

Causal Inference for Asset Pricing*

Valentin Haddad

Zhiguo He

Paul Huebner

Péter Kondor

Erik Loualiche

May 2026

Abstract

This paper provides a guide for using causal inference with asset prices and quantities. Our framework revolves around an elementary assumption about portfolio demand: homogeneous substitution conditional on observables. Under this assumption, standard cross-sectional instrumental variables or difference-in-differences regressions identify the relative demand elasticity between assets with the same observables, the difference between own-price and cross-price elasticity. However, we uncover a missing coefficient problem: cross-sectional estimators mechanically absorb substitution patterns across assets. Recovering substitution is essential to answer many natural counterfactual questions, and requires analyzing the response of portfolios to exogenous time-series variation. The same principles apply to the estimation of multipliers measuring the price impact of supply or demand shocks. Our assumption maps to familiar restrictions on covariance matrices in classical asset pricing models, encompasses demand models such as logit, and accommodates rich substitution patterns even outside of these models. We discuss how to design experiments satisfying this condition and offer diagnostics to validate it.

*Haddad: Anderson School of Management, UCLA and NBER, valentin.haddad@anderson.ucla.edu; He: Graduate School of Business, Stanford University and NBER, hezgh@stanford.edu; Huebner: Stockholm School of Economics, paul.huebner@hhs.se; Kondor: London School of Economics and CEPR, p.kondor@lse.ac.uk; and Loualiche: Carlson School of Management, University of Minnesota, eloualic@umn.edu. We thank Mikhail Chernov, Carter Davis, Zhiyu Fu, William Fuchs, Satoshi Fukuda, Xavier Gabaix, Sergei Glebkin, Paul Goldsmith-Pinkham, Elise Gourier, Kristy Jansen, Ralph Koijen, Jane Li, Lira Mota, Daniel Neuhann, Anna Pavlova, Aleksandra Rzeźnik, Milena Wittwer, Hanbin Yang, Motohiro Yogo, and seminar and conference participants at Minnesota Carlson, Chicago Booth Alumni Insights, UCLA Anderson, Stockholm School of Economics, Demand in Asset Markets Working Group, USC Macro-Finance Workshop, Adam Smith Workshop, Spring Finance Workshop, SFS Cavalcade, FIRS, CFF Asset Pricing and Machine Learning, LBS Summer Symposium, SoFiE Pre-Conference, Gerzensee, NBER Summer Institute, USC, SAFE Asset Pricing Workshop, Chicago Booth, Harvard, Baruch, Yale, University of Zurich, Duke, CEPR Paris Symposium, Imperial College, AFA Annual Meeting, MIT, Toronto Rotman, HEC-McGill Winter Finance, University of Bonn, and Boston College for helpful comments and suggestions.

Introduction

Causal inference methods that leverage plausibly exogenous sources of variation have become essential tools in empirical economics (Angrist and Pischke, 2009). Recently, these methods have gained traction in asset pricing to better understand the demand for financial assets, through both specific experiments like index inclusions or central banks’ asset purchases (Shleifer, 1986; Chang et al., 2014; Krishnamurthy and Vissing-Jorgensen, 2011), and as building blocks for demand systems (Kojen and Yogo, 2019; Haddad et al., 2024c). However, these approaches differ sharply from traditional empirical methods in asset pricing (see, e.g., Cochrane, 2005; Campbell, 2017), which instead prioritize tests of equilibrium relationships such as Euler equations or the CAPM.

We provide a framework for using causal inference in the asset pricing context. We put forward elementary conditions that allow the use of the standard toolbox of causal inference while entertaining a rich set of finance models. Under these conditions, we fully characterize what sources of variation and estimation procedures identify portfolio demand and its equilibrium impact.

Our focus is on how portfolio choice responds to asset prices, or conversely, how asset prices respond to an exogenous shift in demand. Natural experiments offer a window into these questions. If the Fed purchases quasi-randomly some bonds but not others, prices respond, and one can look at how various investors adjust their portfolios. One can also look at the surprise price change following the announcement of a broad purchase program by the Fed, or an intervention such as operation twist, buying more long-term than short-term bonds.¹ But what exactly is learned about investor demand from analyzing these different experiments? Alternatively, one can start from a counterfactual question, such as “what would be the effect of an increase in investor preference for firms with good ESG performance?” and ask which evidence is necessary to inform this counterfactual.

In all generality, all these aspects of demands are necessarily intertwined. Portfolio choice (Markowitz, 1952) — and really basic theory of demand (Debreu, 1959) — implies that the demand for each asset responds not only to the price of one asset, but to the price of every asset. Formally, with N assets, the slope of the demand curve, the demand elasticity, is an $N \times N$ matrix: $\Delta D = \mathcal{E} \Delta P$. Furthermore, prices are connected in equilibrium. In the example above, even if the Fed does not buy a bond, the price of this bond might respond to purchases of its substitutes. Without any restrictions, the only way to make progress is to have N different exogenous sources of variation in prices, a tall order.

¹Selgrad (2023) analyzes variation in which specific bond the Fed purchases in the implementation of quantitative easing programs. Krishnamurthy and Vissing-Jorgensen (2011) focus on the response to the announcement of broad QE programs.

Imposing further economic restrictions on the elasticity matrix is the natural way forward, but the researcher must be careful not to oversimplify. The most extreme simplification — assuming the demand for an asset depends only on its own price — makes the econometrics easy (standard SUTVA), but it fails because it ignores that investors manage portfolios. In classic finance, assets are substitutes if they carry the same risks, effectively acting as bundles like high-risk technology stocks versus low-risk manufacturing stocks. This implies that demand elasticity has two distinct components that must be modeled simultaneously. There is the “stock picker” logic, where if Ford gets expensive relative to GM, you swap them because they have the same risk profile. But there is also the “factor management” logic, where if the value premium decreases and stocks with a high book-to-market ratio get expensive, you might rotate your portfolio out of them and into other factors. We show that leading structural models like [Kojien and Yogo \(2019\)](#) and [Gabaix and Kojien \(2021\)](#) prioritize tractability in ways that limit their ability to capture this dimension. They capture the stock-picking behavior well, but implicitly tether how investors can rotate between broad risk factors. We provide a concrete set of examples illustrating this limitation: these models cannot explain standard portfolio adjustments in response to price changes in high-beta versus low-beta assets. Instead, we enrich the family of possible foundations for portfolio choice while maintaining tractability in estimation.

Our framework relies on a transparent restriction: *homogeneous substitution conditional on observables*. This condition specifies that investors’ demand for any two assets with identical observables responds identically to a price change in any other asset. For instance, if Ford and GM carry the same risk and capital charge, an investor’s demand for them must respond identically to a price shock in First Solar. We show that this condition holds for many foundations: whether investors are managing risk (using factor loadings), navigating balance sheet constraints (using capital charges), or targeting green assets (using ESG scores). This flexibility allows us to capture the diverse motives driving asset demand while keeping the econometrics tractable. Beyond theoretical justifications, we provide a variety of empirical approaches to assess the plausibility of this assumption.

Our assumption simplifies the daunting $N \times N$ elasticity matrix into two naturally interpretable components: *relative elasticity* and *substitution*. The relative elasticity captures the micro response of the stock picker logic: if Ford becomes expensive relative to GM, how much do you swap them? In our baseline, the relative elasticity is a single number; we also show how to make it heterogeneous across assets. The substitution matrix captures reallocation across stocks driven by observables: if assets with a high beta and low capital charge become expensive, how much do you rotate into other assets? This matrix is driven by $K \times K$ parameters, where K is the number of observables. Identifying these two components requires

two distinct empirical strategies.

We formally show that standard cross-sectional identification methods, like instrumental variables or difference-in-differences, successfully recover the relative elasticity. By comparing a treated asset (Ford) to a control asset (GM) with the same observables, the shared substitution effects — e.g., that they both respond to a shock in Tech stocks — wash out. More generally, we prove that controlling for observables in cross-sectional regressions of changes in demand on changes in prices is sufficient to absorb these substitution effects. The researcher, nevertheless, still needs an instrument, since relative prices remain endogenous to demand shifts (e.g., a shock to preference for Ford vs. GM). The resulting regression coefficient consistently estimates the relative elasticity, the difference between own- and cross-price elasticities. Crucially, this means that the seemingly naive cross-sectional regression that appears to ignore substitution works, but it only answers a specific question: how demand responds to *relative* price changes, not absolute ones nor those affecting multiple assets.²

However, the cross-section is blind to substitution. Substitution is a rotation across observables — selling high-beta stocks to buy low-beta ones. Cross-sectional regressions work precisely by controlling for these observables. They mechanically absorb these rotations into the slope coefficients, creating a “missing coefficient” problem: the coefficient on beta cannot distinguish between buying high-beta stocks because their price dropped (substitution) and buying them because you like beta more (a demand shift). This issue echoes the “missing intercept” problem in macroeconomics: general equilibrium effects are absorbed by time fixed effects. But the trap here is subtler: substitution hides in the slopes, not just the intercept. To isolate substitution, the researcher must look to the time series along a few dimensions. Practically, the researcher aggregates data into portfolios along the observables that drive substitution, such as sorting stocks based on their beta. Then, they need to identify exogenous time-series shocks that move the price of these portfolios. By observing how investors tilt their portfolios across these characteristics — swapping high-beta for low-beta stocks — in response to these shocks, the researcher recovers the substitution.

Our framework thus unifies two distinct empirical traditions: cross-sectional instruments (like index inclusions) to measure relative elasticity, and time-series shocks (like aggregate Fed purchases) to measure substitution. By combining these two approaches, we can fully characterize the elasticity matrix and answer counterfactuals that require both logics simultaneously, such as quantifying the equilibrium impact of “operation twist” or a broad ESG mandate. Notably, there is no neat separation of micro and macro. Instead, substitution driven by observables constitutes a “meso” layer: this force is essential to answer questions

²As an aside, the same reasoning rationalizes the demand regression introduced in [Koijen and Yogo \(2019\)](#). However, unlike our partial identification result, parametric restrictions constrain substitution so much in their model that the researcher recovers the full elasticity matrix from this relative estimate alone.

about the cross-section, but requires identification from broad time series.

Our framework applies equally to the measurement of price impacts, or multipliers—the effect of exogenous supply or demand shocks on asset prices. Questions like “how do Fed asset purchases affect bond prices?” reverse the variables of interest: instead of asking how prices move quantities, we ask how quantities move prices. We show that because the multiplier matrix is the inverse of the elasticity matrix of aggregate demand, our assumptions of homogeneous substitution and constant relative elasticity translate directly to this setting. Consequently, the same decomposition holds: cross-sectional regressions identify the relative price impact, while time-series regressions on portfolios recover the aggregate and meso-level multipliers.

We assess the robustness of our approach to a number of potential concerns. Minor deviations from our assumption have a small effect on estimates. We outline how to entertain richer heterogeneity at the cost of stronger assumptions and additional sources of variation. We confirm that our approach is compatible with no-arbitrage conditions and show how to assess empirically whether they play a role. Naturally, settings characterized by deviations from no-arbitrage relations also lend themselves particularly well to our framework.³

In the final part of the paper, we illustrate practically how to apply our framework by studying price impact in the corporate bond market. Consistent with [Chaudhary et al. \(2022\)](#), we find a relative price impact indistinguishable from zero. Going beyond this micro estimate, we recover the spillover structure and show that it is well captured by substitution linear in two observables, duration and credit rating. We use these estimates to evaluate counterfactual designs of central bank asset purchase programs and find that program design — broad versus targeted, as in the SMCCF — materially changes both aggregate price impact and which bonds benefit, including through spillovers to non-purchased bonds.

Together, our results arm the researcher with a guide to use and interpret evidence from natural experiments on price and quantities in asset markets as well as to understand their limits. We spell out: a) which assumptions one needs to defend to use causal inference techniques, b) diagnostics to assess the plausibility of these assumptions, c) which technique and source of variation are appropriate for different economic questions, and d) how to interpret causal estimates.

Related Literature. A long tradition in finance uses plausibly exogenous sources of variation to understand portfolio decisions and the price impact of shifts in demand. Prominent examples include the effect of index inclusion ([Shleifer, 1986](#); [Harris and Gurel, 1986](#); [Chang](#)

³Examples include covered-interest-parity deviations ([Du et al., 2018](#)), the CDS-bond basis ([Bai and Collin-Dufresne, 2019](#)), violations of put-call parity ([Van Binsbergen et al., 2022](#)), or pledgeability premium of corporate bonds ([Chen et al., 2023](#)).

et al., 2014; Pavlova and Sikorskaya, 2022; Greenwood and Sammon, 2024), institutional ownership and fund flows (Gompers and Metrick, 2001; Coval and Stafford, 2007; Lou, 2012; Ben-David et al., 2022; Hartzmark and Solomon, 2022), central bank asset purchases (Krishnamurthy and Vissing-Jorgensen, 2011; Selgrad, 2023; Haddad et al., 2021, 2025), financial constraints (Du et al., 2018; Greenwood and Vissing-Jorgensen, 2018; Haddad and Muir, 2021; Chen et al., 2023), or exchange rates (Evans and Lyons, 2002; Froot and Ramadorai, 2005). This work often incorporates thorough analysis of exogeneity, in particular in the wake of the “credibility revolution” (e.g., Angrist and Pischke, 2009). However, this literature is often more scant in considering a central feature of asset pricing theory, substitution across assets, and whether it affects the validity of inference and the interpretation of estimates. Our framework provides a simple bridge between classical discussions of causal inference and the role of substitution.

Another approach fully specifies and estimates models of portfolio demand and their equilibrium implications. In a seminal article, Kojien and Yogo (2019) derive and estimate a logit model of portfolio choice in the stock market. Subsequent work uses this model either structurally, or as a semi-structural simplification to introduce other mechanisms (e.g. Haddad et al. (2024c)). Applications include quantifying the impact of the rise of passive investing, preferences for sustainable assets (Kojien et al., 2023; Van der Beck, 2021), or the transmission of monetary policy (Lu and Wu, 2023), and have found echo in other settings: the stock market overall Gabaix and Kojien (2021), corporate bonds (Bretscher et al., 2022), treasuries (Jansen et al., 2024; Fang, 2023; Fang and Xiao, 2024), or exchange rates (Kojien and Yogo, 2024; Jiang et al., 2024).

We build on some of the insights from estimation inside of these models. Some important ideas are controlling for common exposures (Kojien and Yogo, 2019), the distinction between micro and macro elasticity (Gabaix and Kojien, 2021; Li and Lin, 2022), heterogeneous substitution (Chaudhary et al., 2022; Aghaee, 2024), substitution along factors (An et al., 2024; An and Huber, 2025; Peng and Wang, 2023), and accounting flexibly for spillovers (Fuchs et al., 2025). Naturally, our simple conditions cannot cover every model; we leave aside considerations of strategic responses (Haddad et al., 2024c), dynamics (Greenwood et al., 2018; Gabaix and Kojien, 2021; Huebner, 2024; He et al., 2025), state-contingent demand shocks (Haddad et al., 2025), intermediary distress (He et al., 2022), or bidding in auctions (Allen et al., 2018). In this context, the contribution of our framework is twofold: it not only provides a unifying formalism to discuss identification across models but also allows discussion of what can be learned from the data before espousing a specific model.

Finally, the role of spillovers is not limited to asset pricing and has been recognized in many other contexts. The industrial organization literature often relaxes the assumption of

independence of irrelevant alternative (IIA) and includes heterogeneous substitution in discrete choice models, such as [Berry et al. \(1995\)](#). While without observables, our assumptions would be closely related to IIA, including the observables entertains heterogeneous substitution. Our setting of portfolio choice in line with finance theory does so without introducing nonlinearities, lending itself to using linear regressions. [Berg et al. \(2021\)](#) discuss spillovers in corporate finance. In macroeconomics, a key concern is the missing intercept problem due to general equilibrium effects, with some recent contributions such as [Chodorow-Reich et al. \(2021\)](#), [Guren et al. \(2021\)](#), [Huber \(2023\)](#), and [Wolf \(2023\)](#).

1 The Challenge of Causal Inference in Asset Pricing

We introduce our estimation target, the demand elasticity matrix, and its central role in answering questions about asset prices. Estimation faces two hurdles. First, rich substitution patterns and price endogeneity make estimation impossible without imposing restrictions. Second, restrictions must be flexible enough to capture well-known asset pricing patterns, something leading structural models fail to do. Our framework in Section 2 solves both problems by relying on simple and flexible assumptions.

1.1 Demand and Demand Elasticity

We focus on a generic setting for identifying the demand for financial assets, with the goal to understand the demand elasticity — that is, how an investor adjusts their portfolio in response to prices.

Demand elasticity: the concept. When an investor maximizes their utility given prices—or simply responds to prices through whatever decision process they follow—this choice characterizes demand as a function of prices.⁴ The demand elasticity is the derivative of this function with respect to prices:

$$\boldsymbol{\mathcal{E}} \equiv \frac{\partial D}{\partial P}. \tag{1}$$

Importantly, $\boldsymbol{\mathcal{E}}$ is a matrix of size $N \times N$, where N is the number of assets. The diagonal elements $\partial D_i / \partial P_i$ measure the own-price elasticities: how demand for Apple responds to the price of Apple. The off-diagonal elements $\partial D_i / \partial P_j$ capture cross-price elasticities: how demand for Apple responds to the price of Nvidia.

⁴[Becker \(1962\)](#) explains the generality of demand as a function of prices.

In the finance context, these cross-price elasticities are essential to a key insight dating back to [Markowitz \(1952\)](#): assets are not distinct goods but instead alternative means of saving with different risk and reward. Investors choose portfolios optimally combining these assets. The most standard example of this approach is mean-variance optimization: an investor chooses their portfolio to maximize $\mathbf{E}(W) - \frac{\gamma}{2} \text{var}(W)$ where W is their future wealth, and γ measures their absolute risk aversion. If assets have constant mean payoffs denoted by vector $\bar{\Pi}$ and covariance matrix denoted by Σ , the vector of demand is:

$$D = \frac{1}{\gamma} \Sigma^{-1} (\bar{\Pi} - P), \quad (2)$$

with the matrix of elasticity \mathcal{E} determined by risk aversion and the covariance between assets: $\mathcal{E} = -\gamma^{-1} \Sigma^{-1}$.⁵ When assets are correlated with each other, they become close substitutes, and their demands respond to each other's prices.

Modern finance research acknowledges many deviations from this simple setting: investors have different beliefs and various cognitive limitations, institutions face many regulations and constraints that influence their portfolio decisions. Still the basic idea of portfolio choice as opposed to asset choice remains. Any model of asset demand implies a matrix of elasticities \mathcal{E} . Conversely, two investors sharing the same elasticity respond the same way to change in prices irrespective of their underlying trading motives.

The demand elasticity of investors is crucial to understand how the equilibrium changes in response to changes in market macrostructure, the broad organization of financial markets ([Haddad and Muir, 2025](#)). For example, the impact of a shift in demand on prices is driven by the inverse of the aggregate demand elasticity. That is, if an investor wakes up and liquidates their portfolio ΔD , the vector of prices moves by $\Delta P = -\mathcal{E}_{agg}^{-1} \Delta D$.⁶ More broadly, armed with investors' demand elasticities, the researcher can answer questions about altering a subset of investors: for example bank regulation, quantitative easing, etc. Appendix B.1 characterizes the set of counterfactual questions that can be answered with demand elasticities in standard competitive equilibrium settings. We also show that demand elasticity remains a relevant concept in models of asymmetric information or even with strategic motives (see also [Haddad et al., 2024c](#)).

⁵If M and Σ respond to changes in prices, the elasticity formula would be different. [Kojien and Yogo \(2019\)](#) characterize demand curves in such a setting. [Kojien et al. \(2023\)](#) add hedging demands.

⁶The equilibrium prices equate the aggregate demand D_{agg} (the sum of demand across investors) to the aggregate supply S , in that $D_{agg}(P) = S$. If demand curves shift by an amount ΔD , the new equilibrium price $P + \Delta P$ satisfies $D_{agg}(P + \Delta P) + \Delta D = S$. With the implicit function theorem, the price impact or multiplier matrix is the inverse of the elasticity matrix of aggregate demand: $\mathcal{M} = -\left(\frac{\partial D_{agg}}{\partial P}\right)^{-1} = -\mathcal{E}_{agg}^{-1}$. Section 4 considers empirical approaches to estimate \mathcal{M} directly as opposed to the demand elasticity of each investor separately.

Demand elasticity: in the data. The objective of this paper is to estimate demand elasticity \mathcal{E} using data on prices and portfolio positions. Portfolio positions change because of prices (movements along the demand curve) and other reasons (shifts in the demand curve). In the mean-variance example, investor beliefs about expected payoffs or volatility could move. Other forces can drive portfolio decisions given prices. We represent all these movements by a component ϵ . Changes in demand follow

$$\Delta D = \mathcal{E} \Delta P + \epsilon \iff \Delta D_i = \sum_j \mathcal{E}_{ij} \Delta P_j + \epsilon_i. \quad (3)$$

One can specify this relation in levels or logs depending on the model of demand. For example, models like CARA preferences are better behaved in levels, while logit demand aligns with logs. In practice, the choice of units should be driven by regularity in the data and the type of model the researcher believes best match this regularity.⁷ We abuse the language of demand estimation slightly and call the matrix \mathcal{E} demand elasticity irrespectively of logs or levels. Also, while we focus on writing specifications in changes to match the standard difference-in-difference framework, similar arguments apply without changes. Regardless of units, if the model is not linear (or log-linear) in prices, we focus on a local approximation of demand; Appendix F discusses the nonlinear case.

1.2 Restrictions are Necessary to Estimate Demand Elasticity

The combination of rich substitution and price endogeneity pose a challenge to the estimation of demand elasticity.

Endogeneity of prices. Equation (3) suggests regressing demand on prices to estimate the demand elasticity \mathcal{E} . This approach is biased when shifts in demand curve ϵ are correlated with prices. Such a correlation is very natural. If positive news come out about a company, the investor would be inclined to buy more of its shares (a positive ϵ) if the price remained unchanged. However, the stock price is also likely to rise in response to the news. Similarly, periods when many investors face financial distress and seek to liquidate their portfolios (a negative ϵ) tend to coincide with depressed asset prices.

A natural solution for this challenge, in line with the causal inference literature, is to use natural experiments: situations where variation in prices is unrelated to shifts in demand curve of the investor. For example, one could focus on the inclusion of a stock in an index as

⁷Appendix H.1 reviews the appropriate units for standard models. See Appendix B for a discussion of the link between elasticity and the multiplier in different units.

an event that increases its price (Shleifer, 1986). Then, measuring the change in portfolio an investor for whom the index status of the stock is irrelevant reveals their demand elasticity.

Substitution. However, it is not that simple: equation (3) highlights that all prices matter for all demands due to substitution between assets. Hence, experiments for a single asset at a time can generally not be used to recover the matrix of elasticity \mathcal{E} , or even the own-price elasticity for this asset.⁸ In the language of causal inference, the stable unit treatment value assumption (SUTVA) is violated: treatment on one unit (for us, an asset) affects other units. A way out would be to include all prices in the demand estimation regression. This is not possible in practice because it requires exogenous sources of variation for each one of the (thousands of) individual prices.⁹

In the face of this challenge, one can deem causal inference hopeless for asset pricing and throw their hands in the air. However, there is a more constructive approach: acknowledge that additional assumptions about the nature of substitution are necessary. After all, this is the second part of Markowitz’ argument: basic economics can inform us about the structure of substitution across assets. In the rest of the paper, we follow this path and put forward simple, flexible conditions guided by these economic principles. Our approach is in the spirit of the causal inference literature: our conditions allow to estimate the elasticity matrix without fully specifying how the equilibrium emerges by using standard econometric methods (Angrist and Pischke, 2009).

These restrictions are flexible enough to model rich substitution patterns in financial markets. Next, we highlight why this flexibility is important.

1.3 Restrictions Should Accommodate Standard Finance Logic

Empirical finance research highlights that price movements are often driven by differential exposure to common factors. In standard finance theories, investors choose their portfolio to manage exposure to common factors. These insights have implications for the structure of the elasticity matrix — specifically cross-asset substitution — that an estimation framework must accommodate. We also show that leading structural models of demand do not.

⁸Consider for simplicity a setting with three assets and no residual shift in demand. In addition, assume that there is no cross-elasticity between assets 1 and 2 and that they have the same own-price elasticity \mathcal{E}_0 . One might want to use an experiment that affects asset 1 but not asset 2 and use it as an instrument to estimate this elasticity. The corresponding instrumental variable would have $Z_1 = 1$ and $Z_2 = 0$. This leads to an IV estimate of the elasticity $\hat{\mathcal{E}} = (\Delta D_1 - \Delta D_2)/(\Delta P_1 - \Delta P_2) = \mathcal{E}_0 + (\mathcal{E}_{13} - \mathcal{E}_{23})\Delta P_3/(\Delta P_1 - \Delta P_2)$, which is a biased estimator of the own-price elasticity \mathcal{E}_0 .

⁹This issue is the well-known challenge for estimating of a fully flexible demand function as in Deaton and Mullbauer (1980).

A workhorse finance model. In standard models, investors focus on factor exposure as a way to manage their risk. Consider the Markowitz model introduced in Section 1.1 with power utility with risk aversion γ .¹⁰ We assume that returns follow a factor structure $R_{it} = \alpha_i + \beta_i F_t + v_{it}$, where $v_{it} \sim \mathcal{N}(0, \sigma_v^2)$, is i.i.d across assets and orthogonal to the factor F_t (as in [Campbell and Viceira \(2002\)](#) or [Kojien and Yogo \(2019\)](#)). We focus on a single factor structure for illustrative purposes; the matrix versions of all the expressions hold with multiple factors. Each asset has a factor loading β_i and α_i is the component of expected return not explained by these factor loadings ($\beta' \alpha = 0$).

In this model, the portfolio share ω_i in each asset takes the following form:

$$\omega_i = \beta_i \mathcal{Y} E(F) + \frac{1}{\gamma \sigma_v^2} \alpha_i, \quad (4)$$

where the scalar \mathcal{Y} depends on risk aversion and the covariance of returns but not prices.¹¹ Two forces shape portfolio decision. First, investors engage in factor management: they invest in a portfolio with weights proportional to the factor loadings β . They buy more or less of this portfolio depending on the risk premium of the factor. Second, when the factor model does not completely explain returns and the alphas are different from zero, investors engage in “arbitrage” by tilting their investment towards assets with positive alphas.¹²

When the vector of prices changes, holding payoffs constant, expected returns change as well — both the alphas and the factor premium. Hence the same two forces also determine the elasticity matrix:

$$\mathcal{E} = \beta \Psi \beta' + \frac{1}{\gamma \sigma_v^2} \mathbf{I}, \quad (5)$$

where β is the $N \times 1$ vector of factor loadings and Ψ is the counterpart of \mathcal{Y} in (4). Cross-asset substitution only depends on the first term. The substitution between asset i and j is determined by their factor loadings β_i and β_j ; we later show that similar substitution patterns can arise for other motives than risk, in which case β_i is replaced by another asset characteristic.¹³ Concretely, if the price of the beta-sorted portfolio increases, this implies a

¹⁰The agent maximizes expected utility $W^{1-\gamma}/(1-\gamma)$ with risk-free rate r_f and lognormally distributed payoffs $\log(\Pi) \sim \mathcal{N}(M, \Sigma)$. Then log returns are $R_i = \log(\Pi_i)/P_i$, which we log-linearize following [Campbell and Viceira \(2002\)](#).

¹¹The covariance matrix of returns is $\Sigma = \beta \Sigma_F \beta' + \sigma_v^2 \mathbf{I}$. One can invert this matrix using the Sherman-Morrison formula and compute portfolio weights and the elasticity. This gives $\Psi = \frac{1}{\gamma \sigma_v^2} (\Sigma_F + \sigma_v^2 (\beta' \beta)^{-1})^{-1}$, and $\mathcal{Y} = \frac{1}{\gamma} (\sigma_v^{-2} \mathbf{I}_K + \Psi \beta' \beta)$.

¹²While arbitrage is often used to describe such trades, they are not the textbooks’ risk-free profit opportunities.

¹³Thus, our framework does not require taking a stand on the age-old debate between factors and characteristics ([Daniel and Titman, 1997](#); [Kelly et al., 2019](#)).

decrease in factor premium, and will alter demand for individual assets more or less depending on their own beta. The second component, driven by “arbitrage,” is diagonal and hence only affects own-price elasticities: if an individual asset is cheaper, the investor increases its position in this specific asset.

Empirical evidence. In the data, we do not know a priori the motives of investors. Yet, the properties of returns are often highly suggestive that motives in the style of factor management are relevant. Treasury bond returns are well explained by two or three common components (Litterman and Scheinkman, 1991). Much work focuses on common factors in stock returns (Fama and French, 1993), although with more debate on whether a low-dimensional representation explains the data. Asset comovement alone does not necessarily imply that substitution occurs along those dimensions.

Evidence from natural experiments where quantity shocks are known is more suggestive of the importance of this type of substitution. For example, Greenwood and Vayanos (2014) show that when the overall supply of treasury changes, the response of bond yields lines up with maturity, even though the composition of the supply changes does not. Krishnamurthy and Vissing-Jorgensen (2011) or Haddad et al. (2024a) document a similar pattern in response to the Fed’s announcement of bond purchases for quantitative easing operations. In the European context, the response across countries of bond yields to asset purchases by central banks lines up well with the riskiness of each country, even though purchases are mostly proportional to country size (e.g., Haddad et al., 2024b). For example, in some interventions, the ECB does not purchase Greek debt, yet its spread decreases more than that of any other country. Finally, Section 5 provides evidence of such substitution across different maturities and quantities of credit risk in corporate bonds.

Limits of leading structural models. To illustrate how existing models cannot capture these forces, we consider the following experiment. In our economy, there is an investor with demand function from equation (4) above as in Campbell and Viceira (2002) and a varying supply of assets. The researcher observes realizations of equilibrium prices and portfolio positions, with all demand parameters held constant (risk-aversion, covariance matrix, and expected payoffs). She estimates her favorite model of portfolio choice from the data; estimation does not face identification challenges because there are no demand shocks in the simulated data.

Then the researcher uses her estimated model to infer the counterfactual change in supply necessary to generate a specific change in prices: a larger observed increase in the price of

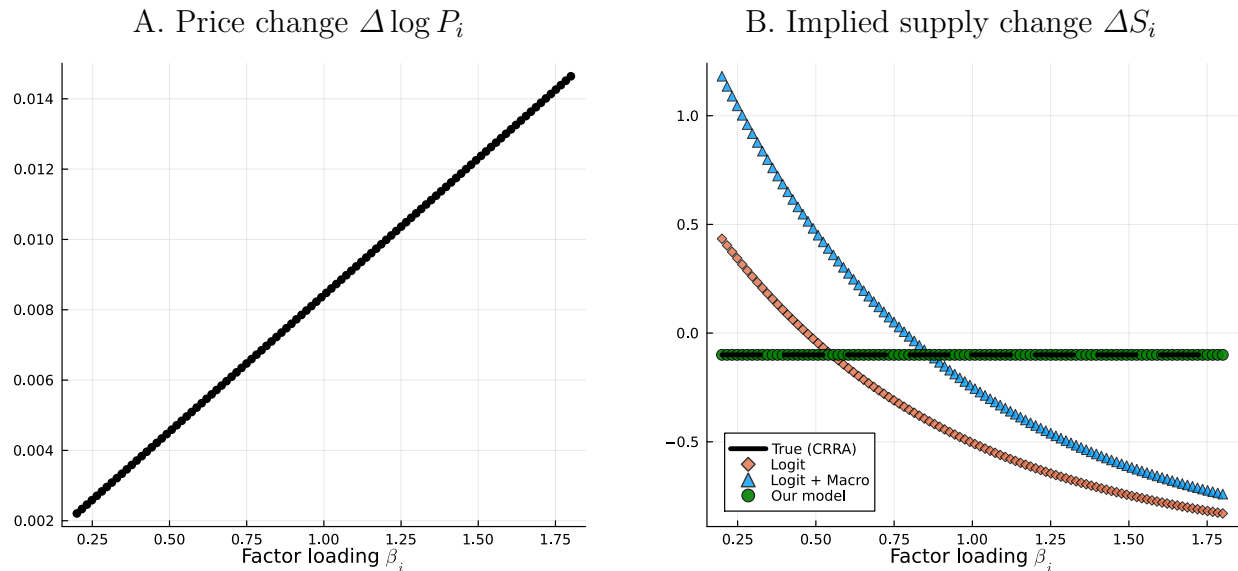


Figure 1: **Supply inference under alternative demand models.** Panel A plots the equilibrium log price change $\Delta \log P_i$ against asset beta β_i when the supply of all assets is reduced uniformly by 10%. High-beta assets experience larger price increases because the reduction in aggregate factor risk lowers the factor premium. Panel B shows the supply change each demand model infers as necessary to generate the price pattern in Panel A. The true shock is a uniform decrease across all assets (dashed black line). Demand models are logit (red), logit enriched with a macro elasticity (blue), and our model (green) introduced in Section 2. Parameters are calibrated to U.S. stock market data (1980–2021); see Appendix I and Table 6 for details on the CRRA economy, estimation, and robustness across alternative specifications.

high-beta assets than that of low-beta assets as in Figure 1, panel A.¹⁴ Concretely, if a researcher had estimated a model of bond demand and then observed the response of prices to the announcement of QE, what pattern of Fed purchases would they predict?

Panel B represents the answer under different models of demand. In our economy, the actual supply shock is an equal decrease in the amount of all assets, the dashed black line. This lowers the amount of factor risk in the economy, which yields a smaller factor premium in equilibrium, and disproportionately affects the price of high-beta assets. Appendix I discusses details of the experiment, parameter choices, and alternative specifications. If researchers instead favor a logit model of demand with beta as a characteristic, as in [Kojen and Yogo \(2019\)](#), they would predict a change in supply that is decreasing with beta, the red line.¹⁵ Because the logit model does not entertain substitution of the form highlighted in (5), it must attribute the differential change in price along betas to a differential change in supply. One

¹⁴One can also interpret this experiment without any reference to an equilibrium. The researcher estimates a model of demand from portfolio responses to exogenous changes in prices. Then they predict the response to a change in price as in panel A of Figure 1.

¹⁵While [Kojen and Yogo \(2019\)](#) start with a framework similar to the theory of this section, they include additional restrictions that specialize the demand function to logit.

might worry that the logit model does not capture aggregate forces, specifically the macro elasticity emphasized by [Gabaix and Koijen \(2024\)](#). We enrich logit with a macro elasticity parameter (as suggested in their appendix G.4) and find no qualitative change in the required supply: on the blue line, the model still needs a differential change in supply to explain the change in prices; we explain precisely why in Section 3.2.1.

Anticipating the next section, we also fit the model of demand that we develop in this paper. Because the model accounts for asset substitution along observable characteristics such as the betas, it is able to infer the actual change in supply from the differential change in price, the green line.

2 A Model of Asset Demand

In this section, we state our model of the elasticity matrix: we give an assumption that is sufficient to make estimation tractable using standard causal inference methods, but flexible enough to entertain many motives for investment. Specifically, we show that it captures a wide range of theoretical foundations, including not only risk-based incentives as in Section 1.3 but also other motives such as regulatory constraints or non-pecuniary objectives.

2.1 Framework

Investors may substitute across assets in diverse ways, as long as these differences are observable.

Assumption A1 (Homogeneous substitution conditional on observables) *Any pair of assets in the estimation sample \mathcal{S} with the same observables shares the same cross-price elasticity with respect to each third asset, within or outside of the estimation sample:*

$$\mathcal{E}_{il} = \mathcal{E}_{jl}, \quad \text{for all } i, j \in \mathcal{S} \text{ such that } X_i = X_j, \text{ and } l \neq i, j, \quad (6)$$

where X_i is the $K \times 1$ vector of observables for asset i . These cross-elasticities are parametrized by a bilinear form \mathcal{E}_{cross} :

$$\mathcal{E}_{il} = \mathcal{E}_{cross}(X_i, X_l) = X_i' \mathcal{E}_X X_l, \quad (7)$$

where \mathcal{E}_X is a $K \times K$ matrix (which may not be symmetric).

Assumption A1 states that for two assets that are comparable along observables, if the price of any third asset—either within or outside the estimation sample—moves, then substitution between the third asset and the two comparable assets is the same. Specifically, take

two firms that share the same characteristics, say Ford and General Motors. When the price of Netflix moves, the investor changes their position as much for Ford assets as for General Motors assets. In contrast, the investor might do something different for a firm with different characteristics, such as Target.

Importantly, assuming that the investor substitutes in the same way with Ford and General Motors is not the same as assuming that the two assets are identical. Features specific to each asset are still allowed to affect how much the investor demands of them, as materialized by the residual ϵ_i in equation (3). For example, in risk-based models, two assets with the same observables can have the same comovement with other assets, which drives substitution, but each of them has a distinct idiosyncratic component to their returns; see Section 2.2.1 below. In practice, the set of observables can include any variable measured prior to the price and quantity changes, including variables derived from past prices. We discuss principles to guide this selection throughout the remainder of the paper.

Analytically, homogeneous substitution implies that cross-price elasticities are a function of the observables, which we write as $\mathcal{E}_{cross}(X_i, X_l)$. To make the model tractable, we further parametrize this function as a bilinear form in observables, $\mathcal{E}_{cross}(X_i, X_l) = X_i' \mathcal{E}_X X_l$. This simple specification encompasses a large space of potential substitution patterns because observables could already include nonlinear transformations of more primitive variables or dummies for their levels.¹⁶ We demonstrate this versatility and practicality in Section 2.2.

Remarkably, this intuitive assumption yields a decomposition of the elasticity matrix into two naturally interpretable components:¹⁷

$$\mathcal{E} = \text{relative elasticity} + \text{substitution} \tag{8}$$

$$= \text{diagonal}(\mathcal{E}_{\text{relative}}) + X \mathcal{E}_X X', \tag{9}$$

where $X = [X_1'; X_2'; \dots; X_N']$ is the $N \times K$ matrix of observables. In this framework, the elasticity matrix is entirely characterized by substitution, encoded by the $K \times K$ matrix \mathcal{E}_X , and the relative elasticity, an asset specific component.¹⁸ These two components generalize the structure of “arbitrage” for the relative elasticity, and “factor management” for the substitution matrix that we observed in equation (5). The substitution matrix encodes cross-price elasticities that depend on observables; we provide an interpretation of the relative elasticity in the next section.

¹⁶Furthermore, it is most often natural to include a constant 1 as part of the observables. If the elasticity satisfies the assumptions for a set of observables X , it also does so for a linear transformation of these observables. For example, one can demean or standardize the observables without loss of generality.

¹⁷See proof in Appendix C.2.

¹⁸The substitution matrix \mathcal{E}_X may not necessarily be symmetric.

To facilitate exposition, we regularize the problem further by assuming that the relative elasticity is constant across assets. In the language of causal inference methods, this corresponds to assuming a homogeneous treatment effect. There are many standard ways to relax this constraint. Section 3.1.2 extends the framework to consider situations where the relative elasticity is not constant and either depends on observable or unobservable sources of variations.

Assumption A2 (Constant relative elasticity) *Assets in the estimation sample have the same value of relative elasticity $\mathcal{E}_{relative}$ with respect to other assets with the same characteristics:*

$$\mathcal{E}_{ii} - \mathcal{E}_{ji} = \mathcal{E}_{relative}, \quad \text{for all } i, j \in \mathcal{S} \text{ such that } X_i = X_j. \quad (10)$$

Assumption A2 ensures a form of symmetry in how the investor responds to the price of assets with the same observables in the sample. It focuses on a specific dimension: the difference between the own-price and cross-price elasticity. We call this difference the relative elasticity. It represents how the demand for one asset relative to another shifts when the price of the asset changes relative to the other, when both assets have the same observables. The next section explains why this quantity is the natural target of cross-sectional regressions.

