

Market Efficiency in the Age of Big Data

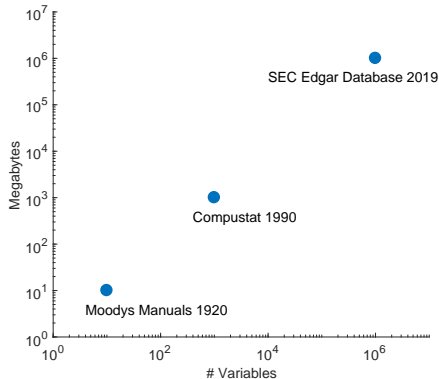
Ian Martin Stefan Nagel

April 2021



Investors' Big Data problem

- Investors face huge number of potential predictors
- $\mathbb{E}[\text{return}] = f(\text{predictors})$
unknown: high-dimensional learning problem



High-dimensional learning in asset pricing

- Standard approaches in asset pricing and market efficiency testing assume rational expectations (RE)
 - ▶ Assumes away learning problem: investors know $f(\cdot)$ in $\mathbb{E}[\text{return}] = f(\text{predictors})$
 - ▶ Motivates in-sample (IS) tests of “market efficiency”:
IS return predictability = risk premium/mispricing
- We show: when investors **learn** about $f(\cdot)$ in **big data** setting, equilibrium asset prices exhibit in-sample predictability
- Combination of learning and big data provides clear motivation for (pseudo-)OOS testing which is lacking in RE framework

Market efficiency

- Fama (1970): A market is efficient if “prices fully reflect all available information”
- Joint hypothesis problem: your “risk premium” is my “pricing inefficiency”
- We make things simple by considering a risk-neutral world
- Then the joint hypothesis problem goes away... but standard tests of market efficiency break down even so

Roadmap

Two steps:

- ① Investors learn about parameters of cash flow generating model and price assets accordingly
- ② Econometrician analyzes equilibrium prices ex post using standard return predictability tests
 - ▶ Properties of IS tests
 - ▶ Properties of OOS tests

Setup

- N assets, $N \times J$ scaled characteristics arranged into a matrix \mathbf{X} , eg,

$$\begin{pmatrix} \text{Size}_{\text{AAPL}} & \text{Leverage}_{\text{AAPL}} & \text{Liquidity}_{\text{AAPL}} & \cdots & \text{CharJ}_{\text{AAPL}} \\ \text{Size}_{\text{AMZN}} & \text{Leverage}_{\text{AMZN}} & \text{Liquidity}_{\text{AMZN}} & \cdots & \text{CharJ}_{\text{AMZN}} \\ \text{Size}_{\text{FB}} & \text{Leverage}_{\text{FB}} & \text{Liquidity}_{\text{FB}} & \cdots & \text{CharJ}_{\text{FB}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{Size}_N & \text{Leverage}_N & \text{Liquidity}_N & \cdots & \text{CharJ}_N \end{pmatrix}$$

- Normally we think of many assets and a limited number of characteristics \implies fixed- J , large- N asymptotics
- This paper: many assets and many characteristics \implies large- N and large- J asymptotics

Setup

- Investors are homogeneous and risk-neutral; interest rate is zero
- Dividend strips: \mathbf{p}_t = prices at t of claims to \mathbf{y}_{t+1}
 - ▶ Think: one period \approx one decade

- Dividend growth $\Delta\mathbf{y}_t$ is predictable based on characteristics \mathbf{X} :

$$\Delta\mathbf{y}_{t+1} = \mathbf{X}\mathbf{g} + \mathbf{e}_{t+1}, \quad \mathbf{e}_{t+1} \sim N(\mathbf{0}, \mathbf{I})$$

- These assumptions are chosen to make life as simple as possible

Setup

- Prices equal expected dividends

$$p_t = \tilde{\mathbb{E}}_t y_{t+1} = y_t + \tilde{\mathbb{E}}_t \Delta y_{t+1} = y_t + \tilde{\mathbb{E}}_t (\mathbf{X} \mathbf{g} + \mathbf{e}_{t+1})$$

