

# Machine Learning Panel Data Regressions with an Application to Nowcasting Price Earnings Ratios

Andrii Babii\*    Ryan T. Ball†    Eric Ghysels‡    Jonas Striaukas§

August 6, 2020

## Abstract

This paper introduces structured machine learning regressions for prediction and nowcasting with panel data consisting of series sampled at different frequencies. Motivated by the empirical problem of predicting corporate earnings for a large cross-section of firms with macroeconomic, financial, and news time series sampled at different frequencies, we focus on the sparse-group LASSO regularization. This type of regularization can take advantage of the mixed frequency time series panel data structures and we find that it empirically outperforms the unstructured machine learning methods. We obtain oracle inequalities for the pooled and fixed effects sparse-group LASSO panel data estimators recognizing that financial and economic data exhibit heavier than Gaussian tails. To that end, we leverage on a novel Fuk-Nagaev concentration inequality for panel data consisting of heavy-tailed  $\tau$ -mixing processes which may be of independent interest in other high-dimensional panel data settings.

*Keywords:* corporate earnings, nowcasting, high-dimensional panels, mixed frequency data, text data, sparse-group LASSO, heavy-tailed  $\tau$ -mixing processes, Fuk-Nagaev inequality.

---

\*University of North Carolina at Chapel Hill - Gardner Hall, CB 3305 Chapel Hill, NC 27599-3305. Email: babii.andrii@gmail.com

†Stephen M. Ross School of Business, University of Michigan, 701 Tappan Street, Ann Arbor, MI 48109. Email: rtball@umich.edu

‡Department of Economics and Kenan-Flagler Business School, University of North Carolina-Chapel Hill. Email: eghysels@unc.edu.

§LIDAM UC Louvain and FRS-FNRS Research Fellow. Email: jonas.striaukas@gmail.com.

# 1 Introduction

The fundamental value of equity shares is determined by the discounted value of future payoffs. Every quarter investors get a glimpse of a firms' potential payoffs with the release of corporate earnings reports. In a data-rich environment, stock analysts have many indicators regarding future cash flows that are available much more frequently. [Ball and Ghysels \(2018\)](#) took a first stab at automating the process using MIDAS regressions. Since their original work, much progress has been made on machine learning regularized mixed frequency regression models. In the current paper, we significantly expand the tools of nowcasting in a data-rich environment by exploiting panel data structures. Panel data regression models are well suited for the firm-level data analysis as both time series and cross-section dimensions can be properly modeled. In such models, time-invariant firm-specific effects are typically modeled in a flexible way which allows capturing heterogeneity in the data. At the same time, machine learning methods are becoming increasingly popular in economics and finance as a flexible way to model relationships between the response and covariates.

In the present paper, we analyze panel data regressions in a high-dimensional setting where the number of time-varying covariates can be very large and potentially exceed the sample size. This may happen when the number of firm-specific characteristics, such as textual analysis news data or firm-level stock returns, is large, and/or the number of aggregates, such as market returns, macro data, etc., is large. In our theoretical treatment, we obtain oracle inequalities for pooled and fixed effects LASSO-type panel data estimators allowing for heavy-tailed  $\tau$ -mixing data. To recognize time series data structures, we rely on a more general sparse-group LASSO (sg-LASSO) regularization with dictionaries, which typically improves upon the unstructured LASSO estimator for time series data in small samples.<sup>1</sup> Importantly, our theory covers LASSO and group-LASSO estimators as special cases.

To recognize that the economic and financial data often have heavier than Gaussian tails, our theoretical treatment relies on a new Fuk-Nagaev panel data concentration inequality. This allows us to characterize the dependence of the performance of LASSO-type estimators on  $N$  (cross-section) and  $T$  (time series), which is especially relevant for modern panel data applications, where both  $N$  and  $T$  can be large; see [Fernández-Val and Weidner \(2018\)](#) for a recent literature review focusing on the low-dimensional panel data case.

Our paper is related to the recent work of [Fosten and Greenaway-McGrevy \(2019\)](#)

---

<sup>1</sup>See [Babii, Ghysels, and Striaukas \(2020b\)](#) for an application of sg-LASSO to ADL-MIDAS model and US GDP nowcasting.

and Khalaf, Kichian, Saunders, and Voia (2020) who focus on nowcasting and/or mixed frequency panel data. In contrast to Khalaf, Kichian, Saunders, and Voia (2020), we explore the time dimension and introduce the LASSO-type regularization. To the best of our knowledge, the theoretical treatment of the sparse-group LASSO regularization for the panel data is not currently available in the literature.<sup>2</sup>

An empirical application to nowcasting firm-specific price/earnings ratios (henceforth P/E ratio) is provided. We focus on the current quarter nowcasts, hence evaluating model-based within quarter predictions for very short horizons. It is widely acknowledged that P/E ratios are a good indicator of the future performance of a particular company and therefore used by analysts and investment professionals to base their decisions on which stocks to pick for their investment portfolios. A typical value investor relies on consensus forecasts of earnings made by a pool of analysts. Hence, we naturally benchmark our proposed machine learning methods against such predictions. Besides, we compare our methods with a forecast combination approach used by Ball and Ghysels (2018) and a simple random walk (RW).

