X_2)(DCX_1)^\top (DCX_1 X_1^\top C^\top D^\top)^{-1}$. Then we use the fact that subtracting $BCX_2$ from both sides of the relationship $X_3 - AX_2 = BU_2$ gives us that $X_3 - (A + BC)X_2 = B(U_2 - CX_2)$. Using our invertability assumptions, this gives us $\widehat{B} = B - E_2 X_1^\top (X_1 X_1^\top)^{-1} X_2 (DCX_1)^\top (DCX_1 X_1^\top C^\top D^\top)^{-1}$. $\qquad\square$

With this, we can analyze the quantities of interest. Let $\hat{\Sigma}_1 = \frac{1}{n} X_1 X_1^\top$. Let $Q \coloneqq DC\hat{\Sigma}_1 C^\top D^\top$.

$$
\begin{aligned}
\mathbb{E}\left[\left\|\widehat{B} - B\right\|_{\mathrm{Fr}}^2 \mid \mathcal{G}\right] &= \frac{1}{n^2} \mathrm{tr}(\mathbb{E}[Q^{-1} DCX_1 X_2^\top (X_1 X_1^\top)^{-1} X_1 E_2^\top E_2 X_1^\top (X_1 X_1^\top)^{-1} X_2 X_1^\top C^\top D^\top Q^{-1}]) \\
&= \frac{\sigma_2^2 d}{n} \mathrm{tr}(\mathbb{E}[Q^{-1} DC\hat{\Sigma}_1 (A + BC)^\top \hat{\Sigma}_1^{-1} (A + BC)\hat{\Sigma}_1 C^\top D^\top Q^{-1}]) \\
&\leq \frac{\sigma_2^2 pd}{n} \kappa_{DC}^2 \left(\frac{\|A + BC\|_{\mathrm{op}}}{\sigma_{\min}(DC)}\right)^2 \mathbb{E}\left[\frac{\kappa_{\hat{\Sigma}_1}^2}{\lambda_{\min}(\hat{\Sigma}_1)}\right].
\end{aligned}
$$

Rearranging and using the definition of $\tau_1$ gives the result.

If $p = d$, then $DC$ is a square, invertible matrix,

$$
\begin{aligned}
\mathbb{E}\left[\left\|\widehat{B} - B\right\|_{\mathrm{Fr}}^2 \mid \mathcal{G}\right] &= \frac{\sigma_2^2 d}{n} \mathrm{tr}[(C^\top D^\top)^{-1}(A + BC)^\top \mathbb{E}\left[\hat{\Sigma}_1^{-1}\right](A + BC)(DC)^{-1}] \\
&\leq \frac{\sigma_2^2 d^2}{n} \left(\frac{\|A + BC\|_{\mathrm{op}}}{\lambda_{\min}(DC)}\right)^2 \left\|\mathbb{E}\left[\hat{\Sigma}_1^{-1}\right]\right\|_{\mathrm{op}}.
\end{aligned}
$$

Rearranging and using the definition of $\tau_2$ gives the result.

## F.6. Proof of Theorem 4

We let $Y(\mathcal{U}, \mathcal{X}) \coloneqq \mathbb{E}[x_2 \mid u_1 \in \mathcal{U}, x_1 \in \mathcal{X}]$, $Z(\mathcal{X}) \coloneqq Z(\mathcal{X}_\alpha)$, $\hat{Y}(\mathcal{U}, \mathcal{X}) \coloneqq \widehat{\mathbb{E}}[x_2 \mid u_1 \in \mathcal{U}, x_1 \in \mathcal{X}]$, and $\hat{Z}(\mathcal{X}) \coloneqq \widehat{Z}(\mathcal{X}_\alpha)$. The proof proceeds by bounding each of the following terms:

$$
\begin{aligned}
\left\|\sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\hat{Z}(\mathcal{X}_\alpha) - \mathbb{E}[x_2 | do(u_1 \coloneqq u)]\right\| &\leq \left\|\sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\hat{Z}(\mathcal{X}_\alpha) - \sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)Z(\mathcal{X}_\alpha)\right\| \\
&+ \left\|\sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)Z(\mathcal{X}_\alpha) - \sum_\alpha Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)Z(\mathcal{X}_\alpha)\right\| \\
&+ \left\|\sum_\alpha Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)Z(\mathcal{X}_\alpha) - \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x]Z(\mathcal{X}_\alpha)\right\| \\
&+ \left\|\sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x]Z(\mathcal{X}_\alpha) - \mathbb{E}[x_2 | do(u_1 \coloneqq u)]\right\|.
\end{aligned}
$$