- Encompasses a range of possible assumptions about expectations $\tilde{\mathbb{E}}$
- Benchmarks to keep in mind...
 - ▶ Rational expectations: investors know \mathbf{g}
 - ▶ OLS: regress past cashflow growth on \mathbf{X} to estimate \mathbf{g}
 - ▶ Random walk: give up on forecasting

Rational expectations: investors know \mathbf{g}

- So $\tilde{\mathbb{E}}_t(\mathbf{X}\mathbf{g} + \mathbf{e}_{t+1}) = \mathbf{X}\mathbf{g}$ and $\mathbf{p}_t = \mathbf{y}_t + \mathbf{X}\mathbf{g}$
- Realized returns $\mathbf{r}_{t+1} = \mathbf{y}_{t+1} - \mathbf{p}_t = \mathbf{e}_{t+1}$
- This is the usual null hypothesis that underlies market efficiency tests, orthogonality conditions, Euler equations
- But it is implausible that investors know \mathbf{g} , especially if J is large
- We focus on the case where investors must learn \mathbf{g}
- Consider, first, two extreme possibilities. . .

Unknown \mathbf{g} : OLS with many predictors

- Investors learn \mathbf{g} by running OLS with as many predictors as assets, $J = N$
- Regression of $\Delta \mathbf{y}_t$ on \mathbf{X} exactly fits $\Delta \mathbf{y}_t$ in sample: $\tilde{\mathbb{E}}_t \Delta \mathbf{y}_{t+1} = \Delta \mathbf{y}_t$
- Prices $\mathbf{p}_t = \mathbf{y}_t + \Delta \mathbf{y}_t$ and returns $\mathbf{r}_{t+1} = \Delta \mathbf{y}_{t+1} - \Delta \mathbf{y}_t$
- Forecast MSE is $\text{var}(\mathbf{e}_{t+1} - \mathbf{e}_t)$, i.e., twice the variance of \mathbf{e}_{t+1}

Unknown \mathbf{g} : The random walk

- Completely give up on prediction: $\tilde{\mathbb{E}}_t \Delta \mathbf{y}_{t+1} = \mathbf{0}$
- Prices $\mathbf{p}_t = \mathbf{y}_t$ and returns $\mathbf{r}_{t+1} = \Delta \mathbf{y}_{t+1}$
- Forecast MSE is $\text{var}(\mathbf{X}\mathbf{g} + \mathbf{e}_{t+1})$

Unknown g : The random walk

- Completely give up on prediction: $\tilde{\mathbb{E}}_t \Delta \mathbf{y}_{t+1} = \mathbf{0}$
- Prices $\mathbf{p}_t = \mathbf{y}_t$ and returns $\mathbf{r}_{t+1} = \Delta \mathbf{y}_{t+1}$
- Forecast MSE is $\text{var}(\mathbf{X}\mathbf{g} + \mathbf{e}_{t+1})$
- If cash-flow growth is hard to predict— $\text{var} \mathbf{X}\mathbf{g} \ll \text{var} \mathbf{e}_{t+1}$ —then the **random walk forecast** may outperform the **OLS forecast**

$$\text{var}(\mathbf{X}\mathbf{g} + \mathbf{e}_{t+1}) \ll \text{var}(\mathbf{e}_{t+1} - \mathbf{e}_t)$$

Bayesian pricing framework: Prior beliefs

- Before seeing data, investors hold informed prior beliefs

$$\mathbf{g} \sim N\left(\mathbf{0}, \frac{\theta}{J}\mathbf{I}\right), \quad \theta > 0$$

