
Causal Inference out of Control: The Steerability of Consumption

Gary Cheng*
Stanford University
chenggar@stanford.edu

Moritz Hardt*
Max Planck Institute for Intelligent Systems, Tübingen
hardt@is.mpg.de

Celestine Mender-Dünner*
Max Planck Institute for Intelligent Systems, Tübingen
cmendler@tuebingen.mpg.de

Abstract

Regulators and academics are increasingly interested in the causal effect that algorithmic actions of a digital platform have on consumption. We introduce a general causal inference problem we call the steerability of consumption that abstracts many settings of interest. Focusing on observational designs, we exhibit a set of assumptions for identifiability that significantly weakens the often unrealistic overlap assumptions of standard designs. The key insight behind our assumptions is to model the dynamics of consumption, viewing the platform as a controller acting on a dynamical system. From this dynamical systems perspective, we show that exogenous variation in consumption and appropriately responsive control actions are sufficient for identifying steerability of consumption. Our results illustrate the fruitful interplay of control theory and causal inference, which we corroborate with examples from econometrics, macroeconomics, and machine learning.

1 Introduction

How much does an increase in advertisement decrease screen time? Do algorithmic recommendations increase consumption of inflammatory content? Does exposure to diverse news sources mitigate political polarization? These are but a few of the questions that firms, researchers, and regulators alike ask about digital platforms. The answer to any such question requires non-trivial causal inference, since past consumption influences both future consumption as well as the algorithmic selection of content. Resolving confounding through randomization in the form of A/B tests is standard in the industry. However, randomization is often not possible. Experiments may be ethically fraught as past experience shows [Kramer et al., 2014, PNAS, 2014], technically challenging to implement, or prohibitively expensive. Moreover, outside investigators may simply not have the power to experimentally intervene in the practices of a platform. Observational causal inference presents an intriguing alternative. Unfortunately, standard observational causal designs would require that the data satisfy an overlap assumption with respect to past consumption. Since user data on digital platforms is often high-dimensional, overlap is unlikely to hold [D’Amour et al., 2017] resulting in invalid inferences.

In this work, we introduce and study a general causal inference problem, called *steerability of consumption*, that unifies numerous questions about the causal effects of algorithmic actions on digital platforms. In our problem formulation, consumption evolves over time according to an unknown dy-

* Authors ordered alphabetically.

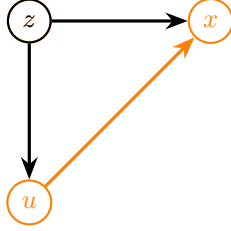


Figure 1: Standard model of the problem

namical system. A platform aims to influence consumption through actions performed at each point in time. Our goal is to identify the effect that the action has on the state of the dynamical system. Taking this control theoretic perspective of the problem, we’re able to exhibit a significantly weaker set of assumptions sufficient for identifiability and estimation. Informally speaking, it suffices that consumption has exogenous variation and that the control action of the platform is non-degenerate.

Our results illustrate the benefits of taking a control-theoretic perspective on causal inference. Rather than omitting the role that time plays in causal inference, we explicitly keep track of a time index. This allows for a more precise problem formulation that admits weaker identifiability assumptions. We illustrate the merits of our problem formulation with examples from econometrics, macroeconomics, and machine learning.

1.1 Our work

At the outset, we’re interested in the causal effect of an action u on a variable x subject to observed confounding by a variable z . We interpret the action u as an algorithmic intervention of a firm in a digital platform, while the variable x captures relevant user features, such as what content the user consumed. Figure 1 presents the standard causal model for the scenario. The confounding variable captures all available past information that influences both the choice of control action u and the variable x . To apply observational causal inference machinery in the standard model, the data generating distribution must assign positive weight to all strata defined by the confounding variable z . In the common case where z is high-dimensional, this assumption is unlikely to hold [D’Amour et al., 2017].

The starting point of our work is therefore a departure from the standard model. Instead, we model a dynamical system where user features x evolves over time. The control action u is a function of previously collected data and updated repeatedly based on the most recent observations of x . For example, video recommendations are updated based on recent views, and ads are chosen based on past clicks. Our model can be seen as postulating that most recent observations are the main driver of confounding. These assumptions result in the dynamical system illustrated in Figure 2.

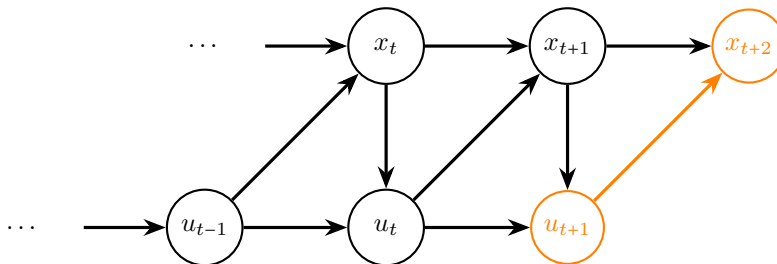


Figure 2: Autoregressive confounding structure

Technically speaking, the graph in Figure 2 could be expressed as an instance of the standard model in Figure 1 where the confounder z corresponds to the set of all past state and action variables. While this general model encapsulates our problem without loss of generality, it ignores the salient structure of the problem, demanding overly strong assumptions. Our main contribution is to instead exploit this structure to make estimating the steerability of consumption from observational data possible under more plausible assumptions. Our technical contributions can be summarized as follows:

1. In Section 2 we discuss our refined structural model for studying the steerability of consumption in detail. The main feature of our model is that it explicitly models temporal dynamics and posits a structure on how previous states and actions interact. We provide several concrete examples of problems that can be approximated by our model.
2. In Section 3, we study identifiability of the steerability of consumption within the context of our model. We demonstrate that exogenous variation in consumption and sufficient expressivity in the control response enable causal identification, circumventing direct interventions on the control action. We also study a linear version of our causal model. In this model, we show that 1) we can weaken the variation in exogenous consumption required for identifiability, and 2) we can weaken assumptions on the expressivity of the control response so long as the auditor observes longer roll-outs of the dynamical system.
3. In Section 4, we present an experiment using actual US interest rate and inflation rate data. Here we provide empirical evidence that overlap assumptions are more likely to be satisfied in our dynamical system model. In Appendix A, we perform empirical investigations on a linear dynamical system perturbed by Gaussian noise to show the connection between the rank condition underlying our theory and identifiability. We find that the conditioning of our observations improves as more variations are accumulated, making the causal inference problem “easier” over time.
4. To complement our identifiability results, in Appendix B, we provide a basic treatment of how practitioners could use our proposed observational model. In particular, we outline a non-parametric estimator of the steerability of consumption, and provide an associated convergence result for the general observational setting. We also propose a Double-ML-type estimator [Chernozhukov et al., 2017] which exploits longer observed roll-outs of the dynamical system. We provide finite sample guarantees for it in the linear setting.

1.2 Background

The fact that digital platforms, their predictions, and their actions non-trivially impact the individuals that interact with the platform has widely been recognized in diverse applications ranging from traffic forecasting, content recommendation to social media [c.f., Shmueli and Tafti, 2020, Thai et al., 2016, Fleder et al., 2010, Adomavicius et al., 2013, Chaney et al., 2018, Krauth et al., 2022].

Recently Hardt et al. [2022] have drawn a formal connection between the extent to which a platform can steer user behavior and the economic concept of power. Their proposal of *performative power* measures the strength of this effect and crucially relies on being able to estimate the steerability of consumption. From this lens, our work provides a roadmap for sufficient conditions for how performative power can be assessed from observational data. More broadly the implications of predictions impacting populations in the context of machine learning have been studied under the umbrella of performative prediction [Perdomo et al., 2020]. Related to our work Mendler-Dünner et al. [2022] focus on estimating the causal effect of predictions on eventual outcomes from observational data. However their problem statement and the static model does not account for time nor take advantage of repeated interactions between the predictor and the population.

Our work draws a connection between the causal question of estimating steerability of consumption and dynamical system theory. The causal model we propose in this paper is reminiscent of dynamical systems in control theory. Abbasi-Yadkori et al. [2011] and Abbasi-Yadkori and Szepesvari [2011] propose methods of controlling linear quadratic control systems with unknown dynamics. Their methods crucially rely on performing system identification—estimating the unknown parameters which describe the dynamics. While our problem can also be stated as a system identification problem, their observation model differs from the one we propose, as they focus only on the linear setting, and they opt to analyze one longer rollout of the dynamical system, instead of shorter, iid rollouts. While we give [Abbasi-Yadkori et al., 2011, Abbasi-Yadkori and Szepesvari, 2011] special attention because of the similarity of their model to ours, system identification, described more generally by Ljung [2010] as the “art and science of building mathematical models of dynamic systems from observed input-output data” is ubiquitous in control theory.

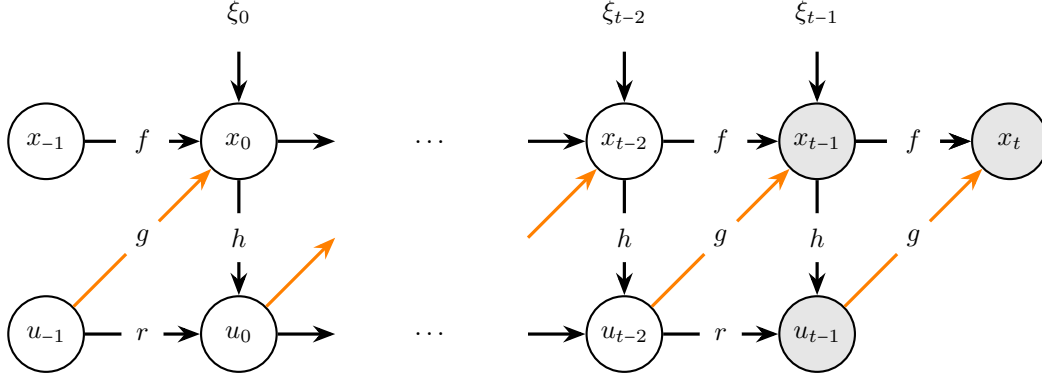


Figure 3: System dynamics with independent state perturbations.

2 Model

As we highlighted in the introduction, a feature of our model is that it unrolls time and makes the temporal dynamics of the problem explicit. Therefore, let $x_t \in \mathbb{R}^d$ and $u_t \in \mathbb{R}^p$ denote the consumption and control action at time step t respectively. We assume for all $t \geq 0$ the dynamics of the system follow

$$\begin{aligned} x_t &= f(x_{t-1}) + g(u_{t-1}) + \xi_t \\ u_t &= h(x_t) + r(u_{t-1}) \end{aligned} \quad (1)$$

with $(u_{-1}, x_{-1}) \sim P_{-1}$ for a given distribution P_{-1} , describing the joint distribution over state and action pairs at time step $t = -1$. x_t is some measure of consumption, and $\xi_t \in \mathbb{R}^d$ will model the exogenous variation in x_t . We do not make any assumption on P_{-1} and ξ_t for now, beyond there existing a density. We focus on time steps $t \geq 0$ without loss of generality. We will use the graph in Figure 3 as a visual representation of our dynamics. The functions $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, $g: \mathbb{R}^p \rightarrow \mathbb{R}^d$, $h: \mathbb{R}^d \rightarrow \mathbb{R}^p$, and $r: \mathbb{R}^p \rightarrow \mathbb{R}^p$ describe how state and control variables affect one another.