In some cases, an asset i does not have a “twin” j with identical observables. This often occurs when the observables are continuous variables, such as the sales of a firm. In this situation, we replace Assumption A2 with its natural extension: $\mathcal{E}_{ii} - \mathcal{E}_{cross}(X_i, X_i) = \mathcal{E}_{relative}$. Combining equation (9) and Assumption A2 gives a simple decomposition of elasticity:

$$\mathcal{E} = \mathcal{E}_{relative}\mathbf{I} + X\mathcal{E}_X X'. \quad (11)$$

In Section 3, we show how to use standard causal inference tools to estimate the two components in this setting: the relative elasticity $\mathcal{E}_{relative}$ (Section 3.1) and the substitution matrix \mathcal{E}_X (Section 3.2). Of course, this tractability is only valuable if the assumptions are plausible in practice; we now turn to this question.

2.2 Applying the Assumptions

The applied researcher who wants to use our framework must motivate its key assumption in their setting using three complementary strategies. First, they can take a stand on the investor’s motives for investment, that is, the underlying theory of portfolio choice. Second, they should choose appropriately which observables to include in X . Finally, they can restrict attention to a specific sample of assets.

The point here is not to argue that our assumption is always true; it is not. Still, we highlight that our assumption holds in a wide range of situations with choices that are intuitive, close to common empirical practice, and line up with standard finance theory. In Section 3.2.3, we discuss diagnostics that allow the researcher to detect when the assumed structure of substitution misses relevant dimensions.

2.2.1 Risk-based motives

The classic theory of portfolio choice relies on the tradeoff between risk and return. Then, the elasticity matrix as in equation (11) is driven by the covariance of asset returns.¹⁹ When this covariance is driven by exposures to a common factor like in Section 1.3, Assumption A1 holds. Each asset’s exposure β_i is the observable X_i driving substitution, as can be immediately observed in equation (5).²⁰ The exact same equation holds with multiple factors, with $\beta_i = X_i$ being the $K \times 1$ vector of factor exposures. Alternatively, under the assumption that the betas are a function of a set of asset characteristics, it is enough to use these characteristics as the observables (Kojien and Yogo, 2019). This corresponds to the instrumented principal components analysis model of Kelly et al. (2019, 2020).

The use of factor models is ubiquitous in empirical finance. Here, the specific notion of factor is for explaining the covariance matrix, not expected returns. For Treasuries, such a covariance is a strong feature of the data. For stocks, factor analysis typically reveals an important “market” factor, followed by a long succession of factors of relatively similar magnitude (Kelly et al., 2019; Kozak et al., 2020; Lopez-Lira and Roussanov, 2020), making it more challenging to capture substitution with a few observables. However, it is also important to note that actual investors might not use a fully sophisticated or rational estimate of covariance matrix for making portfolio decisions. For example, many financial institutions use factor structures such as the Barra model to assess and manage the risk of their portfolios; in this case, these factor loadings are the relevant observables for substitution.

2.2.2 Non-risk drivers of substitution

In practice, portfolio decisions respond to many other dimensions than risk and return. Some investors care about non-pecuniary aspects of the stocks they hold, such as their carbon emissions or ESG characteristics. Mutual funds, pension funds, and endowments often operate under mandates that require a specific mix of assets, while others are guided

¹⁹Beyond the static setting, (He et al., 2025) shows that the same property holds in some dynamic economies.

²⁰Furthermore, Assumption A2 holds if idiosyncratic volatility is constant across assets. In practice, one might be reluctant to assume constant idiosyncratic volatility; Section 3.1.2 shows how to relax this condition.

by broader objectives outlined in their prospectus. When hedge funds take on leveraged positions, haircuts apply and they have to post margins. Banks and insurance companies must ensure that their portfolios satisfy various regulatory targets such as capital adequacy ratios, leverage requirements, or liquidity requirements.

All these dimensions affect which assets these investors choose in the first place, but also how they rebalance their portfolio when prices move. For example, if one of your more environmentally friendly stocks appears overpriced, you might shed this position and replace it with another similarly green position. Hence, these motives can play an important role in the elasticity matrix.

We show formally how this is consistent with our assumptions. Consider a generic representation of such a motive, by adding a quadratic cost and a linear constraint to the mean-variance optimization problem:

$$\max_D \quad D'(M - P) - \frac{\gamma}{2} D' \Sigma D - \frac{\kappa}{2} (D' X^{(1)})^2 \quad (12)$$

$$\text{such that} \quad D' X^{(2)} \leq \Theta. \quad (13)$$

The quadratic cost $\kappa/2 (D' X^{(1)})^2$ captures smooth investment priorities: the more carbon-emitting stocks an investor holds, the less willing she is to hold additional carbon-emitting stocks. The variable $X^{(1)}$ measures the relative contribution of each asset to this total cost — e.g. its carbon emissions — while κ measures the overall willingness to hold carbon emitting stocks. The linear constraint represents hard targets, such as the liquidity ratio that a bank must hold. There, $X^{(2)}$ measures the contribution of each position to the constraint — e.g. its liquidity weight — and Θ is the maximum value capturing the regulatory requirement.

When prices move, the investor balances risk-return considerations with these other non-risk objectives. All these dimensions will shape substitution patterns. We show in Appendix E that the elasticity matrix for this investor satisfies assumptions A1 and A2 with respect to the stacked set of observables $X = [X^{(1)}, X^{(2)}, X^{(3)}]$, where $X^{(3)}$ are the observables necessary to capture the covariance matrix.²¹

Concretely, this implies that the framework accommodates non-risk motives. While the empirical researcher need not know every detail about the origin of these motives, they must take a stand on which observables capture how each asset contributes to them; for example, a firm's carbon emissions for a fund with an ESG mandate.

²¹Clearly, to be able to use our results, the covariance matrix Σ also has to satisfy the assumptions with respect to a set of observables $X^{(3)}$. Specifically, Σ has to be such that the elasticity matrix of the mean-variance problem without constraint and cost function satisfies assumptions A1 and A2.

2.2.3 Experimental design

Instead of imposing structure on all portfolio decisions, researchers can restrict their attention to a subset of assets for which the assumptions are plausible, and for which, as we will see later on, they have appropriate sources of exogenous variation. This corresponds to choosing judiciously the estimation sample \mathcal{S} in Assumptions A1 and A2. For example, one might opt to focus on a narrow set of highly comparable assets, making the assumptions plausible. Of course, this path reduces how much can be identified about the overall elasticity matrix.

Return spreads or homogeneous estimation sample. A widespread practice in asset pricing to isolate local mechanisms is to focus on the behavior of the spread between the price of two comparable assets. In a demand estimation context, [Drechsler et al. \(2024\)](#) assess the effect of monetary policy on mortgage markets by studying the spread between the mortgage rate and the rate of a treasury with the same duration. Another example could be multiple corporate bonds from the same issuer with similar maturity (see [Coppola \(2025\)](#)).

Our framework formalizes the implicit assumptions behind this procedure. This corresponds to imposing the conditions on a subsample \mathcal{S} without observables. In this case, assumption A1 implies that the cross-price elasticity is the same for all assets in the estimation sample, while assumption A2 additionally implies that the own-price elasticity is the same for all assets in the estimation sample:

$$\mathcal{E}_{ii} = \mathcal{E}_{own}, \quad \text{for all } i \in \mathcal{S}, \quad \text{and} \quad \mathcal{E}_{ij} = \mathcal{E}_{cross}, \quad \text{for all } i, j \in \mathcal{S}. \quad (14)$$

In the simple case of a spread this assumes symmetry in the two legs of the spread, a treated and a control.

Moving on to outside assets, which could be a vast set, the substitution between them and the assets in the estimation sample is generally not constant. Still, assumption A1 implies that, for each outside asset, all assets in the estimation sample have the same cross-elasticity. In other words, the demand for any asset in the sample responds in the same way to a change in the prices of each outside assets. Figure 2 illustrates such an elasticity matrix. This setting corresponds to a situation in which observables are constant within the estimation sample, while they can vary arbitrarily across assets outside the estimation sample.

In simple risk-based models like in Section 1.3, in which elasticities are proportional to the inverse of the covariance matrix, this means that all assets in the sample say i and j have the same variance and covariance with each other; and for any outside asset k , the covariance of its return with that of any asset in the sample is constant: $\text{cov}(R_i, R_k) = \text{cov}(R_j, R_k)$. In practice, outside assets are plentiful and this latter condition cannot be fully assessed.

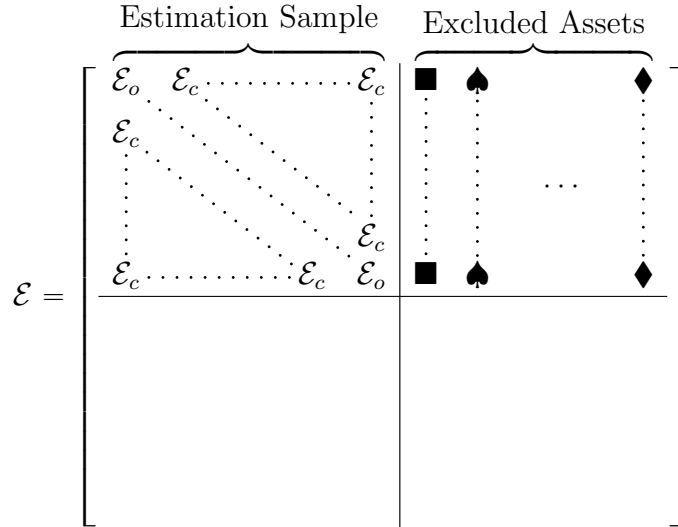


Figure 2: Elasticity matrix satisfying assumptions A1 and A2 for a local experiment.

Still, just as one would assess balance on observables (Athey and Imbens, 2017), researchers should present some corroborating evidence that reports average covariances with a set of broad portfolios for treated and control assets.

Groups of assets. In some settings, the researcher can delineate many such subsamples. Even though homogeneous substitution is not plausible across the whole sample, it is plausible within each of the subsamples. For example, homogeneity might hold for a set of firms in a narrow industry but not across these industries. Practically this corresponds to include group dummies as observables. Chaudhary et al. (2022) explain how omitting group effects in such a situation leads to biased inference. They document the relevance of this bias when measuring the effect of fund flows on corporate bond prices.

Mimicking portfolios as synthetic controls. A variation of this approach particularly well-suited for event-study settings is to construct synthetic controls in the style of hedging portfolios. With a set of treated assets, one can construct portfolios of other assets as the control group for a difference-in-difference study. There are two requirements for this approach to be valid. First, the observables or the factor exposures of the control portfolio must be the same as that of the treated asset. Second, each asset in the control portfolio (as opposed to the combined portfolio returns) must have similar residual volatility as the treated asset.

3 Estimation

We show how to estimate the elasticity matrix under the two previous assumptions, with the help of exogenous sources of variation. As a reminder the the data-generating-process follows equation (3):

$$\Delta D = \mathcal{E} \Delta P + \epsilon, \tag{15}$$

where the elasticity matrix takes the form highlighted in Section 2.1:

$$\mathcal{E} = \mathcal{E}_{\text{relative}} \mathbf{I} + X \mathcal{E}_X X'. \tag{16}$$

We show separately how to estimate the two components.

3.1 Estimating Relative Elasticity

We show that a simple causal inference regression estimates the relative elasticity consistently. Then, we discuss robustness and extensions of the framework.

3.1.1 Identification using standard cross-sectional approaches

Pretend for a moment that demand for asset i responds to the price of asset i only, not to the price of other assets, $\Delta D_i = \mathcal{E} \Delta P_i + \epsilon_i$, where \mathcal{E} is now a scalar. This is the canonical setting of causal inference, and running an instrumental variable estimation is natural in this situation. In this model, one regresses the change in demand for each asset ΔD_i on the change in the price of this asset ΔP_i , using Z_i as an instrument for the price change. This corresponds to the two-stage least square specification:

$$\Delta D_i = \hat{\mathcal{E}} \Delta P_i + \theta' X_i + e_i, \tag{17}$$

$$\Delta P_i = \lambda Z_i + \eta' X_i + u_i, \tag{18}$$

where X_i is the set of observables for each asset. A common case is when the instrument is binary (treatment and control) in which case the estimation is an instrumented difference-in-difference.

The two standard conditions for identification are the relevance and exclusion restrictions. The exclusion restriction captures the idea that the instrument does not affect demand through other channels than the price: $Z_i \perp \epsilon_i | X_i$.²² In other words, the instrument

²²The literature sometimes uses variations of this condition. [Kojien and Yogo \(2019\)](#) focus on

is not correlated with unobservable shifts in the demand curve in the cross-section of assets. For example, even if the experiment leads to general equilibrium effects such as changing the risk-free rate, the exclusion restriction can still be satisfied if the impact of these effects across assets does not correlate with which asset is treated. Relevance is the idea that the instrument Z_i creates variation in prices: $\lambda \neq 0$. In practice, it is not enough for the first stage to be significant at standard confidence levels; it must be strong to avoid issues related to the weak-instrument problem (Stock and Yogo, 2005; Olea and Pflueger, 2013).

Going back to the realistic case with cross-asset substitution, this estimation method seems to ignore the effect of other assets. The next proposition shows that the naive estimator correctly estimates the relative elasticity under our assumptions.

Proposition 1 *Under assumptions A1 and A2, as well as the standard relevance and exclusion restrictions, the two-stage least square estimation of equations (17) and (18) identifies the relative elasticity:*

$$\hat{\mathcal{E}} = \mathcal{E}_{relative}. \tag{19}$$

When the IV estimation is well specified, it identifies the relative elasticity: the difference between the own-price elasticity and the cross-price elasticity for two assets with the same observables. While this result stands in contrast to the intuition of measuring “how demand for each asset responds to its own price,” it is natural. A cross-sectional regression is a comparison across assets in the sample. Even if only the price of the treated asset is shocked, the regression coefficient will still be driven by the response of demand for this asset relative to that for the comparable control asset—hence the relative intensity of the own- and cross-elasticity conditional on observables. In other words, $\hat{\mathcal{E}}$ answers the question: how does the demand for one asset relative to another comparable asset respond to the relative price of these assets?

Proof for the simple case. Appendix C.1 proves Proposition 1. To understand the mechanics of this result, consider a simple case with two assets with the same observables

$E(\exp(\epsilon_i)|X_i, Z_i) = 1$, leading to non-linear estimation. Graves (2025) imposes zero conditional median as opposed to mean, which leads to sequential censored quantile regressions that account for censoring of positions at zero. We favor the linear case for simplicity of exposition and to reflect the most common approach in the empirical literature.

and no demand shifts ϵ . The changes in demands are:

$$\Delta D_1 = \mathcal{E}_{11}\Delta P_1 + \mathcal{E}_{12}\Delta P_2 + \sum_{k>2} \mathcal{E}_{1k}\Delta P_k; \quad (20)$$

$$\Delta D_2 = \mathcal{E}_{22}\Delta P_2 + \mathcal{E}_{21}\Delta P_1 + \sum_{k>2} \mathcal{E}_{2k}\Delta P_k. \quad (21)$$

Assumption A1 implies that the cross-elasticities with respect to other assets ($k > 2$) are identical:

$$\sum_{k>2} \mathcal{E}_{1k}\Delta P_k = \sum_{k>2} \mathcal{E}_{2k}\Delta P_k. \quad (22)$$

When computing the difference $\Delta D_1 - \Delta D_2$, this response to other prices disappears, effectively removing the omitted variable problem due to other assets:

$$\Delta D_1 - \Delta D_2 = (\mathcal{E}_{11} - \mathcal{E}_{21})\Delta P_1 - (\mathcal{E}_{22} - \mathcal{E}_{12})\Delta P_2. \quad (23)$$

Assumption A2 implies that the coefficients on each of the prices are the relative elasticity:

$$\mathcal{E}_{11} - \mathcal{E}_{21} = \mathcal{E}_{22} - \mathcal{E}_{12} = \mathcal{E}_{relative}. \quad (24)$$

Both the response of demand to the own price (measured by \mathcal{E}_{11}) and the response to the price of the other asset (measured by \mathcal{E}_{21}) shape this comparison. Hence, the regression coefficient is the relative elasticity:

$$\hat{\mathcal{E}} = \frac{\Delta D_1 - \Delta D_2}{\Delta P_1 - \Delta P_2} = \mathcal{E}_{relative}. \quad (25)$$

The role of observables. In the richer case where observables differ across assets, it becomes important to control for observables. Two different assets respond differently to the price of other assets. However, assumption A1 ensures that these responses only depend on each asset's observables X_i :

$$\Delta D_i = \mathcal{E}_{ii}\Delta P_i + \sum_{j \neq i} X_i' \mathcal{E}_X X_j \Delta P_j + \epsilon_i \quad (26)$$

$$= (\mathcal{E}_{ii} - X_i' \mathcal{E}_X X_i) \Delta P_i + \sum_j X_i' \mathcal{E}_X X_j \Delta P_j + \epsilon_i \quad (27)$$

$$= \underbrace{(\mathcal{E}_{ii} - X_i' \mathcal{E}_X X_i)}_{\mathcal{E}_{relative}} \Delta P_i + X_i' \underbrace{\sum_j \mathcal{E}_X X_j \Delta P_j}_{\text{constant across assets}} + \epsilon_i. \quad (28)$$

The second term in (28) highlights that substitution, while depending on all other prices, is proportional to X_i . As a result, controlling for X_i in a cross-sectional regression — that is, including a coefficient θ on observables — absorbs the effects of substitution.²³ Furthermore, the first term in (28) shows that the regression (after controlling for X_i) is equivalent to making pairwise comparisons of assets that have the same observables. Hence, following the same reasoning as in the simple case, the estimate $\hat{\mathcal{E}}$ recovers the relative elasticity.

What about equilibrium? One might worry that because prices are an equilibrium outcome, a natural experiment used to create an instrument will generate spillovers across assets, and that this threatens estimation. Say for example that the instrument comes from the Fed quasi-randomly purchasing some bonds but not others, as in [Selgrad \(2023\)](#). Even if the Fed does not buy a specific bond, its price might still respond to purchases of its substitutes. This type of spillovers in the first stage in (18) does not affect our identification result. Instrument exogeneity only requires that treatment status (whether the Fed bought the bond or not) is unrelated to shifts in the investor’s demand curve such as changes in their preferences or their views about the assets: Z_i is orthogonal to the shift in demand ϵ_i , conditional on the observables X_i .

The one economically meaningful situation that threatens identification is when assets are perfect substitutes, such that a no-arbitrage relation perfectly ties their prices together.²⁴ In this case, the relevance condition cannot hold because there is no change in the relative price of treated and control, and identification fails. [Fuchs et al. \(2025\)](#) also discuss how this case creates challenges for demand estimation. Of course, the relevance condition can and should be assessed directly empirically in the first stage. In finite samples, it is not enough for the first stage to be significant at standard confidence levels; it must be strong to avoid weak-instrument problems.²⁵

The potential for this issue in the first stage highlights a key consideration when selecting an appropriate control group for a given treated asset. While a control asset should be similar enough to the treated to satisfy Assumptions A1 and A2, it should not be identical. Because distinct assets naturally contain idiosyncratic risk, their relative prices can fluctuate without triggering arbitrage opportunities. More broadly, non-risk factors such as regulatory constraints can introduce equilibrium wedges that allow similar assets to be priced differently.

²³This reasoning shows that a weaker form of assumption A1 is necessary for Proposition 1 to hold: $\mathcal{E}_{il} = \mathcal{E}_{cross,l}(X_i) = X_i'Y_l$ for arbitrary vectors Y_l . In other words, the dependence to other assets for a given X_i can be arbitrary and does not need to be parametrized by observable characteristics X_l .

²⁴Here, what is important is that these assets are perfect substitutes at the aggregate level (so that their equilibrium prices are exactly the same), not so much that the investor whose demand is estimated treats them as such.

²⁵This can be assessed using the tests of [Stock and Yogo \(2005\)](#) and [Olea and Pflueger \(2013\)](#) for example.

Partial identification. It might be tempting to conclude that the N -dimensional demand curve of equation (3) in which all prices matter for all demands is equivalent to a uni-dimensional demand curve that only depends on the own price and characteristics, as estimated in the regression equation (17). This would be incorrect: the equivalent representation only holds when fixing a specific vector of prices; in other words, while equilibrium quantities demanded satisfy equation (17), the N -dimensional demand curve does not. To put it more formally, Proposition 1 offers a partial identification result, where the naive estimator recovers the relative elasticity but not substitution. As we showed in Section 1.3, models like logit that have the same structure as (17) cannot capture the substitution of traditional finance models. More generally, Section 3.2 shows that you cannot recover the substitution matrix \mathcal{E}_X from the cross-section alone.

3.1.2 Robustness and extensions

Robustness to deviations from the assumptions. In practice, assumptions A1 and A2 are approximations of reality. In Appendix D.1, we assess whether the result of Proposition 1 is robust to small deviations. We show that, as long as the first stage is economically large — that is, the instrument induces a substantial cross-sectional spread in prices — the two-stage least square estimator recovers the relative elasticity up to a bias that is proportional to the distance to the assumptions; small deviations, small bias. For example, a small amount of measurement error in factor loadings is not consequential. Again, the only exception of this rule is with a weak or zero first stage: weak instruments lead to an extreme sensitivity to misspecification (Bound et al., 1995).²⁶

Observed heterogeneity in relative elasticity. We can also entertain deviations from the assumption of constant relative elasticity. Just like Assumption A1 allows cross-elasticities to depend on observables, one can relax Assumption A2 to let the relative elasticity depend on observables. This corresponds to replacing the condition (10) by:

$$\mathcal{E}_{ii} - \mathcal{E}_{cross}(X_i, X_i) = \mathcal{E}_{relative}(X_i) = \mathcal{E}'_r X_i, \quad (29)$$

where \mathcal{E}_r a vector of dimension K that maps observables to the relative elasticity for asset i . For example, if the observable captures the size of a company, this relation allows big stocks to have a different relative elasticity than small stocks, an approach taken, for example, in Haddad et al. (2024c). Another useful application is in the context of the factor models of

²⁶Complementary to the finite-sample weak-IV tests, the relevant diagnostic is the magnitude of the price spread the instrument induces, measured in economic units (Conley et al., 2012; Van Kippersluis and Rietveld, 2018; Andrews et al., 2019).

Section 2.2.1. There, we have seen that idiosyncratic volatility controls the relative elasticity. Therefore, one could include the idiosyncratic volatility of each asset as an observable.

Estimating \mathcal{E}_r simply requires regressing relative demand on the price change and its interactions with observable characteristics, instrumenting appropriately. Appendix Section C.3 proves this result and provides all details on implementation.

Unobserved heterogeneity. A more complicated situation arises when there is unobserved variation in relative elasticity across assets (each asset has its own $\mathcal{E}_{relative,i}$). Nothing comes for free: accommodating this more flexible elasticity matrix \mathcal{E} requires stronger assumptions on the instruments to maintain meaningful identification. Specifically, one needs to assume independence of the instrument with respect to all unobserved sources of heterogeneity — a stronger condition than orthogonality to the demand residual. An example of this is index inclusions: the included and excluded assets are closely related in size but might randomly differ in industries or other characteristics.

Appendix C.6 formalizes conditions to estimate the relative elasticity in this case where Assumption A1 holds but Assumption A2 does not. The two-stage least squares estimation identifies an average of the relative elasticity:²⁷

$$\hat{\mathcal{E}} = \frac{\mathbf{E}_i[\lambda_i \mathcal{E}_{relative,i}]}{\mathbf{E}_i[\lambda_i]} \quad (30)$$

The relative elasticities of assets for which the instrument has a greater impact on prices (large λ_i) are given greater weights. For example, if more illiquid assets have both a higher impact of the instrument λ_i and investors trade them more inelastically (lower \mathcal{E}_{ii}), estimates of relative elasticity will be lower than the unweighted average relative elasticity, and overstate how inelastic the typical asset is.

3.2 Estimating Substitution

Standard cross-sectional causal inference can estimate the relative elasticity, a useful moment for answering micro questions comparing individual assets. There are many other interesting questions concerning more aggregated levels; here, aggregation is across assets. For instance, how do investors rebalance when the price of all small stocks changes relative to all big stocks, or when the price of long-duration bonds changes relative to short-duration ones? At the most aggregated level, how do investors respond when all stocks become more expensive?

²⁷Formally, the estimator yields a local average treatment effect with an added monotonicity condition that the relative impact of the instrument always affects prices in the same direction — all λ_i sharing the same sign. See discussion in Appendix C.6.

This section aims to address these questions. Doing so hinges on estimating cross-price elasticities separately from relative (or own-price) elasticities. Estimating these dimensions of the elasticity matrix, the substitution \mathcal{E}_X , must rely on sources of variation in the time series, one for each characteristic driving substitution plus one for the overall aggregate.

3.2.1 The missing coefficient problem

We show that, unlike for the relative elasticity, observing a single cross-section does not allow identification of the substitution matrix under assumptions A1 and A2. This identification challenge sheds light on how leading structural demand frameworks achieve tractability by implicitly assuming away this missing coefficient problem.

Substitution cannot be estimated using the cross-section. The source of this limitation is the classical challenge of demand estimation: substitution is a response to change in prices along the observables, but demand might also shift along the observables. Formally, the demand shifters ϵ can in general be correlated with the observable X : the investor’s demand for assets with high values of X relative to low values of X can change, even holding prices constant. For example the investors might become more optimistic about the relative prospects of tech firms relative to industrial firms.

Denote $\epsilon_X = (X'X)^{-1}X'\epsilon$ the projection of demand shifts on the observables and $\epsilon_{idio,i} = \epsilon_i - X'_i\epsilon_X$ the residual from this projection. With this notation, we can rewrite (28) as:

$$\Delta D_i = \hat{\mathcal{E}}\Delta P_i + \epsilon_{idio,i} + X'_i \underbrace{(\mathcal{E}_X X' \Delta P + \epsilon_X)}_{\text{coefficient } \theta \text{ (} K \times 1 \text{)}}. \quad (31)$$

Notice that the effect of substitution $\mathcal{E}_X X' \Delta P$ always appears together with the demand shifter ϵ_X . Furthermore, in a single cross-section, irrespective of the number of assets, there is only a single value for these two quantities, and hence they cannot be separated using shocks to one but not the other. Appendix Section C.5 proves this impossibility result.

This observation is the flip side of the result that facilitated the estimation of relative elasticity: the effect of substitution is always proportional to the observables and is absorbed by the coefficients θ in a cross-sectional regression, as shown by the second term in (28). Demand shifts proportional to the observables also shape the coefficient θ . When the observable X is a constant, this issue is the classic missing intercept problem: substitution is absorbed in the intercept of the cross-sectional regression, but so are aggregate shifts in demand. When substitution depends on other variables, this is a “missing coefficient” problem.

While the well-known missing intercept problem implies that cross-sectional regressions

cannot identify aggregate effects (e.g., [Guren et al., 2021](#); [Wolf, 2023](#); [Gabaix and Koijen, 2024](#)), our result reveals a much more pervasive blindness. It demonstrates that even the relative forces driving substitution across assets are absorbed by the slope coefficients, rendering them indistinguishable from demand shifts in standard cross-sectional estimation. In other words, the cross-section is not enough to answer questions about the cross-section.

Leading structural models and the missing coefficient problem. Leading structural models such as the logit demand used in [Koijen and Yogo \(2019\)](#) circumvent the missing coefficient problem. While this feature is an attractive property for tractability, we show it precludes these models from capturing the standard substitution based on risk we explored in Section 1.3. Denote by ω the investor’s portfolio share; in units of log portfolio share on log price, a logit demand system implies an elasticity matrix of the form

$$\mathcal{E}^{logit} = -\alpha \mathbf{I}_N + \alpha \mathbf{1}\omega'. \quad (32)$$

This expression implies that logit satisfies Assumptions A1 and A2, where the observables are initial portfolio weights and a constant.²⁸ Proposition 1 applies, and the cross-sectional estimation identifies the relative elasticity parameter α . Because α is the only parameter governing substitution, this cross-sectional estimate pins down the entire elasticity matrix under the logit specification (225); there is no missing coefficient problem.

However, comparing this elasticity matrix to the one in the standard model of Section 1.3 (equation (5)) reveals that the logit’s restriction on substitution prevents the model from generating the behavior of investors in that workhorse finance theory. As our example demonstrated, a drop in total supply decreases the risk premium, and the required return on riskier assets decreases more than that of less risky assets; otherwise, investors with standard utility functions would not continue to hold the supply of assets. This implies that investors’ substitution matrix is sensitive to risk, i.e., depends on β in our example. The logit model does not offer a degree of freedom to accommodate this dependence (even with an extension including an additional parameter for the macro elasticity as in [Gabaix and Koijen \(2021\)](#)).²⁹

²⁸Formally, the observables for asset i are a constant, $X_i^{[1]} = 1$, and the investor’s initial portfolio weight in the asset, $X_i^{[2]} = \omega_i$. Then the substitution matrix is $\mathcal{E}_X = \begin{pmatrix} 0 & \alpha \\ 0 & 0 \end{pmatrix}$.

²⁹In a logit asset demand system, substitution operates via the outside asset (the risk-free asset), and it is proportional to investors’ pre-determined portfolio shares $\{\omega_i\}$. More precisely, given some counterfactual price vector \mathbf{P} , the resulting demand in terms of portfolio share satisfies $\ln \omega_i(\mathbf{P}) = \ln \omega_0(\mathbf{P}) + \hat{\mathcal{E}} \ln P_i + \theta' X_i$, where X_i potentially includes the risk exposure β_i . Therefore, under a logit asset demand system the movement of risk premium (reflected by a change in the vector \mathbf{P}) only affects the average of portfolio demand (the intercept) but not how it depends on the risk exposure (the coefficient θ). Capturing changes in portfolio dependent on β requires allowing a dependency $\theta(\mathbf{P})$.

For a consistent estimation of the demand elasticity, the extent of this sensitivity still has to be estimated, and the missing coefficient problem arises.

3.2.2 Decomposing substitution

To develop an identification strategy for substitution, we derive a convenient decomposition. In standard finance models like APT (Ross, 1976) or ICAPM (Merton, 1973), a common insight is that one does not need to look at every asset independently but can focus on a few curated portfolios. While we entertain a much richer class of demand functions, a similar property arises in our framework. A few portfolios are enough to capture the effect of substitution.

Proposition 2 *Take an elasticity matrix \mathcal{E} satisfying assumptions A1 and A2, and consider generic changes in demand and price connected by this matrix $\Delta D = \mathcal{E}\Delta P$. Define the change in demand and price aggregated along observables (vectors of size K) and the idiosyncratic counterparts (scalars) as:*

$$\Delta D_X = (X'X)^{-1}X'\Delta D, \quad \Delta P_X = (X'X)^{-1}X'\Delta P, \quad (33)$$

$$\Delta D_{idio,i} = \Delta D_i - X'_i\Delta D_X, \quad \Delta P_{idio,i} = \Delta P_i - X'_i\Delta P_X, \quad (34)$$

such that $\Delta D = \Delta D_{idio} + X\Delta D_X$ and $\Delta P = \Delta P_{idio} + X\Delta P_X$. The response of changes in demand to a change in prices can be decomposed into two sets of components:

$$\text{Micro:} \quad \Delta D_{idio,i} = \mathcal{E}_{relative}\Delta P_{idio,i}, \quad (35)$$

$$\text{Meso-Macro:} \quad \Delta D_X = \check{\mathcal{E}}\Delta P_X, \quad (36)$$

where the $K \times K$ matrix $\check{\mathcal{E}} = \mathcal{E}_{relative}\mathbf{I}_K + \mathcal{E}_X X'X$.

The response of demand to prices can be decomposed into two components. The idiosyncratic part of the response of demand to prices is driven by the relative elasticity discussed above in Section 3.1. The other component of the response corresponds to substitution effects and is summarized by the response of a few broad portfolios. The vectors ΔP_X and ΔD_X represent the return and the demand response for long-short portfolios along the observables X . Specifically, these are the coefficients from a regression of ΔD or ΔP on X . Why can we interpret these coefficients as coming from portfolios? This insight goes back to Fama and MacBeth (1973), who point out that the $(X'X)^{-1}X'$ can be interpreted as a set of portfolio weights based on the variables in X .

Two examples help gain intuition on how the decomposition works. First, consider a setting where there is no variation in observables, that is when the only variable in X is a

constant as in [Gabaix and Koijen \(2021\)](#). Then, ΔD_X and ΔP_X are of dimension 1 and represent the average demand and return across assets: $\Delta D_X = N^{-1} \sum_i \Delta D_i$ and likewise for ΔP_X . Then the scalar $\check{\mathcal{E}}$ in equation (36) is the macro elasticity, how aggregate demand for assets respond to the aggregate return.

Now, in addition to the constant $X^{[1]} = 1$, add heterogeneous substitution along one variable $X^{[2]}$ assumed to have mean zero and unit standard deviation. This adds a additional component to $\Delta D_{X^{[2]}}$ and $\Delta P_{X^{[2]}}$. These quantities measure how the price and the demand for assets with larger values of X change relative to those with lower values. Specifically, $\Delta P_{X^{[2]}} = N^{-1} \sum_i X_i^{[2]} \Delta P_i$ is the change in price of a long-short portfolio sorted on this characteristic; $\Delta D_{X^{[2]}} = N^{-1} \sum_i X_i^{[2]} \Delta D_i$ is the change in portfolio tilt along the characteristic. For example, if X measures the duration of bonds, $\Delta P_{X^{[2]}}$ is the return of a portfolio that goes long high-duration bonds and short low-duration bonds; likewise, $\Delta D_{X^{[2]}}$ measures the portfolio tilt relative to the average duration. This additional component highlights again that heterogeneous substitution cannot be captured with only relative (or micro) and macro elasticity but instead adds an additional intermediate “meso” layer to how demand responds to prices. The formulas in equations (33)-(36) generalize this example to multiple observables.

3.2.3 Using the time series to estimate substitution

The decomposition shows that the effect of substitution is characterized through the behavior of a few time series (as many as the number of observables K), which are the meso and macro aggregators in ΔD_X and ΔP_X . However to estimate substitution, the researcher must still contend with the presence of shifts in demand that are correlated with the price:

$$\Delta D_{X,t} = \check{\mathcal{E}} \Delta P_{X,t} + \epsilon_{X,t} \quad (37)$$

The next proposition shows how to identify $\check{\mathcal{E}}$ with a set of time series instruments.

Proposition 3 *Consider a set of K time series instruments $Z_{X,t}$ for $\Delta P_{X,t}$, such that the exclusion restriction, $Z_X \perp \epsilon_X$, and the relevance condition, $\text{rank}(\text{cov}(Z_X, \Delta P_X)) = K$, both hold. Under assumptions A1 and A2, the set of K two-stage least squares regressions:*

$$\Delta D_{X^{[k]},t} = \sum_l \check{\mathcal{E}}_{k,l} \Delta P_{X^{[l]},t} + e_{X,t}^{[k]}, \quad (38)$$

$$\Delta P_{X^{[k]},t} = \sum_l \Lambda_{k,l} Z_{X^{[l]},t} + u_{X,t}^{[k]}, \quad (39)$$

correctly identifies the matrix $\check{\mathbf{E}}$.³⁰ Combined with an estimate $\hat{\mathbf{E}}$ of $\mathcal{E}_{relative}$, this provides an estimate of the substitution matrix $\mathbf{E}_X = \check{\mathbf{E}} - \mathcal{E}_{relative}\mathbf{I}_K$.

Proposition 3 explains that substitution is identified from regressing the portfolio tilt along each observable ($\Delta D_{X^{[k]}}$) on the change in price of all portfolios sorted on observables (ΔP_X), instrumented by the variables (Z_X). In this regression, the notion of instrument and exogenous variation comes from the time series. That is, the researcher needs a set of variables that moves prices but are orthogonal with demand shifts. The use of such instruments is the solution to the missing coefficients problem introduced in Section 3.2.1, because they separate the effects of substitution on demand from shifts in demand for the observables.

Take the example of bond demand where the researcher focuses on a single observable: duration. Estimating substitution requires running two distinct regressions: one measuring the sensitivity of aggregate demand to the returns of both an aggregate and a duration-tilted portfolio, and another measuring the sensitivity of duration demand to these same returns. Both price variables are required in each regression to capture heterogeneous substitution. For example, if the price of all bonds become more expensive then the investor might be less inclined to invest in a duration-tilted portfolio. To consistently estimate these coefficients, the researcher must employ two separate instruments, one for each endogenous price variable. One could use surprises to government debt issuance or central bank’s overall asset purchases, both in overall quantity (e.g., [Greenwood and Vayanos, 2014](#); [Krishnamurthy and Vissing-Jorgensen, 2011](#)) or changes in the duration of the purchases (e.g., [Hubert de Fraisse, 2022](#)). If one is only interested in the macro multiplier, it is tempting to simplify the analysis and only use a single instrument for the aggregate price and control for the price of duration. However, our analysis highlights that this is a case of a bad control situation which introduces endogeneity ([Angrist and Pischke, 2009](#)).

This result highlights a fundamental tension between estimating relative elasticity and substitution as the set of observables expands. As discussed in Section 3.1, incorporating a rich set of observables can strengthen causal inference for relative elasticities by improving robustness. By contrast, estimating substitution becomes increasingly demanding as dimensionality grows, since each additional observable requires its own source of exogenous variation. In asset classes with a low-dimensional structure, such as Treasuries, this tension is not an issue. In richer environments like equities, however, the researcher must impose struc-

³⁰The estimation equations also take the following matrix form:

$$\Delta D_{X,t} = \check{\mathbf{E}}\Delta P_X + e_X \tag{40}$$

$$\Delta P_{X,t} = \mathbf{\Lambda}Z_{X,t} + u_X. \tag{41}$$

ture to keep the problem tractable. A useful guiding principle is to allow heterogeneity along the dimensions that are central to the question at hand and those that feature prominently in the broader literature. For example, when studying the impact of ESG preferences, it is important to distinguish investors’ willingness to adjust their green tilt in response to changes in the green-brown premium — an economically natural margin that is typically abstracted from in existing estimates.³¹ An alternative strategy explored in the literature is to rely on automatically generated instruments, such as granular instrumental variables (Gabaix and Koijen, 2024), which offer a different way of addressing the challenge of dimensionality.

Falsifying the assumptions. Beyond arguing ex ante for Assumptions A1 and A2, the researcher can subject the framework to falsification tests after estimation. We describe three diagnostics, each flagging a different way the specification might miss a relevant dimension of substitution. A first family of tests, in the spirit of Hausman (1978) and Mundlak (1978), checks whether the estimated relative elasticity is invariant to the source of variation used to identify it. For instance, one can repeat estimation restricting attention to within-issuer variation rather than exploiting the full cross-section of corporate bonds. If the substitution structure is correctly specified, alternative estimators should agree; a meaningful divergence indicates that unmodeled drivers of substitution contaminate the regression. A second test targets a specific candidate: choose an additional characteristic that might drive substitution, add it to X , and re-estimate. If this inclusion alters estimates, that is direct evidence that the candidate characteristic does shape substitution and should be modeled explicitly. A third test addresses the structure of substitution itself: the linear-in-observables form $X\mathcal{E}_X X'$ delivers overidentifying restrictions. Estimating meso responses non-parametrically for separate buckets of X should reproduce the linear pattern implied by the framework, and departures reveal non-linearities in substitution that the specification as written cannot capture. We implement all three diagnostics for the corporate bond market in Section 5.

An alternative approach: assuming symmetric groups. As discussed at the end of Section 2.2.3, one way to ensure assets are comparable is to bin them into groups using dummy variables. Without additional assumptions, however, estimating substitution within this framework would require as many instruments as there are groups. A common alternative to our linear-in-observables approach is to assume symmetry across the groups: the substitution matrix \mathcal{E}_X has constant diagonal and off-diagonal components. In other words, the matrix \mathcal{E}_X (as opposed to \mathcal{E}) satisfies Assumptions A1 and A2 with only a constant as the observable; there is homogeneous substitution *across* group-level portfolios. Chaudhary et al. (2022)

³¹Neither existing demand system estimates (Van der Beck, 2021; Koijen et al., 2023) nor estimates based on classic finance models allow for this force (Berk and Van Binsbergen, 2025).

estimate such a two-tier structure, and the nested logit model is another example.³² In such a setting, a single source of exogenous variation in the cross-section of groups permits estimation of the relative elasticity across groups.