- ▶ Proportionality of prior covariance matrix to \mathbf{I} : can always rotate and rescale \mathbf{X} to make it hold
- ▶ Variance of the elements of \mathbf{g} decline with J : ensures that variance of predictable cash flow growth does not explode when $N, J \rightarrow \infty$
- Investors then learn about \mathbf{g} by observing \mathbf{X} and history $\{\Delta\mathbf{y}_s\}_1^t$, summarized by sample average $\overline{\Delta\mathbf{y}}_t$

Bayesian pricing framework: Posterior mean

- Posterior mean is a ridge regression estimator

$$\tilde{\mathbf{g}}_t = \mathbf{\Gamma}_t (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \overline{\Delta \mathbf{y}}_t$$

i.e., OLS estimator shrunk towards prior mean by the matrix

$$\mathbf{\Gamma}_t = \mathbf{Q} \left(\mathbf{I} + \frac{J}{N\theta t} \mathbf{\Lambda}^{-1} \right)^{-1} \mathbf{Q}'$$

where \mathbf{Q} , $\mathbf{\Lambda}$ come from PC decomposition $\frac{1}{N} \mathbf{X}'\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$

Bayesian pricing framework: Posterior mean

- Posterior mean is a ridge regression estimator

$$\tilde{\mathbf{g}}_t = \mathbf{\Gamma}_t (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \overline{\Delta \mathbf{y}}_t$$

i.e., OLS estimator shrunk towards prior mean by the matrix

$$\mathbf{\Gamma}_t = \mathbf{Q} \left(\mathbf{I} + \frac{J}{N\theta t} \mathbf{\Lambda}^{-1} \right)^{-1} \mathbf{Q}'$$

where \mathbf{Q} , $\mathbf{\Lambda}$ come from PC decomposition $\frac{1}{N} \mathbf{X}'\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$

- Shrinkage strong
 - ▶ if t small (short time dimension)

Bayesian pricing framework: Posterior mean

- Posterior mean is a ridge regression estimator

$$\tilde{\mathbf{g}}_t = \mathbf{\Gamma}_t (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \overline{\Delta \mathbf{y}}_t$$

i.e., OLS estimator shrunk towards prior mean by the matrix

$$\mathbf{\Gamma}_t = \mathbf{Q} \left(\mathbf{I} + \frac{J}{N\theta_t} \mathbf{\Lambda}^{-1} \right)^{-1} \mathbf{Q}'$$

where \mathbf{Q} , $\mathbf{\Lambda}$ come from PC decomposition $\frac{1}{N} \mathbf{X}'\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$

- Shrinkage strong
 - ▶ if t small (short time dimension)
 - ▶ if θ small (prior tightly concentrated around zero)

Bayesian pricing framework: Posterior mean

- Posterior mean is a ridge regression estimator

$$\tilde{\mathbf{g}}_t = \mathbf{\Gamma}_t (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \overline{\Delta \mathbf{y}}_t$$

i.e., OLS estimator shrunk towards prior mean by the matrix

$$\mathbf{\Gamma}_t = \mathbf{Q} \left(\mathbf{I} + \frac{J}{N\theta t} \mathbf{\Lambda}^{-1} \right)^{-1} \mathbf{Q}'$$

where \mathbf{Q} , $\mathbf{\Lambda}$ come from PC decomposition $\frac{1}{N} \mathbf{X}'\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$

- Shrinkage strong
 - ▶ if t small (short time dimension)
 - ▶ if θ small (prior tightly concentrated around zero)
 - ▶ if J/N is large (many predictors)

Bayesian pricing framework: Posterior mean

- Posterior mean is a ridge regression estimator

$$\tilde{\mathbf{g}}_t = \mathbf{\Gamma}_t (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \overline{\Delta \mathbf{y}}_t$$

i.e., OLS estimator shrunk towards prior mean by the matrix

$$\mathbf{\Gamma}_t = \mathbf{Q} \left(\mathbf{I} + \frac{J}{N\theta t} \mathbf{\Lambda}^{-1} \right)^{-1} \mathbf{Q}'$$

where \mathbf{Q} , $\mathbf{\Lambda}$ come from PC decomposition $\frac{1}{N} \mathbf{X}'\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$