Remark Our model can be generalized beyond addition to more sophisticated aggregation functions: i.e., $x_t = H_x(f(x_{t-1}), g(u_{t-1})) + \xi_t$ and $u_t = H_u(h(x_t), r(u_{t-1}))$. Such a generalization could be useful in machine learning settings with strategic behavior. For example, $H_u := \operatorname{argmin}_u L(u, h(x_t), r(u_{t-1}))$ could model an optimization algorithm which uses the previous state and control action to decide the action to take in the next time step. In this work, we focus on the additive setting for expository clarity.

As we mentioned earlier, the auditor is interested in estimating the steerability of consumption from observational data. With respect to the model we outlined, we define it as follows.

Definition 2.1 (Steerability of Consumption). *For the dynamical system specified by Equation (1) and Figure 3, the steerability of consumption with respect to a baseline action u and intervention u' is defined as follows:*

$$\mathcal{S}(u, u') := \mathbb{E}[x_t \mid \operatorname{do}(u_{t-1} := u')] - \mathbb{E}[x_t \mid \operatorname{do}(u_{t-1} := u)].$$

This quantity captures how platform interventions can alter user behavior. We note that because we have a repeating structure in our graph, the steerability of consumption computed using x_t is the same as the associated quantity computed using x_k (e.g., the relationship between x_t and u_{t-1} is the same as x_k and u_{k-1} for any k). Thus, a sufficient condition for identifying $\mathcal{S}(u, u')$ is to identify $\mathbb{E}[x_t \mid \operatorname{do}(u_{t-1} := u)]$.

2.1 Running example

Consider an auditor who is interested in estimating the impact of YouTube's video recommendation on the consumption patterns of its users. Let $y_t \in \mathbb{R}^p$ be some measure of content consumption

(number of hours streamed) for p video categories of interest on YouTube during week t for a given user. Let $z_t \in \mathbb{R}^{d_z}$ comprised of measurements about the platform such as revenue per category, click-through rate per category, unique weekly users, unique advertisers per category, competitors' performance, etc. which could be confounders. We can think of the joint vector $[y_t; z_t] \in \mathbb{R}^d$ as the state variable x_t for $d = p + d_z$. The control 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 controller (YouTube) 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 control action u_{t-1} impacts the average watch habits y_t of users. More specifically, they are interested in the first d coordinates of the steerability of consumption $\mathcal{S}(u, u')$. Recall that the auditor can ignore the other coordinates because they correspond to the confounding variables.

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

$$\begin{aligned} z_t &= f_1(z_{t-1}, y_{t-1}) + g_1(u_{t-1}) + \xi_t^{(1)} \\ y_t &= f_2(z_{t-1}, y_{t-1}) + g_2(u_{t-1}) + \xi_t^{(2)}. \end{aligned}$$

f_1 models how the performance metrics chosen as a target variable by the firm evolve over time, and g_1 models YouTube's ability to control this metric. f_2 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 competitors poach). g_2 , the quantity the auditor wants to estimate, models how much consumption increases as more recommendations get served. $\xi_t^{(1)}, \xi_t^{(2)}$ is the natural variation in user preferences. For example, the Bitcoin price may increase due to changes in economic conditions, leading to many more users watching cryptocurrency videos on YouTube; this change in behavior is independent of past consumption and YouTube's recommendations. We can model the controller similarly as

$$u_t = h(z_t, y_t) + r(u_{t-1}),$$

where h models Youtube's algorithm of how viewer statistics and other metrics affect recommendations in the future. The function r models how YouTube regularizes its recommendations to avoid overfitting to recent activity. For example, a sudden spike in viewership in the "music video" category could be the result of transitive internet trend; overfitting to this viewership spike may be suboptimal for the future. This example serves to illustrate how our model of user dynamics (1) applies to a concrete use-case.

Beyond recommender systems. The steerability of consumption is not a term specific to recommender systems; rather, it is a general term referring to the relationship control actions have on consumption. There are many other settings for which our model is applicable. Microeconomists 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 this quantity is precisely the steerability of consumption. Confounders like product quality can be accounted for in the state variable x_t . A classical problem in macroeconomics is estimating the effect the Federal Interest Rate has on inflation and unemployment. We can use our framework to estimate this effect by modeling the Federal Interest Rate as the control action u_t and inflation and unemployment rates as the state x_t . In this example, GDP and other measures of the global economy could be possible confounders to account for. Many digital advertising platforms (and third-party auditors) are interested in whether personalized advertising increases platform activity. On one hand, advertisements clutter user interfaces, making the user experience less streamlined, but on the other hand, personalized advertisements provide users with more opportunities to engage, giving the platform more influence over user lives. Again, we can use our model to estimate the effect advertisements have on engagement; this time letting x_t be some measure of engagement (e.g., clicks, time online) and u_t be some measure of the type and quantity of ads served, while accounting for confounders like other platform performance measures such as monthly active users.

3 Identifiability from exogenous state variations

Let us start from the general causal graph in Figure 1 and recall classical results from causal identifiability in the presence of an observed confounding variable z [Pearl, 2009]. They tell us that a sufficient condition for identifiability of the causal effect of u on x is *admissibility* and *overlap*.

Definition 3.1 (Admissibility [Pearl, 2009]). *A continuous random variable Z with density p_Z is admissible with respect to treatment U and observation X if it satisfies the adjustment formula:*

$$\mathbb{E}[X \mid \text{do}(U := u)] = \int \mathbb{E}[X \mid U = u, Z = z] p_Z(z) dz. \quad (2)$$

Definition 3.2 (Overlap). *A continuous, \mathbb{R}^d -valued random variable W with density p has overlap if $p(w) > 0$ for all $w \in \mathbb{R}^d$.*

Admissibility is a property of a set of random variables Z that ensures the adjustment formula (2) is a valid tool for estimating causal effects by adjusting for confounding. It states that the causal quantity on the left-hand-side can be expressed as a function of pre-interventional data, and thus can be estimated from such data without bias. Admissibility coupled together with overlap over treatment, confounder pairs—that is $p(u, z) > 0$ for any pair u, z of interest—guarantees that $\mathbb{E}[X \mid U = u, Z = z]$ is well defined for any u, z and thereby allows the auditor to estimate steerability of consumption in a non-parametric way by adjusting for the confounder z using the adjustment formula (2).

To apply the adjustment formula in the context of our refined causal model (Figure 3) we aim to find the smallest set of admissible random variables for adjustment and estimation of the steerability of consumption. Without loss of generality we focus on identifying the steerability of consumption \mathcal{S}_t for $t = 2$. Thus, we are interested in the quantity $\mathbb{E}[x_2 \mid \text{do}(u_1 := u)]$ from which we can derive S_2 . The proof of the following result is found in Appendix C.1.

Proposition 1 (Admissibility in our model). *Given the causal model in Figure 3. The state variable x_1 is admissible with respect to treatment u_1 and observation x_2 .*

Given Proposition 1 and the adjustment formula (2) we know that

$$\mathbb{E}[x_2 \mid \text{do}(u_1 := u)] = \int_z \mathbb{E}[x_2 \mid u_1 = u, x_1 = z] p_{x_1}(z) dz.$$

Being able to compute the right hand side for any u from observational data is a sufficient condition for identifiability. The joint density of (x_1, u_1) being positive everywhere is sufficient for the right hand side to be well defined; thus, the question of identifiability reduces to a question of overlap.

3.1 Key assumptions

We highlight the two assumptions on the dynamical system in (1) that allow us to establish overlap over (x_1, u_1) . The first assumption posits the presence of natural perturbations in the state variable x . In particular, it requires that there is exogenous noise in the system that leads to sufficient variation in the user attribute variables across time.

Assumption 1 (Overlapping Exogenous Noise). *The noise variables in our causal model (Figure 3) are such that $(\xi_0, \xi_1) \mid \{x_{-1} = a, u_{-1} = b\}$ has overlap for all $a \in \mathbb{R}^d$ and $b \in \mathbb{R}^p$.*

We note that the noise variables ξ_0, ξ_1 do not need to be independent of each other or of the previous states, all we need is overlap. The second assumption concerns the controller. It needs to be sufficiently sensitive to the variations in the user attribute variable x , so that the variations provided by Assumption 1 propagate into the control action u .

Assumption 2 (Responsive Control Action). *Let $q_c : \mathbb{R}^d \rightarrow \mathbb{R}^p$ defined as $q_c(y) := r(h(y) + c)$ describe how the current state y affects the next control action given that the previous control action was c . Assume q_c is a surjective, invertible map with an invertible Jacobian for all $c \in \mathbb{R}^p$.*

To illustrate why one might expect the above two assumptions to hold, consider the YouTube recommender system example from Section 2.1. Exogenous perturbations ξ_t can correspond to factors like new trends or economic conditions which often have a very large effect on viewership statistics and the other metrics captured in the user attribute variable x_t , giving us reason to believe that Assumption 1 is satisfied in this example. To argue for surjectivity, recall that r describes the effect

previous control actions have on future recommendations, and h captures the effect the state has on the control action. First, we don't expect the control action to vary wildly between time steps, so we may expect the controller to adopt a simple policy, like discounting the previous action $r(u) = \alpha u$. This action happens to be surjective. Furthermore, in general, the dimensionality of the state is larger than the dimensionality of the control system (the controller has limited actions that it can take to influence the system and there are an arbitrarily large number of metrics and confounders we can add to the state variable). We posit that since $d \geq p$, there exists some setting of the state that maps to every possible control response i.e., h is surjective. Finally, because q_c is the composition of h and r , it is also surjective. The invertible Jacobian condition imposes a form of "monotonicity" on q_c which in the YouTube setting corresponds to: more views in category i should lead to more recommendations in category i .

3.2 General identifiability result

Building on the two assumptions discussed in the previous section, we can now provide our general identifiability results. Proofs can be found in Appendix D.

Theorem 1 (Identifiability from two exogenous state perturbations). *Consider the dynamical system in (1) with any arbitrary P_{-1} and noise variables (ξ_0, ξ_1) satisfying Assumption 1. Let the auditor observe iid samples of (x_1, u_1, x_2) . Then, if Assumption 2 holds, (x_1, u_1) has overlap and the steerability of consumption $\mathcal{S}(u, u')$ is identifiable for any $u, u' \in \mathbb{R}^p$.*

This theorem is stated for the setting where the overlapping exogenous noise occurs in subsequent time steps and the auditor observes the system starting the time step of the second overlapping exogenous noise spike. This result can be generalized further with some additional assumptions. We can handle overlapping exogenous noise spikes which do not occur in subsequent time steps by assuming a variation of Assumption 2 which guarantees the function mapping the effect the first overlapping noise spike on the control action observed is surjective. We can handle settings where the the auditor observes the system after the second overlapping noise spike occurs by assuming all other noise spikes operating on states in other times steps are independent (but need not have overlap). We avoid explicitly detailing these variations for the sake of clarity.