While this symmetry assumption improves tractability, it sacrifices the heterogeneous substitution essential to finance, a feature our model is designed to capture. Consider, for instance, a setting where investors manage portfolio duration, meaning duration itself drives substitution across bonds. In this context, if one groups bonds by duration, Assumptions A1 and A2 are plausible at the asset level (using group dummies as observables), but not at the group level: substitution between 1–5 year bonds and 5–10 year bonds differs from substitution between 1–5 year bonds and 15–20 year bonds.

4 Price Impact

In some contexts, the researcher is not interested separately in the demand of each investor, but instead finds it sufficient to understand the aggregate demand for assets. The price impact of a shift in demand is the inverse of the slope of the aggregate demand curve, $\mathcal{M} = -\mathcal{E}_{agg}^{-1}$ (see Section 1.1).

In this section, we first explain how the direct estimation of price impact is different from that of demand elasticity through the simple case of one asset. Then, we show how to estimate the multiplier matrix \mathcal{M} , in particular how to do so accounting for spillovers to other assets (the off-diagonal elements of \mathcal{M}) by relying on assumptions A1 and A2.

Simple causal inference of price impact. Price impact measures how much prices change in response to an exogenous shift in demand curve. Because in equilibrium aggregate demand does not change if assets are in fixed supply, the empirical setup differs from that of demand estimation. To understand the basic intuition, start with one asset to put aside issues of substitution. While equilibrium demand is fixed, it is possible for demand curves to shift; and we are interested in measuring the impact of such a shift. An idealized example would be an investor waking up in the morning and deciding to buy one share of Apple for no specific reason. Then, the aggregate demand curve for the asset shifts to the right by one unit. In equilibrium, the price must adjust upwards to satisfy market clearing. Figure 3 illustrates this process. Similarly supply shocks can be viewed as the negative of a demand shock and be treated likewise.

In practice, the econometrician starts from a shift of the demand curve Z_i (measured as the actual amount); examples of such shocks from the literature include asset purchases by

³²See [Fang \(2023\)](#) and [Kojien and Yogo \(2024\)](#) for variations of the nested logit structure.

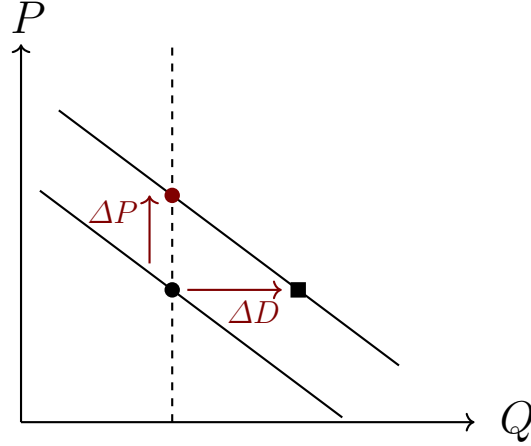


Figure 3: Equilibrium Effect of an Upward Shift in Demand Curve

central banks or rebalancing due to flows in and out of mutual funds (Lou, 2012). Armed with this shift in demand, we run the regression:

$$\Delta P_i = \widehat{\mathcal{M}}Z_i + \epsilon_i, \quad Z_i \perp \epsilon_i. \quad (42)$$

The exclusion restriction $Z_i \perp \epsilon_i$ requires that the change in demand under consideration must be orthogonal to any other demand shifts in the economy. For example, if a group of investors systematically mimicks the Fed’s asset purchases, exogeneity is violated and the regression will be biased; the measured shock undercounts the actual change in demand, and overestimates the price impact.

There is no first stage because Z_i directly measures the magnitude of the shift in the demand curve. This shift does not materialize in equilibrium quantity demanded: prices adjust so that the total quantity demanded stays equal to the fixed supply. Equivalently, this identification condition corresponds to assuming that, if one could measure quantities before prices adjust — the out-of-equilibrium square in Figure 3 — the first-stage coefficient would be one.³³

Handling substitution. Stepping away from the simplistic case of one asset and no spillover, the same issue as for the estimation of demand elasticity arises: all prices are determined together in equilibrium. There is no such thing as “the multiplier” but instead a matrix \mathcal{M} of own-demand and cross-demand multipliers such that $\Delta P = \mathcal{M}\Delta D$, with

³³If the researcher is willing to make stronger assumptions on which group of investors are affected by the demand shock, they can run a first-stage converting an instrument into a demand shock for that group. However, the stronger assumptions correspond to assuming that the demand shock for the subgroup coincides with the aggregate demand shock, what we refer to as assuming a first-stage coefficient of 1.

$\mathcal{M} = -\mathcal{E}_{agg}^{-1}$ (this is a matrix inverse, not element-by-element, hence own-price multiplier depends on both the own-price and cross-price elasticity). Notice that in a slight abuse of notation, ΔD now represents the shift of the aggregate demand curve.³⁴ This implies that in response to a vector of shocks Z , price changes will be:

$$\Delta P = \mathcal{M}Z + \epsilon, \quad (43)$$

where ϵ captures the impact of all other demand shocks. The following lemma suggests that the assumptions we introduced for elasticity also helps with the estimation of the multiplier.

Lemma 4 *The aggregate elasticity \mathcal{E}_{agg} satisfies assumptions A1 and A2 if and only if the multiplier matrix \mathcal{M} does.*

Thanks to this equivalence, we obtain the mirror image of Propositions 1 and 3 for the multiplier matrix: relative multiplier and spillovers can be estimated using classic cross-sectional shocks, and a set of portfolio-level time-series shocks, respectively.

Proposition 5 *Assume that the multiplier matrix \mathcal{M} satisfies assumptions A1 and A2, and hence can be written as $\mathcal{M} = \mathcal{M}_{relative} + X\mathcal{M}_X X'$.*

- Relative multiplier: *If the cross-sectional demand shocks $\{Z_i\}_{i \in \{1, \dots, N\}}$ satisfy the exclusion restriction $Z_i \perp \epsilon_i | X_i$, the cross-sectional regression*

$$\Delta P_i = \widehat{\mathcal{M}}Z_i + \theta' X_i + u_i \quad (44)$$

identifies the relative multiplier $\widehat{\mathcal{M}} = \mathcal{M}_{relative}$, and $\mathcal{M}_{relative} = -\mathcal{E}_{agg,relative}^{-1}$.

- Spillovers: *If the time series demand shocks for each observable $Z_{X,t} = \{Z_{k,t}\}_{k \in \{1, \dots, K\}, t \in \{1, \dots, T\}}$ satisfy the exclusion restriction $Z_{X,t} \perp \epsilon_{X,t}$, the set of time-series regressions*

$$\Delta P_{X,t} = \widetilde{\mathcal{M}}Z_{X,t} + u_{X,t} \quad (45)$$

identifies the substitution matrix, with $\mathcal{M}_X = (\widetilde{\mathcal{M}} - \widehat{\mathcal{M}}I_K)(X'X)^{-1}$.

The relative multiplier measures how the price of one asset relative to another comparable one changes if the relative demand for these assets shifts—how much does the price of Ford change relative to the price of General Motors if the demand for Ford changes relative to the

³⁴That is, ΔD in this section corresponds to the demand shifts aggregated across investors $\epsilon_{agg} = \sum_j \epsilon_j$ in the notation of the previous sections.

demand for General Motors?³⁵ Interestingly, in this case the relative multiplier coincides with the inverse of the relative elasticity, $\mathcal{M}_{relative} = -\mathcal{E}_{relative}^{-1}$.³⁶ In contrast, the cross multipliers captured by the matrix \mathcal{M}_X measure spillovers—how much an increase in demand for one asset affects the price of another asset?

Of course, the researcher should present arguments for the validity of assumptions A1 and A2. The approaches to make this case that we discuss in Section 2 are also relevant here. The next section provides a practical example for the corporate bond market.

5 Empirical Example: Corporate Bond Multipliers

We apply our framework to corporate bonds and estimate the multiplier matrix for that market. First, we choose a source of variation in demand and observables for assumptions A1 and A2, discuss their validity, and conduct pre-estimation diagnostics. Second, we estimate the two components of price impact: relative multiplier and substitution. We conduct a series of specification tests to further evaluate the validity of our setup. Last, we use the estimates to inform counterfactual questions about various designs of central bank interventions in corporate bond markets.

5.1 Identification Strategy and Its Plausibility

Identification of the multiplier relies on combining a source of shifts in demand with assumptions about the structure of the multiplier matrix (A1 and A2). We describe these choices, and discuss their plausibility.

Data and source of shifts in demand. We study the market for corporate bonds, as in [Chaudhary et al. \(2022\)](#). We obtain data on individual corporate bond returns from the WRDS Bond Returns database between 2010 and 2024 and mutual fund bond holdings and flows from the CRSP Survivor-Bias-Free US Mutual Fund Database. As a source of shifts in demand, we use a refinement of flow-induced trading by mutual funds introduced by [Coval and Stafford \(2007\)](#) and [Lou \(2012\)](#). Flow-induced trading aims to capture a mechanical component of how mutual funds adjust their positions in response to flows. The idea is that funds tend to allocate new inflows towards their larger existing positions, independent of new

³⁵A benefit of focusing on multipliers: the estimation does not become degenerate when the price of two assets with identical payoffs remain equal. In this setting, the relative multiplier is 0 and the regression can still identify it.

³⁶This conjunction occurs despite neither own-price and the cross-price elasticities being stable by inversion ($\mathcal{M}_{ij} \neq -1/\mathcal{E}_{ij}$); inverting a matrix is different from inverting each of its elements. [Gabaix and Koijen \(2021\)](#) derive this inversion result for the special case without observables shaping substitution.

arrival of information about these positions. The instrument aggregates these shocks across funds:

$$Z_{it} = \sum_k \frac{A_{k,t-1} \tilde{w}_{i,k,t-1}}{P_{i,t-1} S_{i,t-1}} f_{kt}, \quad (46)$$

where $A_{k,t-1}$ is fund k 's assets under management, $\tilde{w}_{i,k,t-1}$ is the portfolio weight of k on asset i residualized by removing seven common factors from raw weights following [Chaudhry \(2025\)](#), f_{kt} is the flow into fund k at date t , and $P_{i,t-1} S_{i,t-1}$ is the total supply of bond i at $t - 1$ (the product of its price and quantity outstanding).

Exogeneity assumption. Exogeneity requires the changes in demand measured by the instrument to be unrelated to any unobserved source of demand — equation (42). Overall mutual fund flows generically do not satisfy this assumption, but the instrument can. Specifically, equation (46) has the structure of a Bartik shift-share instrument: fund flows f_{kt} are the “shifts” and the residualized portfolio weights $\tilde{w}_{i,k,t-1}$ are the “shares.” [Goldsmith-Pinkham et al. \(2020\)](#) show that exogeneity of the shares is sufficient for identification in this setting, without requiring exogenous shifts. That is, exogeneity is satisfied if the weights $\tilde{w}_{i,k,t-1}$ are uncorrelated with unobserved demand shocks for bond i at time t .

The residualization of [Chaudhry \(2025\)](#) is crucial to argue that the weights are exogenous. If instead, one used actual portfolio shares, the condition could be violated for a couple reasons. First, households might substitute from holding bonds directly to holding them through mutual funds. Bonds with high mutual fund shares would then mechanically be bonds experiencing offsetting decreases in direct household demand, so that shares correlate with an unobserved demand shift. Second, bonds with high mutual fund ownership might also attract demand from other institutional investors, such as insurance companies, pursuing similar strategies. Then portfolio shares would correlate with unobserved demand from these other investors, biasing the estimated price impact.³⁷ Empirically, raw portfolio weights load on common fund characteristics such as growth, size, and income style, suggesting that fund strategies are more likely to overlap with those of other institutions. Removing many common factors isolates idiosyncratic variation in portfolio composition that is less likely to be shared with other investor types, making exogeneity more plausible. For the estimation of spillovers, exogeneity must also hold in the time series after aggregation along the observables we introduce next.

³⁷Because the shares are lagged, violation requires persistent correlation between portfolio composition and unobserved demand—less likely than a contemporaneous one, but not ruled out when allocations and demand patterns are slow-moving (e.g., [Huebner, 2024](#)).

Structure of substitution. Our estimation strategy also relies on homogeneous substitution conditional on a set of observables (assumptions A1 and A2).³⁸ The simplest approach would be to assume unconditional homogeneous substitution. This is unappealing because, if the demand for long-term bonds rises, this will affect the price of other long-term bonds differently from the price of short-term bonds. As discussed in Section 2.2.3, a simple diagnostic for the plausibility of this concern is a test of balance on covariances: do the treated bonds comove in the same way with broad portfolios as control bonds? If they do not, investors would likely substitute these bonds differently.

Figure 4 suggests that they do not. For a given date, we form a long-short portfolio based on whether the instrument Z_{it} is above or below median on that date and compute the beta of this portfolio on a series of broad indices in a two-year range around the date of the sort — the blue dot for that date. Each panel corresponds to a different index: a broad bond index, long-short portfolios based on credit ratings and maturity, and a broad stock index. Bonds with a high instrumented inflow comove more strongly with the credit-rating-sorted portfolio and more weakly with the broad bond index and the duration-sorted portfolio: the blue dots in these panels are systematically away from zero.

It is more plausible to assume homogeneous substitution conditional on duration and credit rating, the two variables that are the most salient characteristics of a corporate bond. Hereafter, we standardize these variables to be mean zero and variance one for each date.³⁹ We revisit our diagnostic for this case: we form portfolios sorted on $Z_{idio,it}$, the residualized instrument with respect to both duration and credit rating. This corresponds to the orange dots in Figure 4, which are much closer to zero, bolstering the plausibility of our assumptions. After the estimation, we provide tests of overidentifying restrictions to assess further whether our framework is misspecified.

Finally, regarding Assumption A2, Appendix Figure 11 shows a similar picture for idiosyncratic volatilities of bonds with high and low values of the instrument. Absent controls for duration and credit rating, these bonds have different idiosyncratic volatilities; after controlling, the volatilities are close to identical.

5.2 Estimates

With these assumptions, we can identify the multiplier matrix by following Proposition 5. For units of change in prices, we use total returns R_{it} , the standard measure in asset pricing

³⁸Note that while we have not emphasized it in the discussion above, the exclusion restriction is also conditional on observables.

³⁹We use credit ratings as a continuous variable by assigning the historical average 10-year default probabilities to each rating category, from Table 65, column D, of [Vazza and Kraemer \(2013\)](#).



Figure 4: **Balance on covariances: factor exposure of long-short portfolios sorted on the instrument.** Each dot represents the factor exposure of a portfolio sorted on the values of the instrument Z_{it} (blue) or the residualized instrument $Z_{idio,it}$ at the corresponding month between 2011:04 and 2023:03. Equation (46) defines Z_{it} , while $Z_{idio,it}$ is obtained by cross-sectionally orthogonalizing the instrument at each date to duration and credit risk (long-run average default probabilities based on S&P ratings). At each date, we form long-short equal-weighted portfolios based on whether the value of the instrument for that bond is above or below the median. We compute the returns of these portfolios over two years centered around t , excluding t , and regress separately these returns on four aggregate factors to obtain the factor loadings represented on the figure. Panel A reports exposures to an aggregate investment-grade corporate bond factor, the total return of the ICE BofA US Corporate Index. Panel B considers a credit-sorted factor, the difference in total return between the ICE BofA US High Yield Index and the ICE BofA US Corporate Index. Panel C uses a duration-sorted factor, the difference in total return between the ICE BofA 15+ Year US Corporate Index and the ICE BofA 1-3 Year US Corporate Index. All bond indices are obtained from FRED. Panel D uses the [Fama and French \(1993\)](#) excess stock market return, obtained from Kenneth French’s data library.

with well-behaved statistical properties.⁴⁰ We estimate the relative multiplier $\widehat{\mathcal{M}}$ with the

⁴⁰It might be tempting to “divide” prices by the observables to remove the role of these variables, e.g. working with changes in yields as opposed to returns. When multiple dimensions drive substitution, no

cross-sectional estimator:

$$R_{it} = \widehat{\mathcal{M}} Z_{it} + \theta_{0,t} + \theta_{DUR,t} DUR_{it} + \theta_{PD,t} PD_{it} + u_{it}. \quad (47)$$

The spillover matrix $\widetilde{\mathcal{M}}$ is obtained from the time-series regressions:

$$R_{agg,t} = \widetilde{\mathcal{M}}_{1,1} Z_{agg,t} + \widetilde{\mathcal{M}}_{1,2} Z_{DUR,t} + \widetilde{\mathcal{M}}_{1,3} Z_{PD,t} + u_{agg,t} \quad (48)$$

$$R_{DUR,t} = \widetilde{\mathcal{M}}_{2,1} Z_{agg,t} + \widetilde{\mathcal{M}}_{2,2} Z_{DUR,t} + \widetilde{\mathcal{M}}_{2,3} Z_{PD,t} + u_{DUR,t} \quad (49)$$

$$R_{PD,t} = \widetilde{\mathcal{M}}_{3,1} Z_{agg,t} + \widetilde{\mathcal{M}}_{3,2} Z_{DUR,t} + \widetilde{\mathcal{M}}_{3,3} Z_{PD,t} + u_{PD,t}, \quad (50)$$

where the subscripts *agg*, *DUR*, and *PD* represent aggregation over all bonds, duration, and credit risk respectively, as in equation (33) (each period, construct $Z_{X,t} = (X'X)^{-1}X'Z_t$).

Table 1: **Relative multiplier $\widehat{\mathcal{M}}$ in corporate bonds**

	Return R_{it}				
	(1)	(2)	(3)	(4)	(5)
<i>Demand shock:</i>					
Z_{it}	0.055 (0.084)			0.492*** (0.128)	2.080*** (0.387)
$Z_{idio,it}$		0.055 (0.084)	0.055 (0.087)		
Date Fixed Effects	Yes	Yes	Yes	Yes	
Duration \times Date Fixed Effects	Yes	Yes			
Credit Risk \times Date Fixed Effects	Yes	Yes			
N	1,041,985	1,041,985	1,041,985	1,041,985	1,041,985
R^2	0.464	0.464	0.293	0.293	0.018

Table 1 reports the results of relative multiplier regressions of bond returns R_{it} on demand shocks Z_{it} and $Z_{idio,it}$ for U.S. non-defaulted corporate bonds. Specifications (1) and (4)–(5) use the flow-induced trading demand shock Z_{it} defined in Equation (46). Specification (1) includes a time fixed effect and controls for a continuous duration variable and a continuous credit risk variable based on average historical default probabilities for each S&P credit rating category, while specification (4) uses only date fixed effects and specification (5) only a common intercept. Specifications (2)–(3) use the demand shock $Z_{idio,it}$ orthogonalized to duration and credit risk each period, with and without controlling for duration and credit risk in the regression. The regressions weigh each date equally. The sample period is 2010:04 to 2024:03. Standard errors are clustered by date and bond issuer.

single change of units can make substitution homogeneous. For example, Appendix J reports counterparts of our estimates for yields, and finds that heterogeneity with respect to duration is mitigated in these units, credit ratings remain relevant.

Relative multiplier. Table 1 estimates the relative multiplier. Column 1 implements the regression of equation (47) to find a relative multiplier $\widehat{\mathcal{M}} \approx 0$, a tightly estimated zero. The estimate suggests that if the demand for one bond relative to another one with the same credit rating and duration increases by 1%, their relative bond price does not change much (the insignificant point estimate corresponds to a response of 5.5bps), as in [Chaudhary et al. \(2022\)](#).

Columns 2 and 3 regress directly the change in price on the residualized instrument $Z_{idio,it}$, with and without controls for the characteristics. This leads to the exact same estimate of the relative multiplier, a mathematical property independent of the specific dataset. This observation highlights that the source of variation for the estimates is variation in the residual component of the instrument $Z_{idio,it}$.

Column 4 estimates the model without controlling for duration and credit risk. This leads to a sizable departure from our baseline estimates, suggesting it is important to account for substitution. Column 5 finds further departure when removing the date fixed effects.

We can also test whether our specification misses some dimensions of substitution. If our model is well specified, including controls for additional observables or changing the level of variation of the data ([Hausman, 1978](#); [Mundlak, 1978](#)) would not alter the estimated relative multiplier. Appendix D.2 confirms that controlling for bond liquidity or focusing on within-issuer variation does not.

Spillover matrix. Table 2 presents estimates of the matrix $\widetilde{\mathcal{M}}$ which encodes the meso- and macro multipliers. Columns 1–3 implement equations (48)–(50). Combining these estimates yields

$$\widetilde{\mathcal{M}} = \begin{pmatrix} 14.430 & 0.425 & 1.127 \\ 7.773 & 4.531 & -0.868 \\ 3.764 & -5.789 & 3.383 \end{pmatrix}, \quad (51)$$

and the implied spillover matrix is $\mathcal{M}_X = (\widetilde{\mathcal{M}} - \widehat{\mathcal{M}}\mathbf{I}_K)(\mathbf{X}'\mathbf{X})^{-1}$.

Column 4 shows that when using the full panel, rather than time-series regressions, one identifies the exact same meso- and macro-multipliers as well as the relative multiplier (coefficient on $Z_{idio,it}$).⁴¹ This confirms that the panel carries no additional information on spillovers, and underscores the separation result inherent to our structure, whereby the cross-section identifies the relative multiplier, while the time-series identifies spillovers.

The first term in the estimated matrix, $(\widetilde{\mathcal{M}})_{11} = 14.430$, is the aggregate multiplier: a

⁴¹This property is not specific to the aggregate returns in column 1, but applies more generally to observable-based portfolios in columns 2 and 3. Appendix Table 9 reports this extended specification.

Table 2: **Macro- and meso multipliers in corporate bonds**

	Return				
	$R_{agg,t}$	$R_{DUR,t}$	$R_{PD,t}$	R_{it}	$R_{agg,t}$
	(1)	(2)	(3)	(4)	(5)
$Z_{agg,t}$	14.430*** (2.707)	7.773*** (1.803)	3.764** (1.392)	14.430*** (2.689)	15.032*** (2.393)
$Z_{DUR,t}$	0.425 (6.713)	4.531 (4.293)	-5.789 (3.948)	0.425 (6.646)	
$Z_{PD,t}$	1.127 (1.866)	-0.868 (1.404)	3.383** (1.025)	1.127 (1.854)	
$Z_{idio,it}$				0.055 (0.087)	
N	168	168	168	1,041,985	168
R^2	0.343	0.173	0.384	0.100	0.342

Table 2 reports the results of macro- and meso multiplier regressions of bond returns on demand shocks for non-defaulted U.S. corporate bonds. Specifications (1)–(3) jointly estimate the matrix $\widetilde{\mathcal{M}}$ from Proposition 5, which together with the relative multiplier $\widehat{\mathcal{M}}$ determines the spillover matrix between observables, \mathcal{M}_X . Specification (4) estimates multipliers mechanically identical to specification (1) using disaggregated, repeated cross-sectional regressions, while adding the relative multiplier $\widehat{\mathcal{M}}$. Specification (5) estimates the macro multiplier in isolation by regressing aggregate bond returns $R_{agg,t}$ on the aggregated instrument $Z_{agg,t}$ in the time series. The $K = 3$ observables are a vector of ones, standardized duration, and standardized credit risk. The sample period is 2010:04 to 2024:03. Robust standard errors are used for specifications (1) to (3) and (5). For specification (4), standard errors are clustered by date and issuer, and regressions are weighted such that each date receives equal weight.

uniform 1% increase in the demand for all bonds leads to a 14% increase in bond prices.⁴² The other terms in the first row characterize the average impact of a shock to the relative demand of long-term vs. short-term bonds, or risky vs. safe bonds. For example, a 1% increase in the relative demand for risky bonds leads to an insignificant 1.13% increase in bond prices. The second and third rows of $\widetilde{\mathcal{M}}$ capture heterogeneous impacts of demand shocks in the cross section.

Figure 5 builds intuition around these numbers. Panel A reports the impact of an aggregate demand shock on the returns of bonds of various duration. The blue line and shaded area represent estimates and standard errors from an estimate constrained to have constant effects. The red line is our main estimate: a line combining the level $(\widetilde{\mathcal{M}})_{11} = 14.430$ with the slope $(\widetilde{\mathcal{M}})_{21} = 7.773$. Long-duration bonds respond more to an aggregate demand shock

⁴²The current literature does not provide estimates at this level of aggregation; existing estimates of portfolio-level multipliers find values between 2.3 and 6 (e.g., [Bouveret and Yu, 2021](#); [Chaudhary et al., 2022](#); [Darmouni et al., 2023](#)); see [Haddad et al. \(2025\)](#) for a review.

than short-duration bonds, with multipliers ranging from 5 to 25. This heterogeneity is statistically significant, the red line is different from the blue line. Panel B shows a similar pattern for the response of bonds with different credit risk. The slope of the red line is now $(\widetilde{\mathcal{M}})_{31} = 3.764$, again statistically and economically significant.

Panels C and D report how demand shocks to long-short portfolios formed on duration and credit risk affect the corresponding portfolio returns, respectively. For duration (panel C), there is a clear monotonic pattern whereby riskier bonds are affected more, but differences are statistically insignificant. For credit risk (panel D), the pattern is similar, and the results are statistically significant.⁴³

Finally, we can also assess the validity of our linear specification for spillovers. We estimate the impact of the three demand shocks separately for bonds in different duration or credit risk buckets, the green dots in Figure 5.⁴⁴ Under the null of our model, the estimates should line up with the red line. Linearity holds well across all four panels, suggesting that these overidentifying restrictions stand.

In Section 3.2.3, we pointed out the potential bias when estimating the macro multiplier without accounting for demand shocks along other dimensions. Excluding the instruments Z_{DUR} and Z_{PD} can result in omitted variable bias, if they correlate with Z_{agg} . Empirically, aggregate flows from mutual funds are not uniform across bonds; they tilt towards short-term bonds (Z_{agg} and Z_{DUR} have a correlation of -0.41). Column 5 implements a regression of aggregate returns on aggregate demand only, similar to [Gabaix and Koijen \(2021\)](#). We find an estimate close to our baseline, suggesting the bias is not a concern in this setting.

5.3 Counterfactual Analysis

The estimated multiplier matrix provides answers to many interesting counterfactual questions, with different parts of \mathcal{M} informing outcomes of different policies. For illustration, consider different implementations of a hypothetical Federal Reserve corporate bond purchase program. From the estimates of the previous section, the full multiplier matrix can be constructed as $\mathcal{M} = \widehat{\mathcal{M}}\mathbf{I} + X\mathcal{M}_X X'$. Applying this matrix to a vector of interventions yields the full cross-section of price responses.

Relative counterfactuals. Imagine the Federal Reserve conducts an asset purchase program in corporate bonds. Instead of buying all bonds, it randomly selects a subset to purchase. This mirrors some aspects of the Fed’s actual implementation of QE in Treasury

⁴³The last two dimensions, the impact of duration shocks on prices of credit and vice versa, are in Appendix Figure 12.

⁴⁴The exclusion restriction underlying these more granular estimates is more stringent because it applies to each bucket individually.

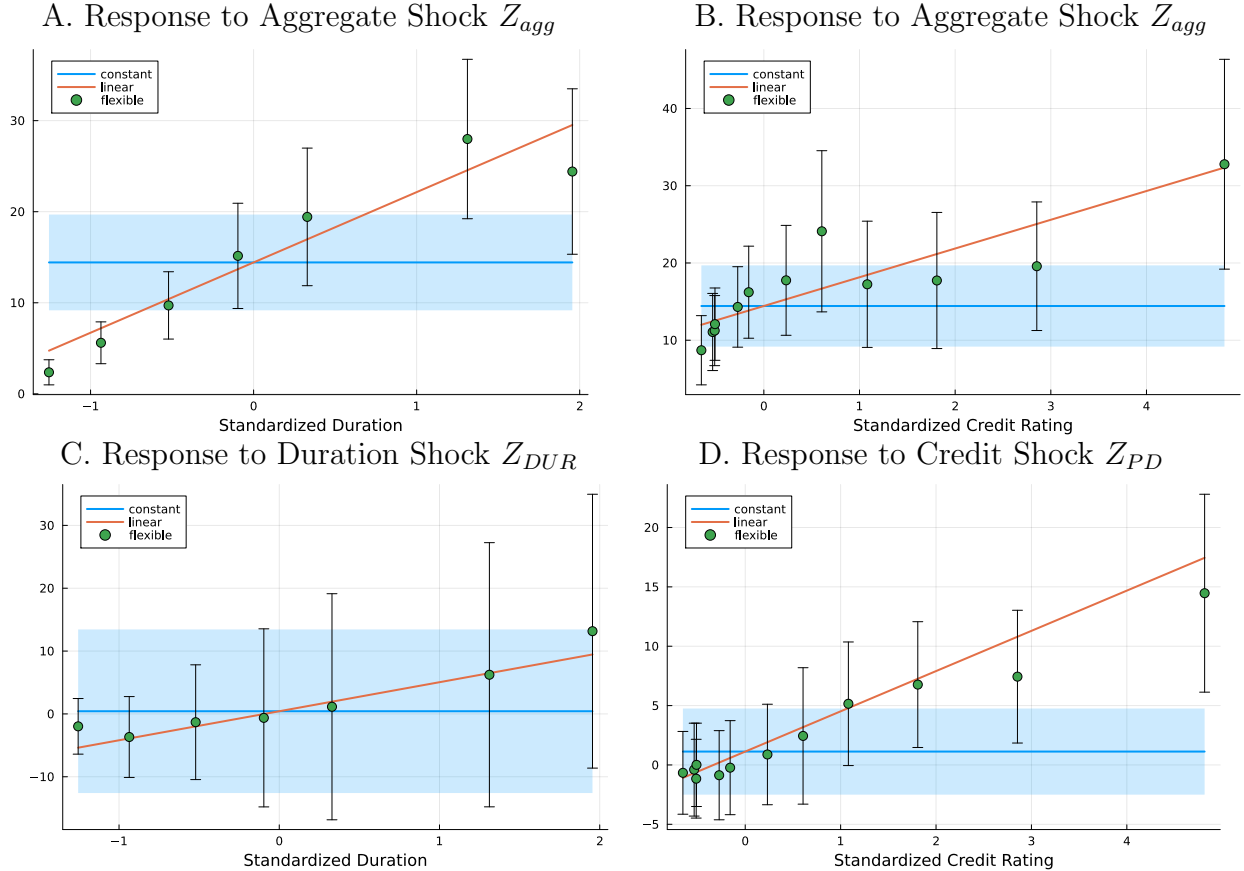


Figure 5: Macro- and meso multipliers in the cross-section. We report the response of portfolios of corporate bonds to aggregate demand shocks Z_{agg} (Panels A and B, in the cross-sections of duration and credit rating) and shocks along duration Z_{DUR} (Panel C) and Z_{PD} (Panel D). Bonds are grouped in seven buckets based on duration: <1 year, 1–3 years, 3–5 years, 5–7 years, 7–10 years, 10–15 years, and 15+ years. Bonds are also grouped by S&P credit rating, with individual notches from A+ through B-, and with AAA/AA and CCC/C ratings pooled at the extremes. The blue lines correspond to the estimates from column (2) of Table 2, which assume identical responses. The orange lines are based on columns (3) and (4), which include linear interaction terms with either duration or credit risk, neutralizing the other. The green dots estimate multipliers separately by duration or credit rating bucket in a pooled panel regression. The sample period is 2010:04 to 2024:03.

securities (Selgrad, 2023). How will these purchases affect the prices of purchased bonds compared to control bonds with the same duration? The answer depends only on the relative multiplier $\widehat{\mathcal{M}}$. With an estimate of $\widehat{\mathcal{M}} \approx 0$, the relative price of the two bonds would not change: substitutability between bonds with the same duration is so strong that buying one bond but not the other moves their prices by the same amount. However, this micro-level comparison does not reveal the total impact of a broad purchase program.

Macro counterfactuals. Consider what happens if the Federal Reserve decides to purchase all corporate bonds. How much would prices change? Unlike in [Gabaix and Koijen \(2021\)](#), the answer is not simply characterized by one macro multiplier; it depends on bond duration and credit rating. For example, consistent with a factor structure based on duration (e.g., [Vayanos and Vila, 2021](#)), long-duration bond prices would move more. If the Fed buys 1% of the supply of all corporate bonds, the price of bond i with observables X_i increases by $X_i' \widetilde{\mathcal{M}}(1, 0, 0)'$ percent.⁴⁵ The price of a BBB rated corporate bond with duration of 1 year increases by 5.3%, while another bond with the same credit rating and a duration of 10 years increases by 18.7%. Effects also increase with credit risk; the price of a 3-year AA-bond increases by 6.4%, while that of a B bond increases by 21.8%.⁴⁶

Meso-macro counterfactuals. The Secondary Market Corporate Credit Facility (SM-CCF), announced in March 2020, restricted purchases to investment-grade bonds with at most five years to maturity — 36% of bonds were eligible. Suppose the Fed buys 1% of the supply of each eligible bond. This intervention raises prices by 2.92% on 1-year BBB bonds and 5.26% on 10-year BBB bonds; for 3-year AA- bonds the impact is 2.08% and for 3-year B bonds 13.27%.⁴⁷

These numbers are mechanically smaller than in the previous counterfactual because the Fed only purchases 36% of the bonds; the cross-sectional patterns differ as well. Strikingly, ineligible bonds with high credit-risk exposure see disproportionate impact: 3-year B bonds reach more than 60% of their macro response without being purchased by the Fed. The Fed’s concentration on short-term bonds implicitly tilts demand against long duration, which through cross-asset substitution transmits along the credit-risk dimension. Our substitution structure captures this pattern through the off-diagonal entry $\widetilde{\mathcal{M}}_{32} = -5.79$. The duration-credit risk spillover channel of [Li \(2025\)](#) offers a foundation for it: when demand tilts away from long-term bonds the term premium increases, which—through life insurers’ negative post-GFC duration gap—boosts their risk-bearing capacity and their demand for risky corporate bonds, compressing credit spreads.

⁴⁵The calculation simplifies in the case where the demand shock is proportional to the observables, $\Delta D = Xd$, then $\mathcal{M}\Delta D = \mathcal{M}Xd = (\widehat{\mathcal{M}}X + X\mathcal{M}_X X'X)d = (\widehat{\mathcal{M}}X + X(\widetilde{\mathcal{M}} - \widehat{\mathcal{M}}\mathbf{I}_K))d = X\widetilde{\mathcal{M}}d$.

⁴⁶In the March 2020 cross-section, a 1-year BBB bond has a standardized credit rating of -0.141 and standardized duration of -1.105 . An 1% aggregate demand shock moves its price by $14.430 + (-1.105) \times 7.773 + (-0.141) \times 3.764 \approx 5.3\%$.

⁴⁷Instead of using the full multiplier matrix \mathcal{M} , an intuitive approximation projects the demand shock on the observables: a positive intercept of 0.36 equal to the share of bonds purchased, a negative duration tilt of -0.34 , and a negative credit-risk tilt of -0.18 . The impact is then given by $X_i' \widetilde{\mathcal{M}}(0.36, -0.34, -0.18)'$ which we find to be within 2bps of the exact counterfactual.

6 Concluding Remarks

This paper provides a framework for using causal inference with asset prices and quantities. Specifically, we provide conditions for estimating the elasticity matrix of asset demand accounting for the natural spillovers that exist between assets when making portfolio choices. Our main assumption is homogeneous substitution conditional on observables: two assets with the same observables are comparable if the demand for them responds in the same way to the price of every other asset. We show that this condition maps naturally to restrictions often imposed in standard asset pricing models, and also provide guidelines to design experiments satisfying this condition and assess its plausibility in the data.

When our conditions hold, the standard cross-sectional difference-in-difference or instrumental variable approach identifies the relative elasticity between comparable assets—that is, the difference between their own-price and cross-price elasticity. Substitution is identified using a set of time series regressions on portfolios sorted on the observables. Identifying both relative elasticity and substitution is crucial to answer many counterfactual questions such as the response to shocks across broad categories of assets. These simple tools and principles offer a straightforward package for researchers wanting to use natural experiments to better understand investment decisions and their equilibrium impact.

Because our conditions are flexible, they can guide empirical design without having to take a strong stance on a specific structural model. Still, these causal estimates should only be a first step towards a deeper understanding of how investors and institutions make portfolio decisions, and how those decisions shape equilibrium prices. For example, in richer models, it becomes interesting to understand how demand elasticities change and respond to outside forces such as an increase in passive investing (Haddad et al., 2024c), the Fed’s support programs (Haddad et al., 2025; Jansen et al., 2024), persistent selling pressures (He et al., 2025), and the broader market macrostructure or other macroeconomics conditions.

References

- Admati, Anat R.**, “A Noisy Rational Expectations Equilibrium for Multi-Asset Securities Markets,” *Econometrica*, 1985, *53* (3), 629–657.
- Aghaee, Alireza**, “The Flattening Demand Curves,” *Working paper*, 2024.
- Allen, Jason, Jakub Kastl, and Milena Wittwer**, “Estimating demand systems with bidding data,” *Available at SSRN 5171755*, 2018.
- An, Yu and Amy W. Huber**, “Demand Propagation Through Traded Risk Factors,” Technical Report 2025.
- , **Yinan Su, and Chen Wang**, “Quantity, Risk, and Return,” *Working Paper*, 2024.
- Anderson, Simon P, André De Palma, and J-F Thisse**, “A representative consumer theory of the logit model,” *International Economic Review*, 1988, pp. 461–466.
- Andrews, Isaiah, James H Stock, and Liyang Sun**, “Weak instruments in instrumental variables regression: Theory and practice,” *Annual Review of Economics*, 2019, *11* (1), 727–753.
- Angrist, Joshua D.**, “Estimating the Labor Market Impact of Voluntary Military Service Using Social Security Data on Military Applicants,” *Econometrica*, 1998, *66* (2), 249–288.
- **and Jörn-Steffen Pischke**, *Mostly Harmless Econometrics: An Empiricist’s Companion*, Princeton University Press, 2009.
- Angrist, Joshua D, Guido W Imbens, and Donald B Rubin**, “Identification of Causal Effects Using Instrumental Variables,” *Journal of the American Statistical Association*, 1996, *91* (434), 444–455.
- Athey, Susan and Guido W Imbens**, “The state of applied econometrics: Causality and policy evaluation,” *Journal of Economic perspectives*, 2017, *31* (2), 3–32.
- Bai, Jennie and Pierre Collin-Dufresne**, “The CDS-bond basis,” *Financial Management*, 2019, *48* (2), 417–439.
- Becker, Gary S.**, “Irrational Behavior and Economic Theory,” *Journal of Political Economy*, 1962, *70* (1), 1–13.
- Ben-David, Itzhak, Jiacui Li, Andrea Rossi, and Yang Song**, “Ratings-driven demand and systematic price fluctuations,” *The Review of Financial Studies*, 2022, *35* (6), 2790–2838.
- Berg, Tobias, Markus Reisinger, and Daniel Streitz**, “Spillover effects in empirical corporate finance,” *Journal of Financial Economics*, 2021, *142* (3), 1109–1127.
- Berk, Jonathan B and Jules H Van Binsbergen**, “The impact of impact investing,” *Journal of Financial Economics*, 2025, *164*, 103972.