- Shrinkage strong
 - ▶ if t small (short time dimension)
 - ▶ if θ small (prior tightly concentrated around zero)
 - ▶ if J/N is large (many predictors)
 - ▶ along unimportant principal components of \mathbf{X} (small **eigenvalues**)

Equilibrium realized returns

Proposition

With assets priced based on $\tilde{\mathbf{g}}_t$, realized returns are

$$\mathbf{r}_{t+1} = \mathbf{y}_{t+1} - \mathbf{p}_t = \mathbf{X}(\mathbf{I} - \mathbf{\Gamma}_t)\mathbf{g} - \mathbf{X}\mathbf{\Gamma}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{e}}_t + \mathbf{e}_{t+1}$$

where $\bar{\mathbf{e}}_t = \frac{1}{t} \sum_{s=1}^t \mathbf{e}_s$

Equilibrium realized returns

Proposition

With assets priced based on $\tilde{\mathbf{g}}_t$, realized returns are

$$\mathbf{r}_{t+1} = \mathbf{y}_{t+1} - \mathbf{p}_t = \mathbf{X}(\mathbf{I} - \mathbf{\Gamma}_t)\mathbf{g} - \mathbf{X}\mathbf{\Gamma}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{e}}_t + \mathbf{e}_{t+1}$$

where $\bar{\mathbf{e}}_t = \frac{1}{t} \sum_{s=1}^t \mathbf{e}_s$

- “underreaction” to \mathbf{X} due to shrinkage

Equilibrium realized returns

Proposition

With assets priced based on $\tilde{\mathbf{g}}_t$, realized returns are

$$\mathbf{r}_{t+1} = \mathbf{y}_{t+1} - \mathbf{p}_t = \mathbf{X}(\mathbf{I} - \mathbf{\Gamma}_t)\mathbf{g} - \mathbf{X}\mathbf{\Gamma}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{e}}_t + \mathbf{e}_{t+1}$$

where $\bar{\mathbf{e}}_t = \frac{1}{t} \sum_{s=1}^t \mathbf{e}_s$

- “underreaction” to \mathbf{X} due to shrinkage
- “overreaction” to estimation error in $\tilde{\mathbf{g}}_t$, dampened by shrinkage

Equilibrium realized returns

Proposition

With assets priced based on $\tilde{\mathbf{g}}_t$, realized returns are

$$\mathbf{r}_{t+1} = \mathbf{y}_{t+1} - \mathbf{p}_t = \mathbf{X}(\mathbf{I} - \mathbf{\Gamma}_t)\mathbf{g} - \mathbf{X}\mathbf{\Gamma}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{e}}_t + \mathbf{e}_{t+1}$$

where $\bar{\mathbf{e}}_t = \frac{1}{t} \sum_{s=1}^t \mathbf{e}_s$

- “underreaction” to \mathbf{X} due to shrinkage
- “overreaction” to estimation error in $\tilde{\mathbf{g}}_t$, dampened by shrinkage
- **unpredictable shock (the only term in RE case)**

Predictive coefficient estimates, h_{t+1}

- Econometrician cross-sectionally regresses (OLS)

$$r_{t+1} = \underbrace{X(I - \Gamma_t)g}_{\text{"underreaction"}} - \underbrace{X\Gamma_t(X'X)^{-1}X'\bar{e}_t}_{\text{"overreaction"}} + \underbrace{e_{t+1}}_{\text{RE}}$$

on characteristics matrix X and obtains predictive coefficients

$$h_{t+1} = (I - \Gamma_t)g - \Gamma_t (X'X)^{-1} X'\bar{e}_t + (X'X)^{-1} X'e_{t+1}$$

- “Kitchen sink” regression approximates what many individual studies have done collectively (“factor zoo”)