In words, this result states that exogenous perturbations on two of the state variables are sufficient for the auditor to identify the steerability of consumption from observations. Interestingly, observing the system after two exogenous perturbations to the state is crucial for identifiability. Observing the system after only one perturbation is not sufficient. We show this necessity by presenting Theorem 2.

Theorem 2 (Unidentifiability of one exogenous state perturbation). *Consider the dynamical system in (1) with initial conditions $x_{-1} = u_{-1} = 0$, and mutually independent noise variables ξ_0 and ξ_1 . Let the auditor receive iid samples of (x_0, u_0, x_1) . Then, for any functions $f : \mathbb{R}^d \rightarrow \mathbb{R}^d, g : \mathbb{R}^p \rightarrow \mathbb{R}^d, h : \mathbb{R}^d \rightarrow \mathbb{R}^p$, the steerability of consumption is unidentifiable; i.e., there exists $(\hat{f}, \hat{g}, \hat{h}) \neq (f, g, h)$ which induces the same distribution of (x_0, u_0, x_1) .*

Intuitively, a single perturbation is not enough for identifiability because the control action is a deterministic function of the user attribute variable. Thus, the auditor is only able to see one value of u_0 for any value of x_0 . The second noise spike impacts the state and control variables differently, providing enough variation for overlap, making the steerability of consumption identifiable. One advantage of our approach is that Assumption 2 is an assumption on the design of the controller (which is typically a deterministic mechanism) and this assumption can be verified with enough knowledge of the control system, while typical overlap assumptions are not, as they model the behaviors of users, which could vary wildly for any number of reasons.

3.3 Identifiability in the linear model

Let us outline how additional structural assumptions can be exploited for identification. We study a linear instantiation of the general model we proposed and show that one noise spike, instead of two, is sufficient for identifiability as long as the auditor observes the system for enough times steps. We provide some intuition why this is the case. Because the system is linear, explicit overlap of (x_t, u_t) is not needed; instead the auditor just needs overlap over enough directions such that their span covers all values that (x_t, u_t) can take on. As we will show, one spike of exogenous state variation and sufficient expressivity in the linear model ensures that observations of (x_0, u_0) , and (x_1, u_1)

provide “one direction” each so that the auditor can use linear combinations of these observations to recover the steerability of consumption.

In this section, we will assume the functions f, g, h, r are linear and defined as

$$\begin{aligned} f(x) &:= Ax & g(x) &:= Bx \\ h(x) &:= Cx & r(x) &:= Dx, \end{aligned} \tag{3}$$

where $A \in \mathbb{R}^{d,d}, B \in \mathbb{R}^{d,p}, C \in \mathbb{R}^{p,d}, D \in \mathbb{R}^{p,p}$. We note that in the linear setting, it is sufficient for the auditor to identify the matrix B , as they can readily compute the steerability of consumption with this information. There are many settings where the dynamical system we want to study is linear, or modeled as linear. For example, a linear controller, such as a proportional-integral (PI) controller, controlling a physical system, such as a quadrotor, can be modeled with linear dynamics [Bouabdallah et al., 2004]. Before we present our results, we begin with an alternative definition of overlap that we will use in this section, termed “full span”.

Definition 3.3 (Full span). *A random vector $x \in \mathbb{R}^d$ has full span if for all vectors $a \in \mathbb{R}^d$ such that $a \neq 0$, $a^T x$ is almost surely not a constant.*

Intuitively, because we are dealing with linear systems, the auditor just needs to be able to see or reconstruct all the directions the data could span, as opposed to observing the full space. Now we present our linear identifiability result in Theorem 3; the proof can be found in Appendix D.3.

Theorem 3 (Simple identifiable DAG). *Consider the dynamical system in (1) with any arbitrary P_{-1} and with functions f, g, h, r defined in (3). Let $\{\xi_t\}_{t \geq 0}$ be a set of mutually independent random vectors, and let ξ_0 have full span. Let the auditor observe iid samples of $(x_0, u_0, \dots, x_{K-1}, u_{K-1}, x_K)$ for any $K \geq 2$. If x_{-1} and u_{-1} are mutually independent of ξ_0 and $[DC, \dots, D^{K-1}C]$ is full row rank, then the steerability of consumption $\mathcal{S}(u, u')$ is identifiable for any $u, u' \in \mathbb{R}^p$. If $x_{-1} = 0$, $u_{-1} = 0$, and $\xi_t = 0$ for $t \geq 1$, then $[DC, \dots, D^{K-1}C]$ being full row rank is necessary for the identifiability of $\mathcal{S}(u, u')$.*

Because Theorem 2 still holds even in the linear setting, Theorem 2 and Theorem 3 together fully characterize the tradeoff between identifiability, number of time steps observed, and rank conditions on the controller dynamics matrices. Summarizing briefly, the auditor cannot identify the steerability of consumption from observations from only observations of (x_0, u_0, x_1) , but identification is possible from observations $(x_0, u_0, \dots, x_{K-1}, u_{K-1}, x_K)$ for $K \geq 2$. Moreover, as K gets larger, the assumptions required on the controller dynamics matrices get weaker. Intuitively this weakening occurs because as the auditor observes more time steps, they have more opportunities to observe all of the “directions” (x_t, u_t) can take on, allowing for more poorly conditioned dynamical systems to be identifiable. We note that in linear control problems, like the aforementioned quadrotor example, the controller dynamics are often known, and thus, the rank conditions sufficient for identifiability are frequently easy to check.

4 Experiments

In this section, we corroborate our theoretical findings with an experiment using actual US interest rate and inflation rate data. Our empirical results suggest that overlap assumptions are more likely to be satisfied in our dynamical system model. We also perform additional synthetic experiments on a linear dynamical system perturbed by Gaussian noise, but we defer these results to Appendix A.

4.1 Federal Interest Rate and Inflation

We apply our model to a time series dataset [Reserve, 2017], containing monthly records of the US Federal Interest Rate and US Inflation rate from October, 1982 to December, 2008 and estimate the causal effect of Federal Interest Rate on the Inflation rate. We will compare the estimates obtained by relying on our modeling assumptions with an estimate obtained from a more naive approach based on the general model from Figure 1 that does not posit a Markovian assumption on how confounding variables interact. We will use an adjustment formula estimator, outlined in Appendix B.1, and we will place a special emphasis on examining how well our overlap assumptions hold as we vary the number previous time steps K that we account for as a confounding variable. We discretize the effective interest rate into two buckets $((0.999, 5.25], (5.25, 11.5])$ and the inflation rate into three

buckets $((1.01, 2.5], (2.5, 3.8], (3.8, 5.9])$. The demarcation of the buckets is chosen such that each bucket marginally contains an equal number of datapoints. For a month $t \in [T]$, u_t will correspond to the discretized effective interest rate, and x_t corresponds to the discretized inflation rate.

Before constructing our estimator, we first specify K . With this quantity, we look at a sliding window over the data $\{(x_t, \dots, x_{t+K}, u_t, \dots, u_{t+K-1})\}_{t=0}^{T-K}$. We will treat these samples as the iid observations the auditor observes. In the context of our adjustment formula estimator, we will use u_{t+K-1} as the treatment variable, $z_{t+K-1} := (x_t, \dots, x_{t+K-1})$ as the confounders, and x_{t+K} as the outcome. We estimate the steerability of consumption using the adjustment formula estimator from Appendix B.1 for a given treatment u as follows

$$\hat{x}(u) := \sum_z \left[\frac{\sum_t x_{t+1} \mathbf{1}\{z_t = z, u_t = u\}}{\sum_t \mathbf{1}\{z_t = z, u_t = u\}} \right] \frac{\sum_t \mathbf{1}\{z_t = z\}}{n},$$

where $n = T - K$ is the number of observations in our dataset. We now provide our findings in Table 1. High interest rate corresponds to an interest rate being in the higher valued bucket (5.31, 11.64]; low interest rate corresponds to the (0.159, 5.31] valued bucket.

	Interest Rate (Intervention u)	Estimated Effect on inflation	Fraction of undefined terms	Probability mass of undefined terms
K=1	High	3.29	0 / 3	0.0%
	Low	2.98	0 / 3	0.0%
K=2	High	3.29	0 / 7	0.0%
	Low	2.99	0 / 7	0.0%
K=3	High	N/A	1 / 14	0.3%
	Low	N/A	2 / 14	1.0%
K=4	High	N/A	5 / 23	3.2%
	Low	N/A	6 / 23	2.9%
K=5	High	N/A	8 / 32	4.2%
	Low	N/A	11 / 32	5.2%
K=6	High	N/A	13 / 41	6.8%
	Low	N/A	16 / 41	7.1%

Table 1: Estimated effects on inflation rate for interest rate interventions.

As K gets larger, there are many z and u pairs where $\sum_t \mathbf{1}\{z_t = z, u_t = u\} = 0$; when this happens, $\hat{x}(u)$ is not well defined. N/A denotes when this occurs. The ‘‘Fraction of undefined terms’’ column corresponds to the number of z values such that $\sum_t \mathbf{1}\{z_t = z, u_t = u\} = 0$ over the total number of values of z where $\frac{\sum_t \mathbf{1}\{z_t = z\}}{n} \neq 0$. The entries of ‘‘Probability mass of undefined terms’’ column is equal to $\sum_z \frac{\sum_t \mathbf{1}\{z_t = z\}}{n} \mathbf{1}\{\sum_t \mathbf{1}\{z_t = z, u_t = u\} = 0\}$. We can see that as K gets larger, the number of undefined, the relative fraction of undefined values, and the mass of said values gets larger.

This experiment highlights how overlap assumptions become harder to satisfy, the more time steps we track as confounders, even when the state and control values are binned-one-dimensional-scalars. We expect these issues get worse if the state (at each time step) is also high dimensional as well, which is often the case in modern machine learning settings. This experiment highlights the pressing need to mitigate overlap issues and showcases how our model can be useful to that end.

5 Discussion

In many settings, our model provides a reasonable first-order approximation to the problem of interest. Certainly, there might be aspects of the system, such as unobserved state variables, or more complex dependencies across time that our model is not able to capture. However, we are interested in understanding how powerful it is to explicitly model time and interactions among confounding variables for tackling causal questions. Thereby, we aim to provide a mostly unexplored direction for designing assumptions for causal inference in data scarce settings and illustrate how to leverage data across multiple time steps for causal estimation.