- Berry, Steven, James Levinsohn, and Ariel Pakes**, “Automobile Prices in Market Equilibrium,” *Econometrica*, 1995, *63* (4), 841–890.
- Binsbergen, Jules H Van, William F Diamond, and Marco Grotteria**, “Risk-free interest rates,” *Journal of Financial Economics*, 2022, *143* (1), 1–29.
- Bound, John, David A Jaeger, and Regina M Baker**, “Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak,” *Journal of the American statistical association*, 1995, *90* (430), 443–450.
- Bouveret, Antoine and Jie Yu**, “Risks and vulnerabilities in the US bond mutual fund industry,” Technical Report 2021.
- Bretscher, Lorenzo, Lukas Schmid, Ishita Sen, and Varun Sharma**, “Institutional corporate bond pricing,” *Swiss Finance Institute Research Paper*, 2022, *21-07*.
- Campbell, John Y.**, *Financial Decisions and Markets: A Course in Asset Pricing*, Princeton University Press, 2017.
- **and Luis Viceira**, *Strategic Asset Allocation: Portfolio Choice for Long-Term Investors*, first ed., Oxford, UK: Oxford University Press, 2002.
- Chang, Yen-Cheng, Harrison Hong, and Inessa Liskovich**, “Regression Discontinuity and the Price Effects of Stock Market Indexing,” *The Review of Financial Studies*, 07 2014, *28* (1), 212–246.
- Chaudhary, Manav, Zhiyu Fu, and Jian Li**, “Corporate bond multipliers: Substitutes matter,” Technical Report 2022.
- Chaudhry, Aditya**, “The Impact of Prices on Analyst Cash Flow Expectations: Reconciling Subjective Beliefs Data with Rational Discount Rate Variation,” *Journal of Financial Economics*, 2025, *forthcoming*.
- Chen, Hui, Zhuo Chen, Zhiguo He, Jinyu Liu, and Rengming Xie**, “Pledgeability and asset prices: Evidence from the Chinese corporate bond markets,” *The Journal of Finance*, 2023, *78* (5), 2563–2620.
- Chodorow-Reich, Gabriel, Plamen T Nenov, and Alp Simsek**, “Stock market wealth and the real economy: A local labor market approach,” *American Economic Review*, 2021, *111* (5), 1613–1657.
- Cochrane, John H**, *Asset pricing*, Princeton University Press, 2005.
- Conley, Timothy G, Christian B Hansen, and Peter E Rossi**, “Plausibly exogenous,” *Review of Economics and Statistics*, 2012, *94* (1), 260–272.
- Coppola, Antonio**, “In safe hands: The financial and real impact of investor composition over the credit cycle,” *The Review of Financial Studies*, 2025, *forthcoming*.

- Coval, Joshua and Erik Stafford**, “Asset fire sales (and purchases) in equity markets,” *Journal of Financial Economics*, 2007, 86 (2), 479–512.
- Daniel, Kent and Sheridan Titman**, “Evidence on the characteristics of cross sectional variation in stock returns,” *the Journal of Finance*, 1997, 52 (1), 1–33.
- Darmouni, Olivier, Kerry Siani, and Kairong Xiao**, “Nonbank Fragility in Credit Markets: Evidence from a Two-Layer Asset Demand System,” Technical Report 2023.
- Davis, Carter**, “The Elasticity of Quantitative Investment,” *The Review of Financial Studies*, 2024, *forthcoming*.
- , **Mahyar Kargar, and Jiacui Li**, “Why Do Portfolio Choice Models Predict Inelastic Demand?,” *Journal of Financial Economics*, 2025, *forthcoming*.
- de Fraise, Antoine Hubert**, “Crowding Out Long-Term Corporate Investment: The Role of Long-Term Government Debt Supply,” 2022.
- Deaton, Angus and John Muellbauer**, “An almost ideal demand system,” *The American economic review*, 1980, 70 (3), 312–326.
- Debreu, Gerard**, *Theory of value: An axiomatic analysis of economic equilibrium*, Vol. 17, Yale University Press, 1959.
- der Beck, Philippe Van**, “Flow-driven ESG returns,” *Swiss Finance Institute Research Paper*, 2021, 21-71.
- Drechsler, Itamar, Alexi Savov, Philipp Schnabl, and Dominik Supera**, “Monetary Policy and the Mortgage Market,” Technical Report, Working paper 2024.
- Du, Wenxin, Alexander Tepper, and Adrien Verdelhan**, “Deviations from covered interest rate parity,” *The Journal of Finance*, 2018, 73 (3), 915–957.
- Duffie, Darrell**, *Dynamic asset pricing theory*, Princeton University Press, 2010.
- Evans, Martin D. D. and Richard K. Lyons**, “Order Flow and Exchange Rate Dynamics,” *Journal of Political Economy*, 2002, 110 (1), 170–180.
- Fama, Eugene F and James D MacBeth**, “Risk, return, and equilibrium: Empirical tests,” *Journal of political economy*, 1973, 81 (3), 607–636.
- **and Kenneth R French**, “Common risk factors in the returns on stocks and bonds,” *Journal of financial economics*, 1993, 33 (1), 3–56.
- Fang, Chuck**, *Monetary policy amplification through bond fund flows*, University of Pennsylvania, 2023.
- **and Kairong Xiao**, “Dissecting Bond Market Transmission of Monetary Policy,” *Available at SSRN 5025417*, 2024.

- Froot, Kenneth A. and Tarun Ramadorai**, “Currency Returns, Intrinsic Value, and Institutional-Investor Flows,” *The Journal of Finance*, 2005, *60* (3), 1535–1566.
- Fuchs, William, Satoshi Fukuda, and Daniel Neuhann**, “Demand-System Asset Pricing: Theoretical Foundations,” *Available at SSRN 4672473*, 2025.
- Gabaix, Xavier and Ralph SJ Koijen**, “In search of the origins of financial fluctuations: The inelastic markets hypothesis,” Technical Report, National Bureau of Economic Research 2021.
- **and** –, “Granular instrumental variables,” *Journal of Political Economy*, 2024, *132* (7), 000–000.
- Goldsmith-Pinkham, Paul, Isaac Sorkin, and Henry Swift**, “Bartik Instruments: What, When, Why, and How,” *American Economic Review*, August 2020, *110* (8), 2586–2624.
- , **Peter Hull, and Michal Kolesár**, “Contamination bias in linear regressions,” *American Economic Review*, 2024, *114* (12), 4015–4051.
- Gompers, Paul A. and Andrew Metrick**, “Institutional Investors and Equity Prices*,” *The Quarterly Journal of Economics*, 02 2001, *116* (1), 229–259.
- Graves, Daniel**, “What Lies Beneath Zero: Censoring, Demand Estimation, and Hidden Beliefs,” Technical Report, Working paper 2025.
- Greenwood, Robin and Dimitri Vayanos**, “Bond supply and excess bond returns,” *The Review of Financial Studies*, 2014, *27* (3), 663–713.
- **and Marco Sammon**, “The Disappearing Index Effect,” *Working paper*, 2024.
- Greenwood, Robin Marc and Annette Vissing-Jorgensen**, “The impact of pensions and insurance on global yield curves,” Technical Report, Harvard Business School 2018.
- Greenwood, Robin, Samuel G Hanson, and Gordon Y Liao**, “Asset Price Dynamics in Partially Segmented Markets,” *The Review of Financial Studies*, 04 2018, *31* (9), 3307–3343.
- Guren, Adam, Alisdair McKay, Emi Nakamura, and Jón Steinsson**, “What Do We Learn from Cross-Regional Empirical Estimates in Macroeconomics?,” *NBER Macroeconomics Annual*, 2021, *35*, 175–223.
- Haddad, Valentin, Alan Moreira, and Tyler Muir**, “When selling becomes viral: Disruptions in debt markets in the COVID-19 crisis and the Fed’s response,” *The Review of Financial Studies*, 2021, *34* (11), 5309–5351.
- , – , **and** –, “Asset purchase rules: How QE transformed the bond market,” Technical Report, Working paper 2024.

- , – , and – , “Asset Purchase Rules in the Euro Area and Their Effect on Bond Markets,” Technical Report, Working paper 2024.
- , – , and – , “Whatever it takes? The impact of conditional policy promises,” *American Economic Review*, 2025, 115 (1), 295–329.
- and **Tyler Muir**, “Do intermediaries matter for aggregate asset prices?,” *The Journal of Finance*, 2021, 76 (6), 2719–2761.
- and – , “Market macrostructure: Institutions and asset prices,” *Annual Review of Financial Economics*, 2025, 17.
- , **Paul Huebner**, and **Erik Loualiche**, “How competitive is the stock market? theory, evidence from portfolios, and implications for the rise of passive investing,” *Working paper*, 2024.
- Harris, Lawrence and Eitan Gurel**, “Price and Volume Effects Associated with Changes in the S&P 500 List: New Evidence for the Existence of Price Pressures,” *The Journal of Finance*, 1986, 41 (4), 815–829.
- Hartzmark, Samuel M and David H Solomon**, “Predictable price pressure,” Technical Report, National Bureau of Economic Research 2022.
- Hausman, J. A.**, “Specification Tests in Econometrics,” *Econometrica*, 1978, 46 (6), 1251–1271.
- Hausman, Jerry and Daniel McFadden**, “Specification Tests for the Multinomial Logit Model,” *Econometrica*, 1984, 52 (5), 1219–1240.
- He, Zhiguo, Paymon Khorrami, and Zhaogang Song**, “Commonality in credit spread changes: Dealer inventory and intermediary distress,” *The Review of Financial Studies*, 2022, 35 (10), 4630–4673.
- , **Peter Kondor**, and **Jessica Shi Li**, “Demand Elasticity in Dynamic Asset Pricing,” Technical Report, Working paper 2025.
- Houweling, Patrick, Albert Mentink, and Ton Vorst**, “Comparing possible proxies of corporate bond liquidity,” *Journal of Banking Finance*, 2005, 29 (6), 1331–1358.
- Huber, Kilian**, “Estimating general equilibrium spillovers of large-scale shocks,” *The Review of Financial Studies*, 2023, 36 (4), 1548–1584.
- Huebner, Paul**, “The Making of Momentum: A Demand-System Perspective,” *Working paper*, 2024.
- Hurwicz, Leonid**, “On the structural form of interdependent systems,” in “Studies in Logic and the Foundations of Mathematics,” Vol. 44, Elsevier, 1966, pp. 232–239.
- Imbens, Guido W and Joshua D Angrist**, “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 1994, 62 (2), 467–475.

- Jansen, Kristy A.E., Wenhao Li, and Lukas Schmid**, “Granular Treasury Demand with Arbitrageurs,” *Available at SSRN 4940397*, 2024.
- Jiang, Zhengyang, Robert J Richmond, and Tony Zhang**, “A portfolio approach to global imbalances,” *The Journal of Finance*, 2024, *79* (3), 2025–2076.
- Kelly, Bryan T, Seth Pruitt, and Yinan Su**, “Characteristics are covariances: A unified model of risk and return,” *Journal of Financial Economics*, 2019, *134* (3), 501–524.
- , – , and – , “Instrumented principal component analysis,” *Available at SSRN 2983919*, 2020.
- Kippersluis, Hans Van and Cornelius A Rietveld**, “Beyond plausibly exogenous,” *The Econometrics Journal*, 2018, *21* (3), 316–331.
- Koijen, Ralph S. J. and Motohiro Yogo**, “A Demand System Approach to Asset Pricing,” *Journal of Political Economy*, 2019, *127* (4), 1475–1515.
- and – , “Exchange Rates and Asset Prices in a Global Demand System,” *Working paper*, 2024.
- Koijen, Ralph S J, Robert J Richmond, and Motohiro Yogo**, “Which Investors Matter for Equity Valuations and Expected Returns?,” *The Review of Economic Studies*, 08 2023, *91* (4), 2387–2424.
- Kozak, Serhiy, Stefan Nagel, and Shrihari Santosh**, “Shrinking the cross-section,” *Journal of Financial Economics*, 2020, *135* (2), 271–292.
- Krishnamurthy, Arvind and Annette Vissing-Jorgensen**, “The effects of quantitative easing on interest rates: channels and implications for policy,” Technical Report, National Bureau of Economic Research 2011.
- Kyle, Albert S.**, “Informed Speculation with Imperfect Competition,” *Review of Economic Studies*, 1989, *56* (3), 317–355.
- Li, Jiacui and Zihan Lin**, “Price Multipliers are Larger at More Aggregate Levels,” *Available at SSRN 4038664*, 2022.
- Li, Ziang**, “Long Rates, Life Insurers, and Credit Spreads,” Technical Report, Working paper 2025.
- Litterman, Robert B and Josè Scheinkman**, “Common Factors Affecting Bond Returns,” *The Journal of Fixed Income*, 1991, *1* (1), 54–61.
- Lopez-Lira, Alejandro and Nikolai L Roussanov**, “Do common factors really explain the cross-section of stock returns?,” *Jacobs Levy Equity Management Center for Quantitative Financial Research Paper*, 2020.
- Lou, Dong**, “A flow-based explanation for return predictability,” *The Review of Financial Studies*, 2012, *25* (12), 3457–3489.

- Lu, Xu and Lingxuan Wu**, “Monetary Transmission and Portfolio Rebalancing: A Cross-Sectional Approach,” Technical Report 2023.
- Markowitz, Harry M**, “Portfolio selection,” *The journal of finance*, 1952, 7, 77–91.
- Merton, Robert C.**, “An Intertemporal Capital Asset Pricing Model,” *Econometrica*, 1973, 41 (5), 867–887.
- Mundlak, Yair**, “On the Pooling of Time Series and Cross Section Data,” *Econometrica*, 1978, 46 (1), 69–85.
- Olea, José Luis Montiel and Carolin Pflueger**, “A robust test for weak instruments,” *Journal of Business & Economic Statistics*, 2013, 31 (3), 358–369.
- Pavlova, Anna and Taisiya Sikorskaya**, “Benchmarking Intensity,” *The Review of Financial Studies*, 08 2022, 36 (3), 859–903.
- Peng, Cameron and Chen Wang**, “Factor Demand and Factor Returns,” *Working paper*, 2023.
- Petajisto, Antti**, “Why do demand curves for stocks slope down?,” *Journal of Financial and Quantitative Analysis*, 2009, 44 (5), 1013–1044.
- Ross, Stephen A**, “The arbitrage theory of capital asset pricing,” *Journal of Economic Theory*, 1976, 13 (3), 341–360.
- Rostek, Marzena and Ji Hee Yoon**, “Imperfect competition in financial markets: Recent developments,” *Journal of Economic Literature*, 2025, 63 (4), 1191–1243.
- Selgrad, Julia**, “Testing the Portfolio Rebalancing Channel of Quantitative Easing,” *Working paper*, 2023.
- Shleifer, Andrei**, “Do Demand Curves for Stocks Slope Down?,” *The Journal of Finance*, 1986, 41 (3), 579–590.
- Stock, James H and Motohiro Yogo**, “Testing for weak instruments in Linear Iv regression,” in “Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg,” Cambridge University Press, 2005, pp. 80–108.
- Vayanos, Dimitri and Jean-Luc Vila**, “A Preferred-Habitat Model of the Term Structure of Interest Rates,” *Econometrica*, 2021, 89 (1), 77–112.
- Vazza, Diane and Nick W Kraemer**, “Default, Transition, and Recovery: 2012 Annual Global Corporate Default Study and Rating Transitions,” RatingsDirect, Standard & Poor’s Ratings Services 3 2013. Accessed: 2026-02-05.
- Welch, Ivo**, “Simply Better Market Betas,” *Critical Finance Review*, 2022, 11 (1), 37–64.
- Wolf, Christian K**, “The missing intercept: A demand equivalence approach,” *American Economic Review*, 2023, 113 (8), 2232–2269.

Appendix

Contents

1	The Challenge of Causal Inference in Asset Pricing	6
1.1	Demand and Demand Elasticity	6
1.2	Restrictions are Necessary to Estimate Demand Elasticity	8
1.3	Restrictions Should Accommodate Standard Finance Logic	9
2	A Model of Asset Demand	13
2.1	Framework	13
2.2	Applying the Assumptions	15
3	Estimation	20
3.1	Estimating Relative Elasticity	20
3.2	Estimating Substitution	25
4	Price Impact	32
5	Empirical Example: Corporate Bond Multipliers	35
5.1	Identification Strategy and Its Plausibility	35
5.2	Estimates	37
5.3	Counterfactual Analysis	42
6	Concluding Remarks	45
A	Elasticity inside models	55
A.1	When does demand elasticity measure a structural relationship?	55
B	Elasticity and Multiplier	57
B.1	Elasticity	57
B.2	From Elasticity to Multiplier	57
C	Proofs and Derivations	57
C.1	Identifying the relative elasticity – Proposition 1	57
C.2	Properties of elasticity under assumptions A1 and A2	58
C.3	Heterogeneous relative elasticities based on observables	61
C.4	Identification beyond the relative elasticity.	62
C.5	Lack of identification of substitution from the cross-section	63
C.6	Estimating a Local Average Treatment Effect	65
C.7	Local Average Treatment Effect relaxing Assumption A1	70

D	Robustness	73
	D.1 Robustness to the assumptions	73
	D.2 Validating the plausibility of assumptions	76
E	Demand beyond risk-based motives for substitution	78
F	A non-linear framework	81
	F.1 Basic concepts	81
	F.2 Relative elasticity vs. substitution and identification.	82
	F.3 Logit, log utility and factor models	83
G	Learning from Prices, Strategic Trading, Dynamic Trading	86
	G.1 Admati (1985) : competitive traders learning from signals and prices with multiple assets	87
	G.2 Multi-asset Kyle (1989) : imperfect competition with common signal	89
	G.3 Dynamic environments	91
H	Estimating elasticity in theoretical models	91
	H.1 Standard models of asset demand	91
	H.2 What about equilibrium spillovers?	93
I	Limits of existing demand models under a factor structure	97
	I.1 Setting	97
	I.2 Model estimation	99
	I.3 Experiments	101
J	Appendix Tables and Figures	106

A Elasticity inside models

A.1 When does demand elasticity measure a structural relationship?

The investor problem. Let us consider the textbook static microfounded model of portfolio decisions. An agent chooses between $N + 1$ assets, where asset 0 is the numeraire, meaning its price is always normalized to 1. The agent believes (rationally or not) that the payoffs (Π_0, \dots, Π_N) of the various assets have joint distribution F . The agent is endowed with quantities $(\bar{Q}_0, \dots, \bar{Q}_N)$ of the various assets. They maximize their expected utility, with utility function u , taking prices $P = (P_1, \dots, P_N)$ as given. This corresponds to:

$$\max_{Q_0, Q_1, \dots, Q_N} \mathbb{E} \left[u \left(\sum_{i=0}^N Q_i \Pi_i \right) \right] \quad (52)$$

$$\text{s.t. } Q_0 + \sum_{i=1}^N Q_i P_i \leq \bar{Q}_0 + \sum_{i=1}^N \bar{Q}_i P_i \quad (53)$$

The optimal positions define a demand function. Focusing only on positions in assets 1 to N , we define this function $D(P) = (Q_1^*(P), \dots, Q_N^*(P))$. Demand is a function from \mathbb{R}^N to \mathbb{R}^N .

In this microeconomic problem, the demand curve entirely characterizes how decisions change when prices change. If one focuses on small changes in prices, changes in positions are characterized by the slope of this demand curve, the Jacobian $\mathcal{E} = \partial D / \partial P = [\partial D_i / \partial P_j]_{i,j=1, \dots, N}$. In a slight abuse of language, we refer to \mathcal{E} as the elasticity of demand. We have immediately $\Delta Q \approx \mathcal{E} \Delta P$ for a small change in the price vector.

One might also be interested in what happens if the distribution of payoffs changes. For example the distribution F might be parametrized by a vector Y , and one can ask how optimal positions change when Y moves. Similarly, preferences might change too: for example the risk aversion of the investor could depend on their mood. We can also encode this in the vector Y . This extends the definition of the demand function to $D(P, Y)$. Note that while in equilibrium, the distribution of returns depends on market clearing, the demand function $D(P, Y)$ describes the agent's optimal portfolio for any arbitrary price vector P and belief system Y , effectively isolating the agent's preferences from the equilibrium constraints. We still refer to the price elasticity of demand as $\mathcal{E} = \partial D / \partial P$.

This elasticity captures a *structural relationship* in the sense of [Hurwicz \(1966\)](#) for any counterfactual about this investor that involves changes in prices but not in the vector Y . The elasticity characterizes how the investor's portfolio changes when prices change locally but the vector Y does not, irrespective of any other details about what determines equilibrium prices.

Inside of an equilibrium. To make things concrete, we flesh out an equilibrium with many such standard agents. Specifically, consider an equilibrium populated by J agents indexed by $j = 1, \dots, J$. Each agent has a demand function $D_j(P, Y_j)$ and an initial endowment \bar{Q}_j . Market clearing requires that the sum of demands equals the aggregate endowment. We

define the aggregate demand for risky assets as $\mathcal{D}(P, Y) = \sum_{j=1}^J D_j(P, Y_j)$ and the aggregate endowment as $\bar{Q}_{\text{mkt}} = \sum_{j=1}^J \bar{Q}_j$. The equilibrium price vector P^* is the solution to:

$$\sum_{j=1}^J D_j(P^*, Y_j) = \bar{Q}_{\text{mkt}} \quad (54)$$

We define the aggregate elasticity of demand as the sum of individual elasticities (Jacobians): $\mathcal{E}_{\text{agg}} = \sum_{j=1}^J \mathcal{E}_j$. Assuming standard conditions for uniqueness and stability, \mathcal{E}_{agg} is negative definite (invertible). We can now compute two structural counterfactuals using the implicit function theorem.

Counterfactual 1: Absorption of a large trade. Consider the arrival of a new trader who wishes to purchase a fixed portfolio Q_{new} of risky assets. This trader is endowed only with the numeraire (asset 0), so they add no new supply of risky assets to the market. Their demand represents a pure removal of shares from the existing pool. Market clearing now requires $\sum D_j(P') + Q_{\text{new}} = \bar{Q}_{\text{mkt}}$. This is equivalent to a negative supply shock of magnitude $\Delta S = -Q_{\text{new}}$ to the risky assets available to the original J agents. Linearizing around the initial equilibrium, the change in prices ΔP required to accommodate this trade satisfies $\mathcal{E}_{\text{agg}} \Delta P \approx -Q_{\text{new}}$. Thus, the price impact is:

$$\Delta P \approx -\mathcal{E}_{\text{agg}}^{-1} Q_{\text{new}} \quad (55)$$

The price response is determined entirely by the aggregate “risk-bearing capacity” of the market (\mathcal{E}_{agg}) and the size of the trade. The existing agents must be induced to reduce their holdings by exactly Q_{new} , and $\mathcal{E}_{\text{agg}}^{-1}$ describes the price concession required to achieve this reduction.

Counterfactual 2: A shock to beliefs, preferences or macrostructure. Consider a change in the characteristics Y_k of a single agent k (e.g., a change in sentiment or risk aversion), while other agents and the total supply remain fixed. The market clearing condition requires the price to adjust to offset agent k 's shift in demand. Differentiating the clearing condition with respect to Y_k :

$$\mathcal{E}_{\text{agg}} \Delta P + \frac{\partial D_k}{\partial Y_k} \Delta Y_k = 0 \implies \Delta P \approx -\mathcal{E}_{\text{agg}}^{-1} \frac{\partial D_k}{\partial Y_k} \Delta Y_k \quad (56)$$

This price change forces all other agents $i \neq k$ to rebalance their portfolios. Since their parameters Y_i are unchanged, their trading is purely a response to the price movement:

$$\Delta Q_i \approx \mathcal{E}_i \Delta P \quad (57)$$

This illustrates the mechanism of how the equilibrium changes: a shock to agent k 's beliefs transmits to agent i 's portfolio strictly through the price channel, governed by agent i 's structural elasticity \mathcal{E}_i .

B Elasticity and Multiplier

B.1 Elasticity

B.2 From Elasticity to Multiplier

As we describe in Section H.1, it is sometimes more suitable to estimate demand elasticities in different units (logarithms, portfolio shares instead of quantities, ...). The inversion result of Section sec:priceimpact applies to these different cases but with slightly adjusted formulas:

$$\mathcal{M}_{\{\log P, \log Q\}} = -\mathcal{E}_{\{\log Q, \log P\}}^{-1}, \quad (58)$$

$$\mathcal{M}_{\{\log P, \log Q\}} = -\left[\mathcal{E}_{\{\log \omega, \log P\}} - (\mathbf{I} - \mathbf{1}\omega')\right]^{-1}, \quad (59)$$

$$\mathcal{M}_{\{\log P, \log Q\}} = -\left[\text{diag}(\omega)^{-1}\mathcal{E}_{\{\omega, \log P\}} - (\mathbf{I} - \mathbf{1}\omega')\right]^{-1}. \quad (60)$$

For example in the case of logit where demand elasticity is measured by regressing the log portfolio share on log price, equation (59) gives us the multiplier in log units: by how many percents do prices move in response to a one percent change in aggregate demand. Similarly, equation (60) is useful for the case of CRRA.

C Proofs and Derivations

C.1 Identifying the relative elasticity – Proposition 1

Start from the general demand equation with demand shocks:

$$\Delta D_i = \mathcal{E}_{ii}\Delta P_i + \sum_{j \neq i} \mathcal{E}_{ij}\Delta P_j + \Delta \bar{D}_i. \quad (61)$$

We recall the two assumptions necessary for identification:

- **Assumption A1.** $X_i = X_j \Rightarrow \mathcal{E}_{il} = \mathcal{E}_{jl} = \mathcal{E}_{\text{cross}}(X_i, X_l) = X_i' \mathcal{E}_X X_l, \quad \forall i, j \in \mathcal{S}, l \neq i, j$, where X_i is a $K \times 1$ vector of observables, and \mathcal{E}_X is a $K \times K$ matrix.
- **Assumption A2.** $\mathcal{E}_{ii} - \mathcal{E}_{\text{cross}}(X_i, X_i) = \mathcal{E}_{jj} - \mathcal{E}_{\text{cross}}(X_j, X_j) = \hat{\mathcal{E}}, \quad \forall i, j \in \mathcal{S}$

Proposition 1 shows that under assumptions A1 and A2 and the exogeneity condition, the IV estimator, conditioning on X_i , identifies coefficient $\hat{\mathcal{E}}$.

Proof. Starting from equation (61), we can rewrite the demand equation as a cross-sectional

regression:

$$\Delta D_i = \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{ij} \Delta P_j + \Delta \bar{D}_i \quad (62)$$

$$= \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{cross}(X_i, X_j) \Delta P_j + \Delta \bar{D}_i \quad (63)$$

$$= (\mathcal{E}_{ii} - \mathcal{E}_{cross}(X_i, X_i)) \Delta P_i + \sum_j \mathcal{E}_{cross}(X_i, X_j) \Delta P_j + \Delta \bar{D}_i \quad (64)$$

$$= \hat{\mathcal{E}} \Delta P_i + \sum_j \mathcal{E}_{cross}(X_i, X_j) \Delta P_j + \Delta \bar{D}_i \quad (65)$$

$$= \hat{\mathcal{E}} \Delta P_i + \sum_j X_i' \mathcal{E}_X X_j \Delta P_j + \Delta \bar{D}_i \quad (66)$$

$$= \hat{\mathcal{E}} \Delta P_i + X_i' \underbrace{\left(\sum_j \mathcal{E}_X X_j \Delta P_j \right)}_{\theta} + \Delta \bar{D}_i \quad (67)$$

$$= \hat{\mathcal{E}} \Delta P_i + \theta' X_i + \Delta \bar{D}_i \quad (68)$$

Equation (64) adds and subtracts $\mathcal{E}_{cross}(X_i, X_i) \Delta P_i$. Equations (65) and (66) use assumptions 2 and 1, respectively. Equation (67) pulls out X_i' from the sum. The remaining part of the sum gets absorbed into θ , a $K \times 1$ vector of cross-sectional constants. These θ are K regression coefficients on the K observables, X_{ik} .

Given the exclusion restriction that $Z_i \perp \Delta \bar{D}_i | X_i$ and the relevance condition that $cov(\Delta P_i, Z_i | X_i) \neq 0$, this is the standard IV setting, and the regression estimates $\hat{\mathcal{E}}$. ■

C.2 Properties of elasticity under assumptions A1 and A2

C.2.1 A matrix representation.

First, we derive a simple matrix representation for an elasticity matrix under our two assumptions.

Lemma 6 *Let \mathcal{E} be an elasticity matrix that satisfies assumptions A1 and A2. Then it can be written as:*

$$\mathcal{E} = \hat{\mathcal{E}} \mathbf{I} + X \mathcal{E}_X X', \quad (69)$$

where $\hat{\mathcal{E}}$ is a scalar equal to the relative elasticity and \mathcal{E}_X is a $K \times K$ matrix.

Proof. Write out the elasticity matrix \mathcal{E} :

$$\mathcal{E} = \begin{pmatrix} \mathcal{E}_{11} & \mathcal{E}_{12} & \dots & \mathcal{E}_{1N} \\ \mathcal{E}_{21} & \mathcal{E}_{22} & \dots & \mathcal{E}_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{E}_{N1} & \mathcal{E}_{N2} & \dots & \mathcal{E}_{NN} \end{pmatrix} = \begin{pmatrix} \hat{\mathcal{E}} + X_1' \mathcal{E}_X X_1 & X_1' \mathcal{E}_X X_2 & \dots & X_1' \mathcal{E}_X X_N \\ X_2' \mathcal{E}_X X_1 & \hat{\mathcal{E}} + X_2' \mathcal{E}_X X_2 & \dots & X_2' \mathcal{E}_X X_N \\ \vdots & \vdots & \ddots & \vdots \\ X_N' \mathcal{E}_X X_1 & X_N' \mathcal{E}_X X_2 & \dots & \hat{\mathcal{E}} + X_N' \mathcal{E}_X X_N \end{pmatrix} \quad (70)$$

The (i, j) element of matrix \mathcal{E} is $[\mathcal{E}]_{ij} = X_i' \mathcal{E}_X X_j = \mathcal{E}_{cross}(X_i, X_j)$, as defined by Assumption 1, for $i \neq j$. The diagonal elements are $[\mathcal{E}]_{ii} = \hat{\mathcal{E}} + X_i' \mathcal{E}_X X_i = \hat{\mathcal{E}} + \mathcal{E}_{cross}(X_i, X_i)$, as in Assumption 2. Since each element in \mathcal{E} directly corresponds to the respective \mathcal{E}_{ij} defined by assumptions A1 and A2, the assumptions are equivalent to the elasticity matrix in (69). ■

C.2.2 Transforming the observables.

The following lemma shows that observables can be recombined in a linear way. In particular they could be demeaned, standardized, or orthogonalized.

Lemma 7 *Let \mathcal{E} be an elasticity matrix that satisfies assumptions A1 and A2 with respect to a set of observables X . If A is a $K \times K$ invertible matrix, \mathcal{E} satisfies assumptions A1 and A2 with respect to the recombined observables $\tilde{X} = XA$.*

Proof. Insert AA^{-1} judiciously into the decomposition of Lemma 6.

$$\mathcal{M} = \widehat{\mathcal{M}}\mathbf{I} + XAA^{-1}\mathcal{M}_X(A')^{-1}A'X' = \widehat{\mathcal{M}}\mathbf{I} + \tilde{X}\mathcal{M}_{\tilde{X}}\tilde{X}', \quad (71)$$

$$\text{with } \mathcal{M}_{\tilde{X}} = A^{-1}\mathcal{M}_X(A')^{-1}. \quad (72)$$

■

For example, if the first observable is a constant and the other ones have mean $\bar{X}_1, \dots, \bar{X}_{K-1}$, the following matrix demeans them:

$$A_{\text{demean}} = \mathbf{I}_K - \begin{pmatrix} 0 & \bar{X}_1 & \cdots & \bar{X}_{K-1} \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{pmatrix}. \quad (73)$$

Importantly, note that there is no reason that orthogonalizing the characteristics makes the substitution matrix $\mathcal{M}_{\tilde{X}}$ diagonal.

C.2.3 Stability by inversion — Proposition ??.

Proposition ?? shows that under assumptions 1 and 2, the multiplier matrix $\mathcal{M} = -\mathcal{E}^{-1}$ also satisfies assumptions 1 and 2, with $\widehat{\mathcal{M}} = -1/\hat{\mathcal{E}}$.

Proof. Start from equation (69), and apply the Woodbury matrix identity:

$$-\mathcal{E}^{-1} = -\left(\hat{\mathcal{E}}\mathbf{I} + X\mathcal{E}_X X'\right)^{-1} \quad (74)$$

$$= -\hat{\mathcal{E}}^{-1}\mathbf{I} + X\left(\hat{\mathcal{E}}^2\mathcal{E}_X^{-1} + \hat{\mathcal{E}}X'X\right)^{-1}X' \quad (75)$$

$$= \widehat{\mathcal{M}}\mathbf{I} + X\mathcal{M}_X X'. \quad (76)$$

This corresponds exactly to assumptions A1 and A2 applied to \mathcal{M} with $\widehat{\mathcal{M}} = -1/\hat{\mathcal{E}}$. ■

C.2.4 Stability by aggregation

We show that that assumptions A1 and A2 are stable by aggregation across investors.

Lemma 8 *Let \mathcal{E}_1 and \mathcal{E}_2 be two elasticity matrices that satisfy assumptions A1 and A2, and (λ_1, λ_2) two scalars. Then the matrix $\lambda_1\mathcal{E}_1 + \lambda_2\mathcal{E}_2$ satisfies assumptions A1 and A2.*

Proof. From lemma 6 we decompose both elasticities which leads to:

$$\lambda_1\mathcal{E}_1 + \lambda_2\mathcal{E}_2 = \left(\lambda_1\hat{\mathcal{E}}_1 + \lambda_2\hat{\mathcal{E}}_2\right)\mathbf{I} + X(\lambda_1\mathcal{E}_{X,1} + \lambda_2\mathcal{E}_{X,2})X' \quad (77)$$

The decomposition and the equivalence from lemma 6 concludes the proof. ■

C.3 Heterogeneous relative elasticities based on observables

We maintain assumption A1. We relax assumption A2 by allowing the relative elasticity to depend linearly on observables:

$$\mathcal{E}_{ii} - \mathcal{E}_{cross}(X_i, X_i) = \mathcal{E}_{relative}(X_i) = \mathcal{E}'_r X_i. \quad (78)$$

From section C.1 and the proof of Proposition 1 we obtain:

$$\Delta D_i = \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{ij} \Delta P_j + \epsilon_i \quad (79)$$

$$= \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{cross}(X_i, X_j) \Delta P_j + \epsilon_i \quad (80)$$

$$= (\mathcal{E}_{ii} - \mathcal{E}_{cross}(X_i, X_i)) \Delta P_i + \sum_j \mathcal{E}_{cross}(X_i, X_j) \Delta P_j + \epsilon_i \quad (81)$$

$$= \mathcal{E}_{relative}(X_i) \Delta P_i + \sum_j X'_i \mathcal{E}_X X_j \Delta P_j + \epsilon_i \quad (82)$$

$$= \underbrace{\mathcal{E}_{relative}(X_i)}_{\mathcal{E}'_r X_i} \Delta P_i + X'_i \underbrace{\left(\sum_j \mathcal{E}_X X_j \Delta P_j \right)}_{\theta} + \epsilon_i \quad (83)$$

$$= \mathcal{E}'_r X_i \Delta P_i + \theta' X_i + \epsilon_i \quad (84)$$

In this case we want to identify the vector \mathcal{E}_r which characterizes the relative elasticity with respect to observables. Identification must rely on a vector of instruments. It is natural to construct those instruments from a single instrument Z_i for the price interacted with the observables. The set of identification conditions is:

$$Z_i X_i \perp \epsilon_i | X_i \quad (85)$$

Under these conditions the two-stage least squares regression proceeds as follows. First, regress each component of $X_i \Delta P_i$ on the vector of instruments $X_i Z_i$ and the observables X_i . The relevance condition is that the matrix of coefficients on the instruments is full-rank. This leads to predicted values of the change in price interacted with observables $\widehat{X_i \Delta P_i}$. Second, regress the change in demand on these predicted values and the observables. The coefficients on the predicted values recovers \mathcal{E}_r . Finally, the relative elasticity for each asset is simply $\mathcal{E}_{relative}(X_i) = \mathcal{E}'_r X_i$.

C.4 Identification beyond the relative elasticity.

We consider aggregation for the generic case of a elasticity matrix \mathcal{E} that satisfies assumptions A1 and A2 for an arbitrary set of observables X . Remember that the first observable is the constant in most cases.

Using Lemma 6, we can represent \mathcal{E} as

$$\mathcal{E} = \widehat{\mathcal{E}}\mathbf{I} + X\mathcal{E}_X X'. \quad (86)$$

To define price and quantity aggregates along the various dimensions of the observables, we regress these vectors on X . We will see that this the natural generalization of the aggregation presented in Proposition ??.

Proposition 9 (Elasticity decomposition with observables in the general case) *Take an elasticity matrix \mathcal{E} satisfying assumptions A1 and A2. Consider generic changes in demand and price connected by this matrix: $\Delta D = \mathcal{E}\Delta P$. Define the change in demand and price aggregated along observables and the idiosyncratic component:*

$$\Delta D_X = (X'X)^{-1}X'\Delta D \quad \Delta P_X = (X'X)^{-1}X'\Delta P. \quad (87)$$

$$\Delta D_{idio,i} = \Delta D_i - X'_i \Delta D_X \quad \Delta P_{idio,i} = \Delta P_i - X'_i \Delta P_X. \quad (88)$$

The response of changes in demand to a change in prices can be decomposed into two sets of components:

$$\text{Micro:} \quad \Delta D_{idio,i} = \widehat{\mathcal{E}}\Delta P_{idio,i} \quad (89)$$

$$\text{Meso-Macro:} \quad \Delta D_X = \check{\mathcal{E}}\Delta P_X, \quad (90)$$

where $\check{\mathcal{E}} = \widehat{\mathcal{E}}\mathbf{I}_K + \mathcal{E}_X X'X$.

Proof. Using the relation $\Delta D = \mathcal{E}\Delta P$ and the decomposition of \mathcal{E} under the two assumptions, we obtain

$$(X'X)^{-1}X'\Delta D = \left(\widehat{\mathcal{E}}(X'X)^{-1}X'\mathbf{I}_N + (X'X)^{-1}X'X\mathcal{E}_X X' \right) \Delta P \quad (91)$$

$$= \left(\widehat{\mathcal{E}}(X'X)^{-1} + \mathcal{E}_X \right) X' \Delta P \quad (92)$$

$$= \left(\widehat{\mathcal{E}}\mathbf{I}_K + \mathcal{E}_X X'X \right) (X'X)^{-1}X' \Delta P \quad (93)$$

$$\Delta D_X = \left(\widehat{\mathcal{E}}\mathbf{I}_K + \mathcal{E}_X X'X \right) \Delta P_X \quad (94)$$

This implies that ΔD_X can be expressed as a linear combination of the K elements of ΔP_X , as opposed to the whole N components of the changes in demand ΔP .

From the definition of the idiosyncratic change in demand:

$$\Delta D_{idio} = \Delta D - X \Delta D_X \quad (95)$$

$$= \hat{\mathcal{E}} \Delta P + X \mathcal{E}_X X' \Delta P - \left(X \hat{\mathcal{E}} + X \mathcal{E}_X X' \right) (X' X)^{-1} X' \Delta P \quad (96)$$

$$= \hat{\mathcal{E}} (\Delta P - X \Delta P_X) + X \mathcal{E}_X X' \Delta P - X \mathcal{E}_X (X' X)^{-1} X' \Delta P \quad (97)$$

$$= \hat{\mathcal{E}} \Delta P_{idio}. \quad (98)$$

Because $\hat{\mathcal{E}}$ is scalar, the idiosyncratic component is determined asset by asset, which concludes the proof. ■

Simple case with no characteristic. We can recover the simpler cases studied in the paper. Proposition ?? corresponds to a single variable X which is constant equal to one. In this case ΔD_X has only one component equal to $\Delta D_{agg} = N^{-1} \sum_i \Delta D_i$, and $\check{\mathcal{E}} = \hat{\mathcal{E}} + N \mathcal{E}_X$ is a scalar equal to the macro elasticity.

With one normalized characteristic. Proposition ?? corresponds to observables that include a constant and a single standardized characteristic that we call X in a slight abuse of notation. There, the regression gives two aggregate prices and quantities: the aggregate component ΔD_{agg} defined as before (the constant of the regression); the meso component $\Delta D_X = N^{-1} \sum_i X_i \Delta D_i$. Then the matrix $\check{\mathcal{E}}$ is 2×2 and equal to

$$\check{\mathcal{E}} = \begin{pmatrix} \hat{\mathcal{E}} + N(\mathcal{E}_X)_{11} & N(\mathcal{E}_X)_{12} \\ N(\mathcal{E}_X)_{21} & \hat{\mathcal{E}} + N(\mathcal{E}_X)_{22} \end{pmatrix} = \begin{pmatrix} \bar{\mathcal{E}}_{agg} & \bar{\mathcal{E}}_X \\ \tilde{\mathcal{E}}_{agg} & \tilde{\mathcal{E}}_X \end{pmatrix}. \quad (99)$$

When observables are group dummies. Consider the case when the observables are dummy variables expressing the belonging to disjoint groups. In this situation, there is no need for a constant. The aggregate demand and price indices have a simple interpretation: they measure the average change in demand and price for each group k :

$$\Delta P_{X,k} = \frac{1}{N_k} \sum_{i \in k} \Delta P_i \quad (100)$$

This implies that demand elasticities have a nested structure. First, there is a relative elasticity within each group $\hat{\mathcal{E}}$ capturing the impact of changes in relative prices within a group. Then individual assets can be replaced by the aggregate portfolio of each group, and there is an elasticity matrix across these aggregate portfolios, $\check{\mathcal{E}}$.