In-sample predictability test: RE null

- Consider the return predictability test statistic

$$T_{re} \equiv \frac{\mathbf{h}'_{t+1} \mathbf{X}' \mathbf{X} \mathbf{h}_{t+1} - J}{\sqrt{2J}}$$

- Standard approach takes RE as null hypothesis, which implies

$$\mathbf{h}_{t+1} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{e}_{t+1}$$

- If so,

$$T_{re} \xrightarrow{d} N(0, 1) \quad \text{as } N, J \rightarrow \infty, J/N \rightarrow \psi > 0$$

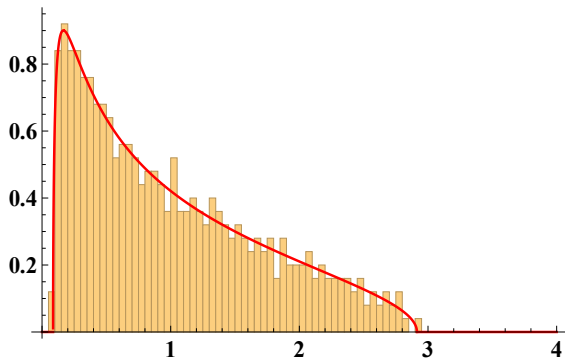
Big data

What happens to the data as $J, N \rightarrow \infty$ with $J/N \rightarrow \psi > 0$?

- There are two possibilities
- Case 1: A few principal components summarize the data
 - ▶ Formally: the eigenvalues of $\frac{1}{N}\mathbf{X}'\mathbf{X}$ tend to zero
 - ▶ Then market efficiency test works as usual, $T_{re} \xrightarrow{d} N(0, 1)$
- Case 2: “Big data”
 - ▶ Formally: the eigenvalues of $\frac{1}{N}\mathbf{X}'\mathbf{X}$ are $> \varepsilon$
 - ▶ This is our case of interest
 - ▶ (Happens if, eg, the entries of \mathbf{X} are iid random variables)

Eigenvalues in an example with iid random X

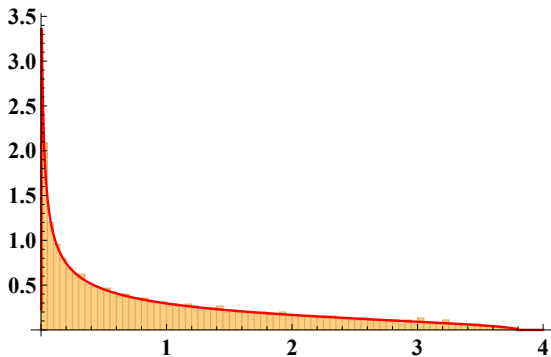
Histogram for $N = 1000$, $J = 500$ & asymptotic limit as $N, J \rightarrow \infty$ with $\frac{J}{N} = 0.5$



- When J is large, random data matrix X has many columns that are roughly orthogonal, hence many eigenvalues close to zero
- But our big data assumption is satisfied (Bai and Yin, 1993)

Eigenvalues in an example with iid random X

Histogram for $N = 1000$, $J = 900$ & asymptotic limit as $N, J \rightarrow \infty$ with $\frac{J}{N} = 0.9$



- When J is large, random data matrix X has many columns that are roughly orthogonal, hence many eigenvalues close to zero
- But our big data assumption is satisfied (Bai and Yin, 1993)

In-sample predictability test with big data

Proposition

The test statistic T_{re} satisfies

$$\frac{T_{re}}{\sqrt{\mu^2 + \sigma^2}} - \frac{\mu - 1}{\sqrt{2(\mu^2 + \sigma^2)}} \sqrt{J} \xrightarrow{d} N(0, 1)$$

where $1 < \mu < 2$ and $1 < \sqrt{\mu^2 + \sigma^2} < 2$ are determined by eigenvalues