### F.6.1. SUPPORTING LEMMAS

We begin with a series of supporting lemmas that will aid us in bounding these terms.

**Lemma F.6.** *Let the conditions of Theorem 1 hold and let $\lambda$ denote the Lebesgue measure for $\mathbb{R}^{d+p}$. For all $A \in \mathcal{N}$ and $B \in \mathcal{M}$, the following implication is true:* $\lambda(A \times B) > 0 \implies (x_1 \in A, u_1 \in B) > 0$.

**Proof**

$$P(x_1 \in A, u_1 \in B) = \int_B \int_A p_{x_1, u_1}(x, u) dx du > 0$$

We know the RHS is positive because the function being integrated is positive by Theorem 1 and the set it's being integrated over has measure greater than 0. $\qquad\square$

**Lemma F.7.** *Let $f : \mathbb{R}^d \to \mathbb{R}^p$ be a L-Lipschitz function. If every element of $\mathcal{N}$ has diameter at most $\varepsilon$ with respect to $\|\cdot\|$, then for all $\mathcal{X} \in \mathcal{N}$, for all $x, y \in \mathcal{X}$, $\|f(x) - f(y)\| \leq L\varepsilon$.*

**Proof**   Follows directly from definitions of diameter and Lipschitz Continuity. $\qquad\square$

**Lemma F.8.** *Consider the data generation model of* (1)*. Let Assumption 3 hold. Let $x_1$ have full support. Then,*

$$\left\| \mathbb{E}[x_2 | do(u_1 := u)] - \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x] Z(\mathcal{X}_\alpha) \right\| \leq L\varepsilon + \left\| \int_{\mathbb{R}^d \setminus \cup \mathcal{N}} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|$$

**Proof**   Let $B := \mathbb{R}^d \setminus \cup \mathcal{N}$ denote the set of points not covered by $\cup \mathcal{N}$. Then, we have the following inequalities:

$$\left\| \mathbb{E}[x_2 | do(u_1 := u)] - \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x] Z(\mathcal{X}_\alpha) \right\|$$

$$\leq \left\| \sum_\alpha \int_{\mathcal{X}_\alpha} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz - \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = r] Z(\mathcal{X}_\alpha) \right\|$$

$$+ \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|$$

$$\leq \sum_\alpha \left\| \int_{\mathcal{X}_\alpha} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz - \mathbb{E}[x_2 | u_1 = u, x_1 = r] Z(\mathcal{X}_\alpha) \right\|$$

$$+ \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|$$

$$\leq \sum_\alpha \int_{\mathcal{X}_\alpha} \left\| \mathbb{E}[x_2 | u_1 = u, x_1 = z] - \mathbb{E}[x_2 | u_1 = u, x_1 = r] \right\| p_{x_1}(z) dz$$

$$+ \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|$$

$$\leq L\varepsilon + \left\| \int_{\mathbb{R}^d \setminus \cup \mathcal{N}} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|.$$

The first and second inequality is from triangle inequality. The third comes from Jensen's inequality. The fourth inequality comes Assumption 3 and Lemma F.7. $\qquad\square$

Lemma F.8 tells us that it suffices to create an estimator that estimates $\sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = r] Z(\mathcal{X}_\alpha)$—supposing that $\cup \mathcal{N}$ is a good approximation of $\mathbb{R}^d$ with respect to $x_1$.