C.5 Lack of identification of substitution from the cross-section

We show that without additional restrictions, substitution cannot be identified from a single cross section. Start with the structural relation: $\Delta D = \mathcal{E} \Delta P + \epsilon$, and impose our assumptions: $\mathcal{E} = \hat{\mathcal{E}} I + X \mathcal{E}_X X'$.

Recall that the demand shift ϵ measures all demand changes that are not caused by a price change. In particular, it can have a relation with the observables X (e.g., if beliefs about relative payoffs of assets with different values of X change) and a non-zero mean (e.g., if the investor demands more assets overall). We can separate ϵ across those components:

$$\epsilon_X = (X'X)^{-1}X'\epsilon \quad (101)$$

$$\epsilon_{idio,i} = \epsilon_i - X'_i\epsilon_X, \quad (102)$$

with all components of ϵ_X potentially different from 0 even in the limit of many assets (large N) and imposing no constraints on all other model quantities.

Plugging into demand, we obtain:

$$\Delta D_i = \hat{\mathcal{E}}\Delta P_i + X'_i\mathcal{E}_X X'\Delta P + X'_i\epsilon_X + \epsilon_{idio,i} \quad (103)$$

$$= \hat{\mathcal{E}}\Delta P_i + \epsilon_{idio,i} + X'_i \underbrace{(\mathcal{E}_X X'\Delta P + \epsilon_X)}_{K \times 1} \quad (104)$$

Proposition 10 *No free parameter of the matrix \mathcal{E}_X can be identified from a single cross section, even under the restriction that some of coefficients of \mathcal{E}_X are 0.*

Proof. Denote \mathcal{E}_X^{true} and ϵ_X^{true} the true values of \mathcal{E}_X and ϵ_X . For any other guess \mathcal{E}_X^{false} , the model with $\mathcal{E}_X = \mathcal{E}_X^{false}$ and $\epsilon_X = \epsilon_X^{true} + (\mathcal{E}_X^{true} - \mathcal{E}_X^{false})X'\Delta P$ is observationally equivalent to the true model. As long as such guesses exist, that is, as long as \mathcal{E}_X has at least one free parameter, this concludes the proof of no identification. ■

Simplest case: only a constant. Consider the simplest possible case, where X is a constant. Writing cross-sectional means at date t with a bar, and noting that \mathcal{E}_X is a scalar in this case, we have:

$$\Delta D_{i,t} = \hat{\mathcal{E}}\Delta P_{i,t} + \mathcal{E}_X \overline{\Delta P}_t + \bar{\epsilon}_t + \epsilon_{idio,i,t} \quad (105)$$

$$= (\mathcal{E}_X \overline{\Delta P}_t + \bar{\epsilon}_t) + \hat{\mathcal{E}}\Delta P_{i,t} + \epsilon_{idio,i,t} \quad (106)$$

Both substitution $\mathcal{E}_X \overline{\Delta P}_t$ and the aggregate demand shift $\bar{\epsilon}_t$ contribute to the constant of a cross-sectional regression, and there is no way to separate them. This is a version of the missing intercept problem.

Obtaining partial identification with symmetry across observables. One way to obtain some partial identification of substitution is to impose that the same parameter drives substitution across many observables. Then, if one also assumes that the number of observables grows as the number of assets increases in the cross section, one can identify part of the substitution matrix.

We illustrate this approach when the observables are dummy variables belonging to a given group. A priori, substitution across groups could follow any matrix \mathcal{E}_X . But one might want to assume that all groups substitute symmetrically, that is, $\mathcal{E}_X = \mathcal{E}_{own-g}I + \mathcal{E}_{cross-g}11'$, with \mathcal{E}_{own-g} an own-group elasticity and $\mathcal{E}_{cross-g}$ a cross-group elasticity. In this case, we

have:

$$\Delta D_i = \hat{\mathcal{E}} \Delta P_i + \mathcal{E}_{own-g} N_g \Delta P_g + \mathcal{E}_{cross-g} N \overline{\Delta P} \text{ if } i \in g \quad (107)$$

A cross-sectional regression with an instrument for ΔP_g across groups allow to recover \mathcal{E}_{own-g} . Notice that $\mathcal{E}_{cross-g}$ remains unidentified. This approach corresponds to repeating the relative elasticity estimation at a higher level of aggregation: the matrix \mathcal{E}_X (in contrast to \mathcal{E}) satisfies assumptions A1 and A2 with only a constant as observable.

The key assumption here is not that the observables correspond to disjoint groups, but instead that there is a common substitution parameter that affects each of the observables separately. Its plausibility depends on context: for example in a simple mean variance setting, it does not apply if the groups are based on levels of factor loadings.

The nested logit model assumes such a symmetry across groups (the “nests”) and additionally imposes that the missing intercept $\mathcal{E}_{cross-g}$ is pinned down by the other parameters.

C.6 Estimating a Local Average Treatment Effect

This section removes Assumption A2, deriving a Local Average Treatment Effect result for relative elasticities—a counterpart to Proposition 1—under stronger assumptions on instrument exogeneity.

The structural equations (data-generating process) are

$$\Delta D_i = \sum_{j=1}^N \mathcal{E}_{ij} \Delta P_j + \epsilon_i \quad (108)$$

$$\Delta P_i = \sum_{j=1}^N \Lambda_{ij} Z_j + \mu_i, \quad (109)$$

where equations (108) and (109) represent the structural demand equation and the structural impact of the instrument under an unrestricted spillover matrix Λ , respectively. The elasticity matrix \mathcal{E} satisfies Assumption A1 but not A2, taking the form

$$\mathcal{E} = \text{diag}(\mathcal{E}_{relative,i}) + X \mathcal{E}_X X', \quad (110)$$

such that relative elasticities $\mathcal{E}_{relative,i} = \mathcal{E}_{ii} - X_i' \mathcal{E}_X X_i$ are heterogeneous across assets. Without loss of generality, the $N \times K$ matrix of observables X is of rank K , and we assume that K remains fixed as we take limits. This ensures in particular that the operator $\mathbf{Q}_X = \mathbf{I} - X(X'X)^{-1}X'$, which constructs residuals after controlling for X , is well-defined.

The econometrician estimates the same two-stage least squares (2SLS) regression as in Proposition 1:

$$\Delta D_i = \hat{\mathcal{E}} \Delta P_i + X_i' \theta + e_i, \quad (111)$$

instrumenting ΔP_i with Z_i while controlling for X_i .

We assume strict instrument exogeneity and independence, meaning that Z_i is randomly assigned and independent of all other components:⁴⁸

1. $Z_i \perp Z_j$ for all $i \neq j$, with $\mathbf{E}[Z_i] = 0$, $\text{var}(Z_i) = \sigma_Z^2$ and $\mathbf{E}[Z_i^4] < \infty$ for all i .⁴⁹
2. $Z_i \perp (\mu_j, \epsilon_j)$ for all i, j .
3. $Z_i \perp (\Lambda_{kl}, X_j, \mathcal{E}_{relative,j})$ for all i, j, k, l .

The standard exclusion restriction follows directly from these conditions. We also assume relevance, meaning that the first stage retains power after partialling out observables:

$$\text{rank}(\text{cov}(\Delta P_i, Z_i | X)) = 1. \quad (112)$$

Equivalently, this requires $\text{tr}(\mathbf{Q}_X \mathbf{\Lambda}) \neq 0$, ensuring that the instruments generate price variation that is not fully absorbed by observables.⁵⁰

The following proposition gives an interpretation of the 2SLS regression coefficient as a local average treatment effect (LATE).

Proposition 11 (LATE of relative elasticities) *Under instrument exogeneity, independence, relevance, and Assumption A1, the two-stage least squares estimator from regression (111) identifies:*

$$\hat{\mathcal{E}} \xrightarrow{p} \sum_{i=1}^N \omega_i \mathcal{E}_{relative,i}, \quad \text{with } \omega_i = \frac{[\mathbf{Q}_X \mathbf{\Lambda}']_{ii}}{\text{tr}(\mathbf{Q}_X \mathbf{\Lambda})}, \quad (113)$$

where $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii} = \sum_j [\mathbf{Q}_X]_{ij} \Lambda_{ij}$ and $\mathbf{Q}_X = \mathbf{I} - X(X'X)^{-1}X'$ is the residualization operator. The weights ω_i sum to one and are determined by the first-stage spillover matrix $\mathbf{\Lambda}$ after projecting out observables.

Moreover, if $\mathbf{\Lambda}$ has the same bilinear structure as \mathcal{E} —that is, $\Lambda_{ij} = X_i' \mathbf{M}_X X_j$ for all $i \neq j$, for some $K \times K$ matrix \mathbf{M}_X —and the monotonicity condition $\lambda_i := \Lambda_{ii} - X_i' \mathbf{M}_X X_i \geq 0$ holds for all i , then the weights are non-negative and the estimand is a convex combination of relative elasticities.

We discuss the proposition before turning to the proof.

⁴⁸Full statistical independence is sufficient but stronger than necessary: the expectation calculations in the proof use only first- and second-moment conditions. Independence across instruments is used in the concentration argument underlying the probability limit.

⁴⁹Zero mean is without loss of generality: because the first element of X_i is 1, an intercept is included in the controls, and \mathbf{Q}_X annihilates any common instrument mean. Identical variance is assumed for notational convenience. With heterogeneous variances $\text{var}(Z_i) = \sigma_{Z,i}^2$, the weights in Proposition 11 become $\omega_i = [\mathbf{Q}_X \text{diag}(\sigma_{Z,j}^2) \mathbf{\Lambda}']_{ii} / \text{tr}(\mathbf{Q}_X \text{diag}(\sigma_{Z,j}^2) \mathbf{\Lambda})$, and all results carry through, including non-negative weights under bilinear $\mathbf{\Lambda}$ with monotonicity. Intuitively, instruments with larger variance generate more first-stage price variation and receive proportionally greater weight in the estimand. The finite fourth moment ensures concentration of the quadratic forms in the 2SLS estimand (Step 4 of the proof).

⁵⁰This condition permits $\mathbf{\Lambda}$ (and $\mathbf{Q}_X \mathbf{\Lambda}$) to be rank-deficient: some assets may be redundant, in the sense that their price variation is fully explained by other assets and observables. Such assets have $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii} = 0$ and receive zero weight $\omega_i = 0$ in the estimand (113). Identification requires only that *some* assets have nonzero first-stage variation after partialling out X .

Role of controls. The controls X_i serve two purposes. First, they absorb the cross-elasticity spillovers $X_i' \mathcal{E}_X X' \Delta P$, which are common across assets sharing the same characteristics. Second, the residualization \mathbf{Q}_X removes the component of first-stage variation that is explained by observables, ensuring that the IV variation is “idiosyncratic” relative to X . Without controls, the IV estimand would include an additional bias term from the uncontrolled cross-elasticity channel.

No contamination. With heterogeneous treatment effects and controls, 2SLS estimands do not generally equal a convex weighted average of treatment effects: heterogeneity in how units respond to both the treatment and the controls can distort the Frisch–Waugh–Lovell residualization (Goldsmith-Pinkham et al., 2024). In our setting, absent additional structure, demand for asset i depends on the entire vector of prices through heterogeneous cross-elasticities \mathcal{E}_{ij} , creating an N -treatment problem where contamination is generic.

Assumption A1 resolves this by imposing homogeneous substitution conditional on observables. Because all cross-price effects take the form $X_i' \mathcal{E}_X X_j$, the entire cross-price channel collapses to the common term $X_i' \theta$ in equation (117). Since θ is constant across assets in the cross-section, it is fully absorbed by the controls and annihilated by \mathbf{Q}_X , regardless of how the own-price elasticities $\mathcal{E}_{relative,i}$ vary.

As a result, once we condition on X , the model reduces to a single-treatment IV problem in the sense of Angrist (1998). Heterogeneity remains only in the own-price elasticities $\mathcal{E}_{relative,i}$, allowing the 2SLS estimand to be a first-stage-weighted average free of contamination.

Classic LATE and interpretation of weights. In the standard LATE framework (Imbens and Angrist, 1994; Angrist et al., 1996), IV with heterogeneous treatment effects identifies a first-stage-weighted average. When the first stage is diagonal ($A_{ij} = 0$ for $i \neq j$), the weights reduce to $\omega_i \propto A_{ii}$: assets whose prices are most responsive to their own instrument receive greater weight. In this case the result holds even without Assumption A1; see Appendix Section C.7.

However, with a full spillover matrix $\mathbf{\Lambda}$, the weights become $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii}$ —the total instrument-induced price response of asset i , combining direct and indirect channels, after projecting out the component explained by observables.

This distinction matters when spillovers are heterogeneous. An asset with a modest direct effect A_{ii} can receive large weight if it is strongly exposed to other assets’ instruments through off-diagonal elements of $\mathbf{\Lambda}$. Conversely, an asset with a large direct effect may receive little weight if spillovers from correlated assets are absorbed by the controls X .

Economically, the source of heterogeneous first-stage exposure here is much richer than in the diagonal case: it reflects equilibrium price transmission across assets, rather than simple, direct sensitivity to an instrument. For example, if more illiquid assets both experience larger residualized price responses $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii}$ and larger own-price elasticities $\mathcal{E}_{relative,i}$, the IV estimand will overstate how inelastic the typical asset is.

Sign of weights. In general, while the weights ω_i sum to one, some can be negative, meaning that the estimand is not guaranteed to be a convex combination of relative elasticities.

A negative weight $\omega_i < 0$ arises when $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii} < 0$: the net effect of instruments on asset i 's price is opposite in sign to the average effect across assets.

The second part of Proposition 11 gives sufficient conditions for non-negative weights. When $\mathbf{\Lambda}$ has the same bilinear structure as \mathcal{E} , the annihilator \mathbf{Q}_X reduces the weights to $[\mathbf{Q}_X]_{ii} \lambda_i$ (Step 5 of the proof), where $\lambda_i = \Lambda_{ii} - X_i' \mathcal{M}_X X_i$ is the “relative” first-stage coefficient—the analog of $\mathcal{E}_{relative,i}$ for the first stage. The monotonicity condition $\lambda_i \geq 0$ then requires that each asset’s own first-stage response weakly exceeds the cross-asset spillover between two assets sharing its characteristics.

Proof. *Step 1: Rewriting the demand equation.* Starting from the structural demand equation (108) and the elasticity decomposition (110):

$$\Delta D_i = \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{ij} \Delta P_j + \epsilon_i \quad (114)$$

$$= (\mathcal{E}_{relative,i} + X_i' \mathcal{E}_X X_i) \Delta P_i + \sum_{j \neq i} X_i' \mathcal{E}_X X_j \Delta P_j + \epsilon_i \quad (115)$$

$$= \mathcal{E}_{relative,i} \Delta P_i + X_i' \mathcal{E}_X \sum_{j=1}^N X_j \Delta P_j + \epsilon_i \quad (116)$$

$$= \mathcal{E}_{relative,i} \Delta P_i + X_i' \theta + \epsilon_i, \quad (117)$$

where $\theta = \mathcal{E}_X \sum_j X_j \Delta P_j$ is a $K \times 1$ vector that is constant across assets in the cross-section, but varies across realizations. The third equality adds and subtracts $X_i' \mathcal{E}_X X_i \Delta P_i$ and uses Assumption A1 to write all cross-elasticity terms as $X_i' \mathcal{E}_X X_j$, then factors out $X_i' \mathcal{E}_X$ from the full sum over j .

Step 2: The 2SLS estimand via Frisch–Waugh–Lovell. By the FWL theorem, the 2SLS coefficient on ΔP_i in regression (111) equals the coefficient from the residualized regression:

$$\hat{\mathcal{E}} = \frac{\sum_i Z_i^\perp \cdot \Delta D_i}{\sum_i Z_i^\perp \cdot \Delta P_i}, \quad (118)$$

where $Z^\perp = \mathbf{Q}_X Z$ denotes the residual of Z after projecting out X .

Substituting (117) into the numerator:

$$\begin{aligned} \text{Numerator} &= \sum_i Z_i^\perp (\mathcal{E}_{relative,i} \Delta P_i + X_i' \theta + \epsilon_i) \\ &= \sum_i \mathcal{E}_{relative,i} Z_i^\perp \Delta P_i + \theta' \underbrace{\left(\sum_i X_i Z_i^\perp \right)}_{= X' \mathbf{Q}_X Z = 0} + \sum_i Z_i^\perp \epsilon_i \\ &= \sum_i \mathcal{E}_{relative,i} Z_i^\perp \Delta P_i + \sum_i Z_i^\perp \epsilon_i. \end{aligned} \quad (119)$$

The $X' \theta$ terms vanish because \mathbf{Q}_X is the projector orthogonal to the column space of X .

Step 3: Expectations over instruments. Conditional on all structural parameters $\Theta = (\mathcal{E}_{relative,i}, \Lambda_{kl}, X_i, \mu_i, \epsilon_i)_{i,k,l}$, only Z is random. We compute the expected numerator and denominator.

Since $Z^\perp = \mathbf{Q}_X Z$, the i -th element is $Z_i^\perp = \sum_k [\mathbf{Q}_X]_{ik} Z_k$. Using the structural price change $\Delta P_i = \sum_j \Lambda_{ij} Z_j + \mu_i$, independence of Z from μ , and $\mathbf{E}[Z_k Z_j] = \sigma_Z^2 \mathbf{1}_{k=j}$ (since instruments are zero-mean and independent):

$$\begin{aligned} \mathbf{E}[Z_i^\perp \Delta P_i | \Theta] &= \sum_k \sum_j [\mathbf{Q}_X]_{ik} \Lambda_{ij} \mathbf{E}[Z_k Z_j] \\ &= \sigma_Z^2 \sum_j [\mathbf{Q}_X]_{ij} \Lambda_{ij} = \sigma_Z^2 [\mathbf{Q}_X \mathbf{\Lambda}']_{ii}, \end{aligned} \quad (120)$$

where the last equality uses the definition $[\mathbf{Q}_X \mathbf{\Lambda}']_{ii} = \sum_j [\mathbf{Q}_X]_{ij} [\mathbf{\Lambda}']_{ji} = \sum_j [\mathbf{Q}_X]_{ij} \Lambda_{ij}$.

For the denominator, summing (120) over i :

$$\mathbf{E}[\text{Denominator} | \Theta] = \sigma_Z^2 \sum_i [\mathbf{Q}_X \mathbf{\Lambda}']_{ii} = \sigma_Z^2 \text{tr}(\mathbf{Q}_X \mathbf{\Lambda}). \quad (121)$$

For the numerator, the ϵ term vanishes since $\epsilon_i \in \Theta$ and the instruments are zero-mean: $\mathbf{E}[Z_i^\perp \epsilon_i | \Theta] = \epsilon_i \sum_k [\mathbf{Q}_X]_{ik} \mathbf{E}[Z_k] = 0$. From (119) and (120):

$$\mathbf{E}[\text{Numerator} | \Theta] = \sigma_Z^2 \sum_i \mathcal{E}_{relative,i} [\mathbf{Q}_X \mathbf{\Lambda}']_{ii}. \quad (122)$$

Step 4: The conditional estimand. The numerator and denominator are quadratic forms in Z . As $N \rightarrow \infty$, since instruments are independently assigned across assets and first-stage exposure is dispersed across many assets, cross-sectional averages satisfy a law of large numbers. Both the numerator and denominator concentrate around their conditional expectations.^{51,52}

Since $\mathbf{E}[\text{Denom} | \Theta] = \sigma_Z^2 \text{tr}(\mathbf{Q}_X \mathbf{\Lambda}) \neq 0$ by relevance, the continuous mapping theorem implies

$$\hat{\mathcal{E}} \xrightarrow{p} \frac{\sum_i \mathcal{E}_{relative,i} [\mathbf{Q}_X \mathbf{\Lambda}']_{ii}}{\text{tr}(\mathbf{Q}_X \mathbf{\Lambda})} = \sum_i \omega_i \mathcal{E}_{relative,i}, \quad (123)$$

with weights $\omega_i = [\mathbf{Q}_X \mathbf{\Lambda}']_{ii} / \text{tr}(\mathbf{Q}_X \mathbf{\Lambda})$ that depend on $(\mathbf{\Lambda}, X)$ but not on $\mathcal{E}_{relative,i}$.

Step 5: Non-negative weights under bilinear $\mathbf{\Lambda}$ with monotonicity. Suppose $\Lambda_{ij} = X_i' \mathbf{M}_X X_j$ for all $i \neq j$, so that $\mathbf{\Lambda} = \text{diag}(\lambda_i) + X \mathbf{M}_X X'$ with $\lambda_i = \Lambda_{ii} - X_i' \mathbf{M}_X X_i$. Then $\mathbf{\Lambda}' =$

⁵¹This requires that first-stage exposure is not concentrated in a vanishing fraction of assets as N grows. Intuitively, as additional assets are added to the cross-section, sufficiently many must contribute independent residualized first-stage variation after partialling out X ; no single asset may account for a non-negligible share of total variation.

⁵²Independence of instruments across assets ensures that cross-terms vanish when taking expectations of quadratic forms such as $\sum_i Z_i^\perp \Delta P_i$, so that averaging across assets yields concentration.

$\text{diag}(\lambda_i) + X\mathbf{M}'_X X'$, and since $\mathbf{Q}_X X = 0$:

$$\mathbf{Q}_X \boldsymbol{\Lambda}' = \mathbf{Q}_X \text{diag}(\lambda_i). \quad (124)$$

Thus $[\mathbf{Q}_X \boldsymbol{\Lambda}']_{ii} = [\mathbf{Q}_X]_{ii} \lambda_i$. Since \mathbf{Q}_X is a symmetric idempotent matrix, its diagonal elements satisfy $[\mathbf{Q}_X]_{ii} \in [0, 1]$. Under the monotonicity condition $\lambda_i \geq 0$ for all i , we have $[\mathbf{Q}_X \boldsymbol{\Lambda}']_{ii} \geq 0$ for all i , so $\omega_i \geq 0$.

■

C.7 Local Average Treatment Effect relaxing Assumption A1

The data-generating process under heterogeneous treatment effects is:

$$\Delta D_i = \mathcal{E}_{ii} \Delta P_i + \sum_{j \neq i} \mathcal{E}_{ij} \Delta P_j + \epsilon_i \quad (125)$$

$$\Delta P_i = \lambda_i Z_i + \mu_i \quad (126)$$

This setup removes Assumption A1 of homogeneous substitution conditional on observables, but restricts the data-generating process to feature no spillovers in how instruments transmit to prices.

The instrument Z_i , with constant variance $\text{var}(Z_i) = \text{var}(Z), \forall i$, is randomly assigned and independent of everything else:

$$Z_i \perp\!\!\!\perp Z_j \quad \forall i \neq j \quad (127)$$

$$Z_i \perp\!\!\!\perp \mathcal{E}_{kl} \quad \forall i, k, l \quad (128)$$

$$Z_i \perp\!\!\!\perp \lambda_j \quad \forall i, j \quad (129)$$

$$Z_i \perp\!\!\!\perp \mu_j \quad \forall i, j \quad (130)$$

$$Z_i \perp\!\!\!\perp \epsilon_j \quad \forall i, j \quad (131)$$

After substituting (126) into (125), we derive the estimate from the demand equation

$$\Delta D_i = \mathcal{E}_{ii} \lambda_i Z_i + \sum_{j \neq i} \mathcal{E}_{ij} \lambda_j Z_j + \mathcal{E}_{ii} \mu_i + \sum_{j \neq i} \mathcal{E}_{ij} \mu_j + \epsilon_i \quad (132)$$

We now state and prove the following Proposition on the LATE regression:

Proposition 12 *Assume that the data-generating process of the first stage follows:*

$$\Delta P_i = \lambda_i Z_i + \mu_i, \quad \text{with } Z_i \text{ independent of } (\mu_i, \lambda_i), \quad (133)$$

and that the instrument is independent of own- and cross-price elasticities as well as the demand residual

$$(\mathcal{E}_{ii}, \mathcal{E}_{ij}, \epsilon_i) | Z_i \sim (\mathcal{E}_{ii}, \mathcal{E}_{ij}, \epsilon_i). \quad (134)$$

Then, the two-stage least square estimation of equations (17) and (18) without observables

identifies the local average of the relative elasticity:

$$\hat{\mathcal{E}} = \frac{\mathbf{E}_i \{ \lambda_i (\mathcal{E}_{ii} - \mathbf{E}_j(\mathcal{E}_{ji})) \}}{\mathbf{E}_i(\lambda_i)}. \quad (135)$$

Proof. Without loss of generality, define a centered instrument \tilde{Z}_i as

$$\tilde{Z}_i \equiv Z_i - \frac{1}{N} \sum_j Z_j, \quad (136)$$

such that we have the following properties:

$$\sum_{j \neq i} \tilde{Z}_j = -\tilde{Z}_i \quad (137)$$

$$\begin{aligned} \text{cov}(\tilde{Z}_i, \tilde{Z}_j) &= \underbrace{\text{cov}(Z_i, Z_j)}_{=0} - \frac{1}{N} \underbrace{\sum_k \text{cov}(Z_k, Z_j)}_{=\text{var}(Z)} - \frac{1}{N} \underbrace{\sum_l \text{cov}(Z_i, Z_l)}_{=\text{var}(Z)} + \frac{1}{N^2} \underbrace{\text{cov}\left(\sum_k Z_k, \sum_l Z_l\right)}_{=N \text{var}(Z)} \\ &= -\frac{1}{N} \text{var}(Z), \quad \forall j \neq i \end{aligned} \quad (138)$$

$$= -\frac{1}{N} \text{var}(Z), \quad \forall j \neq i \quad (139)$$

$$\text{var}(\tilde{Z}) = \text{var}\left(Z_i - \frac{1}{N} \sum_j Z_j\right) \quad (140)$$

$$= \underbrace{\text{var}(Z_i)}_{=\text{var}(Z)} - \frac{2}{N} \underbrace{\text{cov}\left(Z_i, \sum_j Z_j\right)}_{=\text{var}(Z)} + \frac{1}{N^2} \underbrace{\text{var}\left(\sum_j Z_j\right)}_{=N \text{var}(Z)} \quad (141)$$

$$= \text{var}(Z) \left(1 - \frac{2}{N} + \frac{1}{N}\right) = \frac{N-1}{N} \text{var}(Z) \quad (142)$$

$$\text{cov}(Z_j, \tilde{Z}_i) = \text{cov}\left(Z_j, Z_i - \frac{1}{N} \sum_k Z_k\right) \quad (143)$$

$$= \text{cov}(Z_j, Z_i) - \frac{1}{N} \underbrace{\text{cov}\left(Z_j, \sum_k Z_k\right)}_{=\text{var}(Z)} \quad (144)$$

$$= \begin{cases} \text{var}(Z) - \frac{1}{N} \text{var}(Z) = \frac{N-1}{N} \text{var}(Z) = \text{var}(\tilde{Z}) & \text{if } j = i \\ 0 - \frac{1}{N} \text{var}(Z) = -\frac{\text{var}(Z)}{N} = -\frac{\text{var}(\tilde{Z})}{N-1} & \text{if } j \neq i \end{cases} \quad (145)$$

We are interested in $\text{cov}(\Delta D_i, \tilde{Z}_i)$ and $\text{cov}(\Delta P_i, \tilde{Z}_i)$. Since \tilde{Z}_i is mean-zero, by the law of iterated expectations we have:

$$\text{cov}(\Delta D_i, \tilde{Z}_i) = \mathbb{E} \left[\Delta D_i \tilde{Z}_i \right] = \mathbb{E} \left[\mathbb{E} \left[\Delta D_i \tilde{Z}_i | \Theta \right] \right], \quad (146)$$

where Θ is a set that contains all \mathcal{E}_{ij} and λ_i .

The first-stage covariance is:

$$\text{cov} \left(\Delta P_i, \tilde{Z}_i \right) = \mathbb{E} \left[\mathbb{E} \left[\Delta P_i \tilde{Z}_i | \Theta \right] \right] = \mathbb{E} \left[\lambda_i \text{var}(\tilde{Z}) \right] = \text{var}(\tilde{Z}) \mathbb{E} [\lambda_i]. \quad (147)$$

For the demand side, starting from equation (132):

$$\mathbb{E} \left[\Delta D_i \tilde{Z}_i | \Theta \right] = \mathcal{E}_{ii} \lambda_i \mathbb{E}[Z_i \tilde{Z}_i] + \sum_{j \neq i} \mathcal{E}_{ij} \lambda_j \mathbb{E}[Z_j \tilde{Z}_i] \quad (148)$$

$$= \text{var}(\tilde{Z}) \left(\mathcal{E}_{ii} \lambda_i - \frac{1}{N-1} \sum_{j \neq i} \mathcal{E}_{ij} \lambda_j \right), \quad (149)$$

where we used the independence of Z from (μ, ϵ) .

Taking expectations over Θ :

$$\text{cov}(\Delta D_i, \tilde{Z}_i) = \text{var}(\tilde{Z}) \left(\mathbb{E}[\mathcal{E}_{ii} \lambda_i] - \frac{1}{N-1} \mathbb{E} \left[\sum_{j \neq i} \mathcal{E}_{ij} \lambda_j \right] \right). \quad (150)$$

For the second term, apply an index swap:

$$\mathbb{E} \left[\sum_{j \neq i} \mathcal{E}_{ij} \lambda_j \right] = \frac{1}{N} \sum_i \sum_{j \neq i} \mathcal{E}_{ij} \lambda_j \stackrel{i \leftrightarrow j}{=} \frac{1}{N} \sum_i \sum_{j \neq i} \mathcal{E}_{ji} \lambda_i = \mathbb{E} \left[\lambda_i \sum_{j \neq i} \mathcal{E}_{ji} \right]. \quad (151)$$

Substituting back into (150):

$$\text{cov}(\Delta D_i, \tilde{Z}_i) = \text{var}(\tilde{Z}) \mathbb{E} \left\{ \lambda_i \left(\mathcal{E}_{ii} - \frac{1}{N-1} \sum_{j \neq i} \mathcal{E}_{ji} \right) \right\}. \quad (152)$$

The IV estimator is the ratio of (150) to (147):

$$\hat{\mathcal{E}} = \frac{\mathbb{E} \left\{ \lambda_i \left(\mathcal{E}_{ii} - \frac{1}{N-1} \sum_{j \neq i} \mathcal{E}_{ji} \right) \right\}}{\mathbb{E}[\lambda_i]} \quad (153)$$

$$= \frac{\mathbb{E} \left\{ \lambda_i (\mathcal{E}_{ii} - \mathbb{E}_{j \neq i}[\mathcal{E}_{ji}]) \right\}}{\mathbb{E}[\lambda_i]}, \quad (154)$$

where $\mathbb{E}_{j \neq i}[\mathcal{E}_{ji}] = \frac{1}{N-1} \sum_{j \neq i} \mathcal{E}_{ji}$ is the average cross-elasticity in column i , excluding the diagonal. This is the local average treatment effect: the relative elasticity $(\mathcal{E}_{ii} - \mathbb{E}_{j \neq i}[\mathcal{E}_{ji}])$ weighted by first-stage intensity λ_i . ■

D Robustness

D.1 Robustness to the assumptions

We assess the robustness of the identification result of Proposition 1 with respect to deviations from the assumptions. For simplicity, we focus on a case with two assets and deviations from Assumption A2. The argument generalizes to other types of deviations.

We first show that the conclusions are robust when the first stage is economically large. Then we illustrate potential issues in presence of weak instruments in the context of a model.

D.1.1 Robustness to deviations from Assumption A2

Consider a setting with two assets and an arbitrary elasticity matrix. In response to an exogenous shock, the relative change in demand is:

$$\Delta D_1 - \Delta D_2 = \mathcal{E}_{11}\Delta P_1 + \mathcal{E}_{12}\Delta P_2 - \mathcal{E}_{22}\Delta P_2 - \mathcal{E}_{21}\Delta P_1 \quad (155)$$

$$= (\mathcal{E}_{11} - \mathcal{E}_{21})\Delta P_1 + (\mathcal{E}_{22} - \mathcal{E}_{12})\Delta P_2 \quad (156)$$

We denote the relative elasticities by $\mathcal{E}_{rel,1} = \mathcal{E}_{11} - \mathcal{E}_{21}$ and $\mathcal{E}_{rel,2} = \mathcal{E}_{22} - \mathcal{E}_{12}$. Rearranging the terms leads to

$$\Delta D_1 - \Delta D_2 = \frac{\mathcal{E}_{rel,1} + \mathcal{E}_{rel,2}}{2} (\Delta P_1 - \Delta P_2) + \frac{\mathcal{E}_{rel,1} - \mathcal{E}_{rel,2}}{2} (\Delta P_1 + \Delta P_2) \quad (157)$$

Dividing by the relative change in price in response to the shock, we obtain the estimator:

$$\frac{\Delta D_1 - \Delta D_2}{\Delta P_1 - \Delta P_2} = \frac{\mathcal{E}_{rel,1} + \mathcal{E}_{rel,2}}{2} + \frac{\mathcal{E}_{rel,1} - \mathcal{E}_{rel,2}}{2} \frac{\Delta P_1 + \Delta P_2}{\Delta P_1 - \Delta P_2} \quad (158)$$

The first term is the average relative elasticity. The second term is the potential bias: the heterogeneity of relative elasticities times the ratio of sum of changes in prices to their difference. The denominator $\Delta P_1 - \Delta P_2$ is the first stage of the regression, and the relevance condition is $\Delta P_1 - \Delta P_2 \neq 0$.

It is straightforward to see that if the relevance condition is satisfied, the identification result of relative elasticity is robust to small deviations from assumption A2. Formally, consider a family of experiments (or models) indexed by a variable x , such that the experiment for $x = 0$ satisfies assumption A2 but it does not otherwise. If elasticities and changes in prices are continuous in x , then the IV estimator is also continuous in x as long as $\Delta P_1 - \Delta P_2 \neq 0$.

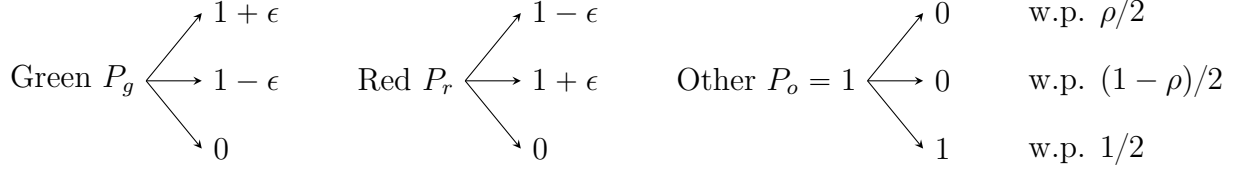
When the relevance condition is not satisfied, the last term of equation (158) becomes infinite, and the bias is large relative to the actual coefficient. The same issue arises with weak instruments in the neighborhood of this limit, when $\Delta P_1 - \Delta P_2$ is small. We illustrate this in an example next. In practice, one should assess the magnitude of the price spread induced by the instrument and confirm it is economically large, not merely non-zero.

D.1.2 An example with discontinuous estimates and a weak first stage

Setting. We consider a variation from the model of Section H.2, where the probability of the different states is asymmetric. Furthermore, we measure elasticities in the “wrong” units,

portfolio share on price rather than portfolio share on log price. These two features create deviations from the Assumption A2. The experiment is a shock to the endowment of the green asset E_g .

The setting is the same as before except for a change in the probability of the states:



The optimal portfolio shares are:

$$\omega_g(P_g, P_r) = \frac{P_g((\epsilon^2 - 1)P_g + P_r(4\rho\epsilon + (\epsilon - 1)^2))}{4(\epsilon^2 + 1)P_gP_r + 2(\epsilon^2 - 1)P_g^2 + 2(\epsilon^2 - 1)P_r^2} \quad (159)$$

$$\omega_r(P_g, P_r) = \frac{P_r(P_g((\epsilon + 1)^2 - 4\rho\epsilon) + (\epsilon^2 - 1)P_r)}{4(\epsilon^2 + 1)P_gP_r + 2(\epsilon^2 - 1)P_g^2 + 2(\epsilon^2 - 1)P_r^2} \quad (160)$$

Equilibrium and elasticities. Assume that the endowments are $E_g = E_r = 1/2$ and $E_o = 1$. Then equilibrium prices are

$$P_g = 1 - \epsilon(1 - 2\rho), \quad P_r = 1 + \epsilon(1 - 2\rho). \quad (161)$$

At this equilibrium, the elasticity matrix of the portfolio shares with respect to the level of prices for the green and red asset is:

$$\mathcal{E}_{gg} = \frac{\partial \omega_g}{\partial P_g} = \frac{(\epsilon^2 - 1)((2\rho - 1)\epsilon - 1)}{32(\rho - 1)\rho\epsilon^2}, \quad (162)$$

$$\mathcal{E}_{rr} = \frac{\partial \omega_r}{\partial P_r} = -\frac{(\epsilon^2 - 1)((2\rho - 1)\epsilon + 1)}{32(\rho - 1)\rho\epsilon^2}, \quad (163)$$

$$\mathcal{E}_{gr} = \frac{\partial \omega_g}{\partial P_r} = \frac{(\epsilon^2 - 1)((2\rho - 1)\epsilon + 1)}{32(\rho - 1)\rho\epsilon^2}, \quad (164)$$

$$\mathcal{E}_{rg} = \frac{\partial \omega_r}{\partial P_g} = -\frac{(\epsilon^2 - 1)((2\rho - 1)\epsilon - 1)}{32(\rho - 1)\rho\epsilon^2}. \quad (165)$$

This leads to the two relative elasticities:

$$\mathcal{E}_{rel,g} = \mathcal{E}_{gg} - \mathcal{E}_{rg} = \frac{(\epsilon^2 - 1)((2\rho - 1)\epsilon - 1)}{16(\rho - 1)\rho\epsilon^2}, \quad (166)$$

$$\mathcal{E}_{rel,r} = \mathcal{E}_{rr} - \mathcal{E}_{gr} = \frac{(\epsilon^2 - 1)(-(2\rho - 1)\epsilon - 1)}{16(\rho - 1)\rho\epsilon^2} \quad (167)$$

Then the terms from the difference-in-difference estimator are:

$$\frac{\mathcal{E}_{rel,g} + \mathcal{E}_{rel,r}}{2} = \frac{1 - \epsilon^2}{16(\rho - 1)\rho\epsilon^2} \quad (168)$$

$$\frac{\mathcal{E}_{rel,g} - \mathcal{E}_{rel,r}}{2} = \frac{(\epsilon^2 - 1)(2\rho - 1)}{16(\rho - 1)\rho\epsilon^2} \quad (169)$$

Note that constant relative elasticity, assumption A2, holds only if $\rho = \frac{1}{2}$. We work in a neighborhood of assumption A2, where $\rho \sim \frac{1}{2}$. To weaken the first stage and make the assets perfect substitute, we take ϵ to zero. These expressions become approximately

$$\frac{\mathcal{E}_{rel,g} + \mathcal{E}_{rel,r}}{2} \approx \frac{-1}{4\epsilon^2} \quad (170)$$

$$\frac{\mathcal{E}_{rel,g} - \mathcal{E}_{rel,r}}{2} \approx \frac{2\rho - 1}{4\epsilon^2} \quad (171)$$

Equilibrium prices as a function of the endowments are

$$P_g(E_o, E_g, E_r) = \frac{E_o((\epsilon^2 - 1)E_g - E_r(4\rho\epsilon + (\epsilon - 1)^2))}{-2(\epsilon^2 + 1)E_gE_r + (\epsilon^2 - 1)E_g^2 + (\epsilon^2 - 1)E_r^2} \quad (172)$$

$$P_r(E_o, E_g, E_r) = \frac{E_o((\epsilon^2 - 1)E_r - E_g((\epsilon + 1)^2 - 4\rho\epsilon))}{-2(\epsilon^2 + 1)E_gE_r + (\epsilon^2 - 1)E_g^2 + (\epsilon^2 - 1)E_r^2}. \quad (173)$$

Around the initial equilibrium, the changes in prices are:

$$\Delta P_g = \frac{\partial P_g}{\partial E_g} = -4\epsilon\rho - (\epsilon - 1)^2 \quad (174)$$

$$\Delta P_r = \frac{\partial P_r}{\partial E_g} = \epsilon^2 - 1 \quad (175)$$

The term controlling the bias is:

$$\frac{\Delta P_g + \Delta P_r}{\Delta P_g - \Delta P_r} = \frac{2\epsilon\rho - \epsilon + 1}{\epsilon(2\rho + \epsilon - 1)} \quad (176)$$

Putting it all together. We will study what happens around $\rho = 1/2$, so, echoing our general setup, we call $x = 2\rho - 1$. When $x = 0$, assumption A2 is satisfied, and the estimator is unbiased.