Our high-frequency regressors include traditional macro and financial series as well as non-standard series generated by the textual analysis. We consider structured pooled and fixed effects sg-LASSO panel data regressions with mixed frequency data (sg-LASSO-MIDAS). The fixed effects estimator yields sparser models compared to pooled regressions with the Revenue growth and the first lag of the dependent variable are selected throughout the out-of-sample period. BAA less AAA bond yield spread, firm-level volatility, and news textual analysis Aggregate Event Sentiment index are also selected very frequently. Our results show the superior performance of sg-LASSO-MIDAS over analysts' predictions, forecast combination method, and firm-specific time series regression models. Besides, the sg-LASSO-MIDAS regressions perform better than unstructured panel data regressions with the elastic net regularization.

Regarding the textual news data, it is worth emphasizing that the time series of news data is sparse since for many days are without firms-specific news and we impute zero values. The nice property of our mixed frequency data treatment with dictionaries, imputing zeros also implies that non-zero entries get weights with a decaying pattern for distant past values in comparison to the most recent daily news data. As a result, our ML approach is particularly useful to model news data which

---

<sup>2</sup>The panel data regressions with the LASSO penalty is used in microeconometrics since Koenker (2004); see also Lamarche (2010), Kock (2013), Belloni, Chernozhukov, Hansen, and Kozbur (2016), Lu and Su (2016), Kock (2016), Harding and Lamarche (2019), Chiang, Rodrigue, and Sasaki (2019) and Chernozhukov, Hausman, and Newey (2019) among others. The group LASSO is considered in Su, Shi, and Phillips (2016), Lu and Su (2016), and Farrell (2015) among others.

is sparse in nature.

The paper is organized as follows. Section 2 introduces the models and estimators. Oracle inequalities for sparse group LASSO panel data regressions appear in Section 3. Section 4 covers Fuk-Nagaev inequalities for panel data. Results of our empirical application analyzing price earnings ratios for a panel of individual firms are reported in Section 5. Technical material appears in the Appendix and conclusions appear in Section 6.

## 2 Methodology

In this section we describe briefly the methodological approach, while in Section 3 we provide more details and the supporting theoretical results. We focus on the pooled and the fixed effects panel regressions with the sparse-group LASSO (sg-LASSO) regularization. The best linear predictor for firm  $i = 1, \dots, N$  in the panel data setting is

$$\alpha_i + x_{it}^\top \beta,$$

where  $\alpha_i, i = 1, \dots, N$  are fixed intercepts. We consider predictive regressions with homogeneous and heterogeneous entity-specific intercepts.

### 2.1 Pooled sg-LASSO

In the pooled regressions, we ignore the cross-sectional heterogeneity, assuming that  $\alpha_i = \alpha$  for all  $i = 1, \dots, N$  and the pooled sg-LASSO estimator is a solution to

$$\min_{(a,b) \in \mathbf{R} \times \mathbf{R}^p} \|\mathbf{y} - a\mathbf{1} - \mathbf{X}b\|_{NT}^2 + 2\lambda\Omega(b),$$

where

$$\Omega(b) = \gamma|b|_1 + (1 - \gamma)\|b\|_{2,1}$$

is a penalty function and  $\gamma \in [0, 1]$  is the relative weight of the LASSO and the group LASSO penalties.

Intuitively, the low- or high-frequency lags of a single covariate define a group which might be a dense signal provided that this covariate is relevant for prediction. The dense signals are not well-captured by the unstructured LASSO estimator of Tibshirani (1996). Indeed, the lags of a single covariate are temporally related, hence, taking this group structure into account might improve upon the predictive performance of the unstructured LASSO estimator in small samples; see Section 2.3 for more details how the dense time series signal is mapped into the sparse-group structure with dictionaries.

## 2.2 Fixed effects sg-LASSO

In contrast to the pooled regressions, in the fixed effects regressions we estimate the heterogeneous slope parameters  $\alpha_i, i = 1, \dots, N$ , and use them subsequently to construct the best linear predictors. The fixed effects sg-LASSO estimator is a solution to

$$\min_{(a,b) \in \mathbf{R}^{N+p}} \|\mathbf{y} - Ba - \mathbf{X}b\|_{NT}^2 + 2\lambda\Omega(b),$$

where  $B = I_N \otimes \iota$  and  $\Omega$  is the sparse-group LASSO penalty. Note that we consider the fixed effects as a dense signal and leave them unpenalized.<sup>3</sup> The sparse group structure is defined by low- and high-frequency lags similarly to the pooled regressions as explained in the following subsection.