**Lemma F.9.** *Consider the data generating process from* (1)*. Let $x_1, u_1$ have full support. Let Assumption 4 hold, then*

$$\left\| \mathbb{E}[Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_{\alpha(x)})] - \mathbb{E}[x_2 | u_1 = u, x_1 = x] \right\| \leq \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \mathbb{E}[\|x_2\| | u_1 = u, x_1 = x].$$

**Proof** Fix any $u \in \cup \mathcal{M}, x \in \cup \mathcal{N}$. Let $Z \coloneqq (x_1, u_1)$, $z \coloneqq (x, u)$, and $A \coloneqq \mathcal{X}_{\alpha(x)} \times \mathcal{U}_{\beta(u)}$. Observe that $\mathbb{E}[Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_{\alpha(x)})] = \mathbb{E}[x_2 \mid Z \in A]$. Note that these conditional expectations exist because $Z$ has full support and by construction $A$ has positive Lebesgue measure. The following holds

$$
\|\mathbb{E}[x_2 \mid Z \in A] - \mathbb{E}[x_2 \mid Z = z]\| = \left\| \int_{\mathbb{R}^d} x \left[ \frac{P(Z \in A | x_2 = x)}{P(Z \in A)} - \frac{p(Z = z | x_2 = x)}{p(Z = z)} \right] p_{x_2}(x) dx \right\|
$$

$$
\leq \int_{\mathbb{R}^d} \|x\| \left| \frac{P(Z \in A | x_2 = x)}{P(Z \in A)} - \frac{p(Z = z | x_2 = x)}{p(Z = z)} \right| p_{x_2}(x) dx
$$

$$
\leq \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \int_{\mathbb{R}^d} \|x\| \frac{p(Z = z | x_2 = x)}{p(Z = z)} p_{x_2}(x) dx
$$

$$
= \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \mathbb{E}[\|x_2\| \mid Z = z].
$$

The first inequality is an application of Jensen's inequality. The second inequality is an application of Assumption 4 and the fact that the diameter of $A$ is no more than $\varepsilon$. $\qquad \square$

### F.6.2. APPLYING LEMMAS TO BOUND TERMS

Armed with these lemmas we can proceed with bounding each of the aforementioned terms.

**First term:** Recall that the following holds for a $\tau^2$-subgaussian random variable $X$

$$
P(|X - \mathbb{E}[X]| > \delta|\mathbb{E}[X]|) \leq 2 \exp\left( \frac{-\delta^2 \mathbb{E}[X]^2}{2\tau^2} \right).
$$

For any $\alpha$, $\hat{Z}(\mathcal{X}_\alpha)$ is $\frac{1}{4n}$ subgaussian. This means we need $n = \frac{1}{2\delta^2 Z(\mathcal{X}_\alpha)^2} \log(4|\mathcal{N}|/\rho)$ samples to get $\mathcal{U}_{\beta(u)} \hat{Z}(\mathcal{X}_\alpha)$ within error of $\delta Z(\mathcal{X}_\alpha)$ of $Z(\mathcal{X}_\alpha)$ with probability $\rho/(2|\mathcal{N}|)$. Using union bound, we have that with probability with at least $1 - \rho/2$,

$$
\left\| \sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) \hat{Z}(\mathcal{X}_\alpha) - \sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) Z(\mathcal{X}_\alpha) \right\|
$$

$$
\leq \sum_\alpha \|\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\| |\hat{Z}(\mathcal{X}_\alpha) - Z(\mathcal{X}_\alpha)|
$$

$$
\leq \delta \sum_\alpha \|\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\| Z(\mathcal{X}_\alpha)
$$

$$
\leq \delta \sum_\alpha \|\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) - Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\| Z(\mathcal{X}_\alpha) + \delta \sum_\alpha \|Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\| Z(\mathcal{X}_\alpha)
$$

$$
\leq \delta\gamma + \delta \sum_\alpha \|Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)\| Z(\mathcal{X}_\alpha)
$$

$$
\leq \delta\gamma + \delta R + \delta \mathbb{E}[\|g(u_1)\| \mid u_1 \in \mathcal{U}_{\beta(u)}, x_1 \in \mathcal{X}_\alpha]
$$

$$
\leq \delta\gamma + 2\delta R
$$

where the first inequality comes from triangle inequality. The second inequality comes from subgaussianity. The third inequality is from triangle inequality. The fourth inequality is from the bound of the **Second term** below. The fifth and sixth inequalities are from triangle inequality, compactness, and from the fact $\mathbb{E}[\xi_t] = 0$.