We plug all the expressions above in equation (158):

$$\frac{\Delta D_1 - \Delta D_2}{\Delta P_1 - \Delta P_2} = \frac{\mathcal{E}_{rel,1} + \mathcal{E}_{rel,2}}{2} + \frac{\mathcal{E}_{rel,1} - \mathcal{E}_{rel,2}}{2} \frac{\Delta P_1 + \Delta P_2}{\Delta P_1 - \Delta P_2} \quad (177)$$

$$\approx \frac{-1}{4\epsilon^2} + \frac{x}{4\epsilon^2} \frac{1}{\epsilon(x + \epsilon)} \quad (178)$$

Both the average relative elasticity and the difference in relative elasticity go to infinity at the same pace ($1/\epsilon^2$). However, the weak first stage amplifies the bias by another order of

magnitude. To visualize this issue, it is more natural to compute the relative bias of the estimator:

$$\frac{\frac{\Delta D_1 - \Delta D_2}{\Delta P_1 - \Delta P_2} - \frac{\mathcal{E}_{rel,1} + \mathcal{E}_{rel,2}}{2}}{\frac{\mathcal{E}_{rel,1} + \mathcal{E}_{rel,2}}{2}} \approx -\frac{x}{\epsilon(x + \epsilon)} \quad (179)$$

The bias term present when $x \neq 0$ is an order of magnitude large than the correct estimate in the limit of a weak instrument.

D.2 Validating the plausibility of assumptions

D.2.1 Plausibility of relative multiplier regressions

Complementary to the balance-on-covariance diagnostics in Figure 4, one approach to assessing the plausibility of Assumption A1 is to progressively restrict the sample to more homogeneous subsets, so that the assumption of homogeneous substitution conditional on observables is increasingly likely to hold, and check whether the estimated relative multiplier $\widehat{\mathcal{M}}$ changes. This is the same idea as for experimental design in Section 2.2.3: the more similar the bonds being compared, the less scope for violations of Assumption A1.

In the discrete choice literature, this logic underlies the [Hausman and McFadden \(1984\)](#) test for the Independence of Irrelevant Alternatives, which compares parameter estimates from the full choice set against those from a restricted subset, or nest. In corporate bonds, a natural and narrow nest is the bond issuer: bonds from the same issuer share many characteristics, so within-issuer variation provides a more controlled comparison. We implement this test via the [Mundlak \(1978\)](#) regression approach.

Table 3 provides the results, testing whether the relative multiplier estimated only from within-issuer variation (specification (2)) differs from the pooled estimate (specification (1), or specification (1) from Table 1). The relevant coefficient is the one on $\bar{Z}_{issuer,it}$, which is defined as the across-bond average realization of the instrument Z_{it} for each issuer-date pair, and which is small and statistically insignificant. This is even though adding bond issuer \times date fixed effects absorbs a lot of additional variation, as seen from the increase in the R^2 coefficient between specifications (1) and (2), and tightens standard errors. The results suggest that relative multipliers estimated within issuers are not significantly different from estimates from our main specification. Given the tension between estimating relative multipliers and spillovers—that each additional observable used to strengthen causal inference for relative multipliers requires its own source of exogenous variation for spillovers—the within-issuer result is reassuring: it suggests that the pooled estimate is valid, without requiring a separate time-series instrument for each issuer to operationalize spillover estimation at that level of disaggregation.

A similar idea can be applied to non-nested models for tests of omitted observables: does adding additional potentially relevant observables affect the estimate of the relative multiplier $\widehat{\mathcal{M}}$? After duration and credit risk, another natural characteristic of a corporate bond is its liquidity. Table 4 shows that adding a proxy for bond liquidity, log outstanding amount, does not significantly change the relative multiplier.

We want to emphasize that we view these diagnostics as an informative tool for researchers

Table 3: Hausman test for Relative multiplier $\widehat{\mathcal{M}}$

	Return R_{it}		
	(1)	(2)	(3)
Z_{it}	0.055 (0.084)	0.067 (0.056)	0.067 (0.056)
$\bar{Z}_{issuer,it}$			-0.048 (0.122)
Date Fixed Effects	Yes		Yes
$X_{it}^1 \times$ Date Fixed Effects	Yes	Yes	Yes
$X_{it}^2 \times$ Date Fixed Effects	Yes	Yes	Yes
Issuer \times Date Fixed Effects		Yes	
$\bar{X}_{issuer,it}^1 \times$ Date Fixed Effects			Yes
$\bar{X}_{issuer,it}^2 \times$ Date Fixed Effects			Yes
N	1,041,985	1,041,985	1,041,985
R^2	0.464	0.783	0.466

Table 3 reports the results of a Hausman (1978) test for the relative multiplier of bond returns R_{it} on flow-induced trading demand shock Z_{it} defined in (46) for U.S. non-defaulted corporate bonds. Specification (1) is the same as specification (3) in Table 1, controlling for date fixed effects and a continuous duration variable X_{it}^1 and credit risk variable based on average historical default probabilities for each S&P credit rating category X_{it}^2 for each date. Specification (2) zooms in further by replacing date fixed effects with bond issuer \times date fixed effects. Specification (3) adds the Hausman (1978) test via a Mundlak (1978) regression, testing whether the coefficients in (1) and (2) are different through the coefficient on the issuer-average of the instrument, $\bar{Z}_{issuer,it}$. The regressions weigh each date equally. The sample period is 2010:04 to 2024:03. Standard errors are clustered by date and bond issuer.

to support the plausibility of Assumption A1, but not as definite tests, as they are neither necessary nor sufficient for A1 to hold. In part, this is because rejection can occur because Assumption A1 is violated given a set of observables, so that zooming in further or adding observables means less violation and hence a more consistent estimator, but it can also occur because of heterogeneous treatment effects as described in Section 3.1.2.⁵³ For instance, if the relative multiplier differs across bond issuers, the within-issuer estimator in specification (2) captures a different weighted average of issuer-level multipliers than the pooled estimator in specification (1), causing the Mundlak coefficient to reject even absent any violation of A1. On the flip side, a failure to reject in these tests is not sufficient, given the inherent limitation of this approach of pitting one model against another. There may yet be more omitted observables, or even unobservables, that drive substitution and spillovers, the potential existence of which implies that failure to reject does not definitively show that A1 is satisfied. That said, the within-issuer comparison in Table 3 is particularly informative in this regard, as it absorbs all issuer-level unobservables, so that only within-issuer variation

⁵³Goldsmith-Pinkham et al. (2020) emphasize a similar diagnostic result in the context of heterogeneity with Bartik instruments.

Table 4: Coefficient stability test for Relative multiplier \widehat{M} : bond liquidity

	Return R_{it}	
	(1)	(2)
Z_{it}	0.055 (0.084)	0.055 (0.084)
$D \times Z_{it}$		-0.022 (0.017)
date Fixed Effects	Yes	
$X_{it}^1 \times$ date Fixed Effects	Yes	
$X_{it}^2 \times$ date Fixed Effects	Yes	
$D \times$ Date Fixed Effects		Yes
$X_{it}^1 \times D \times$ Date Fixed Effects		Yes
$X_{it}^2 \times D \times$ Date Fixed Effects		Yes
$D \times X_{it}^3 \times D \times$ Date Fixed Effects		Yes
N	1,041,985	2,083,970
R^2	0.464	0.467

Table 4 reports the results of a coefficient stability test for the relative multiplier of bond returns R_{it} on flow-induced trading demand shock Z_{it} defined in (46) for U.S. non-defaulted corporate bonds. Specification (1) is the same as specification (3) in Table 1, controlling for date fixed effects and a continuous duration variable X_{it}^1 and credit risk variable X_{it}^2 for each date. To test whether adding a liquidity control changes the coefficient on Z_{it} , the dataset is stacked: two copies of the data are created, indexed by $D \in \{0, 1\}$, with all fixed effects and controls interacted with D so that each copy estimates its own specification. The liquidity control, standardized log bond amount outstanding (e.g., Houweling et al., 2005), enters only the $D = 1$ copy. The coefficient on the interaction $D \times Z_{it}$ directly tests whether the relative multiplier differs across the two specifications. The regressions weigh each date equally. The sample period is 2010:04 to 2024:03. Standard errors are clustered by date and bond issuer.

in unobservables across bonds poses a remaining threat.

E Demand beyond risk-based motives for substitution

Consider the problem of investors combining risk-based mean-variance demand where the covariance matrix Σ is characterized by a set of characteristics $X^{(3)}$ with a cost of holding assets that is quadratic in another set of characteristics $X^{(1)}$ and a portfolio constraint linear in yet another set characteristics $X^{(2)}$. Section 2.2.2 is a special case of this that assumes that $X^{(1)}$ and $X^{(2)}$ each only contain one observable: carbon intensity and a bank's liquidity ratio. The proposition below generalizes this.

Proposition 13 (Mean-variance demand with quadratic cost and linear constraint)

Assume that investors choose their demand according to the problem

$$\max_D D'(M - P) - \frac{\gamma}{2} D' \Sigma D - \frac{\kappa}{2} D' X^{(1)} X^{(1)'} D \quad (180)$$

$$\text{such that } D' X^{(2)} \leq \Theta, \quad (181)$$

where D , M , and P are the $N \times 1$ vectors of investor demand, expected payoffs, and prices, γ is risk aversion, Σ the $N \times N$ covariance matrix, κ controls the quadratic cost function, $X^{(1)}$ and $X^{(2)}$ are the $N \times K_1$ and $N \times K_2$ matrices of stock characteristics, and Θ is a $1 \times K_2$ vector that controls the linear constraint.

Further assume that the risk-based component of investor demand satisfies assumptions A1 and A2, i.e.,

$$-\frac{1}{\gamma} \Sigma^{-1} = \hat{\mathcal{E}}^{(3)} I + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'}, \quad (182)$$

where $X^{(3)}$ is another $N \times K_3$ matrix of stock characteristics, and $\mathcal{E}_{X^{(3)}}$ the $K_3 \times K_3$ matrix of substitution between observables.

Then the resulting demand curve satisfies assumptions A1 and A2 conditional on the stacked observables $\mathbb{X} = [X^{(1)}, X^{(2)}, X^{(3)}]$.

Proof. By Lemma 6, to proof the proposition, we need to show that the elasticity matrix \mathcal{E} can be expressed as

$$\mathcal{E} = \hat{\mathcal{E}} I + \mathbb{X} \mathcal{E}_{\mathbb{X}} \mathbb{X}'. \quad (183)$$

Start by putting together equations (180) and (181) in the Lagrangian

$$\mathcal{L}(D, \lambda) = D'(M - P) - \frac{\gamma}{2} D' \Sigma D - \frac{\kappa}{2} D' X^{(1)} X^{(1)'} D - \lambda(D' X^{(2)} - \Theta), \quad (184)$$

where λ is the $K_2 \times 1$ Lagrange multiplier on the linear constraint.

Setting the first-order condition with respect to D to zero, and solving for D , yields

$$D = \left(\underbrace{\gamma \Sigma + \kappa X^{(1)} X^{(1)'}}_{\equiv \Omega} \right)^{-1} (M - P - X^{(2)} \lambda) \quad (185)$$

$$= \Omega^{-1} (M - P - X^{(2)} \lambda), \quad (186)$$

where

$$\Omega = \gamma \Sigma + \kappa X^{(1)} X^{(1)'}. \quad (187)$$

Plugging into the linear constraint to solve for λ :

$$(M - P - X^{(2)}\lambda)' \Omega^{-1} X^{(2)} = \Theta \quad (188)$$

$$\implies (M - P)' \Omega^{-1} X^{(2)} - \lambda' X^{(2)'} \Omega^{-1} X^{(2)} = \Theta \quad (189)$$

$$\implies \lambda = [X^{(2)'} \Omega^{-1} X^{(2)}]^{-1} [X^{(2)'} \Omega^{-1} (M - P) - \Theta']^+ \quad (190)$$

Plugging the Lagrange multipliers back into optimal investor demand gives:

$$D = \Omega^{-1} (M - P) - \Omega^{-1} X^{(2)} [X^{(2)'} \Omega^{-1} X^{(2)}]^{-1} [X^{(2)'} \Omega^{-1} (M - P) - \Theta']^+ \quad (191)$$

The elasticity matrix therefore is:

$$\frac{dD}{dP} = -\Omega^{-1} + \Omega^{-1} X^{(2)} S_b [S_b' X^{(2)'} \Omega^{-1} X^{(2)} S_b]^{-1} S_b' X^{(2)'} \Omega^{-1} \quad (192)$$

Here, S_b is the binding constraint selection matrix, which for the first-order condition selects the columns of $X^{(2)}$ for which constraints are binding.

Start now with the part for when the inequality constraints are all non-binding:

$$-\Omega^{-1} = \hat{\mathcal{E}}^{(3)} I + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} \quad (193)$$

$$+ \left(\hat{\mathcal{E}}^{(3)} I + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} \right) X^{(1)} \underbrace{\left[\frac{1}{\kappa} I - X^{(1)'} (\gamma \Sigma)^{-1} X^{(1)} \right]^{-1}}_{\equiv \mathcal{H}} X^{(1)'} \left(\hat{\mathcal{E}}^{(3)} I + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} \right) \quad (194)$$

$$= \hat{\mathcal{E}}^{(3)} I + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} + \left(\hat{\mathcal{E}}^{(3)} \right)^2 X^{(1)} \mathcal{H} X^{(1)'} + \hat{\mathcal{E}}^{(3)} X^{(1)} \mathcal{H} X^{(1)'} X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} \quad (195)$$

$$+ \hat{\mathcal{E}}^{(3)} X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} X^{(1)} \mathcal{H} X^{(1)'} + X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} X^{(1)} \mathcal{H} X^{(1)'} X^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} \quad (196)$$

$$= \hat{\mathcal{E}}^{(3)} I + \underbrace{[X^{(1)}, X^{(3)}]}_{\equiv X^{(1,3)}} \underbrace{\begin{bmatrix} \left(\hat{\mathcal{E}}^{(3)} \right)^2 \mathcal{H} & \hat{\mathcal{E}}^{(3)} \mathcal{H} X^{(1)'} X^{(3)} \mathcal{E}_{X^{(3)}} \\ \hat{\mathcal{E}}^{(3)} \mathcal{E}_{X^{(3)}} X^{(3)'} X^{(1)} \mathcal{H} & \mathcal{E}_{X^{(3)}} + \mathcal{E}_{X^{(3)}} X^{(3)'} X^{(1)} \mathcal{H} X^{(1)'} X^{(3)} \mathcal{E}_{X^{(3)}} \end{bmatrix}}_{\equiv \mathcal{F}} [X^{(1)}, X^{(3)}]' \quad (197)$$

$$= \hat{\mathcal{E}}^{(3)} I + X^{(1,3)} \mathcal{F} X^{(1,3)'} \quad (198)$$

When the linear constraints are not binding, the elasticity matrix satisfies assumptions A1 and A2 conditional on the stacked observables $[X^{(1)}, X^{(3)}]$.

For the case that some constraints are not binding, define:

$$\mathcal{G} \equiv S_b [S_b' X^{(2)'} \Omega^{-1} X^{(2)} S_b]^{-1} S_b' \quad (199)$$

The elasticity matrix is

$$\frac{dD}{dP} = \widehat{\mathcal{E}}^{(3)} I + X^{(1,3)} \mathcal{F} X^{(1,3)'} + \left(\widehat{\mathcal{E}}^{(3)} I + X^{(1,3)} \mathcal{F} X^{(1,3)'} \right) X^{(2)} \mathcal{G} X^{(2)'} \left(\widehat{\mathcal{E}}^{(3)} I + X^{(1,3)} \mathcal{F} X^{(1,3)'} \right) \quad (200)$$

$$= \widehat{\mathcal{E}} I + \mathbb{X} \mathcal{E}_{\mathbb{X}} \mathbb{X}', \quad (201)$$

where

$$\widehat{\mathcal{E}} = \widehat{\mathcal{E}}^{(3)} \quad (202)$$

$$\mathbb{X} = [X^{(1)}, X^{(3)}, X^{(2)}] \quad (203)$$

$$\mathcal{E}_{\mathbb{X}} = \begin{bmatrix} \mathcal{F} + \mathcal{F} X^{(1,3)'} X^{(2)} \mathcal{G} X^{(2)'} X^{(1,3)} \mathcal{F} & \widehat{\mathcal{E}}^{(3)} \mathcal{F} X^{(1,3)'} X^{(2)} \mathcal{G} \\ \widehat{\mathcal{E}}^{(3)} \mathcal{G} X^{(2)'} X^{(1,3)} \mathcal{F} & \left(\widehat{\mathcal{E}}^{(3)} \right)^2 \mathcal{G} \end{bmatrix}. \quad (204)$$

The elasticity matrix satisfies assumptions A1 and A2 conditional on the stacked observables \mathbb{X} .

■

F A non-linear framework

We derive properties for a family of non-linear demand functions which satisfy locally our assumption of homogeneous substitution conditional on observables. Doing so provides more general intuition behind our results in linear structures.

Because the non-linear structural models considered in [Kojen and Yogo \(2019\)](#) also belong to this family of demand functions, this framework also allows us to better understand the connection of our results with properties of those models. In particular, we explain the restrictions imposed by the logit form relative to arbitrary factor models, simple factor models (with constant variance and expected payoffs), and more general demand functions.

F.1 Basic concepts

We consider a setting with an investor, N assets indexed by i , and K observables for each asset. We start with a general demand function defined as a mapping from the vector of (log) prices \mathbf{p} and the $N \times K$ matrix of observables \mathbf{x} to a vector of positions D (portfolio shares in our applications):

$$D(\mathbf{p}, \mathbf{x}) : \mathbb{R}^N \times \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^N$$

It will be helpful to define the following property.

Definition 14 (HCO functions) *A function $F : \mathbb{R}^N \times \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^N$ is homogenous-conditional-on-observables (HCO) if $\forall i, [F(\mathbf{p}, \mathbf{x})]_i = f(p_i, x_i; \mathbf{p}, \mathbf{x})$ for a function $f : (\mathbb{R} \times \mathbb{R}^K) \times (\mathbb{R}^N \times \mathbb{R}^{N \times K}) \rightarrow \mathbb{R}$ for each i .*

That is, for a fixed overall price and observables vector, the value for each element is given by the same (scalar-valued) function of its own price and observables.

Then, in the spirit of the discussion of properties of factor models of [Kojien and Yogo \(2019\)](#), we can restrict attention to a subset of general demand functions as follows.

Definition 15 (HCO demand) *A demand function is a homogenous-conditional-on-observables demand if it is a HCO function.*

With HCO demand functions, individual positions can be written as:

$$[D(\mathbf{p}, \mathbf{x})]_i = d(p_i, x_i; \mathbf{p}, \mathbf{x}).$$

This notation emphasizes the dual role of prices and observables. On the one hand, the same function $d(\cdot, \cdot; \mathbf{p}, \mathbf{x})$ describes how the demand of each asset depends on its own price and own observables only. On the other hand, this mapping varies with the vector (\mathbf{p}, \mathbf{x}) . Thinking of this vector as the state of the economy, a HCO demand describes a mapping which is possibly state-dependent, but identical across assets.

Naturally, the choice of observables is what gives meaningful restrictions to this definition. For example, if the observables \mathbf{x} includes each asset's "name", i , then all demand functions are also HCO demand.

An example of HCO demand is logit:

$$[D^{logit}(\mathbf{p}, \mathbf{x})]_i = \frac{\exp(-\alpha p_i + \beta' x_i)}{1 + \sum_{j=1}^N \exp(-\alpha p_j + \beta' x_j)},$$

because the numerator is a function of p_i and x_i only, while the denominator is a fixed function (i.e. that does not depend on i) of \mathbf{p} and \mathbf{x} .

F.2 Relative elasticity vs. substitution and identification.

HCO demand leads to a natural decomposition of the elasticity matrix between a relative elasticity and a substitution matrices. HCO demand implies an elasticity matrix:

$$\mathcal{E} = \frac{\partial D}{\partial \mathbf{p}} = \underbrace{\text{diag} \left(\frac{\partial d}{\partial p_i} (p_i, x_i; \mathbf{p}, \mathbf{x}) \right)}_{\text{relative elasticity, } N \times N} + \underbrace{\begin{bmatrix} \frac{\partial d}{\partial \mathbf{p}} (p_1, x_1; \mathbf{p}, \mathbf{x})' \\ \vdots \\ \frac{\partial d}{\partial \mathbf{p}} (p_N, x_N; \mathbf{p}, \mathbf{x})' \end{bmatrix}}_{\substack{1 \times N \\ \text{substitution, } N \times N}}$$

where the derivative in the second term is with respect to the third argument of d , not a total derivative. We call the first term relative elasticity and the second one substitution. To understand why the first one is a relative elasticity, notice that if two assets have the same price and observables, this term is equal to the difference between their own-price and

cross-price elasticity:

$$\begin{aligned}\mathcal{E}_{ii} - \mathcal{E}_{ji} &= \underbrace{\left(\frac{\partial d}{\partial p_i}(p_i, x_i; \mathbf{p}, \mathbf{x}) + \left[\frac{\partial d}{\partial \mathbf{p}}(p_i, x_i; \mathbf{p}, \mathbf{x}) \right]_i \right)}_{\mathcal{E}_{ii}} - \underbrace{\left[\frac{\partial d}{\partial \mathbf{p}}(p_j, x_j; \mathbf{p}, \mathbf{x}) \right]_i}_{\mathcal{E}_{ji}} \\ &= \frac{\partial d}{\partial p_i}(p_i, x_i; \mathbf{p}, \mathbf{x}) \text{ if } (p_i, x_i) = (p_j, x_j)\end{aligned}$$

The substitution matrix captures how investor reallocate between assets when their price change. Any HCO demand satisfies homogeneous substitution conditional on all observables and the price:

$$\mathcal{E}_{il} = \mathcal{E}_{jl} = \left[\frac{\partial d}{\partial \mathbf{p}}(p_i, x_i, \mathbf{p}, \mathbf{x}) \right]_l \text{ if } (p_i, x_i) = (p_j, x_j)$$

Indeed, this corresponds to assumption A1 in the text when the observables \mathbf{x} are variables that the econometrician can measure.

Identification. The cross-section can allow to identify relative elasticity by comparing demand for two assets with the same observables but nearby prices. However, because (\mathbf{p}, \mathbf{x}) are fixed in a given cross-section, identification of substitution is generally impossible with the cross-section.

This limitation of the cross-section can be overcome by imposing additional restrictions. For example, when there is no substitution, that is the demand function does not depend on the price vector per se (its third argument), we have: $d(p_i, x_i, \mathbf{p}, \mathbf{x}) = d(p_i, x_i; \mathbf{x})$. Elasticity is relative elasticity, and hence can be estimated from the cross-section alone.

Another case is logit. Even though this model has non-zero substitution, it can be estimated from the cross-section because parameters determining substitution can be identified by measuring relative elasticity. Specifically the relative elasticity vector is $-\alpha\omega$ and the substitution matrix is $\alpha\omega\omega'$ where $\omega = D(\mathbf{p}, \mathbf{x})$ is the realized vector of portfolio weights. This calculation also highlights that the structure of substitution is very restricted in logit: the substitution matrix of rank 1, and the effects must be proportional to portfolio weights. To better understand how limiting these restrictions are, we compare logit to other demand models in the following section.

F.3 Logit, log utility and factor models

Our main interest in this section is to what extent the demand of an investor with a standard utility function and particular views on the dynamics of expected returns might be represented with logit demand.

For this purpose, we consider the demand of a log investor with log-normal returns (like in Section H.1) as the simplest example of a standard utility.

$$D^{log}(\mathbf{p}, \mathbf{x}) = \Sigma(\mathbf{p}, \mathbf{x})^{-1} \mu(\mathbf{p}, \mathbf{x})$$

where $\mu(\mathbf{p}, \mathbf{x})$ is the expected return vector and $\Sigma(\mathbf{p}, \mathbf{x})$ is the return covariance matrix.

For the representation of views of the investor on the structure of asset returns, we define first a class of factor models where factor loadings might be state-dependent.

Definition 16 (General factor models) *A general M-factor model is defined by functions $\mu(\mathbf{p}, \mathbf{x}) : \mathbb{R}^N \times \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^N$ and $\Sigma(\mathbf{p}, \mathbf{x}) : \mathbb{R}^N \times \mathbb{R}^{N \times K} \rightarrow \mathbb{R}^{N \times N}$ describing the expected return vector and covariance matrix respectively such that*

$$\Sigma(\mathbf{p}, \mathbf{x}) = \text{diag}(\underbrace{\sigma_\epsilon^2(\mathbf{p}, \mathbf{x})}_{N \times 1}) + \underbrace{\beta(\mathbf{p}, \mathbf{x})}_{N \times M} \underbrace{\Sigma_F}_{M \times M} \beta(\mathbf{p}, \mathbf{x})'$$

, where $\mu(\mathbf{p}, \mathbf{x})$, $\sigma_\epsilon^2(\mathbf{p}, \mathbf{x})$ and each column of $\beta(\mathbf{p}, \mathbf{x})$ are HCO functions and Σ_F is a covariance matrix.

A general M-factor model lets the functions mapping prices and observables to expected payoffs, to factor loadings and to idiosyncratic risk for each asset to freely vary with the state of the economy (\mathbf{p}, \mathbf{x}) . As we show below, this is a rich enough set that for any HCO demand $D^{HCO}(\mathbf{p}, \mathbf{x})$ one can always find a particular factor model that the log investor's demand exactly corresponds to the chosen HCO demand: $D^{HCO}(\mathbf{p}, \mathbf{x}) = D^{log}(\mathbf{p}, \mathbf{x})$. This is a generalized version of Corollary 1 of [Kojien and Yogo \(2019\)](#) who specialize the function $D^{HCO}(\mathbf{p}, \mathbf{x})$ to be the logit demand.

Proposition 17 *Fix two functions $\sigma_\epsilon^2(\mathbf{p}, \mathbf{x})$ and an $N \times 1$ $\beta(\mathbf{p}, \mathbf{x})$ which are HCO. For any HCO demand $D^{HCO}(\mathbf{p}, \mathbf{x})$, there exists a general 1-factor model with the corresponding covariance matrix $\Sigma(\mathbf{p}, \mathbf{x}) = \text{diag}(\sigma_\epsilon^2(\mathbf{p}, \mathbf{x})) + \sigma_F^2 \underbrace{\beta(\mathbf{p}, \mathbf{x}) \beta(\mathbf{p}, \mathbf{x})'}_{N \times 1}$, such that log utility demand with this factor model yields the same demand function.*

Proof. Choose $\mu(\mathbf{p}, \mathbf{x}) = \Sigma(\mathbf{p}, \mathbf{x})D^{HCO}(\mathbf{p}, \mathbf{x})$ then clearly $D^{log}(\mathbf{p}, \mathbf{x}) = D^{HCO}(\mathbf{p}, \mathbf{x})$. Therefore, we have to show only that $\mu(\mathbf{p}, \mathbf{x})$ is HCO. For this, note that

$$\begin{aligned} [\Sigma(\mathbf{p}, \mathbf{x})D^{HCO}(\mathbf{p}, \mathbf{x})]_i = & \\ & [\sigma_\epsilon^2(\mathbf{p}, \mathbf{x})]_i d^{HCO}(p_i, \mathbf{x}_i; \mathbf{p}, \mathbf{x}) + \left[\beta(\mathbf{p}, \mathbf{x}) \underbrace{\beta'(\mathbf{p}, \mathbf{x}) D^{HCO}(\mathbf{p}, \mathbf{x})}_{\text{scalar}} \right]_i \end{aligned}$$

which, given that $\beta(\mathbf{p}, \mathbf{x})$ is HCO gives the proof. ■

This result shows that for any HCO demand, and a fortiori for logit demand, there exist factor models that microfound it. However, the reverse is clearly not true: an arbitrary factor model does not give rise to logit demand. To move closer to common finance intuition, we consider a restricted class of models often used to think about portfolio choice with stable variance and expected payoffs.

Definition 18 (Stable factor model) *A stable M-factor model (based on observables) is a factor model where expected payoff idiosyncratic risk and factor loadings depend on observables only: $\mu(\mathbf{p}, \mathbf{x}) = M(\mathbf{x}) - \mathbf{p}$, $\sigma_\epsilon^2(\mathbf{p}, \mathbf{x}) = \sigma_\epsilon^2(\mathbf{x})$, and $\beta(\mathbf{p}, \mathbf{x}) = \beta(\mathbf{x})$ where $M(\mathbf{x})$ and $\sigma_\epsilon^2(\mathbf{x})$ and the columns of $\beta(\mathbf{x})$ are all HCO.*

One might wonder if there are stable factor models that the demand of a log investor can be represented with the logit form. This cannot hold overall due to different functional forms: the stable factor model is linear in log prices. To make the comparison more meaningful we ask if such an equivalence holds locally. To do so, we define first-order equivalence.

Definition 19 *Two demand functions $D^1(\mathbf{p}, \mathbf{x})$ and $D^2(\mathbf{p}, \mathbf{x})$ are first-order equivalent around a point $(\mathbf{p}_0, \mathbf{x}_0)$ if they have same value and Jacobian with respect to the price matrix at that point:*

$$\begin{aligned} D^1(\mathbf{p}_0, \mathbf{x}_0) &= D^2(\mathbf{p}_0, \mathbf{x}_0) \\ \underbrace{\frac{\partial D^1}{\partial \mathbf{p}}(\mathbf{p}_0, \mathbf{x}_0)}_{N \times N} &= \frac{\partial D^2}{\partial \mathbf{p}}(\mathbf{p}_0, \mathbf{x}_0) \end{aligned}$$

We obtain that in general the answer is negative.

Proposition 20 *Out of the set of all stable M -factor models, there is only one under which a logit demand model with $\alpha > 0$ can be first-order equivalent to the demand of the log investor. This specific model has 1 factor with identical factor loadings and idiosyncratic variance inversely proportional to demand. Under any other stable factor model, logit demand is not a valid approximation of the demand of log investor.*

Proof. For a stable factor model, we have $D^{\log}(\mathbf{p}, \mathbf{x}) = \Sigma(\mathbf{x})^{-1} (M(\mathbf{x}) - p)$, so $\partial D^{\log} / \partial p = -\Sigma(\mathbf{x})^{-1}$. For logit we have: $\partial D^{\logit} / \partial p = -\alpha \text{diag}(\omega) (I - \mathbf{1}\omega') = -\alpha \text{diag}(\omega) + \alpha \omega \omega'$, where ω is the investor's portfolio share vector. We can invert it with the Sherman-Morrison formula and identify with Σ :

$$\begin{aligned} -(\partial D^{\logit} / \partial p)^{-1} &= \alpha^{-1} (I - \mathbf{1}\omega')^{-1} \text{diag}(\omega)^{-1} \\ &= \alpha^{-1} \left(I + \frac{1}{1 - \omega' \mathbf{1}} \mathbf{1}\omega' \right) \text{diag}(\omega)^{-1} \\ \Sigma &= \underbrace{\alpha^{-1} \text{diag}(\omega)^{-1}}_{\text{idiosyncratic risk}} + \underbrace{\alpha^{-1} \frac{1}{1 - \omega' \mathbf{1}} \mathbf{1}\mathbf{1}'}_{\text{single factor}} \end{aligned}$$

Comparing this to the covariance matrix under a generic stable factor model,

$$\Sigma(\mathbf{p}, \mathbf{x}) = \underbrace{\text{diag}(\sigma_\epsilon^2(\mathbf{x}))}_{N \times 1} + \underbrace{\beta(\mathbf{x})}_{N \times M} \underbrace{\Sigma_F}_{M \times M} \beta(\mathbf{x})',$$

concludes the statement. ■

The proof also illustrates immediately that logit can never be the approximation of a stable multi-factor model that cannot be reduced to a single factor. In such a model, the substitution matrix is of rank equal to the number of factors. Intuitively investors substitute along portfolios corresponding to the various risk factors, differently for assets with different loading on those factors. More broadly, HCO demands include models where, following a price increase for a given position, the investor would substitute disproportionately with assets with similar observables.

A two-asset example. We illustrate that this limitation arises even in the simplest possible 2×2 example with the same variance. Fix the vectors \mathbf{p} and $d(p_i; \mathbf{p}) = \omega_i$ that we are looking for first-order equivalence around. Such position could come from a factor model with covariance matrix for any correlation ρ and variance σ^2 :

$$\Sigma = \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$$

If we have a valid approximating logit model, it would feature:

$$-(\partial D^{\text{logit}} / \partial \mathbf{p})^{-1} = \alpha^{-1} \begin{pmatrix} \omega_1^{-1} + \frac{1}{1-\omega_1-\omega_2} & \frac{1}{1-\omega_1-\omega_2} \\ \frac{1}{1-\omega_1-\omega_2} & \omega_2^{-1} + \frac{1}{1-\omega_1-\omega_2} \end{pmatrix}$$

Clearly the two matrix can never be identical if $\omega_1 \neq \omega_2$ because the diagonal terms must be equal. Even if we assume that our point of approximation has a given value $\omega_1 = \omega_2 = \bar{\omega}$, the models are identical only if it matches both on-diagonal and off-diagonal elements:

$$\begin{aligned} \sigma^2 &= \alpha^{-1} \left(\bar{\omega}^{-1} + \frac{1}{1-2\bar{\omega}} \right) \\ \sigma^2 \rho &= \alpha^{-1} \frac{1}{1-2\bar{\omega}} \end{aligned}$$

We can already see the issue: there is only one degree of freedom in logit (α) but two degrees of freedom for the covariance matrix (σ^2 and ρ). Let us construct the corresponding contradiction. Subtracting the second equation from the first one gives:

$$\sigma^2 (1 - \rho) \bar{\omega} = \alpha^{-1}$$

Plugging this expression for α^{-1} in the second equation leads to:

$$\begin{aligned} \sigma^2 \rho &= \sigma^2 (1 - \rho) \frac{\bar{\omega}}{1 - 2\bar{\omega}} \\ \frac{\rho}{1 - \rho} &= \frac{\bar{\omega}}{1 - 2\bar{\omega}} \end{aligned}$$

The right-hand-side is fixed, this is our point of approximation. The left-hand-side could take any value as ρ is a free parameter.

G Learning from Prices, Strategic Trading, Dynamic Trading

In the main text, we showed that the linear approximation of demand function

$$\mathbf{D}_t = \mathcal{E} \mathbf{P}_t + \boldsymbol{\epsilon}_t \tag{205}$$

where $\mathbf{D}_t \in \mathbb{R}^n$ is the vector of (net) demands for n assets, $\mathbf{P}_t \in \mathbb{R}^n$ is the price vector, $\mathcal{E} \in \mathbb{R}^{n \times n}$ is a constant matrix, and $\boldsymbol{\epsilon} \in \mathbb{R}^n$ is a price-independent intercept, is consistent with investors trading for risk-sharing purposes, potentially in the presence of portfolio constraints and a factor structure of returns. Here we present two examples building on standard models to demonstrate that this form is also consistent with: (i) asymmetric information models where investors learn from prices, and (ii) models based on strategic trading with imperfect competition. However, we also emphasize some caveats. We explain that we implicitly require a sufficiently stable economic environment that (205) remains meaningful to estimate. Also, in the last part of this appendix we discuss that in dynamic environments demand functions might be considerably more complex to be represented by a mapping like (205).

G.1 Admati (1985): competitive traders learning from signals and prices with multiple assets

This section presents a multi-asset noisy rational expectations model following Admati (1985).

Environment. There are n risky assets with payoff vector $\mathbf{m} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_m)$ and price vector $\mathbf{p} \in \mathbb{R}^n$. A continuum of competitive, risk-averse traders indexed by $i \in [0, 1]$ have CARA utility with absolute risk aversion $\gamma > 0$.

Each trader observes a private signal

$$\mathbf{s}_i = \mathbf{m} + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon),$$

i.i.d. across i and independent of \mathbf{m} . Traders also observe the equilibrium price \mathbf{p} and use it to update beliefs. Noise traders supply an exogenous quantity $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_u)$.

Individual demand. Trader i is a price-taker and chooses position $\mathbf{d}_i \in \mathbb{R}^n$ to maximize CARA expected utility. Standard CARA-normal optimization yields:

$$\mathbf{d}_i = \frac{1}{\gamma} \boldsymbol{\Sigma}_{m|s,p}^{-1} (\mathbb{E}[\mathbf{m} | \mathbf{s}_i, \mathbf{p}] - \mathbf{p}),$$

where $\boldsymbol{\Sigma}_{m|s,p} = \text{Var}(\mathbf{m} | \mathbf{s}_i, \mathbf{p})$ is the posterior covariance.

Bayesian updating. In equilibrium, the price is informationally equivalent to a signal $\boldsymbol{\eta} = \mathbf{m} + \boldsymbol{\xi}$ for some noise $\boldsymbol{\xi}$ with covariance $\boldsymbol{\Sigma}_\eta$. The posterior precision combines prior and signal precisions:

$$\boldsymbol{\Sigma}_{m|s,p}^{-1} = \boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\epsilon^{-1} + \boldsymbol{\Sigma}_\eta^{-1}.$$

The posterior mean is a precision-weighted average:

$$\mathbb{E}[\mathbf{m} | \mathbf{s}_i, \mathbf{p}] = \boldsymbol{\Sigma}_{m|s,p} (\boldsymbol{\Sigma}_\epsilon^{-1} \mathbf{s}_i + \boldsymbol{\Sigma}_\eta^{-1} \boldsymbol{\eta}).$$

Aggregate demand. Integrating over traders and using $\int_0^1 \mathbf{s}_i di = \mathbf{m}$ by the law of large numbers:

$$\mathbf{D} = \int_0^1 \mathbf{d}_i di = \frac{1}{\gamma} \boldsymbol{\Sigma}_\varepsilon^{-1}(\mathbf{m} - \mathbf{p}) + \frac{1}{\gamma} \boldsymbol{\Sigma}_\eta^{-1}(\boldsymbol{\eta} - \mathbf{p}) + \frac{1}{\gamma} \boldsymbol{\Sigma}_m^{-1}(-\mathbf{p}).$$

Equilibrium. Market clearing $\mathbf{D} = \mathbf{u}$ determines the equilibrium. In the unique linear equilibrium, the price signal precision satisfies:

$$\boldsymbol{\Sigma}_\eta^{-1} = \frac{1}{\gamma^2} \boldsymbol{\Sigma}_\varepsilon^{-1} \boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_\varepsilon^{-1}.$$

The precision of \mathbf{m} given prices alone is:

$$\boldsymbol{\Sigma}_{m|p}^{-1} = \boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\eta^{-1} = \boldsymbol{\Sigma}_m^{-1} + \frac{1}{\gamma^2} \boldsymbol{\Sigma}_\varepsilon^{-1} \boldsymbol{\Sigma}_u^{-1} \boldsymbol{\Sigma}_\varepsilon^{-1}.$$

Prices become fully revealing as $\gamma \rightarrow 0$ or $\boldsymbol{\Sigma}_u \rightarrow \mathbf{0}$.