References

- Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *COLT*, 2011.
- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Gediminas Adomavicius, Jesse C. Bockstedt, Shawn P. Curley, and Jingjing Zhang. Do recommender systems manipulate consumer preferences? a study of anchoring effects. *Information Systems Research*, 24(4):956–975, 2013.
- S. Bouabdallah, A. Noth, and R. Siegwart. Pid vs lq control techniques applied to an indoor micro quadrotor. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (IEEE Cat. No.04CH37566)*, volume 3, pages 2451–2456 vol.3, 2004. doi: 10.1109/IROS.2004.1389776.
- Allison June-Barlow Chaney, Brandon M Stewart, and Barbara E. Engelhardt. How algorithmic confounding in recommendation systems increases homogeneity and decreases utility. *Proceedings of the 12th ACM Conference on Recommender Systems*, 2018.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James M. Robins. Double/debiased machine learning for treatment and structural parameters. *Econometrics: Econometric & Statistical Methods - Special Topics eJournal*, 2017.
- Alexander D’Amour, Peng Ding, Avi Feller, Lihua Lei, and Jasjeet S. Sekhon. Overlap in observational studies with high-dimensional covariates. *Journal of Econometrics*, 2017.
- Daniel Fleder, Kartik Hosanagar, and andreas buja. Recommender systems and their effects on consumers: the fragmentation debate. 06 2010.
- Moritz Hardt, Meena Jagadeesan, and Celestine Mender-Dünner. Performative power. *ArXiv*, abs/2203.17232, 2022.
- Tamas Jambor, Jun Wang, and Neal Lathia. Using control theory for stable and efficient recommender systems. *Proceedings of the 21st international conference on World Wide Web*, 2012.
- Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M. Kakade, and Michael I. Jordan. A short note on concentration inequalities for random vectors with subgaussian norm. *ArXiv*, abs/1902.03736, 2019.
- Adam D. I. Kramer, Jamie Guillory, and Jeffrey T. Hancock. Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences of the United States of America*, 111 24:8788–90, 2014.
- Karl Krauth, Yixin Wang, and M.I. Jordan. Breaking feedback loops in recommender systems with causal inference. *ArXiv*, abs/2207.01616, 2022.
- Lennart Ljung. Perspectives on system identification. *Annual Reviews in Control*, 34(1):1–12, 2010. ISSN 1367-5788. doi: <https://doi.org/10.1016/j.arcontrol.2009.12.001>. URL <https://www.sciencedirect.com/science/article/pii/S1367578810000027>.
- Celestine Mender-Dünner, Frances Ding, and Yixin Wang. Anticipating performativity by predicting from predictions. *ArXiv*, abs/2208.07331, 2022.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd edition, 2009. ISBN 052189560X.
- Juan Perdomo, Tijana Zrnic, Celestine Mender-Dünner, and Moritz Hardt. Performative prediction. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 7599–7609. PMLR, 2020.
- PNAS. Editorial expression of concern: Experimental evidence of massivescale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, 111(29):10779–10779, 2014. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1412469111>.

Federal Reserve. Federal reserve interest rates, 1954-present. <https://www.kaggle.com/datasets/federalreserve/interest-rates>, 2017.

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

Jérôme Thai, Nicolas Laurent-Brouty, and Alexandre M. Bayen. Negative externalities of gps-enabled routing applications: A game theoretical approach. In *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, pages 595–601, 2016. doi: 10.1109/ITSC.2016.7795614.

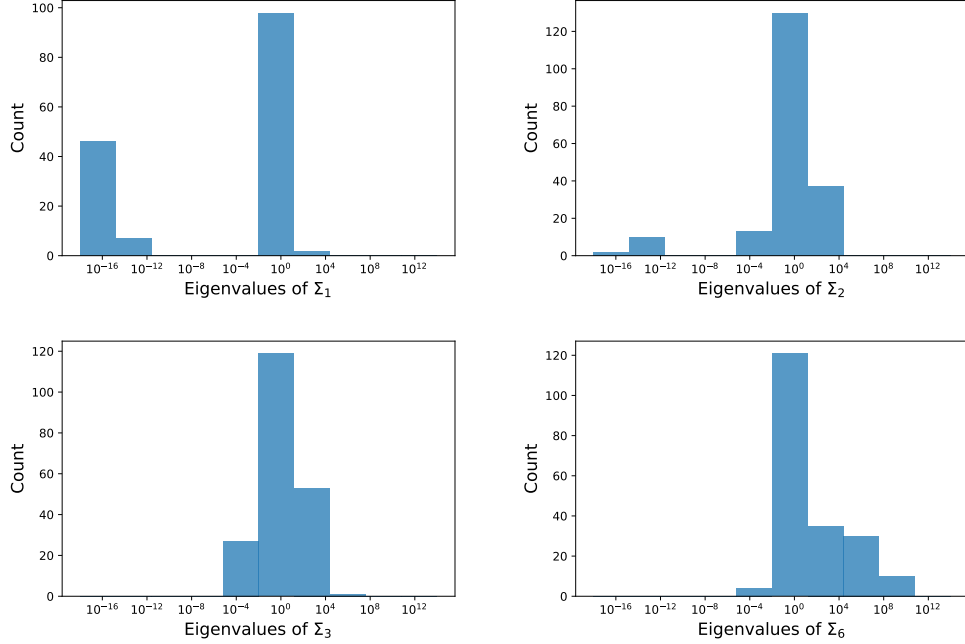


Figure 4: Histograms of eigenvalues of Σ_t defined in Appendix A

A Synthetic Experiments

We analyze the dynamical system (1) with linear dynamics (3) with independent Gaussian noise on the states to show how the noise leads to identifiability. In particular we show how the conditioning of the problem evolves across time steps and how the presence of noise makes the problem easier.

We instantiate the model as follows: We let $\xi_t \stackrel{d}{=} \mathcal{N}(0, I)$ for all t , starting from $x_{-1} = u_{-1} = 0$. 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 $B = WW^T/n$. We repeat this process to generate 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 $r = 80$.

For this system, we can explicitly write down how the covariance matrix of (x_t, u_t) , denoted Σ_t , evolves. Namely, from the dynamics

$$\begin{bmatrix} x_t \\ u_t \end{bmatrix} = J \begin{bmatrix} x_{t-1} \\ u_{t-1} \end{bmatrix} + M \varepsilon_t \quad \text{where} \quad J := \begin{bmatrix} A & B \\ CA & CB + D \end{bmatrix} \quad M := \begin{bmatrix} I \\ C \end{bmatrix}.$$

we can deduce that

$$\Sigma_t = J \Sigma_{t-1} J^T + M M^T = \sum_{k=0}^{t-1} (J^k) M M^T (J^k)^T.$$

We now plot the histogram of the eigenvalues of Σ_t for a random system that we generated. The important quantity to look out for is whether Σ_t is full rank. Directions in the null space of Σ_t (i.e., with zero eigenvalue) are directions that are unobserved by the auditor. Thus, when Σ_t becomes full rank, the system becomes identifiable because (x_t, u_t) becomes fully-spanning.

We note that in this system, $\text{rank } DC = 80 = r < d$ and $\text{rank } C = 80 = r < d$. Even though the setup of this system is not the same as Theorem 3, we believe linearity gives us a correspondence between the settings. Assuming for a moment that there is an equivalence, Theorem 3 tells us that observing (x_1, u_1) is not sufficient for identifiability. This is consistent with Figure 4 as there are still eigenvalues with value zero. However, since in this system $\text{rank}[DC, D^2C] = 100 = d$, Theorem 3 tells us observing (x_2, u_2) (i.e., observing three noise spikes) is sufficient for identifiability. This is

also consistent with Figure 4, as all eigenvalues are bounded away from zero now. Moreover, we see that the eigenvalues of get larger as more time passes: the eigenvalue mass of Σ_6 is further to the right of the eigenvalue mass of Σ_3 in Figure 4.

Our experiments corroborate our theoretical results which say that three noise spikes for this low rank system is necessary and sufficient for identifiability. Furthermore, our experiment also suggests that more noise spikes over more time steps make the observations better conditioned, likely making estimating the steerability of consumption easier for the auditor in practice.

B Estimators

To complement our identifiability results, we provide a basic treatment of how practitioners could estimate steerability of consumption from observational data under our proposed model. First, we outline a non-parametric estimator based on the adjustment formula and provide an associated convergence result for the non-parametric setting. This estimator assumes the auditor has iid observations of triplets (x_t, u_t, x_{t+1}) and explicitly requires overlap of (x_t, u_t) . Then, we propose a Double-ML-type estimator [Chernozhukov et al., 2017] that assumes the auditor has iid observations of $(x_t, u_t, x_{t+1}, u_{t+1}, x_{t+2})$ and analyze it in the linear model. This method exploits the longer roll-outs in the observational data and has the advantage of always being well defined, even when the overlap conditions needed for theoretical guarantees do not hold. For this section, without loss of generality, we will assume the dynamical system (1) starts at time index $t = -1$ and the auditor observes some subset of the variables $(x_0, u_0, x_1, u_1, x_2)$, similar to previous sections.

B.1 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 steerability of consumption. In order to do this, we first discretize a compact subset of the user state $\mathcal{X} \subset \mathbb{R}^d$ and control spaces $\mathcal{U} \subset \mathbb{R}^p$ with a cover, defined as follows.

Definition B.1 (Cover). *For a space \mathcal{X} , let $\mathcal{N} = \{U_\alpha \subset \mathcal{X} : \alpha \in A\}$ be an indexed family of sets U_α . \mathcal{N} is a cover of \mathcal{X} if $\mathcal{X} \subset \cup_{U \in \mathcal{N}} U$. \mathcal{N} is an ε -cover with respect to $\|\cdot\|_2$ if for all $U \in \mathcal{N}$ for all $x, y \in U$, $\|x - y\|_2 \leq \varepsilon$. \mathcal{N} is a disjoint cover if, additionally, for all $U, V \in \mathcal{N}$, $U \cap V = \emptyset$.*

Definition B.2 (Cover Representatives). *For a space \mathcal{X} and a disjoint cover \mathcal{N} , $\mathcal{R} \subset \mathcal{X}$ is a cover representative of \mathcal{N} if for all $A \in \mathcal{N}$, there exists $x \in \mathcal{R}$ such that $x \in A$ and for all $y \in \mathcal{R}$ where $y \neq x$, $y \notin A$.*

With this discretization, we can estimate the integral from Definition 3.1 using a finite sum. The specifics of how the estimator is constructed can be found in Algorithm 1, but we first introduce some notation. We let $y^{(k)}$ denote the k th sample of the random variable y . For a cover \mathcal{N} , we let $\mathcal{N}(u) \in \mathcal{N}$ denote the set where $u \in \mathcal{N}(u)$. We define $n_{u,x} := \sum_{k \in [n]} \mathbf{1}\{u_1^{(k)} \in \mathcal{N}_{\mathcal{U}}(u), x_1^{(k)} \in \mathcal{N}_{\mathcal{X}}(x)\}$ to be the number of observations which fall into $\mathcal{N}_{\mathcal{U}}(u)$ and $\mathcal{N}_{\mathcal{X}}(x)$.