**Second term:** For any $\alpha$, $\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)$ is $\frac{\sigma^2}{n_{u,x}}$ subgaussian, which means its $\frac{d\sigma^2}{n_{u,x}}$ norm-subgaussian by Lemma 1 from (Jin et al., 2019). Thus, the following inequality holds

$$
P(\|\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) - \mathbb{E}[\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)]\| \geq t) \leq 2 \exp\left( -\frac{t^2 n_{u,x}}{2d\sigma^2} \right).
$$

This means we need $n_{u,x} = \frac{2d\sigma^2}{\gamma^2} \log(4|\mathcal{N}|/\rho)$ samples to get $\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)$ with error $\gamma$ of $\mathbb{E}[\hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha)]$ with probability $\rho/(2|\mathcal{N}|)$. Moreover, because the conditions of Lemma F.6 are met, we know these requirements will hold for all $n_{u,x}$ for large enough $n$. Using union bound, we have that with probability with at least $1 - \rho/2$,

$$\left\| \sum_\alpha \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) Z(\mathcal{X}_\alpha) - \sum_\alpha Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) Z(\mathcal{X}_\alpha) \right\| \leq \sum_\alpha \left\| \hat{Y}(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) - Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) \right\| Z(\mathcal{X}_\alpha) \leq \gamma$$

The first inequality comes from Jensen's inequality. The second comes from subgaussianity.

**Third term:**

$$\left\| \sum_\alpha Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) Z(\mathcal{X}_\alpha) - \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x] Z(\mathcal{X}_\alpha) \right\|$$
$$\leq \sum_\alpha \left\| Y(\mathcal{U}_{\beta(u)}, \mathcal{X}_\alpha) - \mathbb{E}[x_2 | u_1 = u, x_1 = x]) \right\| Z(\mathcal{X}_\alpha)$$
$$\leq \frac{2\eta}{1-\eta} \sum_\alpha \mathbb{E}[\|x_2\| | u_1 = u, x_1 = x] Z(\mathcal{X}_\alpha)$$
$$\leq \frac{2\eta}{1-\eta} (2R + c_1)$$

The first inequality comes from Jensen's inequality. The second comes from Lemma F.9. The third inequality comes from triangle inequality.

**Fourth term:** Recalling that $B \coloneqq \mathbb{R}^d \setminus \cup\mathcal{N}$.

$$\left\| \sum_\alpha \mathbb{E}[x_2 | u_1 = u, x_1 = x] Z(\mathcal{X}_\alpha) - \mathbb{E}[x_2 | do(u_1 \coloneqq u)] \right\|$$
$$\leq L\varepsilon + \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\|$$
$$\leq L\varepsilon + \mathbb{E}[\|f(x_1)\| \mathbf{1}\{x_1 \in B\}] + P_{x_1}(B) R.$$

The first inequality comes from Lemma F.8. The second inequality comes from $\mathbb{E}\xi_2 = 0$, triangle inequality, Jensen's inequality, and the definition of $R$.

Union bounding over the two events and bounding the first and second terms and combining all the inequalities gives the result.

# G. Relaxing Assumption 1

In context of Theorem 1, we can replace Assumption 1 with the following weaker assumption

**Assumption 6.** *Let $\xi_T$ be such that $\mathbb{E}[\xi_T] = \mathbb{E}[\xi_T | u_1 = u, x_1 = z]$ for all $u \in \mathbb{R}^p$ and $z \in \mathbb{R}^d$.*

This "no-correlation" type assumption is required for showing admissibility (Proposition 1), and it only needs to apply to the exogenous variation affecting the time step $T$ we are interested in estimating $\mathrm{PE}_T$. Having said that, Assumption 1 is necessary for Theorem 2. Mutual independence is crucial for our proof technique.