Identifying \mathcal{E} and $\boldsymbol{\epsilon}$. Writing aggregate demand as $\mathbf{D} = \mathcal{E}\mathbf{p} + \boldsymbol{\epsilon}$:

$$\mathcal{E} = -\frac{1}{\gamma} (\boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\varepsilon^{-1} + \boldsymbol{\Sigma}_\eta^{-1}) = -\frac{1}{\gamma} \boldsymbol{\Sigma}_{m|s,p}^{-1},$$

$$\boldsymbol{\epsilon} = \frac{1}{\gamma} (\boldsymbol{\Sigma}_\varepsilon^{-1} \mathbf{m} + \boldsymbol{\Sigma}_\eta^{-1} \boldsymbol{\eta}).$$

Cross-asset demand elasticities arise from learning: even if private signals are independent across assets ($\boldsymbol{\Sigma}_\varepsilon$ diagonal), correlated noise trading ($\boldsymbol{\Sigma}_u$ non-diagonal) creates cross-asset information linkages through $\boldsymbol{\Sigma}_\eta^{-1}$.

Identification of \mathcal{E} . To identify the demand elasticity matrix \mathcal{E} empirically, one needs instruments that shift prices but are uncorrelated with the intercept $\boldsymbol{\epsilon}$. In this model, $\boldsymbol{\epsilon}$ depends on the fundamental \mathbf{m} and the price signal $\boldsymbol{\eta}$ —both of which carry information about payoffs. A valid instrument must therefore be a *non-informational* shock to prices. The natural candidate is **noise trading \mathbf{u}** : shocks to the supply from noise traders move equilibrium prices but are, by assumption, orthogonal to fundamentals \mathbf{m} and hence to $\boldsymbol{\epsilon}$. In practice, this could correspond to index rebalancing flows, liquidity-driven sales (e.g., mutual fund redemptions), or other demand shocks uncorrelated with asset payoffs.

Importantly, while a shock to noise trading is a valid instrument, a shock that investors *learn* to originate from non-fundamental sources might *not* be valid in this model. For example, if the Federal Reserve unexpectedly announces a random purchase of certain assets, investors would know that the resulting price movement is uninformative about \mathbf{m} . However, such knowledge constitutes additional information not present in the model setup: here, investors assign probabilities to the sources of price changes based on their priors ($\boldsymbol{\Sigma}_m$, $\boldsymbol{\Sigma}_\varepsilon$, $\boldsymbol{\Sigma}_u$) and cannot perfectly disentangle fundamental from non-fundamental shocks. If they could identify the shock's source with certainty, their inference problem—and hence the equilibrium demand intercept $\boldsymbol{\epsilon}$ —would change.

This points to a crucial, even if implicit, assumption behind our analysis. We assume that *the environment is sufficiently stable* that the $\mathbf{D}_t = \mathcal{E} \mathbf{P}_t + \boldsymbol{\epsilon}_t$ relationship continue to hold at least approximately, at least for the relevant period. In this particular framework, this implies that the parameters of the model are stable.

G.2 Multi-asset Kyle (1989): imperfect competition with common signal

This section presents a multi-asset version of Kyle (1989)'s model of imperfect competition with demand schedules. The setup and solution method follow the unified framework developed by Rostek and Yoon (2025), who provide a comprehensive treatment of imperfectly competitive financial markets encompassing static, dynamic, centralized, and decentralized settings. Our common-signal specification is a special case of their general framework with symmetric traders and no inference from prices.

Environment. There are n risky assets with payoff vector $\mathbf{m} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_m)$ and prices \mathbf{p} . There are $N \geq 3$ risk-averse speculators with CARA utility (risk aversion $\gamma > 0$), all observing the *same* signal:

$$\mathbf{s} = \mathbf{m} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon).$$

Noise traders supply $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_u)$. Traders submit demand schedules (not quantities), and the market clears without a market maker. Because signals are common, prices reveal no information beyond \mathbf{s} .

Posterior beliefs. All traders share the same posterior:

$$\mathbb{E}[\mathbf{m} \mid \mathbf{s}] = \boldsymbol{\Sigma}_m (\boldsymbol{\Sigma}_m + \boldsymbol{\Sigma}_\epsilon)^{-1} \mathbf{s}, \quad \boldsymbol{\Sigma}_{m|\mathbf{s}} = (\boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\epsilon^{-1})^{-1}.$$

Equilibrium characterization. Following Rostek and Yoon (2025), equilibrium is characterized by two conditions: (i) each trader optimizes given their price impact, and (ii) price impacts are mutually consistent (a fixed point condition). In the symmetric equilibrium, each speculator k submits a demand schedule:

$$\mathbf{d}_k(\mathbf{p}, \mathbf{s}) = \mathbf{B}\mathbf{s} - \boldsymbol{\Gamma}\mathbf{p},$$

where $\boldsymbol{\Gamma}$ is the price sensitivity matrix. The individual price impact $\boldsymbol{\Lambda}_I$ (the slope of each trader's inverse residual supply) satisfies the fixed point condition:

$$\boldsymbol{\Lambda}_I = [(N - 1)\boldsymbol{\Gamma}]^{-1}.$$

Optimization. Each speculator k maximizes expected utility taking price impact into account. The first-order condition equates marginal utility to marginal payment:

$$\mathbb{E}[\mathbf{m} \mid \mathbf{s}] - \gamma \boldsymbol{\Sigma}_{m|\mathbf{s}} \mathbf{d}_k = \mathbf{p} + \boldsymbol{\Lambda}_I \mathbf{d}_k,$$

which gives the best response:

$$\mathbf{d}_k = (\gamma \boldsymbol{\Sigma}_{m|s} + \boldsymbol{\Lambda}_I)^{-1} (\mathbb{E}[\mathbf{m} | \mathbf{s}] - \mathbf{p}).$$

Hence $\boldsymbol{\Gamma} = (\gamma \boldsymbol{\Sigma}_{m|s} + \boldsymbol{\Lambda}_I)^{-1}$.

Solving the fixed point. Substituting $\boldsymbol{\Gamma}$ into the fixed point condition:

$$\boldsymbol{\Lambda}_I = \frac{\gamma \boldsymbol{\Sigma}_{m|s} + \boldsymbol{\Lambda}_I}{N-1} \Rightarrow (N-2)\boldsymbol{\Lambda}_I = \gamma \boldsymbol{\Sigma}_{m|s} \Rightarrow \boxed{\boldsymbol{\Lambda}_I = \frac{\gamma}{N-2} \boldsymbol{\Sigma}_{m|s}}$$

Equilibrium coefficients. Substituting back:

$$\boxed{\boldsymbol{\Gamma} = \frac{N-2}{\gamma(N-1)} \boldsymbol{\Sigma}_{m|s}^{-1} = \frac{N-2}{\gamma(N-1)} (\boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\varepsilon^{-1})}$$

For the signal sensitivity \mathbf{B} , matching coefficients on \mathbf{s} in the best response (using $\mathbb{E}[\mathbf{m} | \mathbf{s}] = \boldsymbol{\Sigma}_m (\boldsymbol{\Sigma}_m + \boldsymbol{\Sigma}_\varepsilon)^{-1} \mathbf{s}$):

$$\boxed{\mathbf{B} = \boldsymbol{\Gamma} \boldsymbol{\Sigma}_m (\boldsymbol{\Sigma}_m + \boldsymbol{\Sigma}_\varepsilon)^{-1} = \frac{N-2}{\gamma(N-1)} \boldsymbol{\Sigma}_\varepsilon^{-1}}$$

where the last equality uses the matrix identity $\boldsymbol{\Sigma}_{m|s}^{-1} \boldsymbol{\Sigma}_m (\boldsymbol{\Sigma}_m + \boldsymbol{\Sigma}_\varepsilon)^{-1} = \boldsymbol{\Sigma}_\varepsilon^{-1}$.

Market clearing and aggregate price impact. From market clearing $\sum_{k=1}^N \mathbf{d}_k = \mathbf{u}$, the equilibrium price satisfies:

$$\mathbf{p} = \boldsymbol{\Lambda} [N\mathbf{B}\mathbf{s} - \mathbf{u}], \quad \boldsymbol{\Lambda} = (N\boldsymbol{\Gamma})^{-1} = \frac{\gamma(N-1)}{N(N-2)} \boldsymbol{\Sigma}_{m|s}.$$

The relationship between individual and aggregate price impact is $\boldsymbol{\Lambda}_I = \frac{N}{N-1} \boldsymbol{\Lambda}$.

Identifying \mathcal{E} and $\boldsymbol{\epsilon}$. Aggregate demand is $\mathbf{D} = N\mathbf{d}_k = N\mathbf{B}\mathbf{s} - N\boldsymbol{\Gamma}\mathbf{p}$. Thus:

$$\begin{aligned} \mathcal{E} &= -N\boldsymbol{\Gamma} = -\frac{N(N-2)}{\gamma(N-1)} (\boldsymbol{\Sigma}_m^{-1} + \boldsymbol{\Sigma}_\varepsilon^{-1}), \\ \boldsymbol{\epsilon} &= N\mathbf{B}\mathbf{s} = \frac{N(N-2)}{\gamma(N-1)} \boldsymbol{\Sigma}_\varepsilon^{-1} \mathbf{s}. \end{aligned}$$

Unlike the competitive Hellwig/Admati model, cross-asset elasticities here arise purely from the covariance structure of payoffs and signals—not from learning, since all traders share the same information.

Identification of \mathcal{E} . To identify \mathcal{E} in this setting, one again needs instruments that shift prices but are uncorrelated with $\boldsymbol{\epsilon}$. Here, $\boldsymbol{\epsilon}$ depends on the common signal $\mathbf{s} = \mathbf{m} + \boldsymbol{\varepsilon}$, which contains information about fundamentals. As in the Hellwig/Admati model, **noise**

trading \mathbf{u} is a valid instrument: it shifts equilibrium prices (through market clearing) but is orthogonal to \mathbf{s} and hence to $\boldsymbol{\epsilon}$. An interesting feature of this model is that $\boldsymbol{\Sigma}_u$ does not appear in the equilibrium demand coefficients $(\mathbf{B}, \boldsymbol{\Gamma})$ —noise trading affects price *levels* but not how traders respond to signals or prices. Nevertheless, cross-sectional or time-series variation in \mathbf{u} (e.g., ETF creation/redemption activity, futures rolls, or settlement-driven flows) provides the exogenous price variation needed to trace out demand curves. Shocks correlated with the signal—such as public news announcements that also enter \mathbf{s} —would be invalid instruments.

G.3 Dynamic environments

In all our motivating models, here and in the main text, the distribution of investors’ relevant pay-offs is exogenously given. In standard dynamic settings as in [Merton \(1973\)](#), this is not the case. The key issue is that in dynamic models, investors’ demand depends not only on current prices but on the entire distribution of future returns—including volatilities, covariances, and correlations with future investment opportunities. Then, the specification where the current price is the main determinant of demand, $\mathbf{D} = \boldsymbol{\mathcal{E}}\mathbf{P} + \boldsymbol{\epsilon}$ might not be appropriate. Simply, demand curves have higher dimensions in a dynamic environment. As [He et al. \(2025\)](#) argues in these environment an instrument based on a supply shock is less trivially sufficient for identification than in static models.

H Estimating elasticity in theoretical models

We take a brief detour through theoretical models. We first show how commonly used models, once considered in appropriate units, relate to our identification assumptions. Then, we provide an example explaining why simple equilibrium considerations do not affect our identification results.

H.1 Standard models of asset demand

We discuss how standard models of asset demand relate with the identification assumptions. For each model, we derive the appropriate units and parameter restrictions under which the demand regression is well specified.⁵⁴ Table 5 summarizes the results.

Constant absolute risk aversion. In the mean-variance model (CARA) described above, we have seen the direct mapping between covariance matrix and the elasticity matrix when considering a relation between the level of demand the level of prices: $\boldsymbol{\mathcal{E}} = \partial\mathbf{D}/\partial\mathbf{P} = \gamma^{-1}\boldsymbol{\Sigma}^{-1}$. Section 2.2.1 show that Assumptions A1 and A2 are satisfied if the covariance matrix has a factor structure with factor loadings which depend on the observables.

⁵⁴[Petajisto \(2009\)](#), [Davis et al. \(2025\)](#), and [Davis \(2024\)](#) quantify elasticities in these models.

	CARA	CRRA	Logit
Regression units “demand” LHS	demand D_i	portfolio shares $P_i D_i / W$	log portfolio shares $\log(P_i D_i / W)$
Regression units “price” RHS	price P_i	log price $\log P_i$	log price $\log P_i$
Own Price Elasticity \mathcal{E}_{ii}	$\frac{R_f}{\gamma} [\Sigma^{-1}]_{ii}$	$\frac{1}{\gamma} [\Sigma^{-1}]_{ii}$	$\alpha(1 - \omega_i)$
Cross price Elasticity \mathcal{E}_{ij}	$\frac{R_f}{\gamma} [\Sigma^{-1}]_{ij}$	$\frac{1}{\gamma} [\Sigma^{-1}]_{ij}$	$-\alpha\omega_j$
Relative Elasticity $\hat{\mathcal{E}} = \mathcal{E}_{ii} - \mathcal{E}_{ji}$	$\frac{R_f}{\gamma} ([\Sigma^{-1}]_{ii} - [\Sigma^{-1}]_{ji})$	$\frac{1}{\gamma} ([\Sigma^{-1}]_{ii} - [\Sigma^{-1}]_{ji})$	α

Table 5: Three standard models of asset demands.

Constant relative risk aversion. Preferences with constant relative risk aversion (CRRA) are the workhorse model of macro-finance. Utility in this case is given by $u(C) = C^{1-\gamma}/(1-\gamma)$, with now γ being the constant relative risk-aversion. Assume that the risk-free rate is r_f and that there are N assets with payoffs $X = \{X_i\}_i$ at time 1, with prices $\{P_i\}$. Hence, asset returns are $R_i = X_i/P_i$.

To solve for the optimal demands, we assume that the payoffs follow a lognormal distribution: $\log X \sim \mathcal{N}(M, \Sigma)$ and log-linearize portfolio returns following [Campbell and Viceira \(2002\)](#).⁵⁵ For an investor with wealth W , the optimal demand is:

$$D_i = \frac{1}{\gamma} \frac{W}{P_i} \left[\Sigma^{-1} \left(M - \log P - r_f + \frac{1}{2} \text{diag}(\Sigma) \right) \right]_i \quad (207)$$

This implies that when considering the relation between portfolio weights, $\omega_i = P_i D_i / W$, and log prices, the elasticity matrix is (this relation is exact in continuous time; see [Duffie, 2010](#); [He et al., 2025](#)):

$$\mathcal{E} = \frac{\partial \omega}{\partial \log P} = -\frac{1}{\gamma} \Sigma^{-1}. \quad (208)$$

This is the same elasticity as the CARA case, albeit with different units: portfolio weights on log prices. Therefore, our earlier discussion relating properties of the covariance matrix (here of log returns) and the identification assumptions apply to this case as well.

Logit. [Kojien and Yogo \(2019\)](#) introduce a model of portfolio demand of the logit form. They show existence of factor models giving rise to this demand for an investor with log utility. Unlike the model we just studied, these factor models feature a covariance matrix and expected returns that depend nonlinearly on prices, and hence have a different elasticity

⁵⁵We log-linearize the return of portfolio ω , $r_p = \log R_p$ as:

$$r_p - r_f = \log(\omega' \exp(\mathbf{r} - r_f)) \simeq \omega'(\mathbf{r} - r_f) + \frac{1}{2} \omega' \text{diag}(\Sigma) - \frac{1}{2} \omega' \Sigma \omega. \quad (206)$$

matrix; Appendix F details this distinction. The logit model is also commonly used in industrial organization, as well as in many fields in economics including trade and spatial equilibrium models. There, it is most often motivated by aggregation of a consumer discrete choice model, but can also apply to an individual’s choice of consumption basket.⁵⁶

For an investor with initial wealth W , the expenditure shares or portfolio weights are:⁵⁷

$$\omega_i = \frac{P_i D_i}{W} = \frac{\exp(-\alpha p_i + \theta' X_i + \epsilon_i)}{1 + \sum_l \exp(-\alpha p_l + \theta' X_l + \epsilon_l)}, \quad (210)$$

where p_i is the log of the price of asset i , X_i observable demand shifters, and ϵ_i the unobserved component of demand.

When considering the relation between log portfolio weights and log prices, the elasticity matrix is:

$$\mathcal{E} = \frac{\partial \log \omega}{\partial \log P} = -\alpha (\mathbf{I} - \mathbf{1}\omega'), \quad (211)$$

where ω is the vector of portfolio weights given in (210). Note that \mathcal{E} in general is not symmetric in this case. The coefficient α is the only demand parameter that determines the matrix of demand elasticity, as opposed to the whole covariance matrix in the CARA and CRRA cases. Further, this matrix always satisfies assumptions A1 ($\mathcal{E}_{jk} = \mathcal{E}_{ik} = \alpha \omega_k$) and A2 ($\mathcal{E}_{ii} - \mathcal{E}_{ji} = \alpha$), with α being the relative elasticity of demand.

Same relative elasticity vs. same elasticity matrix. Note that the fact that the simple risk-based models and the logit model can both satisfy the two assumptions only implies that they lead to the same estimation of the relative elasticity. Even when these models have the same value of relative elasticity, they exhibit different elasticity matrices. Figure 6 illustrates this nuance, with three distinct elasticity matrices that share the same relative elasticity.

H.2 What about equilibrium spillovers?

The reader might be surprised that, so far, we have not discussed the concept of equilibrium, which is usually central in asset pricing. This is not because we assume that the world is not in equilibrium: equation (??) is a change in demand in equilibrium. Instead, we can do so because identifying specific sources of variation in prices — the instrument Z_i — and assuming that an investor’s demand satisfies Assumptions A1 and A2 is enough to estimate this investor’s demand elasticity without understanding the entire structure of the equilibrium.

In this section, we work out a simple equilibrium model to illustrate this insight. The setting is inspired by [Fuchs et al. \(2025\)](#) who point out that endogenous cross-asset spillovers

⁵⁶[Anderson et al. \(1988\)](#) derives the utility that leads to logit shares as optimal demand.

⁵⁷If there is not outside asset, the model of expenditure shares becomes:

$$\omega_i = \frac{P_i D_i}{W} = \frac{\exp(-\alpha p_i + \theta' X_i + \epsilon_i)}{\sum_l \exp(-\alpha p_l + \theta' X_l + \epsilon_l)}. \quad (209)$$

$$\begin{bmatrix} \hat{\mathcal{E}} & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots \\ \vdots & \dots & \dots & \dots \\ 0 & \dots & 0 & \hat{\mathcal{E}} \end{bmatrix} \qquad \begin{bmatrix} \frac{1}{1-\rho}\hat{\mathcal{E}} & \frac{\rho}{1-\rho}\hat{\mathcal{E}} & \dots & \frac{\rho}{1-\rho}\hat{\mathcal{E}} \\ \frac{\rho}{1-\rho}\hat{\mathcal{E}} & \dots & \dots & \dots \\ \vdots & \dots & \dots & \dots \\ \frac{\rho}{1-\rho}\hat{\mathcal{E}} & \dots & \frac{\rho}{1-\rho}\hat{\mathcal{E}} & \frac{1}{1-\rho}\hat{\mathcal{E}} \end{bmatrix}$$

(a) Diagonal matrix.

(b) Symmetric matrix.

$$\begin{bmatrix} (1-\omega_1)\hat{\mathcal{E}} & -\omega_2\hat{\mathcal{E}} & \dots & -\omega_N\hat{\mathcal{E}} \\ -\omega_1\hat{\mathcal{E}} & (1-\omega_2)\hat{\mathcal{E}} & \dots & \dots \\ \vdots & \dots & \dots & \dots \\ -\omega_1\hat{\mathcal{E}} & \dots & -\omega_{N-1}\hat{\mathcal{E}} & (1-\omega_N)\hat{\mathcal{E}} \end{bmatrix}$$

(c) Logit-style matrix.

Figure 6: Different elasticity matrices with the same relative elasticity $\hat{\mathcal{E}}$

can imply a low measured own-price elasticity even if the true own price elasticity is near infinite. This result considers a different regression from the causal inference framework of this paper. We show that the issue arises because the omitted variable bias that we have pointed out in Section ?? is present. We also explain that, because the example satisfies our assumptions A1 and A2, a standard difference-in-difference regression is unbiased, once we recognize that it recovers the relative elasticity.

Setting. The economy is populated by a representative agent with log utility. There are three assets with different payoffs in three possible states of the world, with payoffs as follows:

$$\begin{array}{ccc} \begin{array}{l} \text{Green } P_g \end{array} \begin{array}{l} \nearrow 1 + \epsilon \\ \rightarrow 1 - \epsilon \\ \searrow 0 \end{array} & \begin{array}{l} \text{Red } P_r \end{array} \begin{array}{l} \nearrow 1 - \epsilon \\ \rightarrow 1 + \epsilon \\ \searrow 0 \end{array} & \begin{array}{l} \text{Other } P_o = 1 \end{array} \begin{array}{l} \nearrow 0 \\ \rightarrow 0 \\ \searrow 1 \end{array} \quad \begin{array}{l} \text{w.p. } 1/4 \\ \text{w.p. } 1/4 \\ \text{w.p. } 1/2 \end{array} \end{array}$$

The “other” asset acts as a numéraire, whose price is normalized to 1. Denote the prices of the green and red assets P_g and P_r . These two assets become closer substitutes as ϵ goes towards 0. Indeed, in the limit, any price difference between P_g and P_r represents an arbitrage opportunity. The representative agent has endowments E_g , E_r , and E_o , which implies that their wealth is $W = P_g E_g + P_r E_r + E_o$.

Demand and equilibrium. We can first derive the demand function, that is, the optimal portfolio share as a function of prices:

$$\omega_g(P_g, P_r) = \frac{P_g}{2} \frac{\epsilon^2(P_r + P_g) + (P_r - P_g)}{\epsilon^2(P_r + P_g)^2 - (P_r - P_g)^2}, \quad \omega_r(P_g, P_r) = \frac{P_r}{2} \frac{\epsilon^2(P_r + P_g) + (P_g - P_r)}{\epsilon^2(P_r + P_g)^2 - (P_r - P_g)^2}. \quad (212)$$

Market-clearing for the two assets, $\omega_g W = P_g E_g$ and $\omega_r W = P_r E_r$ lead to equilibrium prices as functions of the endowments:

$$P_g(E_o, E_g, E_r) = E_o \frac{\epsilon^2(E_g - E_r) - (E_g + E_r)}{\epsilon^2(E_g - E_r)^2 - (E_g + E_r)^2}, \quad P_r(E_o, E_g, E_r) = E_o \frac{\epsilon^2(E_r - E_g) - (E_g + E_r)}{\epsilon^2(E_g - E_r)^2 - (E_g + E_r)^2}. \quad (213)$$

As an initial equilibrium, we assume that the endowments are $E_g = E_r = 1/2$ and $E_o = 1$. It is then immediate that $P_r = P_g = 1$.

Demand elasticities. We can compute the demand elasticities: how individual demand would respond to a change in prices. Because utility is CRRA, we measure the sensitivity of portfolio shares to log prices (in line with Section H.1) around the initial equilibrium values of prices:

$$\mathcal{E}_{own} = \frac{\partial \omega_g}{\partial \log P_g} = \frac{1}{8} - \frac{1}{8\epsilon^2}; \quad (214)$$

$$\mathcal{E}_{cross} = \frac{\partial \omega_g}{\partial \log P_r} = -\frac{1}{8} + \frac{1}{8\epsilon^2}. \quad (215)$$

The expressions for ω_r are identical. The relative elasticity is $\mathcal{E}_{own} - \mathcal{E}_{cross} = 1/4 - 1/(4\epsilon^2)$.

These measures show that when the two assets are near-identical, $\epsilon \rightarrow 0$, any deviation from parity would lead to a near-infinite increase in demand for the cheaper asset, and near-infinite decrease in demand for the more expensive one. This is the standard arbitrage argument.

Running regressions. We are interested in whether various regressions around a supply shock for one of the assets can identify these elasticities. A shift in supply of the green asset leads to the price changes

$$\frac{\partial \log P_g}{\partial E_g} = -(1 + \epsilon^2), \quad \frac{\partial \log P_r}{\partial E_g} = -(1 - \epsilon^2), \quad (216)$$

around the equilibrium.

[Fuchs et al. \(2025\)](#) correctly point out that regressing the demand for the green asset on the change in its price using such a supply shock does not recover the own price elasticity. This regression corresponds to taking the ratio of the change in portfolio to the change in

price across equilibria:

$$\frac{d\omega_g/dE_g}{d\log P_g/dE_g} = -\frac{1}{4} \frac{1 - \varepsilon^2}{1 + \varepsilon^2} \neq \mathcal{E}_{own}. \quad (217)$$

In particular, when $\varepsilon \rightarrow 0$, the regression coefficient on the left-hand-side converges to $-1/4$, while the own-price elasticity goes to infinity. Unpacking the total derivative explains the source of the bias:

$$\frac{d\omega_g/dE_g}{d\log P_g/dE_g} = \frac{\frac{\partial\omega_g}{\partial\log P_g} \frac{\partial\log P_g}{\partial E_g} + \frac{\partial\omega_g}{\partial\log P_r} \frac{\partial\log P_r}{\partial E_g}}{\frac{d\log P_g}{dE_g}} = \mathcal{E}_{own} + \mathcal{E}_{cross} \frac{\partial\log P_r/\partial E_g}{\partial\log P_g/\partial E_g} \quad (218)$$

The change in demand for the green asset is not driven only by the change in its own price but also by the change in price of its substitute the red asset, because $\mathcal{E}_{cross} \neq 0$. In the language of regressions, the price of the red asset is acting as a correlated omitted variable. Intuitively, the induced price drop of the green asset would lead to a large increase of its demand if the price of the red asset remained high. However, in equilibrium the price of the red asset drops too, resulting in only a moderate change in the demand for the green asset.

However, the canonical causal inference framework corresponds to a standard difference-in-difference regression in this setting with two assets. The regression coefficient is the ratio of the difference of change in portfolio weight to the difference of change in price, as in equation (??):

$$\frac{d\omega_g/dE_g - d\omega_r/dE_g}{d\log P_g/dE_g - d\log P_r/dE_g} = \frac{1}{4} - \frac{1}{4\varepsilon^2} = \mathcal{E}_{own} - \mathcal{E}_{cross}. \quad (219)$$

The difference-in-difference coefficient correctly identifies the relative elasticity. Indeed, this setting satisfies Assumptions A1 and A2: because the elasticity matrix is symmetric, substitution is homogeneous and the relative elasticity is constant. Also, because the endowment shock does not have a direct effect on log investors' choice of portfolio shares, the standard exogeneity condition is satisfied. Note that the identified relative elasticity is unbounded when ε converges to 0, in line with the economic intuition that the relative demand for near arbitrage assets should react strongly to a change in their relative price. The estimator leads to this limit because the relative change in portfolio remains finite, while the relative change in price goes to zero in this limit due to arbitrage.

This example highlights an important conceptual point: Demand elasticities are well defined regardless of the source of the change in prices. Precisely, the demand curve maps the price vector to the quantity vector through any source of change in price that are not accompanied by changes in other drivers of demand. As a result, it is not surprising that equilibrium spillovers are not a problem for the identification of demand elasticities per se. Instead, the econometrician has to be careful that prices of other assets might introduce omitted variable bias. Assumptions A1 and A2 ensure that this is not the case for a standard difference-in-difference estimator, which is then an unbiased estimator of relative elasticity.

I Limits of existing demand models under a factor structure

This appendix describes implementation details and parameter choices of the static CRRA demand model under a factor structure introduced in Section 1.3. We conduct three counterfactual experiments varying design and parameter choices. In all cases, existing demand models fail to capture the factor structure; in all cases, our framework matches it.

I.1 Setting

We consider N risky assets indexed by i with log-normal payoffs $\log X \sim \mathcal{N}(M, \Sigma)$, in fixed supply given by a vector S . Let P denote the equilibrium price vector. The covariance matrix follows a one-factor model with asset-specific loadings β_i , as in equation (4):

$$\Sigma = \sigma_v^2 I + \sigma_F^2 \beta \beta', \quad (220)$$

$$\log X_i = M_i + \beta_i F + v_i, \quad (221)$$

where $F \sim \mathcal{N}(0, \sigma_F^2)$, $v_i \sim \mathcal{N}(0, \sigma_v^2)$, and all terms are independent. The expected log payoff M_i maps to the notation in Section 1.3 via $\alpha_i = M_i - \log P_i - \beta_i \mathbb{E}[F]$, the component of expected log return not explained by factor exposure.

A risk-free asset is in perfectly elastic supply at log rate $r_f = 0$.

There are two investors:

- A **CRRA investor** with risk aversion γ , endowed with E units of the risky assets and E_{r_f} units of the risk-free asset.
- An **investor with deep pockets** (in risk-free assets) who purchases A units of risky assets, parameterizing the supply variation in the experiment of Section 1.3. The residual supply for the CRRA investor is $S = E - A$.

Everyone is a price taker, and equilibrium clears the market for risky assets.

CRRA demand. Let W denote the wealth of the CRRA investor, D their vector of risky asset positions, and $\omega_i = D_i P_i / W$ the portfolio weight in asset i . The investor maximizes expected CRRA utility over terminal wealth with log-normally distributed payoffs. Following the log-linearization of [Campbell and Viceira \(2002\)](#), optimal static demand is:

$$\omega^{\text{CRRA}} = \frac{1}{\gamma} \Sigma^{-1} \left(M - \log P - r_f + \frac{1}{2} \text{diag}(\Sigma) \right), \quad (222)$$

where the $+\frac{1}{2} \text{diag}(\Sigma)$ term is the Jensen's inequality correction for log-normal returns. Unlike the dynamic [Campbell and Viceira \(2002\)](#) framework, this is a one-period problem and contains no intertemporal hedging demand. Initial wealth is $W = E_{r_f} + E'P$.

Equilibrium. Equilibrium requires optimal demand to equal supply: $D = S$, i.e. $\omega_i = S_i P_i / W$. Substituting into (222) and rearranging:

$$\log P = M - r_f + \frac{1}{2} \text{diag}(\Sigma) - \gamma \Sigma \omega, \quad \omega = \frac{S \circ P}{W}, \quad W = E_{r_f} + E'P. \quad (223)$$

This system defines a fixed-point map $g(\omega) = S \circ P(\omega) / W(\omega)$, where $P(\omega) = \exp(M - r_f + \frac{1}{2} \text{diag}(\Sigma) - \gamma \Sigma \omega)$ and $W(\omega) = E_{r_f} + E'P(\omega)$. We numerically solve for fixed points $\omega = g(\omega)$.

Parameter choices. All parameter choices are summarized in Table 6. We choose $N = 100$ to create a rich cross-section while maintaining parsimony, and explore robustness versions at $N = 50, 200$. $E = 1$ and $E_{r_f} = 100$ are normalizations for the numbers of shares in each of the 100 risky assets and one risk-free asset. The risk-free rate is normalized to 0.

The common expected log payoff M_i does not directly affect expected returns as log prices (modulo wealth effects) move one-to-one with it. It mainly controls the asset allocation between the risk-free asset and risk assets. We choose $M_i = 1$ in the baseline, so that the CRRA investor's portfolio share in the risk-free asset (≈ 0.27) closely matches the ratio of around $25/(25+60) \approx 0.29$ of the supply of U.S. treasuries (consolidated with Fed holdings) to U.S. treasuries and total U.S. stock market capitalization as of early 2026. For robustness, we produce results under two alternative parameter choices, $M = 0.1$ and $M = 2$, corresponding to risk-free asset portfolio shares of around 0.46 and 0.12, respectively.

We choose β_i to be equally spaced on $[0.2, 1.8]$. Equal spacing is for parsimony, and to show most clearly how results for experiments differ across β_i . The distribution of betas is centered around 1 in line with the factor corresponding to a market factor, such that the average asset naturally has a beta of 1. The endpoints of the uniform spacing are chosen such that the cross-sectional standard deviation across betas matches the time-series average of the cross-sectional standard deviation of betas, based on the beta estimates from Welch (2022) between 1980–2021. The volatility of the market factor, $\sigma_F = 0.15$, is chosen to match the time-series volatility of the market excess return between 1980–2021 (≈ 0.1551), with data taken from the Kenneth French data library.⁵⁸ We choose $\sigma_v = 0.4$ for idiosyncratic volatility to match the volatility of CRSP excess returns in 1980–2021 (excluding bottom quintile micro cap stocks). Specifically, average equal-weighted total volatility is 0.44, which translates to an idiosyncratic volatility of 0.41 for the average stock at $\beta_i = 1$.

In our baseline parametrization, we choose a constant relative risk aversion of $\gamma = 5$, and produce results under alternative parameters of $\gamma = 2$ and $\gamma = 10$, toward the middle, lower, and upper end of commonly accepted parameter choices in asset pricing.

⁵⁸Since $r_f = 0$ in the model, the market volatility corresponds directly to the volatility of excess market returns. For comparison, the volatility of total market returns over the same period is 0.1546.

Table 6: **Parameter choices for the CRRA demand simulation**

Symbol	Name	Baseline	Alternative choices
N	Number of assets	100	50, 200
E_i	Risky asset endowment	1	—
E_{rf}	Risk-free endowment	100	—
r_f	Risk-free rate	0	—
M_i	Expected log payoff	1	0.1, 2
β_i	Market betas	[0.2, 1.8], equally spaced	—
σ_F	Market volatility	0.15	—
σ_v	Idiosyncratic volatility	0.4	—
γ	Risk aversion	5	2, 10

Table 6 summarizes the parameter choices for the CRRA demand simulation described in Appendix I. The return process is calibrated to U.S. stock market data between 1980 and 2021. Details are provided in the text.

I.2 Model estimation

We next describe how we generate data that we can estimate models from, precisely what models we estimate, and how we conduct the counterfactuals.

Step 1: Simulate data from the CRRA model. Fix baseline supply $S_0 = E$ and solve the CRRA equilibrium (223) to obtain baseline P_0, ω_0, W_0 . Then simulate $T = 50$ periods of supply perturbations under $\sigma_S = 0.05$:⁵⁹

$$S_t = S_0 \circ \exp(\sigma_S \cdot \eta_t - \frac{1}{2}\sigma_S^2), \quad \eta_t \sim \mathcal{N}(0, I), \quad (224)$$

and solve the equilibrium at each S_t to obtain repeated cross-sections $\{\omega_t, \log P_t\}_{t=1}^T$.

For all experiments, we include a robustness version that treats wealth as fixed (e.g., [Kojien and Yogo, 2019](#)) at the initial level of W_0 . Specifically, the market-clearing condition becomes $\omega_i = S_i P_i / W_0$ with W_0 held constant across all supply perturbations—both during simulation and when computing counterfactual prices.

Step 2a: Estimate logit model. The logit model defines portfolio weights as:

$$\omega_i^{\text{Logit}} = \frac{\exp(\alpha \log P_i + \varepsilon_i)}{1 + \sum_j \exp(\alpha \log P_j + \varepsilon_j)}, \quad (225)$$

with the risk-free share $\omega_0 = 1 / (1 + \sum_j \exp(\alpha \log P_j + \varepsilon_j))$. The parameters of the model are α (elasticity) and ε_i (asset-specific tastes), potentially related to β_i .

⁵⁹These parameters are not included in Table 6 as they are not primitives of the model but rather control the precision of estimation.

We estimate the logit model by logarithmizing both sides and estimating in levels (Kojien and Yogo, 2019):

$$\log \omega_{i,t} = \alpha \log P_{i,t} + \mu_i + a_t + b_t \beta_i + u_{i,t}, \quad (226)$$

where α is the elasticity parameter estimated, μ_i are asset fixed effects, a_t time fixed effects (absorbing the log outside share), and b_t time-varying slopes on β_i .

One may argue that any failure of logit here may just be the result of misspecification in function form, especially from forcing a log on portfolio shares. We therefore estimate an alternative log-linear specification in differences,⁶⁰ which is identical to the relative elasticity estimation in Step 2c, with $\alpha_{lin} = \hat{\mathcal{E}}$, but without the substitution component of our model. Accordingly, the counterfactual under this specification uses the same linear form as our model, $\omega_0 + \mathcal{E} \Delta \log P$, but with the substitution matrix set to $\hat{\mathcal{E}} I_N$.

Step 2b: Estimate logit + macro model. As suggested in Gabaix and Kojien (2024), we enrich logit with a macro elasticity. Define the total risky portfolio share $\omega_{\text{TOT},t} \equiv \sum_i \omega_{i,t} = 1 - \omega_{0,t}$ and the aggregate risky-asset price change $\bar{P}_{\text{agg},t} = \frac{1}{N} \sum_i P_{i,t}$. We estimate

$$\Delta \log \omega_{\text{TOT},t} = \Theta \Delta \log \bar{P}_{\text{agg},t} + u_t. \quad (227)$$

Further, Gabaix and Kojien (2021) emphasize value-weighting. Accordingly, we estimate alternative specifications that replace the equal-weighted average with a value-weighted price index $\bar{P}_{\text{agg},t}^{vw} = \sum_i w_i P_{i,t}$, where $w_i = P_{0,i} / \sum_j P_{0,j}$.

Step 2c: Estimate our model. We estimate the relative elasticity $\hat{\mathcal{E}}$ from a cross-sectional regression, controlling for β_i as the observable, interacted with time-fixed effects, as this is a repeated cross-section:

$$\Delta \omega_{i,t} = \hat{\mathcal{E}} \Delta \log P_{i,t} + \theta_{1,t} \beta_i + \theta_{0,t} + e_{i,t}. \quad (228)$$

Denote by X_i a cross-sectionally demeaned and normalized (to unit norm) β_i . We estimate substitution from the time-series by running $K = 2$ bivariate regressions, estimating four substitution parameters:

$$\begin{pmatrix} \frac{1}{\sqrt{N}} \mathbf{1}' \Delta \omega_t \\ X' \Delta \omega_t \end{pmatrix} = \check{\mathcal{E}} \begin{pmatrix} \frac{1}{\sqrt{N}} \mathbf{1}' \Delta \log P_t \\ X' \Delta \log P_t \end{pmatrix} + \begin{pmatrix} e_t^{(1)} \\ e_t^{(2)} \end{pmatrix}, \quad (229)$$

where $\check{\mathcal{E}}$ is 2×2 , combining substitution and the relative elasticity as in Proposition 2.

⁶⁰Differences here are defined as perturbations from the baseline, i.e., $\Delta Y_{i,t} \equiv Y_{i,t} - Y_{i,0}$ for any variable Y between simulation t and the baseline 0.

Step 3: Run experiments I to III. Finally, we conduct three experiments based on the estimated models. Experiment I computes counterfactual prices from the CRRA model and inverts each estimated model for the implied supply shocks. Experiments II and III fix the supply shocks and compare the predicted price changes across models.

Experiment I evaluates the effects of a uniform supply shock of 10% to all risky assets, i.e., setting $S_1 = 0.9S_0$. We solve for the CRRA model-implied counterfactual prices P_1 , and W_1 based on equilibrium conditions in (223). The experiment then asks what size of supply shocks are required for each asset under different models—logit, logit + macro, and our model—to justify the counterfactual price P_1 from the CRRA model. Panel A of Figure 1 visualizes the log price change $\Delta \log P$ we require each demand model to match. For logit, we apply counterfactual prices P_1 , the estimates $\hat{\alpha}$ and $\hat{\epsilon}_i$ to Equation (225), which defines counterfactual portfolio shares. In our model, the counterfactual portfolio share is simply the initial portfolio share plus $\mathcal{E} \Delta \log P$, where \mathcal{E} is the estimated elasticity matrix combining the relative elasticity with substitution, i.e.,

$$\mathcal{E} = \hat{\epsilon}I + \left[\frac{1}{\sqrt{N}} \mathbf{1}, X \right] (\check{\mathcal{E}} - \hat{\epsilon}I_2) \left[\frac{1}{\sqrt{N}} \mathbf{1}, X \right]'. \quad (230)$$

The model combining logit and a macro elasticity computes counterfactual shares as the product of the predicted total risky asset portfolio share times within-risky asset portfolio shares, $\omega_i^{\text{Logit+Macro}}(P_1) = \omega_{\text{TOT}}^{\text{macro}}(P_1) \cdot \tilde{\omega}_i(P_1)$, where

$$\omega_{\text{TOT}}^{\text{macro}}(P_1) = \omega_{\text{TOT},0} \cdot \exp(\hat{\Theta}(\log \bar{P}_{\text{agg},1} - \log \bar{P}_{\text{agg},0})), \quad (231)$$

$$\tilde{\omega}_i(P_1) = \frac{\tilde{\omega}_{0,i} \cdot \exp(\hat{\alpha}(\log P_{1,i} - \log P_{0,i}))}{\sum_j \tilde{\omega}_{0,j} \cdot \exp(\hat{\alpha}(\log P_{1,j} - \log P_{0,j}))}, \quad (232)$$

with $\tilde{\omega}_{0,i}$ the within-risky assets portfolio share of asset i in the baseline equilibrium.