To summarize, this estimator first discretizes the space of observations and then uses sample averages to estimate the terms in the adjustment formula Equation (2). We will need some mild assumptions to prove a guarantee on the estimator. Our first assumption controls how much previous user state and control 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 $z, z' \in \mathcal{X}$ and $u, u' \in \mathcal{U}$ it holds that*

$$\|\mathbb{E}[x_2|u_1 = u, x_1 = z] - \mathbb{E}[x_2|u_1 = u', x_1 = z']\| \leq L \left\| \begin{bmatrix} u \\ z \end{bmatrix} - \begin{bmatrix} u' \\ z' \end{bmatrix} \right\|.$$

We also need to control how far the empirical estimate $\mathbb{E}[Y^n(u, x)]$ 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. *For all $w, w' \in \mathcal{X}$ and $u, u' \in \mathcal{U}$, for which*

$$\left\| \begin{bmatrix} u \\ w \end{bmatrix} - \begin{bmatrix} u' \\ w' \end{bmatrix} \right\| \leq \varepsilon.$$

Algorithm 1 Adjustment formula estimator

Require: Cover granularity $\varepsilon > 0$, error tolerance $\delta \in (0, 1)$, $\gamma > 0$, failure probability tolerance $\rho \in (0, 1)$, compact subset $\mathcal{U} \subset \mathbb{R}^p$, compact subset $\mathcal{X} \subset \mathbb{R}^d$, control action $u \in \mathcal{U}$ of interest.

Step 1. Form a disjoint $\varepsilon/2$ -cover of \mathcal{U} such that every element has Lebesgue measure greater than 0 and label it $\mathcal{N}_{\mathcal{U}}$; let $\mathcal{R}_{\mathcal{U}}$ denote its respective cover representative. Form a disjoint $\varepsilon/2$ -cover of \mathcal{X} such that every element has Lebesgue measure greater than 0 and label it $\mathcal{N}_{\mathcal{X}}$; let $\mathcal{R}_{\mathcal{X}}$ denote its respective cover representative.

Step 2. Collect enough samples such that $n_{u,x} \geq \frac{2d\sigma^2}{\gamma^2} \log(4|\mathcal{R}_{\mathcal{X}}|/\rho)$ for all $x \in \mathcal{X}$ and the total number of samples n satisfies $n \geq \max_{x \in \mathcal{R}_{\mathcal{X}}} \frac{1}{2\delta^2 P(x_1 \in \mathcal{N}_{\mathcal{X}}(x))^2} \log(4|\mathcal{R}_{\mathcal{X}}|/\rho)$. Take data to form dataset $\mathcal{D}^n = \{(x_1^{(k)}, u_1^{(k)}, x_2^{(k)})\}_{k=1}^n, u \in \mathbb{R}^p$.

Step 3. Let

$$Y^n(u, x) := \frac{\sum_{k \in [n]} x_2^{(k)} \mathbf{1}\{u_1^{(k)} \in \mathcal{N}_{\mathcal{U}}(u), x_1^{(k)} \in \mathcal{N}_{\mathcal{X}}(x)\}}{\sum_{k \in [n]} \mathbf{1}\{u_1^{(k)} \in \mathcal{N}_{\mathcal{U}}(u), x_1^{(k)} \in \mathcal{N}_{\mathcal{X}}(x)\}}$$

$$Z^n(x) := \frac{1}{n} \sum_{k=1}^n \mathbf{1}\{x_1^{(k)} \in \mathcal{N}_{\mathcal{X}}(x)\}.$$

return the following estimator:

$$\hat{x}_2 := \sum_{x \in \mathcal{R}_{\mathcal{X}}} Y^n(u, x) Z^n(x).$$

For any $x \in \mathcal{X}$, the following condition on the density p holds for some $\eta(\varepsilon) \in (0, 1)$ such that $\lim_{\varepsilon \rightarrow 0} \eta(\varepsilon) = 0$:

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

This assumption ensures that the conditional distribution is “stable” in any ε -neighborhood. 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 iid samples of (x_1, u_1, x_2) . Suppose x_2 is σ^2 -subgaussian conditioned on u_1 and x_1 . Let $\mathbb{E}[\xi_2 | x_1 = x, u_1 = w] = 0$, $\mathbb{E}[\|\xi_2\| | x_1 = x, u_1 = w] \leq c_1$ for all $x \in \mathbb{R}^d$ and $w \in \mathbb{R}^p$. Let f and g be continuous functions, and define D such that $\sup_{x \in \mathcal{X}, w \in \mathcal{U}} \max(\|f(x)\|, \|g(w)\|) \leq D$. Let the conditions of Theorem 1 hold, Assumption 3 hold with L and Assumption 4 hold with η . For any compact subset $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{U} \subset \mathbb{R}^p$, for any specified $u \in \mathcal{U}$, for n large enough, Algorithm 1 executed with parameters $\varepsilon, \delta, \gamma, \rho, \mathcal{U}, \mathcal{X}, u$ returns \hat{x}_2 which has the following guarantee with probability at least $1 - \rho$:

$$\|\hat{x}_2 - \mathbb{E}[x_2 | do(u_1 := u)]\| \leq \delta\gamma + 2\delta D + \gamma + \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} (2D + c_1)$$

$$+ L\varepsilon + \mathbb{E}[\|f(x_1)\| \mathbf{1}\{x_1 \notin \mathcal{X}\}] + (1 - P_{x_1}(\mathcal{X}))D.$$

The proof of Theorem 4 can be found in Appendix E.1. We quickly go through all the terms in the bound, and verify that they can all be made arbitrarily small (with sufficient samples). δ and γ can be made smaller, so long as the auditor receives proportionally enough samples. The auditor can create a finer discretization to make ε smaller and therefore η 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., \mathcal{X} —becomes larger.

B.2 Double-ML-type estimator for separable model

Now we present an estimator which leverages the structure of the data generation procedure and performs a regression between residuals, reminiscent of Double Machine Learning [Chernozhukov et al., 2017]. The procedure requires hypotheses (i.e., sets of functions $\mathcal{F}^{d \rightarrow p}$, $\mathcal{F}^{d \rightarrow d}$, and $\mathcal{F}^{p \rightarrow d}$) of

Algorithm 2 Double ML estimator

Require: Loss functions ℓ_1, ℓ_2, ℓ_3 , Hypothesis classes $\mathcal{F}^{d \rightarrow p} \subset \{f \mid f: \mathbb{R}^d \rightarrow \mathbb{R}^p\}$, $\mathcal{F}^{d \rightarrow d} \subset \{f \mid f: \mathbb{R}^d \rightarrow \mathbb{R}^d\}$, and $\mathcal{F}^{p \rightarrow d} \subset \{f \mid f: \mathbb{R}^p \rightarrow \mathbb{R}^d\}$, Dataset $\mathcal{D}^n = \{(x_2^{(k)}, u_1^{(k)}, x_1^{(k)})\}_{k=1}^n$

Step 1.

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{F}^{d \rightarrow p}} \frac{1}{n} \sum_{k=1}^n \ell_1(u_0^{(k)}, h(x_0^{(k)}))$$

Step 2.

$$\hat{s} = \operatorname{argmin}_{s \in \mathcal{F}^{d \rightarrow d}} \frac{1}{n} \sum_{k=1}^n \ell_2(x_1^{(k)}, s(x_0^{(k)}))$$

Step 3.

$$\hat{g} = \operatorname{argmin}_{g \in \mathcal{F}^{p \rightarrow d}} \frac{1}{n} \sum_{k=1}^n \ell_3(x_2^{(k)} - \hat{s}(x_1^{(k)}), g(u_1^{(k)} - \hat{h}(x_1^{(k)})))$$

return $\hat{g}: \mathbb{R}^p \rightarrow \mathbb{R}^d$

what the relationships between control and state are. The details of the estimator are presented in Algorithm 2.

We now present some intuition about the estimator by analyzing it in the linear setting from Section 3.3. Our results can be generalized further, but we feel these generalizations only obfuscate the main message; thus, we focus on this simplified setting. In particular, for the remainder for the section, assume data is generated according to the dynamical system Equation (1) with functions f, g, h, r defined in Equation (3). Accordingly, in the context of Algorithm 2, we will consider linear function classes i.e., $\mathcal{F}^{d \rightarrow p} = \mathbb{R}^{p,d}$, $\mathcal{F}^{d \rightarrow d} = \mathbb{R}^{d,d}$, and $\mathcal{F}^{p \rightarrow d} = \mathbb{R}^{d,p}$. And we will assume ℓ_1, ℓ_2, ℓ_3 corresponds to square loss i.e., $\ell_j(a, b) = \frac{1}{2} \|a - b\|_2^2$. Our results can be generalized further, but we feel these generalizations only obfuscate the main message; thus, we focus on this simplified setting. Our result will need the following assumption to hold.

Observe that the following relationship holds between control and state variables in the linear setting.

$$x_2 - (A + BC)x_1 = B(u_1 - Cx_1).$$

Thus, if we had good estimates of $H := (A + BC)$ and C denoted by \hat{H} and \hat{C} respectively, we could regress $x_2 - \hat{H}x_1$ against $u_1 - \hat{C}x_1$ to get an estimate of the steerability of consumption B . In fact, this is exactly what Algorithm 2 is doing in the linear setting. Let $x_t^{(k)}, u_t^{(k)}, \xi_t^{(k)}$ denote the k th observations of x_t, u_t , and ξ_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 n samples of x_t, u_t , and ξ_t respectively. Algorithm 2 in the linear setting corresponds to first estimating H and C with the values \hat{H} and \hat{C} respectively before computing an estimate of the steerability of consumption \hat{B} . Explicitly, Algorithm 2 is performing these three computations

$$\begin{aligned} \hat{C} &= \operatorname{argmin}_{C \in \mathbb{R}^{p,d}} \frac{1}{2n} \|U_0 - CX_0\|_{\text{Fr}}^2 = U_0 X_0^T (X_0 X_0^T)^{-1} \\ \hat{H} &= \operatorname{argmin}_{H \in \mathbb{R}^{d,d}} \frac{1}{2n} \|X_1 - HX_0\|_{\text{Fr}}^2 = X_1 X_0^T (X_0 X_0^T)^{-1} \\ \hat{B} &= \operatorname{argmin}_{B \in \mathbb{R}^{d,p}} \frac{1}{2n} \|X_2 - \hat{H}X_1 - B(U_1 - \hat{C}X_1)\|_{\text{Fr}}^2 \\ &= (X_2 - \hat{H}X_1)(U_1 - \hat{C}X_1)^T ((U_1 - \hat{C}X_1)(U_1 - \hat{C}X_1)^T)^{-1}. \end{aligned}$$

Before we present our convergence result, we will need the following assumption to be satisfied.