- Therefore,

$$T_{re} \approx \sqrt{\mu^2 + \sigma^2} N(0, 1) + \frac{\mu - 1}{\sqrt{2}} \sqrt{J}$$

- In a big data world, we are almost certain to reject the RE null

Interpretation as a trading strategy

- Consider a characteristics-based trading strategy with weights

$$\mathbf{w}_{IS,t} = \mathbf{X}\mathbf{h}_{t+1}, \quad r_{IS,t+1} = \mathbf{w}'_{IS,t}\mathbf{r}_{t+1}$$

(“in-sample” because \mathbf{h}_{t+1} estimated using returns \mathbf{r}_{t+1})

- We can rewrite $r_{IS,t+1} = \mathbf{h}'_{t+1}\mathbf{X}'\mathbf{r}_{t+1}$ as

$$r_{IS,t+1} = \mathbf{h}'_{t+1}\mathbf{X}'\mathbf{X}\mathbf{h}_{t+1}$$

- Econometrician's test is equivalent to checking whether the trading strategy does well

Conclusions so far

- Asset returns under high-dimensional learning are very different from asset returns under RE, or in a “small data” world
- IS return predictability need not be consequence of risk premia or behavioral biases
- Not an econometric issue: the RE null is simply false, because learning + big data makes returns predictable in sample even without risk premia or behavioral biases
- Existence of a “factor zoo” based on IS predictability evidence not surprising in high-dimensional setting

(Absence of) out-of-sample return predictability

Proposition

Consider an out-of-sample strategy with predicted returns as portfolio weights, $r_{OOS,t+1} = \mathbf{r}'_{t+1} \mathbf{X} \mathbf{h}_{s+1}$ where $t \neq s$. Then $\mathbb{E} [\mathbf{r}'_{t+1} \mathbf{X} \mathbf{h}_{s+1}] = 0$

- **Forward** case $t > s$ is natural: Investors are Bayesian so the econometrician cannot “beat” investors
- **Backward** case $t < s$ is more surprising. Not a tradable strategy, but interesting for research
 - ▶ Suggests backwards OOS tests (e.g., Linnainmaa and Roberts 2018) and cross-validation (e.g., Kozak, Nagel and Santosh 2020; Bryzgalova, Pelger, and Zhu 2020) could be appropriate for Bayesian learning setting

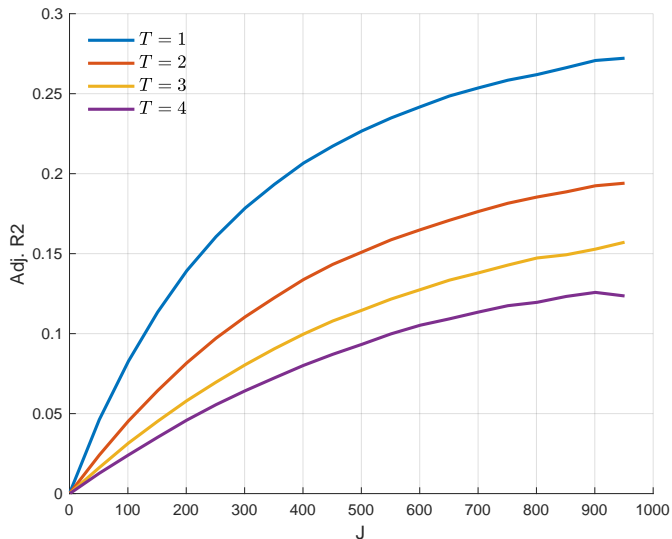
Finite-sample analysis: Simulations

- Simulate cash-flows, prices, returns for $N = 1000$ assets
- To generate data, we set $\theta = 1$ in