Experiment II evaluates the impact of an asymmetric supply shock with opposing signs for high-versus-low β assets. Specifically, it corresponds to reducing supply by 20% for assets in the highest- β quartile, and increasing supply by 20% for assets in the lowest- β quartile, analogous to an “operation-twist” type of supply shock. Panel A of Figure 7 visualizes the experiment. We then invert all estimated demand functions to obtain the predicted price change of each estimated demand model given the structure of the supply shock, and compare it to the one generated under the CRRA data-generating process. This is the inverse procedure to how Experiment I is calculated, and generally requires numerical inversion. Experiment III is similar, but reduces supply by 20% for assets with below-average β , leaving the remaining assets unchanged, as shown in Panel B of Figure 7.

I.3 Experiments

We now discuss the results from the three experiments. Across all of them, our model correctly captures the factor structure in counterfactuals, while logit and logit + macro struggle to match it.

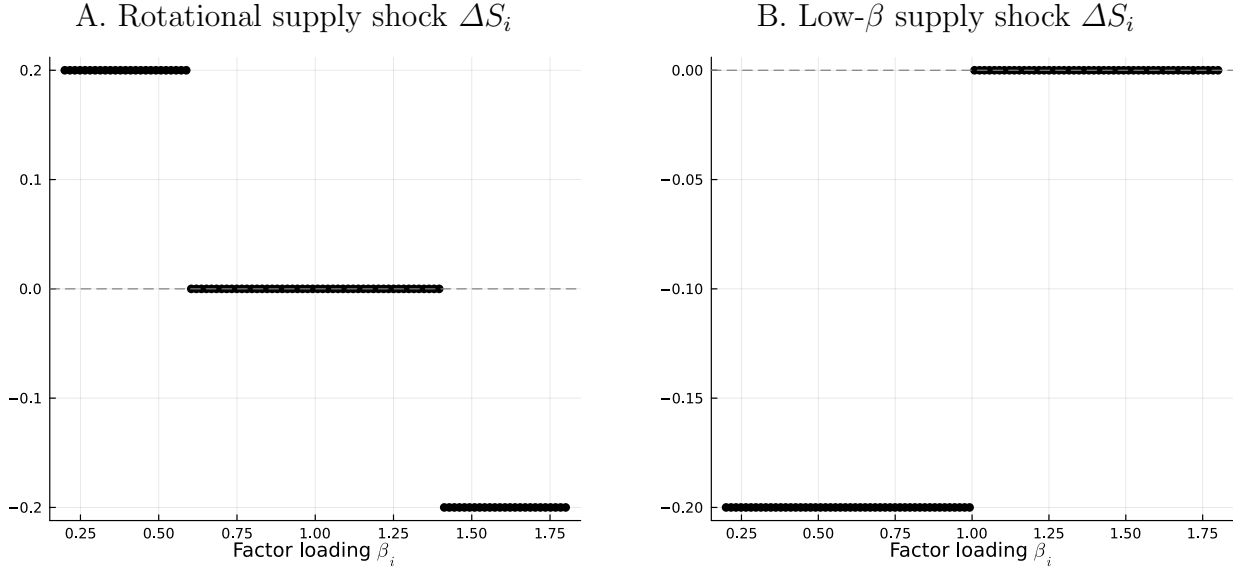


Figure 7: **Supply shocks in experiments II and III.** Figure 7 shows the supply perturbations ΔS_i as a function of asset beta β_i for experiments II and III. Panel A displays the rotational shock, which reduces supply by 20% for top-quartile β assets and increases it by 20% for bottom-quartile β assets. Panel B displays the low- β shock, which reduces supply by 20% for assets with below-average β and leaves remaining assets unchanged.

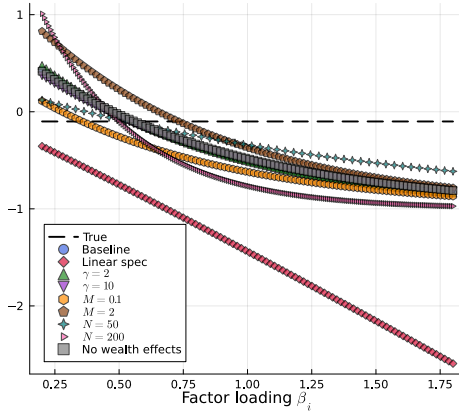
Experiment I: uniform supply shock A uniform supply shock of 10% of supply moves the price of low- β assets less than that of high- β assets, as shown in Panel A of Figure 1. How large a demand shock is needed to create that same movement in prices according to different empirical models? Panel B provides the answer for each model. Our model correctly infers the 10% shock that is uniform across the β distribution. In contrast, both logit and logit + macro incorrectly conclude that the only way to move prices of high- β assets more than those of low- β assets is through a demand shock that involves buying a lot more of high- β than low- β assets.

This conclusion is not specific to the chosen parameters or other choices as part of the experiment. While exact quantitative results matter based on parameter values, all alternative parameter choices— $\gamma = 2, 10$, $M = 0.1, 2$, $N = 50, 200$ —and other aspects of the experiment—estimating a log-linear specification instead of logit, excluding wealth effects, using value-weighting—qualitatively lead to the same conclusion as for the baseline version of the experiment, as shown in Panels A (logit) and Panel B (logit + macro) of Figure 8.

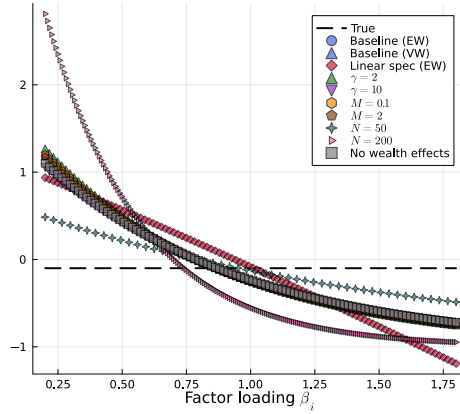
Experiment II: Rotational supply shock Experiment II is an “operation twist” type of supply shock, a rotation that involves reducing supply by 20% for top-quartile β assets while adding 20% supply to bottom-quartile β assets, as seen in Panel A of Figure 7.

Panel A of Figure 9 shows that generally, the price impact of the rotational supply shock is again, just like in experiment I, close to linearly increasing in β under the data-generating process, with two distinct jumps. First, the linear increase derives from the same mechanism as before; supply reduction in high- β assets with supply increase in low- β assets lowers

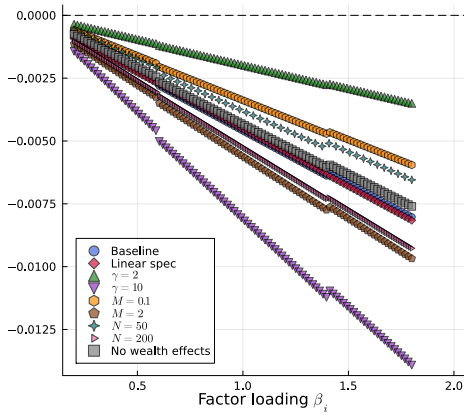
A. Implied supply change ΔS_i from logit



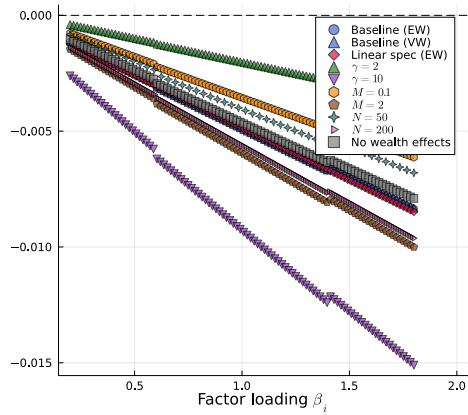
B. Implied ΔS_i from logit + macro



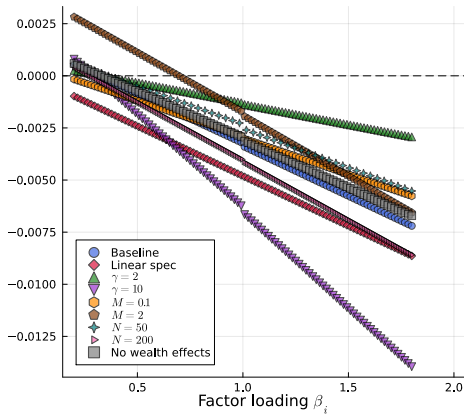
C. Logit price error (rotation)



D. Logit + macro price error (rotation)



E. Logit price error (low- β)



F. Logit + macro price error (low- β)

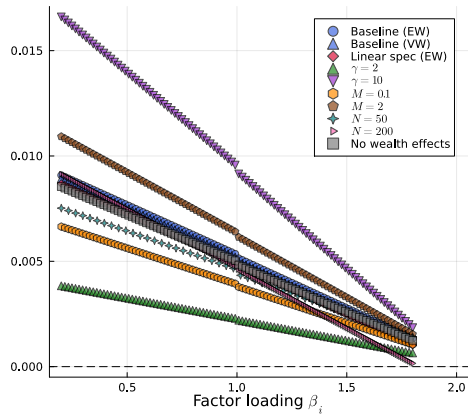


Figure 8: **Robustness of logit and logit + macro demand models across experiments.** Panels A–B plot the implied supply change, as in Figure 1, against asset beta β_i for the uniform 10% supply reduction of experiment I, under different robustness specifications for logit and logit + macro; the dashed black line shows the true supply change. Panels C–F plot the price prediction error $\Delta \log P_i^{\text{model}} - \Delta \log P_i^{\text{true}}$ relative to baseline prices P_0 for the rotation shock of experiment II (C–D) and the low- β shock of experiment III (E–F) under logit and logit + macro; the dashed line marks zero error. Each series corresponds to a different parameterization (γ , M , N ; see Table 6 for details), or specification variant: using log-linear specifications, adding a value-weighted price index (macro only), or shutting down wealth effects.

the risk premium associated with market risk, the response to which is proportional to β_i . Second, the distinct jumps are best understood by comparing an asset just above or below. For example, the supply of assets with a beta of just above the cutoff of 1.4 is reduced in the supply shock, while that of assets just below it is not.⁶¹ Accordingly, their relative price moves based on the inverse of the relative elasticity: the discrete jump. Our model again exactly matches this behavior.

In contrast, logit fails to match the general pattern. That is precisely because while it does generate the discontinuities in the right places, it fails to capture the linear increase derived from the factor risk premium. The same is true for logit + macro. Since the supply shock is “long-short,” there is (almost) no aggregate shock for the macro multiplier to create aggregate price variation.

As before, results are robust across parameter and implementation choices, as seen in Panels C (logit) and D (logit + macro) from Figure 8. Since the true, CRRA-implied price change generally varies with parameters in this experiment (and the third one), we plot the deviations from the CRRA-implied price changes, rather than raw price changes.

A. Price change $\Delta \log P$ from rotation shock

B. Price change $\Delta \log P$ from low- β shock

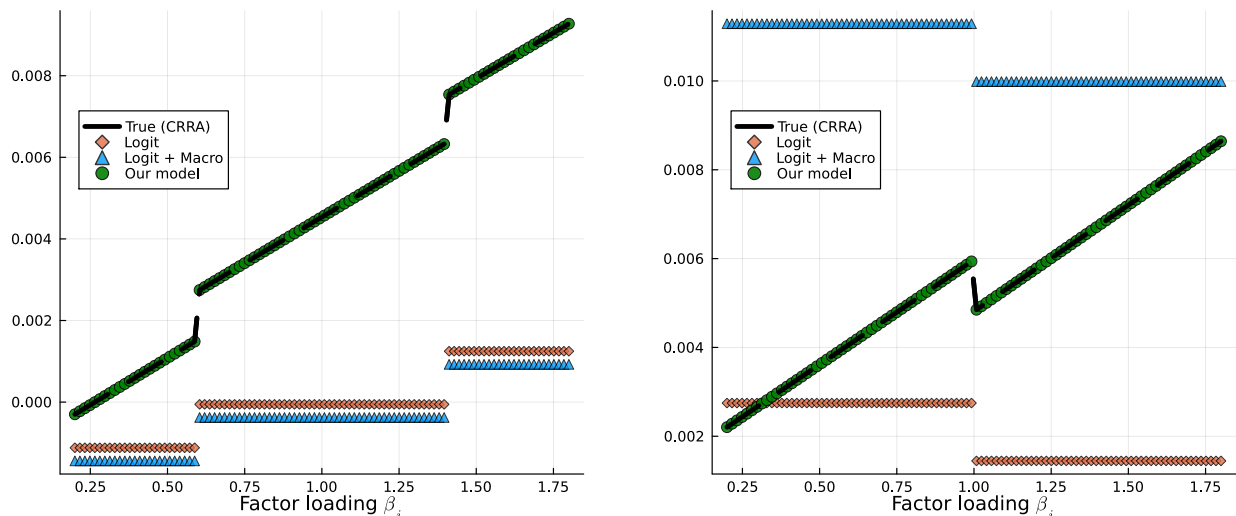


Figure 9: Price responses in experiments II and III. Figure 9 shows the equilibrium log price change $\Delta \log P_i$ relative to baseline prices P_0 as a function of asset beta β_i for experiments II and III. The dashed black line is the true CRRA response; colored markers show predictions from the logit (orange), logit + macro (blue), and our model (green). Panel A displays the response to the rotational shock of experiment II. Panel B displays the response to the low- β shock of experiment III. The supply shocks are detailed in Figure 7.

Experiment III: low- β supply shock The final experiment corresponds to a supply reduction of 20% in assets with a β below 1, as seen in Panel B of Figure 7.

Panel B of Figure 9 shows that the price change under the CRRA data-generating process takes a similar form as in the previous experiment, combining a linear increase based on β_i

⁶¹Such discrete cutoffs are similar in spirit to regression discontinuity designs.

with a discontinuity at $\beta_i = 1$, where assets slightly below move more because their supply is reduced, while that for assets with a beta just above one is not. And again, our model matches this counterfactual experiment.

As in the previous experiment, logit again only succeeds in capturing the discontinuity at $\beta_i = 1$, while failing to capture the response to changing factor risk premia. In contrast, logit + macro, while also creating the same discontinuity, now moves all prices by more than the largest true price change across assets. The reason is that for the macro part, the supply shock looks like a large aggregate shock; it involves reducing supply by 20% for half of the assets, almost like a 10% uniform supply shock. This, of course, misses that the shock comes from low β assets, and failing to account for that means overestimating the true magnitude of the shock. Since logit + macro overestimates the true magnitude of the shock, it naturally also overstates its impact on prices.

Again, results are robust across parameter and implementation choices, this time seen in Panels E (logit) and F (logit + macro) from Figure 8.

J Appendix Tables and Figures

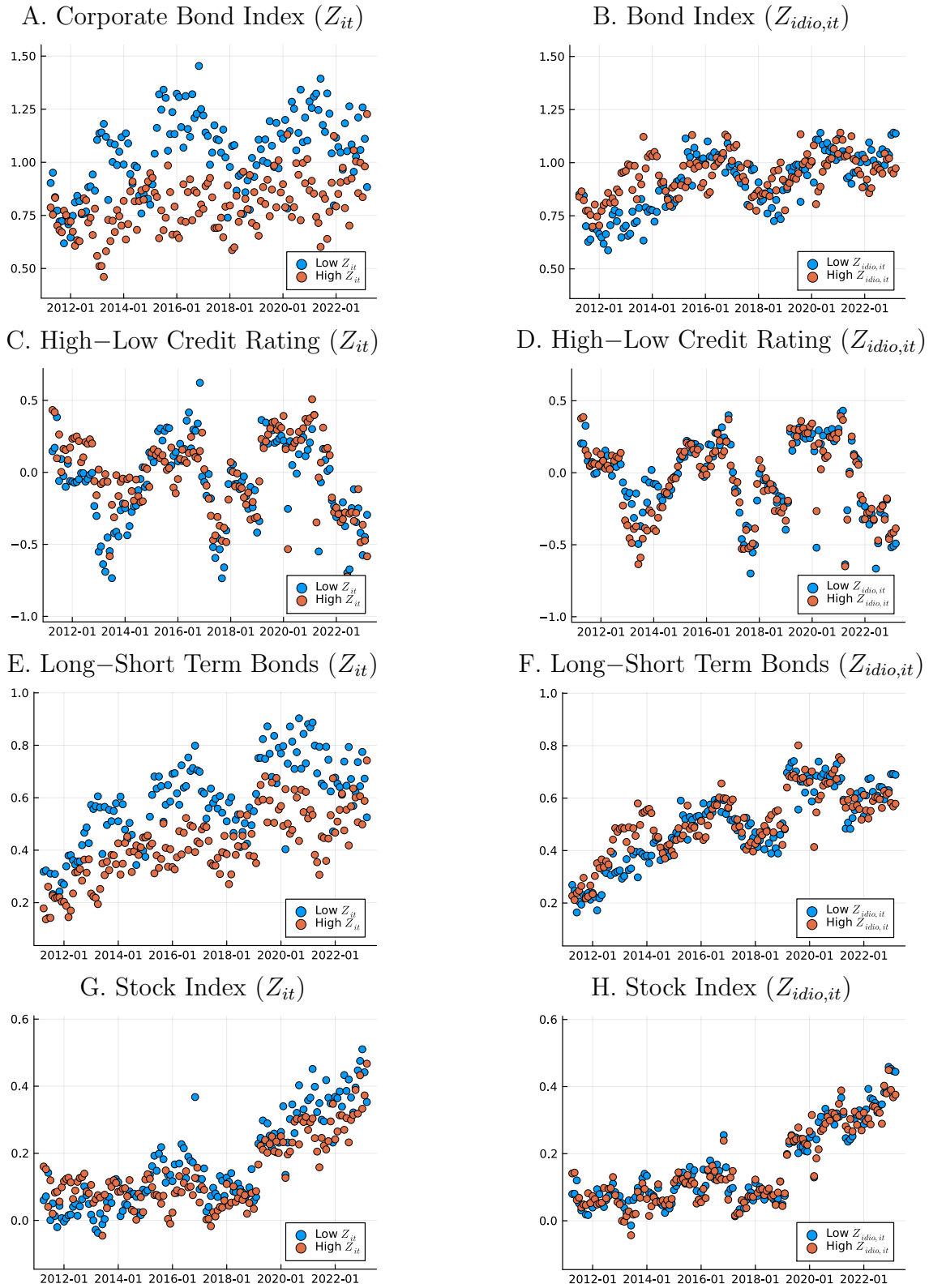


Figure 10: **Balance on covariances: exposure of portfolios sorted on demand shocks to various factors.** Figure 10 follows the exact definitions from Figure 4, but instead of showing the exposure of long-short portfolios to various factors, it shows the exposure for the long (orange) and short (blue) legs separately, sorted based on Z_{it} in the left panels and $Z_{idio,it}$ in the right panels.

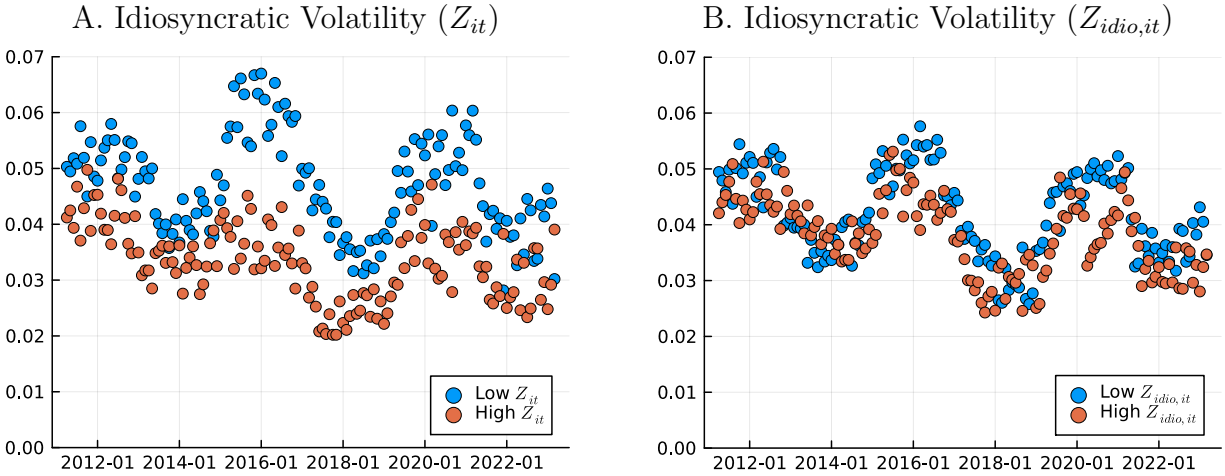


Figure 11: Balance on variances: average idiosyncratic volatility sorted on demand shocks. Figure 11 reports average idiosyncratic volatilities per group sorted on both the raw demand shock Z_{it} (blue) and the demand shock $Z_{idio,it}$ (orange) that is cross-sectionally orthogonalized to duration and credit risk, as proxied by long-run average default probabilities based on S&P ratings, at each point in time. At each date, we compute idiosyncratic volatilities for each corporate bond over a two-year window centered around t , excluding t , with respect to four factors: the ICE BofA US Corporate Index Total Return, the difference between the ICE BofA US High Yield Index Total Return and the ICE BofA US Corporate Index Total Return, the difference between the ICE BofA 15+ Year US Corporate Index Total Return and the ICE BofA 1-3 Year US Corporate Index Total Return, and the [Fama and French \(1993\)](#) excess stock market return. We present the equal-weighted average idiosyncratic volatility among bonds with above or below median demand shock Z_{it} (or $Z_{idio,it}$). The data for factors is from FRED and the Kenneth French data library. The time series is from 2011:04 to 2023:03.

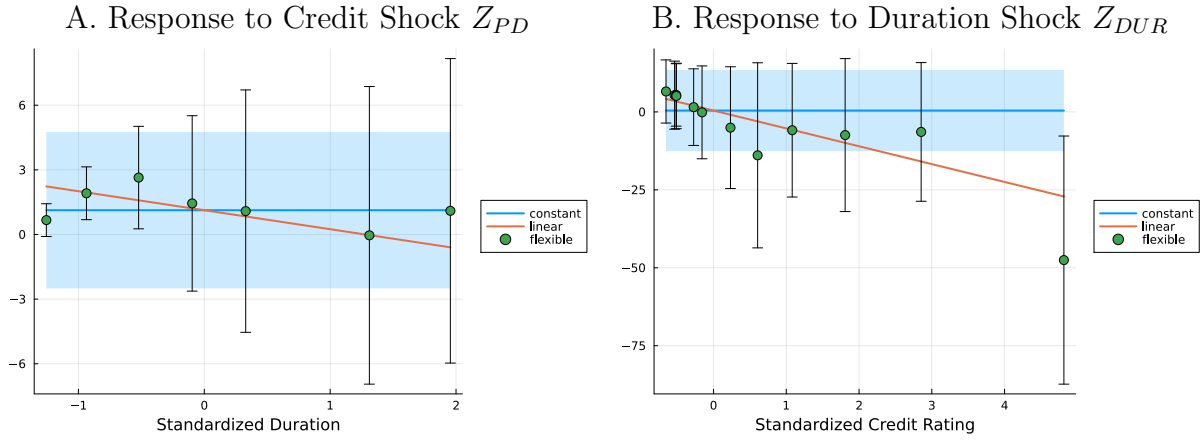


Figure 12: Off-diagonal macro- and meso multipliers in the cross-section. Figure 12 reports the off-diagonal elements omitted from Figure 5. That is, the response of bonds sorted on duration to a credit shock Z_{PD} (Panel A), and the response of bonds sorted on credit rating to a duration shock Z_{DUR} (Panel B). Bonds are grouped in seven buckets based on duration: <1 year, 1–3 years, 3–5 years, 5–7 years, 7–10 years, 10–15 years, and 15+ years. Bonds are grouped in seven buckets based on duration: <1 year, 1–3 years, 3–5 years, 5–7 years, 7–10 years, 10–15 years, and 15+ years. Bonds are also grouped by S&P credit rating, with individual notches from A+ through B-, and with AAA/AA and CCC/C ratings pooled at the extremes. The blue lines correspond to the estimates from column (2) of Table 2, which assume identical responses. The orange lines are based on columns (3) and (4), which include linear interaction terms with either duration or credit risk, neutralizing the other. The green dots estimate multipliers separately by duration or credit rating bucket in a pooled panel regression. The sample period is 2010:04 to 2024:03.

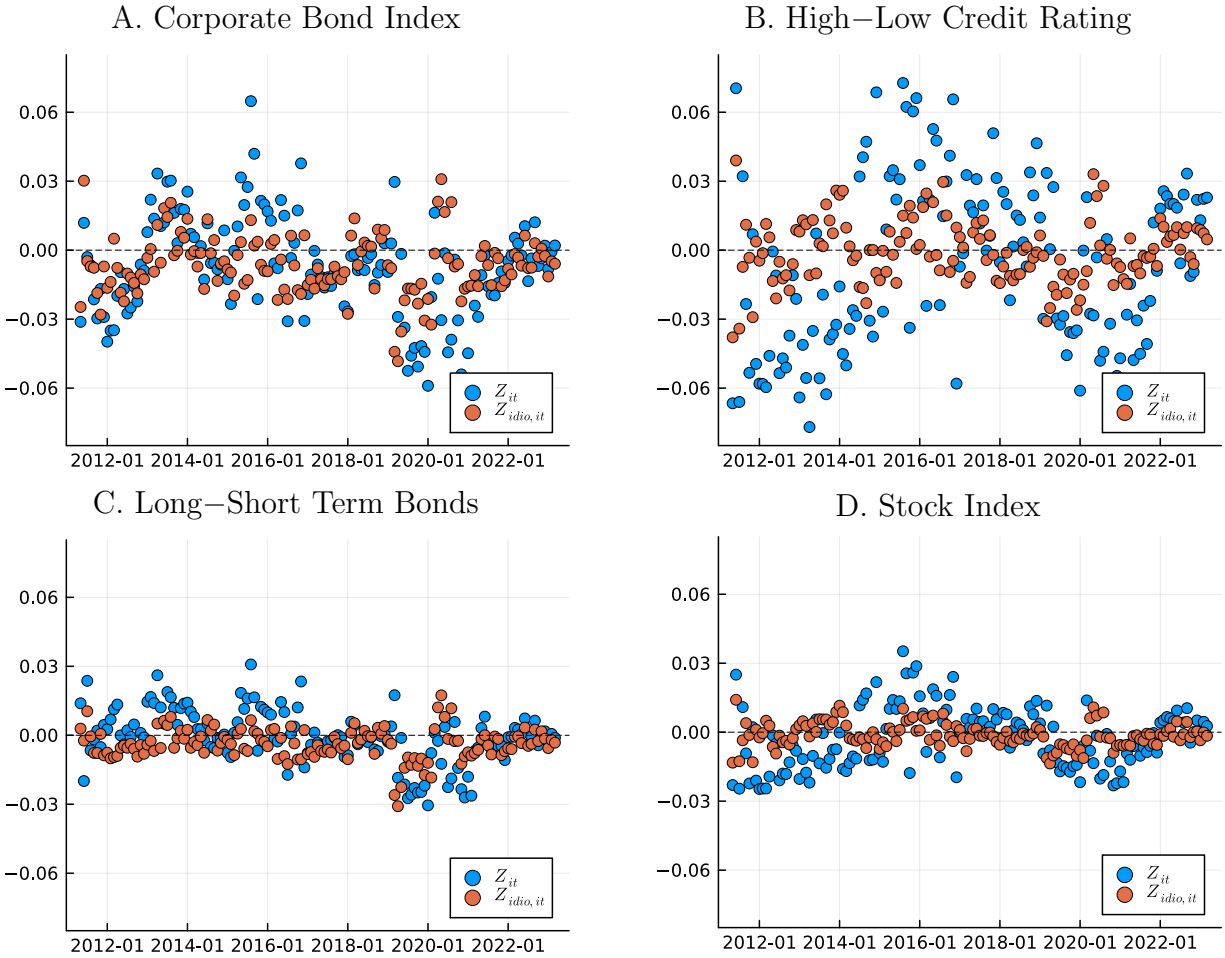


Figure 13: Balance on covariances in yield changes: exposure of long-short portfolios sorted on demand shocks to various factors. Figure 13 reports regression coefficients from balance-on-covariance regressions based on both the raw demand shock Z_{it} (blue) and the demand shock $Z_{idio,it}$ (orange) that is cross-sectionally orthogonalized to duration and credit risk, as proxied by long-run average default probabilities based on S&P ratings, at each point in time. At each date, we form long-short equal-weighted portfolios based on whether Z_{it} (or $Z_{idio,it}$) is above or below the median. We compute the yield changes of these portfolios over two years centered around t , excluding t , and regress these yield changes on four aggregate factors. Panel A shows the time-series of coefficients for regressions on an aggregate investment-grade corporate bond factor, the ICE BofA US Corporate Index Total Return. Panel B uses the difference between aggregate high-yield and investment-grade corporate bond factors, the ICE BofA US High Yield Index Total Return and the ICE BofA US Corporate Index Total Return. Panel C uses the difference between the ICE BofA 15+ Year US Corporate Index Total Return and the ICE BofA 1-3 Year US Corporate Index Total Return. Panel D uses the [Fama and French \(1993\)](#) excess stock market return. The data for factors in panels A to C is from FRED, while the data for the excess market return in Panel D is from the Kenneth French data library. The time series is from 2011:05 to 2023:03.

Table 7: **Relative yield multiplier \widehat{M} in corporate bonds**

	Yield change ΔY_{it}				
	(1)	(2)	(3)	(4)	(5)
<i>Demand shock:</i>					
Z_{it}	-0.058 (0.030)			-0.281*** (0.060)	-0.584*** (0.134)
$Z_{idio,it}$		-0.058 (0.030)	-0.058 (0.030)		
Date Fixed Effects	Yes	Yes	Yes	Yes	
Duration \times Date Fixed Effects	Yes	Yes			
Credit Risk \times Date Fixed Effects	Yes	Yes			
N	1,016,830	1,016,830	1,016,830	1,016,830	1,016,830
R^2	0.099	0.099	0.040	0.042	0.006

Table 7 reports the results of relative multiplier regressions of yield changes ΔY_{it} on demand shocks Z_{it} and $Z_{idio,it}$ for non-defaulted U.S. corporate bonds. Specifications (1) and (4)–(5) use the flow-induced trading demand shock Z_{it} defined in Equation (46). Specification (1) includes a time fixed effect and controls for a continuous duration variable and a continuous credit risk variable based on average historical default probabilities for each S&P credit rating category, while specification (4) uses only date fixed effects and specification (5) only a common intercept. Specifications (2)–(3) use the demand shock $Z_{idio,it}$ orthogonalized to duration and credit risk each period, with and without controlling for duration and credit risk in the regression. The regressions weigh each date equally. The sample period is 2010:05 to 2024:03. Standard errors are clustered by date and issuer.

Table 8: Macro- and meso yield multipliers in corporate bonds

	Yield Change				
	$\Delta Y_{agg,t}$	$\Delta Y_{DUR,t}$	$\Delta Y_{PD,t}$	ΔY_{it}	$\Delta Y_{agg,t}$
	(1)	(2)	(3)	(4)	(5)
$Z_{agg,t}$	-2.506*** (0.526)	0.146 (0.150)	-1.475 (0.839)	-2.506*** (0.523)	-2.961*** (0.565)
$Z_{DUR,t}$	1.142 (1.485)	-0.943* (0.391)	3.040 (2.507)	1.142 (1.471)	
$Z_{PD,t}$	-0.498 (0.330)	0.227* (0.109)	-0.939 (0.491)	-0.498 (0.329)	
$Z_{idio,it}$				-0.058 (0.030)	
N	167	167	167	1,016,830	167
R^2	0.452	0.205	0.277	0.018	0.443

Table 8 reports the results of macro- and meso multiplier regressions of yield changes on demand shocks for non-defaulted U.S. corporate bonds. Specifications (1)–(3) jointly estimate the matrix $\widehat{\mathbf{M}}$ from Proposition 5, which together with the relative multiplier $\widehat{\mathcal{M}}$ determines the spillover matrix between observables, \mathbf{M}_X . Specification (4) estimates multipliers mechanically identical to specification (1) using disaggregated, repeated cross-sectional regressions, while adding the relative multiplier $\widehat{\mathcal{M}}$. Specification (5) estimates the macro multiplier in isolation by regressing aggregate bond yield changes $\Delta Y_{agg,t}$ on the aggregated instrument $Z_{agg,t}$ in the time series. The $K = 3$ observables are a vector of ones, standardized duration, and standardized credit risk. The sample period is 2010:04 to 2024:03. Robust standard errors are used for specifications (1) to (3) and (5). For specification (4), standard errors are clustered by date and issuer, and regressions are weighted such that each date receives equal weight.

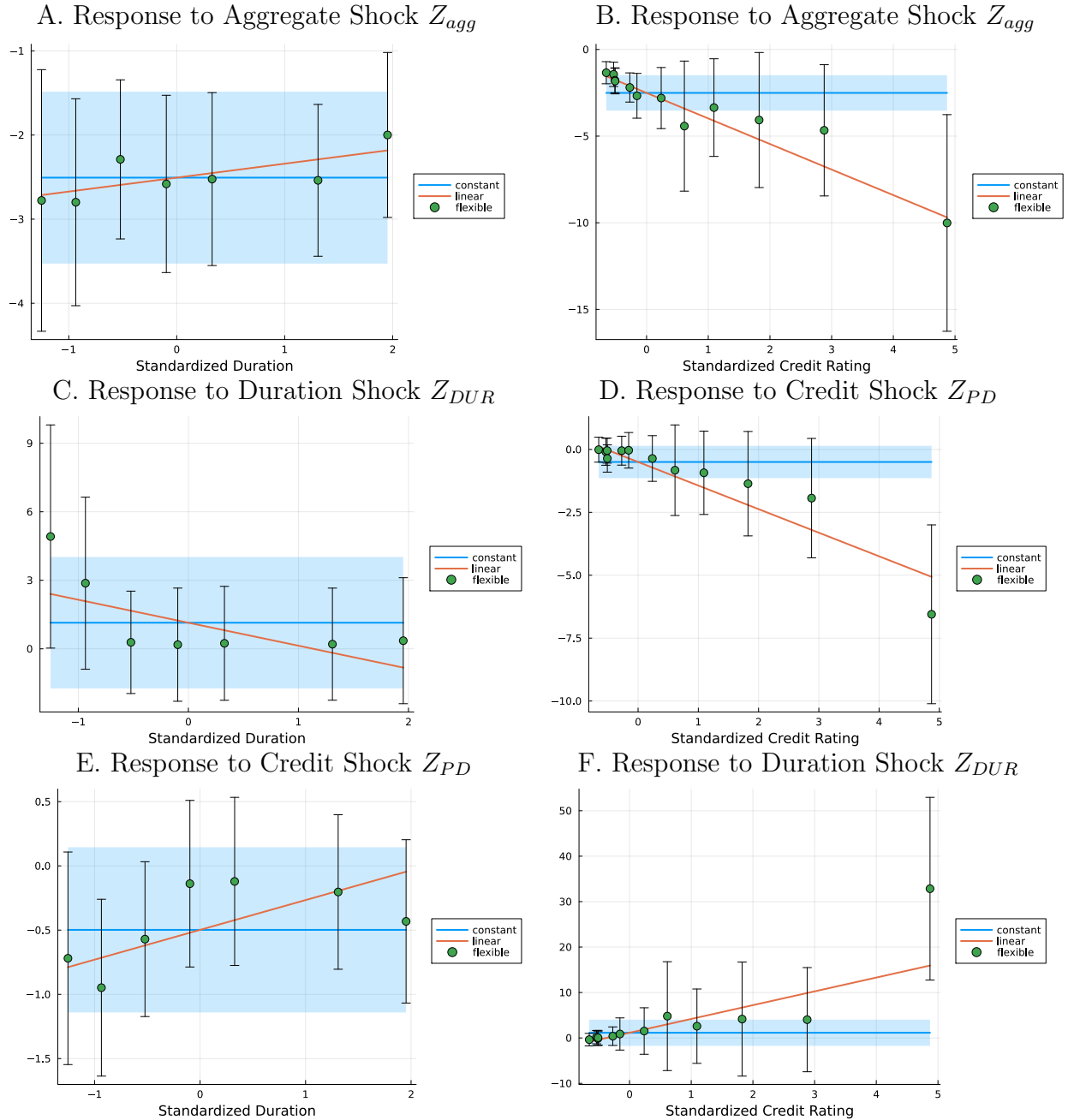


Figure 14: Macro- and meso yield multipliers in the cross-section. Figure 14 reports the yield response of portfolios of corporate bonds to aggregate demand shocks Z_{agg} (Panels A and B, in the cross-sections of duration and credit rating) and shocks along duration Z_{DUR} (Panels C and F) and Z_{PD} (Panels D and E). Bonds are grouped in seven buckets based on duration: <1 year, 1–3 years, 3–5 years, 5–7 years, 7–10 years, 10–15 years, and 15+ years. Bonds are also grouped by S&P credit rating, with individual notches from A+ through B-, and with AAA/AA and CCC/C ratings pooled at the extremes. The blue lines correspond to the estimates from column (2) of Table 8, which assume identical responses. The orange lines are based on columns (3) and (4), which include linear interaction terms with either duration or credit risk, neutralizing the other. The green dots estimate multipliers separately by duration or credit rating bucket in a pooled panel regression. The sample period is 2010:04 to 2024:03.

Table 9: Macro- and meso multipliers in corporate bonds with ZCA whitening

	Return					
	$R_{agg,t}$	$R_{DUR,t}$	$R_{PD,t}$	R_{it}		$R_{agg,t}$
	(1)	(2)	(3)	(4)	(5)	(6)
$Z_{agg,t}$	14.315*** (2.669)	7.018*** (1.815)	3.548** (1.324)	14.315*** (2.651)	14.315*** (2.645)	15.032*** (2.393)
$Z_{DUR,t}$	-0.187 (6.729)	5.541 (4.234)	-6.100 (3.901)	-0.187 (6.663)	-0.187 (6.658)	
$Z_{PD,t}$	1.162 (2.156)	-0.489 (1.443)	2.466 (1.305)	1.162 (2.138)	1.162 (2.137)	
$Z_{idio,it}$				0.055 (0.087)	0.055 (0.084)	
Duration DUR_{it}					-0.002 (0.001)	
Credit Risk PD_{it}					-0.002 (0.001)	
$Z_{agg,t} \times DUR_{it}$					7.018*** (1.793)	
$Z_{agg,t} \times PD_{it}$					3.548** (1.322)	
$Z_{DUR,t} \times DUR_{it}$					5.541 (4.180)	
$Z_{DUR,t} \times PD_{it}$					-6.100 (3.885)	
$Z_{PD,t} \times DUR_{it}$					-0.489 (1.429)	
$Z_{PD,t} \times PD_{it}$					2.466 (1.305)	
N	168	168	168	1,041,985	1,041,985	168
R^2	0.343	0.125	0.386	0.100	0.138	0.342

Table 9 reports the results of macro- and meso multiplier regressions of bond returns on demand shocks for non-defaulted U.S. corporate bonds. It is identical to Table 2, other than that it uses ZCA whitening of standardized duration and credit risk at each date to create orthonormal observables X . It also adds another specification (5), showing that all estimates using disaggregated data are mechanically identical to those obtained from time-series regressions in specifications (1)–(3). The sample period is 2010:04 to 2024:03. Robust standard errors are used for specifications (1) to (3) and (6). For specifications (4) and (5), standard errors are clustered by date and issuer, and regressions are weighted such that each date receives equal weight.