Assumption 5 (ρ -Bounded System Dynamics). *The linear dynamical system specified by (1) and (3) has ρ -Bounded System Dynamics if*

$$\frac{\|A + BC\|_{\text{op}}}{\sigma_{\min}(DC)} \leq \rho.$$

Let's investigate the consequences of this assumption. First, consider the quantity $\|x_1 - \xi_1\|_2 / \|u_1\|_2$; this is the ratio between the magnitude of the state and control action after one time step of dynamics ignoring noise and assuming the system starts from equilibrium $x_{-1} = u_{-1} = 0$. Because $\|x_1 - \xi_1\|_2 / \|u_1\|_2 \leq \frac{\|A+BC\|_{\text{op}}\|x_0\|_2}{\lambda_{\min}(DC)\|x_0\|_2}$, having Bounded System Dynamics ensures that we have a bound on the relative magnitude between state and control action. We will also use the notation $\kappa_A := \sigma_{\max}(A)/\sigma_{\min}(A)$ to denote the condition number of a matrix A and $\hat{\Sigma}_0 := \frac{1}{n} \sum_{k=1}^n \xi_0^{(k)} (\xi_0^{(k)})^T$ to denote the sample covariance of the first noise spike. We now provide our convergence result for the Double ML-type estimator in Theorem 5; the proof can be found in Appendix E.2.

Theorem 5. *Consider the dynamical system in (1) with $x_{-1} = u_{-1} = 0$ and with functions f, g, h, r defined in (3). Let the auditor observe n iid samples of $(x_0, u_0, x_1, u_1, x_2)$ and let DC be full column rank. Suppose that $\mathbb{E}\|\xi_1\|_2^2 = \sigma_1^2 d$, and the data generation model has ρ -Bounded System Dynamics. Let \mathcal{G} denote the event where $X_0 X_0^T$ is invertible. If $\mathbb{E}\left[\frac{\kappa_{\hat{\Sigma}_0}^2}{\lambda_{\min}(\hat{\Sigma}_0)}\right] \leq \tau_0$, then*

$$\frac{1}{pd} \mathbb{E}\left[\|\hat{B} - B\|_{\text{Fr}}^2 \mid \mathcal{G}\right] \leq \frac{\sigma_1^2 \rho^2 \kappa_{DC}^2 \tau_0}{n}.$$

In the case where $p = d$, if $\|\mathbb{E}[\hat{\Sigma}_0^{-1}]\|_{\text{op}} \leq \tau_1$ we have that

$$\frac{1}{d^2} \mathbb{E}\left[\|\hat{B} - B\|_{\text{Fr}}^2 \mid \mathcal{G}\right] \leq \frac{\sigma_1^2 \rho^2 \tau_1}{n}.$$

We note that rank condition on DC in this result is consistent with the rank condition from the identifiability result in the linear setting (Theorem 3); both results require DC to be full rank for identifiability. The τ_0, τ_1 conditions are a bit technical, but they essentially just require ξ_0 to be well behaved. We provide a simple Gaussian noise example to better illustrate.

Example 1 (Gaussian noise): Suppose ξ_0 and ξ_1 are drawn iid from $\mathcal{N}(0, \sigma_0^2 I_d)$ and $p = d$. We have $\mathbb{E}\|\xi_1\|_2^2 = \sigma_1^2 d$. For $n \geq d$, $\hat{\Sigma}_0$ is almost surely invertible. $(X_0 X_0^T / \sigma_0^2)^{-1}$ has an inverse Wishart distribution and thus, $\mathbb{E}[\hat{\Sigma}_0^{-1}] = \frac{n}{(n-d-1)\sigma_0^2} I_d$ for $n > d + 1$. Theorem 5 gives us

$$\frac{1}{d^2} \mathbb{E}\|\hat{B} - B\|_{\text{Fr}}^2 \leq \frac{\sigma_1^2 \rho^2}{(n-d-1)\sigma_0^2}.$$

◇

C Proofs of Section 2

C.1 Proof of Proposition 1

Recall that the do action alters the data generation model by deleting incoming edges into u_1 .

$$\begin{aligned} \mathbb{E}[x_2 | do(u_1 := u)] &= \mathbb{E}[f(x_1) + g(u) + \xi_t] \\ &= \int \mathbb{E}[f(z) + g(u) + \xi_t] p_{x_1}(z) dz \\ &= \int \mathbb{E}[f(x_1) + g(u_1) + \xi_t | u_1 = u, x_1 = z] p_{x_1}(z) dz \\ &= \int \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \end{aligned}$$

D Proofs of Section 3

Lemma D.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

D.1 Proof of Theorem 1

Overlap We will first show that (x_1, u_1) has joint 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 overlap 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 (ξ'_0, ξ_1) still has overlap. Using this modified notation, we have

$$\begin{aligned} x_0 &= \xi'_0 \\ u_0 &= h(\xi'_0) + c \\ x_1 &= f(\xi'_0) + g(h(\xi'_0) + c) + \xi_1 \\ u_1 &= h(x_1) + r(h(\xi'_0) + c) \end{aligned}$$

We first show that x_1 has overlap. Recall ξ_1 has positive density over \mathbb{R}^d . Because addition by a constant is an invertible, differentiable function, Lemma D.1 implies that $f(\xi'_0) + g(h(\xi'_0) + c) + \xi_1 | \xi'_0$ has positive density over \mathbb{R}^d . Since ξ'_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 overlap, it suffices to show that $u_1 | x_1$ has overlap over $\mathbb{R}^p \times \mathbb{R}^d$. 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 D.1 tells us that $u_1 | x_1$ would have positive density everywhere. Furthermore, because q_c satisfies Assumption 2, the conditions of Lemma D.1 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 overlap, the denominator is positive. Since ξ'_0 has overlap, $p_{\xi'_0}(a) > 0$ 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.

Identifiability Because p_{u_1, x_1} is positive everywhere, $\mathbb{E}[x_2 | u_1 = u, x_1 = z]$ is well defined. Additionally, because x_1 has density, $\int_z \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz$ is well defined as well. Finally because our model is admissible as stated in Proposition 1, $\mathbb{E}[x_2 | do(u_1 := u)] = \int_z \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz$. 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, $\mathbb{E}[x_2 | do(u_1 := u)]$ can be computed from the distribution of observations (x_1, u_1, x_2) . Because this quantity identifiable, the steerability of consumption $\mathcal{S}(u, u')$ is also identifiable for any $u, u' \in \mathbb{R}^d$.

D.2 Proof of Theorem 2

Proof Define a measurable function $\Delta : \mathbb{R}^p \rightarrow \mathbb{R}^d$ such that $\Delta \neq 0$. 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 (ξ_0, ξ_1) , let $(\hat{\xi}_0, \hat{\xi}_1)$ an identically distributed copy. Let x_1, x_0, u_0 be sampled according to the dynamics specified by (1) using the functions f, g, h and noise variables (ξ_0, ξ_1) . Let $\hat{x}_1, \hat{x}_0, \hat{u}_0$ be sampled according to the dynamics specified by (1) using the functions $\hat{f}, \hat{g}, \hat{h}$ in place of f, g, h and noise

variables $(\hat{\xi}_0, \hat{\xi}_1)$ in place of (ξ_0, ξ_1) . We see that

$$\begin{aligned}
x_0 &\stackrel{d}{=} \xi_0 \stackrel{d}{=} \hat{x}_0 \\
u_0 &\stackrel{d}{=} h(\xi_0) \stackrel{d}{=} \hat{h}(\xi_0) \stackrel{d}{=} \hat{u}_0 \\
x_1 &\stackrel{d}{=} f(\xi_0) + g(h(\xi_0)) + \xi_1 \\
&\stackrel{d}{=} f(\xi_0) + \Delta(h(\xi_0)) + g(h(\xi_0)) - \Delta(h(\xi_0)) + \xi_1 \\
&\stackrel{d}{=} \hat{f}(\xi_0) + \hat{g}(\hat{h}(\xi_0)) + \xi_1 \stackrel{d}{=} \hat{x}_1.
\end{aligned}$$

□

D.3 Proof of Theorem 3

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 D.2. *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. □

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

$$\Theta_x := \begin{bmatrix} A^T \\ B^T \end{bmatrix} \quad \Theta_u := \begin{bmatrix} C^T \\ D^T \end{bmatrix} \quad \xi_x := \begin{bmatrix} \xi_0^T \\ \vdots \\ 0 \end{bmatrix} \quad \xi_u := 0$$

Let $\hat{\Theta}_x, \hat{\Theta}_u$ be defined with respect to $\hat{A}, \hat{B}, \hat{C}, \hat{D}$. Let $\hat{\xi}_x \stackrel{d}{=} \xi_x$ and $\hat{\xi}_u \stackrel{d}{=} \xi_u$. Let (A, B, C, D) and ξ_x, ξ_u induce P_K and let $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$ and $\hat{\xi}_x, \hat{\xi}_u$ induce \hat{P}_K . Let $x := [x_0, \dots, x_K]^T$ and $u := [u_0, \dots, u_{K-1}]^T$ be observations from P_K and let \hat{x} and \hat{u} defined with hat variables be observations from \hat{P}_K . 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

$$\begin{aligned}
x - \xi_x &\stackrel{d}{=} Q_x \Theta_x & u - \xi_u &\stackrel{d}{=} Q_u \Theta_u \\
\hat{x} - \hat{\xi}_x &\stackrel{d}{=} \hat{Q}_x \hat{\Theta}_x & \hat{u} - \hat{\xi}_u &\stackrel{d}{=} \hat{Q}_u \hat{\Theta}_u
\end{aligned}$$

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

Lemma D.3. *If $x - \xi_x \stackrel{d}{=} Q_x \Theta_x$, then $x - \xi_x \stackrel{d}{=} Q_x \Theta_x \stackrel{d}{=} \hat{Q}_x \hat{\Theta}_x \stackrel{d}{=} \hat{x} - \hat{\xi}_x$. Similarly, if $u - \xi_u \stackrel{d}{=} Q_u \Theta_u$, then $u - \xi_u \stackrel{d}{=} Q_u \Theta_u \stackrel{d}{=} \hat{Q}_u \hat{\Theta}_u \stackrel{d}{=} \hat{u} - \hat{\xi}_u$.*

Proof Recall that the random variables in the vector z corresponds to nodes in the causal directed acyclic graph shown in Figure 3. Define $\sigma : \mathbb{Z} \rightarrow \mathbb{Z}$ such that $z_{\sigma(i)}$ is in sorted DAG order with

respect to the DAG in Figure 3 (i.e, the parents of $z_{\sigma(i)}$ have σ indices smaller than $\sigma(i)$ and its children have σ indices larger than $\sigma(i)$). We proceed inductively to show that $z_{\sigma(i)} \stackrel{d}{=} \hat{z}_{\sigma(i)}$.