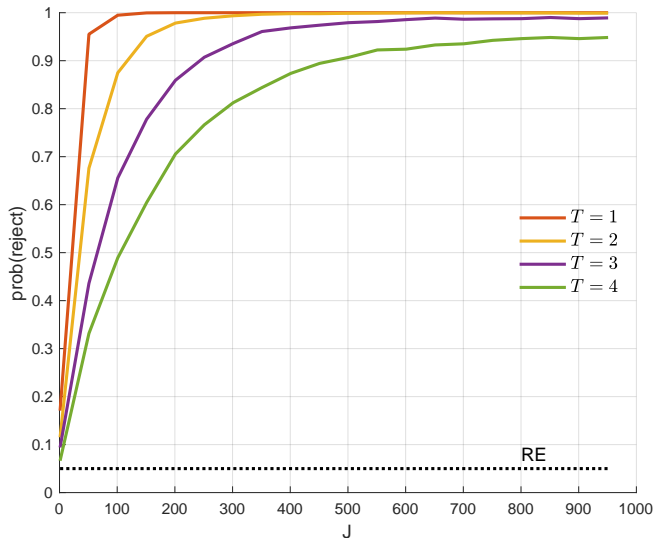
$$\Sigma_g = \frac{\theta}{J} \mathbf{I}$$

- $\theta =$ ratio of forecastable/residual cash-flow growth variance
 - ▶ Based on analyst expectations, Chen, Karceski, and Lakonishok (2003) find forecastable/residual cash-flow growth variance of 0.4 at 10yr horizon
- Econometrician regresses r_{T+1} on \mathbf{X} after investors have learned about \mathbf{g} for T periods

Adjusted R^2



Rejection probability of no-return-predictability null



Variations in the paper

- So far, shrinkage was purely due to objectively correct informative prior beliefs of investors
- If (time-varying?) cost to observe predictor variables, this may induce excess shrinkage \implies positive OOS returns
- Similar results when investors deal with big data by using Lasso rather than ridge regressions

Empirical illustration: IS vs OOS predictability

- Suppose returns from earlier augmented with risk premium/mispricing component $\mathbf{X}\boldsymbol{\gamma}$

$$\mathbf{r}_{t+1} = \mathbf{X}\boldsymbol{\gamma} + \mathbf{X}(\mathbf{I} - \boldsymbol{\Gamma}_t)\mathbf{g} - \mathbf{X}\boldsymbol{\Gamma}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{e}}_t + \mathbf{e}_{t+1}$$

- OOS returns measure importance of risk premium/mispricing:

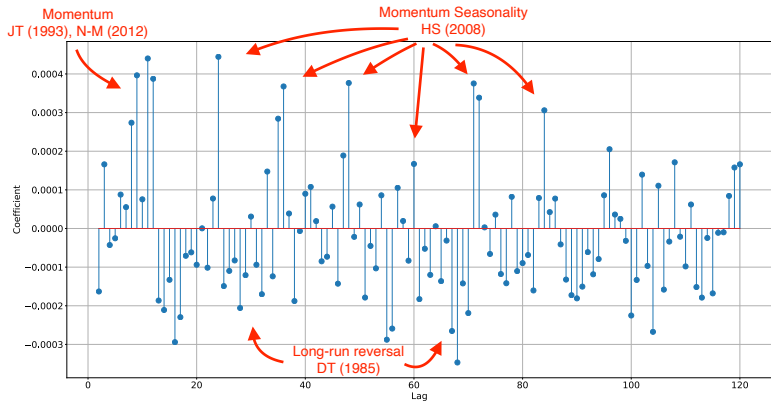
$$\boldsymbol{\gamma}'\mathbf{X}'\mathbf{X}\boldsymbol{\gamma} = \mathbb{E}[r_{OOS,t+1}]$$

Empirical illustration: IS vs OOS predictability

- Use past returns of each stock (available, in principle, for decades) to predict returns in month t with
 - ▶ Returns in months $t - 2, \dots, t - 120$
 - ▶ Squared returns in months $t - 2, \dots, t - 120$
- All U.S. stocks on CRSP, except market cap $<$ 20th NYSE percentile or price $<$ \$1 at the end of month $t - 1$
- All predictors cross-sectionally demeaned and standardized to unit S.D. each month
- Ridge regression with leave-one-year-out cross-validation to choose penalty parameter value