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

Inductive step: suppose $z_{\sigma(j)} \stackrel{d}{=} \hat{z}_{\sigma(j)}$ jointly over all j . We know that $z_{\sigma(j+1)} \stackrel{d}{=} L(\{z_{\sigma(i)}\}_{i < j}, \xi, \hat{\Theta}_x, \hat{\Theta}_u)$ where L is linear in $\{z_{\sigma(i)}\}_{i < j}$, linear in ξ , linear with respect to $\hat{\Theta}_x$, and linear in $\hat{\Theta}_u$. By the inductive hypothesis we know that $z_{\sigma(j+1)} \stackrel{d}{=} L(\{z_{\sigma(i)}\}_{i < j}, \xi, \hat{\Theta}_x, \hat{\Theta}_u) \stackrel{d}{=} L(\{\hat{z}_{\sigma(i)}\}_{i < j}, \hat{\xi}, \hat{\Theta}_x, \hat{\Theta}_u) \stackrel{d}{=} \hat{z}_{\sigma(j+1)}$, 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 $\hat{Q} \stackrel{d}{=} Q$. This proves the desired result. \square

Now that we have established our supporting lemmas, we can now prove our desired result.

Claim 1: Necessity and sufficiency when $x_{-1} = u_{-1} = \xi_t = 0$ for $t \geq 1$. Let $X_t \in \mathbb{R}^{d,d}$ and $U_t \in \mathbb{R}^{p,d}$ be defined such that $x_t = X_t \varepsilon_0$ and $u_t = U_t \varepsilon_0$. Further define the following random matrix:

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

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

$$\begin{aligned} x - \xi_x &\stackrel{d}{=} \hat{x} - \xi_x \stackrel{d}{=} \hat{Q}_x \hat{\Theta}_x \stackrel{d}{=} Q_x \hat{\Theta}_x \\ u - \xi_u &\stackrel{d}{=} \hat{u} - \xi_u \stackrel{d}{=} \hat{Q}_u \hat{\Theta}_u \stackrel{d}{=} Q_u \hat{\Theta}_u. \end{aligned} \tag{4}$$

Finally, defining the fixed matrices $X := [X_0, \dots, X_K]^T$, $U := [U_0, \dots, U_{K-1}]^T$, and

$$Q_X := \begin{bmatrix} X_{-1}^T & U_{-1}^T \\ X_0^T & U_0^T \\ \vdots & \vdots \\ X_{K-1}^T & U_{K-1}^T \end{bmatrix},$$

we can rewrite (4) as

$$\begin{bmatrix} \xi_0^T \\ \vdots \\ \xi_0^T \end{bmatrix} \odot X = \begin{bmatrix} \xi_0^T \\ \vdots \\ \xi_0^T \end{bmatrix} \odot Q_X \hat{\Theta}_x. \tag{5}$$

Using Lemma D.2 we know the above equality holds if and only if the following holds

$$X = Q_X \hat{\Theta}_x. \tag{6}$$

Lemma D.3 tells us B is identifiable if and only if the entries of Θ_x corresponding to B is unique (6). 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 D.3 to show that $\hat{P}_K = P_K$; 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 and 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 $\hat{\Theta}_x = \Theta_x + v1^T$ is also a

solution to (6); 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

$$Q_X = \begin{bmatrix} I & C^T \\ X_1^T & (CX_1 + DU_0)^T \\ \vdots & \vdots \\ X_{K-1}^T & (CX_{K-1} + DU_{T-2})^T \end{bmatrix} = M_1 \begin{bmatrix} I & C^T \\ 0 & (DU_0)^T \\ \vdots & \vdots \\ 0 & (DU_{T-2})^T \end{bmatrix} = M_2 \begin{bmatrix} I & C^T \\ 0 & (DC)^T \\ \vdots & \vdots \\ 0 & (D^{K-1}C)^T \end{bmatrix}.$$

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

$$\tilde{Q}_X := \begin{bmatrix} I & C^T \\ 0 & (DC)^T \\ \vdots & \vdots \\ 0 & (D^{K-1}C)^T \end{bmatrix}.$$

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 $[DC, \dots, D^{K-1}C]$ being full row rank. Tracing back all the if and only if statements gives the result.

Claim 2: Sufficiency even when $x_{-1} \neq 0$ and $u_{-1} \neq 0$. In this setting, the proof for Claim 1 holds up to Equation (5). Equation (5) changes to the following

$$\begin{bmatrix} \xi_0^T \\ \vdots \\ \xi_0^T \end{bmatrix} \odot X + w_1(x_{-1}, u_{-1}, \xi_1, \xi_2) = \begin{bmatrix} \xi_0^T \\ \vdots \\ \xi_0^T \end{bmatrix} \odot Q_X \hat{\Theta}_x + w_2(x_{-1}, u_{-1}, \xi_1, \xi_2).$$

By Lemma D.2, we know that these equalities hold if and only if Equation (6) holds, $w_1(x_{-1}, u_{-1}, \xi_1, \xi_2) = w_2(x_{-1}, u_{-1}, \xi_1, \xi_2)$ holds. $\text{null}(Q_X) \perp \text{span}(\mathcal{S})$ 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 (6). The rest of the argument in Claim 1 follows identically.

E Proofs of Appendix B

E.1 Proof of Theorem 4

The proof proceeds by bounding each of the following terms.

$$\begin{aligned} & \left\| \sum_{x \in \mathcal{R}_X} Y^n(u, x) Z^n(x) - \mathbb{E}[x_2 | do(u_1 := u)] \right\| \leq \\ & \left\| \sum_{x \in \mathcal{R}_X} Y^n(u, x) Z^n(x) - \sum_{x \in \mathcal{R}_X} Y^n(u, x) P(x_1 \in \mathcal{N}_X(x)) \right\| \\ & + \left\| \sum_{x \in \mathcal{R}_X} Y^n(u, x) P(x_1 \in \mathcal{N}_X(x)) - \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)] P(x_1 \in \mathcal{N}_X(x)) \right\| \\ & + \left\| \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)] P(x_1 \in \mathcal{N}_X(x)) - \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 = u, x_1 = x] P(x_1 \in \mathcal{N}_X(x)) \right\| \\ & + \left\| \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 = u, x_1 = x] P(x_1 \in \mathcal{N}_X(x)) - \mathbb{E}[x_2 | do(u_1 := u)] \right\| \end{aligned}$$

E.1.1 Supporting lemmas

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

Lemma E.1. *Let the conditions of Theorem 1 hold and let λ denote the Lebesgue measure for \mathbb{R}^{d+p} . Let $\mathcal{N}_{\mathcal{X}}$ be a cover of \mathcal{X} and $\mathcal{N}_{\mathcal{U}}$ be a cover of \mathcal{U} . If for all $A \in \mathcal{N}_{\mathcal{X}}$ and $B \in \mathcal{N}_{\mathcal{U}}$, $\lambda(A \times B) > 0$, then $P(x_1 \in A, u_1 \in B) > 0$ for all $A \in \mathcal{N}_{\mathcal{X}}, B \in \mathcal{N}_{\mathcal{U}}$.*

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. \square

Lemma E.2. *Let $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} \subset \mathbb{R}^p$, and let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a L -Lipschitz function. Let \mathcal{N} be any ε -cover of \mathcal{X} for $\varepsilon > 0$. Then for all $U \in \mathcal{N}$, for all $x, y \in U$, $\|f(x) - f(y)\| \leq L\varepsilon$.*

Proof Follows directly from definitions of cover and Lipschitz Continuity. \square

Lemma E.3. *Consider the data generation model of (1). Let Assumption 3 hold. Let x_1 have overlap. Let $\mathcal{R}_{\mathcal{X}}$ denote the ε -cover representatives of $\mathcal{X} \subset \mathbb{R}^d$. Then,*

$$\begin{aligned} & \left\| \mathbb{E}[x_2 | do(u_1 := u)] - \sum_{r \in \mathcal{R}_{\mathcal{X}}} \mathbb{E}[x_2 | u_1 = u, x_1 = r] P(x_1 \in \mathcal{N}_{\mathcal{X}}(r)) \right\| \\ & \leq L\varepsilon + \left\| \int_{\mathbb{R}^d \setminus \mathcal{X}} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \end{aligned}$$

Proof

Let $B := \mathbb{R}^d \setminus \cup_{A \in \mathcal{N}_{\mathcal{X}}} A$ denote the set of points not covered by the epsilon cover.

$$\begin{aligned} & \left\| \mathbb{E}[x_2 | do(u_1 := u)] - \sum_{r \in \mathcal{R}_{\mathcal{X}}} \mathbb{E}[x_2 | u_1 = u, x_1 = r] P(x_1 \in \mathcal{N}_{\mathcal{X}}(r)) \right\| \\ & \leq \left\| \sum_{r \in \mathcal{R}_{\mathcal{X}}} \int_{\mathcal{N}_{\mathcal{X}}(r)} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz - \sum_{r \in \mathcal{R}_{\mathcal{X}}} \mathbb{E}[x_2 | u_1 = u, x_1 = r] P(x_1 \in \mathcal{N}_{\mathcal{X}}(r)) \right\| \\ & \quad + \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \\ & \leq \sum_{r \in \mathcal{R}_{\mathcal{X}}} \left\| \int_{\mathcal{N}_{\mathcal{X}}(r)} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz - \mathbb{E}[x_2 | u_1 = u, x_1 = r] P(x_1 \in \mathcal{N}_{\mathcal{X}}(r)) \right\| \\ & \quad + \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \\ & \leq \sum_{r \in \mathcal{R}_{\mathcal{X}}} \int_{\mathcal{N}_{\mathcal{X}}(r)} \|\mathbb{E}[x_2 | u_1 = u, x_1 = z] - \mathbb{E}[x_2 | u_1 = u, x_1 = r]\| p_{x_1}(z) dz \\ & \quad + \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \\ & \leq L\varepsilon + \left\| \int_{\mathbb{R}^d \setminus \mathcal{X}} \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \end{aligned}$$

The first and second inequality is from triangle inequality. The third comes from Jensen's inequality. The fourth inequality comes Assumption 3, Lemma E.2, and the fact that $\mathcal{X} \subset \cup_{A \in \mathcal{N}_{\mathcal{X}}} A$. \square

Lemma E.3 tells us that it suffices to create an estimator that estimates $\sum_{r \in \mathcal{R}_{\mathcal{X}}} \mathbb{E}[x_2 | u_1 = u, x_1 = r] P(x_1 \in \mathcal{N}_{\mathcal{X}}(r))$ —supposing that \mathcal{X} is a good approximation of \mathbb{R}^d with respect to x_1 .