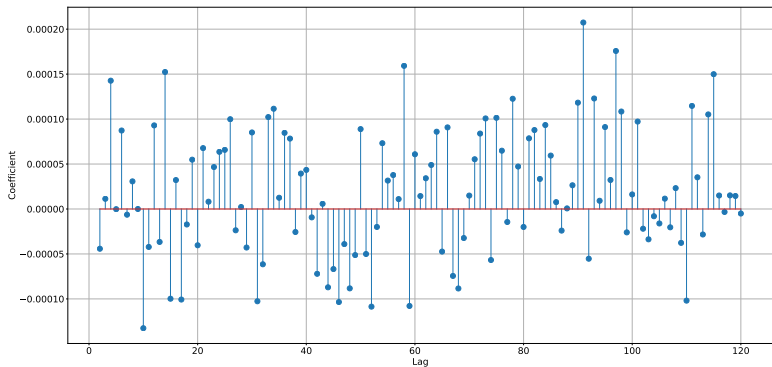
In-sample: Past return coefficients

Sample period: 1971-2018



In-sample: Past squared return coefficients

Sample period: 1971-2018



Estimating risk premia/mispricing in presence of learning

- Recall: Estimate $\gamma'X'X\gamma$ from sample version of

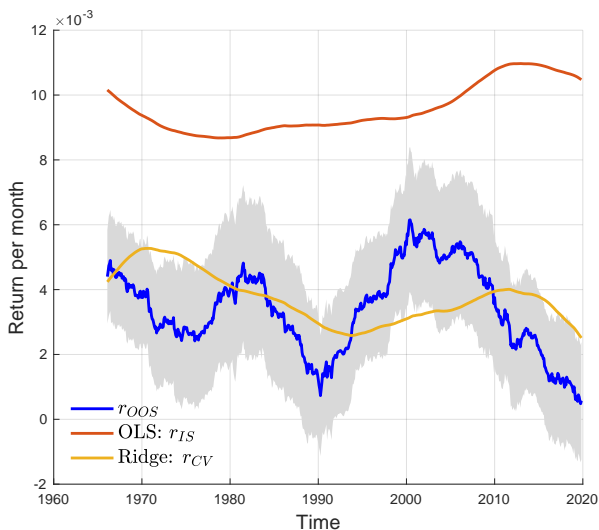
$$\gamma'X'X\gamma = \mathbb{E}[r_{OOS,t+1}]$$

- h_t estimated with OLS in backwards 20-year moving window up to month t and used to form

$$r_{OOS,t+1} = r'_{t+1}Xh_t$$

- $r_{OOS,t+1}$ averaged in 10-year moving windows
- Compare with two other returns in backwards 20-year window
 - ▶ In-sample return based on OLS estimates r_{IS}
 - ▶ Return on validation folds for cross-validated ridge regression r_{CV}

Estimating risk premia/mispricing in presence of learning



IS vs OOS returns

- In an RE model, expected IS and OOS portfolio returns would both equal $\gamma'X'X\gamma$
- If investors learn, this is still true for the OOS portfolio return
- But the IS return is distorted by learning-induced components that are not predictable OOS
- Seems that the learning case is relevant
- IS predictability does not carry over to OOS predictability and hence does not reflect risk premia demanded by investors ex ante, or persistent belief distortions

Implications: Market Efficiency in the Age of Big Data

- In Big Data setting, RE (investors know g) is implausible
- Learning (about g) has strong effects on asset prices
- Risk premia & bias theories should focus on explaining OOS, not IS, return predictability
- Investor learning provides clear motivation for (pseudo-)OOS testing which is lacking in RE framework