Lemma E.4. *Consider the data generating process from (1). Let x_1, u_1 have overlap. Let Assumption 4 hold, then*

$$\|\mathbb{E}[Y^n(u, x)] - \mathbb{E}[x_2 | u_1 = u, x_1 = x]\| \leq \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \mathbb{E}[\|x_2\| | u_1 = u, x_1 = x].$$

Proof Fix any $u \in \mathcal{U}, x \in \mathcal{X}$. Let $Z := (x_1, u_1)$, $z := (x, u)$, and $A := \mathcal{N}_{\mathcal{X}}(x) \times \mathcal{N}_{\mathcal{U}}(u)$. Observe that $\mathbb{E}[Y^n(u, x)] = \mathbb{E}[x_2 \mid Z \in A]$. Note that these conditional expectations exist because Z has overlap and by construction A has positive Lebesgue measure.

$$\begin{aligned}
\|\mathbb{E}[x_2 \mid Z \in A] - \mathbb{E}[x_2 \mid Z = z]\| &= \left\| \int_{\mathbb{R}^d} x \left[\frac{P(Z \in A \mid x_2 = x)}{P(Z \in A)} - \frac{p(Z = z \mid x_2 = x)}{p(Z = z)} \right] p(x_2 = x) dx \right\| \\
&\leq \int_{\mathbb{R}^d} \|x\| \left| \frac{P(Z \in A \mid x_2 = x)}{P(Z \in A)} - \frac{p(Z = z \mid x_2 = x)}{p(Z = z)} \right| p(x_2 = x) dx \\
&= \int_{\mathbb{R}^d} \|x\| \left| \frac{\int_A p(Z = a \mid x_2 = x) da}{\int_A p(Z = a) da} - \frac{p(Z = z \mid x_2 = x)}{p(Z = z)} \right| p(x_2 = x) dx \\
&\leq \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \int_{\mathbb{R}^d} \|x\| \frac{p(Z = z \mid x_2 = x)}{p(Z = z)} p(x_2 = x) dx \\
&= \frac{2\eta(\varepsilon)}{1 - \eta(\varepsilon)} \mathbb{E}[\|x_2\| \mid Z = z].
\end{aligned}$$

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 ε . \square

E.1.2 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 τ^2 -subgaussian random variable X .

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

$Z^n(x)$ is $\frac{1}{4n}$ subgaussian. This means we need $n = \frac{1}{2\delta^2 P(x_1 \in \mathcal{N}_{\mathcal{X}}(x))^2} \log(4|\mathcal{R}_{\mathcal{X}}|/\rho)$ samples to get $Z^n(x)$ within error of $\delta P(x_1 \in \mathcal{N}_{\mathcal{X}}(x))$ of $P(x_1 \in \mathcal{N}_{\mathcal{X}}(x))$ with probability $\rho/(2|\mathcal{R}_{\mathcal{X}}|)$. Using union bound, we have that with probability with at least $1 - \rho/2$,

$$\begin{aligned}
&\left\| \sum_{x \in \mathcal{R}_{\mathcal{X}}} Y^n(u, x) Z^n(x) - \sum_{x \in \mathcal{R}_{\mathcal{X}}} Y^n(u, x) P(x_1 \in \mathcal{N}_{\mathcal{X}}(x)) \right\| \leq \sum_{x \in \mathcal{R}_{\mathcal{X}}} \|Y^n(u, x)\| |Z^n(x) - P(x_1 \in \mathcal{N}_{\mathcal{X}}(x))| \\
&\leq \delta \sum_{x \in \mathcal{R}_{\mathcal{X}}} \|Y^n(u, x)\| P(x_1 \in \mathcal{N}_{\mathcal{X}}(x)) \\
&\leq \delta \sum_{x \in \mathcal{R}_{\mathcal{X}}} \|Y^n(u, x) - \mathbb{E}[x_2 \mid u_1 \in \mathcal{N}_{\mathcal{U}}(u), x_1 \in \mathcal{N}_{\mathcal{X}}(x)]\| P(x_1 \in \mathcal{N}_{\mathcal{X}}(x)) \\
&\quad + \delta \sum_{x \in \mathcal{R}_{\mathcal{X}}} \|\mathbb{E}[x_2 \mid u_1 \in \mathcal{N}_{\mathcal{U}}(u), x_1 \in \mathcal{N}_{\mathcal{X}}(x)]\| P(x_1 \in \mathcal{N}_{\mathcal{X}}(x)) \\
&\leq \delta\gamma + \delta \sum_{x \in \mathcal{R}_{\mathcal{X}}} \|\mathbb{E}[x_2 \mid u_1 \in \mathcal{N}_{\mathcal{U}}(u), x_1 \in \mathcal{N}_{\mathcal{X}}(x)]\| P(x_1 \in \mathcal{N}_{\mathcal{X}}(x)) \\
&\leq \delta\gamma + \delta D + \delta \mathbb{E}[\|g(u_1)\| \mid u_1 \in \mathcal{N}_{\mathcal{U}}(u), x_1 \in \mathcal{N}_{\mathcal{X}}(x)] \\
&\leq \delta\gamma + 2\delta D
\end{aligned}$$

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: $Y^n(u, x)$ 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

$$P(\|Y^n(u, x) - \mathbb{E}[Y^n(u, x)]\| \geq t) \leq 2 \exp\left(-\frac{t^2 n_{u,x}}{2d\sigma^2}\right)$$

This means we need $n_{u,x} = \frac{2d\sigma^2}{\gamma^2} \log(4|\mathcal{R}_X|/\rho)$ samples to get $Y^n(u,x)$ with error γ of $\mathbb{E}[Y^n(u,x)]$ with probability $\rho/(2|\mathcal{R}_X|)$. Moreover, because the conditions of Lemma E.1 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 $1 - \rho/2$,

$$\begin{aligned} & \left\| \sum_{x \in \mathcal{R}_X} Y^n(u,x) P(x_1 \in \mathcal{N}_X(x)) - \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)] P(x_1 \in \mathcal{N}_X(x)) \right\| \\ & \leq \sum_{x \in \mathcal{R}_X} \|Y^n(u,x) - \mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)]\| P(x_1 \in \mathcal{N}_X(x)) \\ & \leq \gamma \end{aligned}$$

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

Third term:

$$\begin{aligned} & \left\| \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)] P(x_1 \in \mathcal{N}_X(x)) - \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 = u, x_1 = x] P(x_1 \in \mathcal{N}_X(x)) \right\| \\ & \leq \sum_{x \in \mathcal{R}_X} \|\mathbb{E}[x_2 | u_1 \in \mathcal{N}_U(u), x_1 \in \mathcal{N}_X(x)] - \mathbb{E}[x_2 | u_1 = u, x_1 = x]\| P(x_1 \in \mathcal{N}_X(x)) \\ & \leq \frac{2\eta}{1-\eta} \sum_{x \in \mathcal{R}_X} \mathbb{E}[\|x_2\| | u_1 = u, x_1 = x] P(x_1 \in \mathcal{N}_X(x)) \\ & \leq \frac{2\eta}{1-\eta} (2D + c_1) \end{aligned}$$

The first inequality comes from Jensen's inequality. The second comes from Lemma E.4. The third inequality comes from triangle inequality.

Fourth term:

$$\begin{aligned} & \left\| \sum_{x \in \mathcal{R}_X} \mathbb{E}[x_2 | u_1 = u, x_1 = x] P(x_1 \in \mathcal{N}_X(x)) - \mathbb{E}[x_2 | do(u_1 := u)] \right\| \leq L\varepsilon + \left\| \int_B \mathbb{E}[x_2 | u_1 = u, x_1 = z] p_{x_1}(z) dz \right\| \\ & \leq L\varepsilon + \mathbb{E}[\|f(x_1)\| \mathbf{1}\{x_1 \in B\}] + P_{x_1}(B)D. \end{aligned}$$

The first inequality comes from Lemma E.3. The second inequality comes from $\mathbb{E}\xi_2 = 0$, triangle inequality, Jensen's inequality, and the definition of D .

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

E.2 Proof of Theorem 5

We first introduce a helpful supporting lemma.

Lemma E.5. *Suppose n samples are drawn iid from P_2 . If $X_0 X_0^T$ is invertible, then $\widehat{C} = C$ and $\widehat{H} = A + BC + E_1 X_0^T (X_0 X_0^T)^{-1}$. If $X_0 X_0^T$ is invertible and $DC X_0 X_0^T C^T D^T$ is invertible, then $\widehat{B} = B - E_1 X_0^T (X_0 X_0^T)^{-1} X_1 (DC X_0)^T (DC X_0 X_0^T C^T D^T)^{-1}$.*

Proof Substituting $U_0 = CX_0$ and $X_1 = (A + BC)X_0 + E_1$ 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 $\widehat{H} = A + BC + E_1 X_0^T (X_0 X_0^T)^{-1}$ by the first result. We observe that $U_1 = CX_1 + DU_0 = CX_1 + DCX_0$ to get that $\widehat{B} = (X_2 - \widehat{H}X_1)(DCX_0)^T (DCX_0 X_0^T C^T D^T)^{-1}$. Then we use the fact that subtracting BCX_1 from both sides of the relationship $X_2 - AX_1 = BU_1$ gives us that

$X_2 - (A + BC)X_1 = B(U_1 - CX_1)$. Using our invertability assumptions, this gives us $\widehat{B} = B - E_1 X_0^T (X_0 X_0^T)^{-1} X_1 (DC X_0)^T (DC X_0 X_0^T C^T D^T)^{-1}$. \square

With this, we can analyze the quantities of interest. Let $\hat{\Sigma}_0 = \frac{1}{n} X_0 X_0^T$. Let $Q := DC \hat{\Sigma}_0 C^T D^T$.

$$\begin{aligned} \mathbb{E} \left[\|\widehat{B} - B\|_{\text{Fr}}^2 \mid \mathcal{G} \right] &= \frac{1}{n^2} \text{tr}(\mathbb{E}[Q^{-1} DC X_0 X_1^T (X_0 X_0^T)^{-1} X_0 E_1^T E_1 X_0^T (X_0 X_0^T)^{-1} X_1 X_0^T C^T D^T Q^{-1}]) \\ &= \frac{\sigma_1^2 d}{n} \text{tr}(\mathbb{E}[Q^{-1} DC \hat{\Sigma}_0 (A + BC)^T \hat{\Sigma}_0^{-1} (A + BC) \hat{\Sigma}_0 C^T D^T Q^{-1}]) \\ &\leq \frac{\sigma_1^2 p d}{n} \kappa_{DC}^2 \left(\frac{\|A + BC\|_{\text{op}}}{\sigma_{\min}(DC)} \right)^2 \mathbb{E} \left[\frac{\kappa_{\hat{\Sigma}_0}^2}{\lambda_{\min}(\hat{\Sigma}_0)} \right]. \end{aligned}$$

Rearranging and using the definition of τ_0 gives the result.

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

$$\begin{aligned} \mathbb{E} \left[\|\widehat{B} - B\|_{\text{Fr}}^2 \mid \mathcal{G} \right] &= \frac{\sigma_1^2 d}{n} \text{tr}[(C^T D^T)^{-1} (A + BC)^T \mathbb{E}[\hat{\Sigma}_0^{-1}] (A + BC) (DC)^{-1}] \\ &\leq \frac{\sigma_1^2 d^2}{n} \left(\frac{\|A + BC\|_{\text{op}}}{\lambda_{\min}(DC)} \right)^2 \|\mathbb{E}[\hat{\Sigma}_0^{-1}]\|_{\text{op}}. \end{aligned}$$

Rearranging and using the definition of τ_1 gives the result.