## 2.3 Mixed Frequency Data Panels

Motivated by our empirical application, we allow the high-dimensional set of covariates to be sampled at a higher frequency than the outcome variable. Let  $K$  be the total number of time-varying covariates  $\{x_{i,t-j/m,k}, i \in [N], t \in [T], j \in [m], k \in [K]\}$  possibly measured at some higher frequency with  $m$  observations for every low-frequency period  $t$  and consider the following MIDAS panel data regression

$$y_{it} = \alpha_i + \sum_{k=1}^K \psi(L^{1/m}; \beta_k) x_{it,k} + u_{it},$$

where  $\psi(L^{1/m}; \beta_k) x_{it,k} = 1/m \sum_{j=1}^m \beta_{j,k} x_{i,t-j/m,k}$  is the high-frequency lag polynomial. For  $m = 1$ , we retain the standard panel data regression model, while  $m > 1$  signifies that the high-frequency lags of  $x_{i,t,k}$  are also included. For large  $m$ , there is a proliferation of the total number of the estimated parameters which reduces the finite-sample predictive performance.

The MIDAS literature offers various parametrization of weights, see Ghysels, Santa-Clara, and Valkanov (2006); Ghysels, Sinko, and Valkanov (2006). More recently, Babii, Ghysels, and Striaukas (2020b) proposed a new approach based on dictionaries linear in parameters to approximate the MIDAS weight function which is particularly useful for high-dimensional MIDAS regression models. The sparse-group LASSO allows for the data-driven approximation to the MIDAS weight functions from the dictionary promoting sparsity between groups (covariate selection) and within groups (MIDAS weight approximation).

---

<sup>3</sup>The pooled and the fixed effects estimators can be efficiently computed using a variant of coordinate descent algorithm proposed by Simon, Friedman, Hastie, and Tibshirani (2013).

More concretely, we focus on a class of weights well-approximated by some collection of functions  $w^L(s) = (w_1(s), \dots, w_L(s))^\top$ , called the *dictionary*,

$$\omega_L(s; \beta_k) = (w^L(s))^\top \beta_k.$$

Instead of estimating a large number of parameters pertaining to the high-frequency lag polynomial  $\psi(L^{1/m}; \beta_k)x_{it,k} = 1/m \sum_{j=1}^m \beta_{j,k}x_{i,t-j/m,k}$ , we estimate a lower-dimensional parameter  $\beta_k$  in

$$\frac{1}{m} \sum_{j=1}^m \omega_L(j/m; \beta_k)x_{i,t-j/m,k}.$$

The attractive feature of the sparse-group LASSO estimator is that it can learn the MIDAS weight function from the dictionary in a nonlinear data-driven way and, at the same time, selects covariates defined as groups of time series lags. In practice, dictionaries lead to the design matrix  $\mathbf{X}$  structured appropriately; see [Babii, Ghysels, and Striaukas \(2020b\)](#) for more details on how to construct the design matrix.

Note that such weights depend linearly on the parameter  $\beta_k$  which allows for the efficient estimation of the high-dimensional MIDAS panel regression model, cf., [Khalaf, Kichian, Saunders, and Voia \(2020\)](#) for the low-dimensional non-linear case. A suitable dictionary for our purposes are Legendre polynomials, which are in the class of orthogonal polynomials and have very good approximating properties.<sup>4</sup> In practice, orthogonal polynomials typically outperform non-orthogonal counterparts, e.g. Almon polynomials, or unrestricted lags in small samples; see [Babii, Ghysels, and Striaukas \(2020b\)](#) for further details and Monte Carlo simulation study supporting this choice.

## 2.4 Tuning parameter

We consider several approaches to select the tuning parameter  $\lambda$ . First, we adapt the  $k$ -fold cross-validation to the panel data setting. To that end, we resample the data by blocks respecting the time-series dimension and creating folds based on individual firms instead of the pooled sample. We use 5-fold cross-validation as the sample size of the dataset we consider in our empirical application is relatively small. We also consider the following three information criteria: BIC, AIC, and corrected AIC (AICc). Assuming that  $y_{it}|x_{it}$  are i.i.d. draws from  $N(\alpha_i + x_{it}^\top \beta, \sigma^2)$ ,

---

<sup>4</sup>More precisely, we can approximate any continuous weight function in the  $L_\infty[0, 1]$  norm, and more generally, any square-integrable function in the  $L_2[0, 1]$  norm, hence, the discontinuous MIDAS weights are not ruled-out.

the log-likelihood of the sample is

$$\mathcal{L}(\alpha, \beta, \sigma^2) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{t=1}^T (y_{it} - \alpha_i - x_{it}^\top \beta)^2.$$

Then, for the pooled model, the BIC criterion is

$$\text{BIC} = \frac{\|\mathbf{y} - \hat{\alpha}\nu - \mathbf{X}\hat{\beta}\|_{NT}^2}{\hat{\sigma}^2} + \frac{\log(NT)}{NT} \times \hat{df},$$

where  $df$  denotes the degrees of freedom. The degrees of freedom are estimated as  $\hat{df} = |\hat{\beta}|_0 + 1$  for the pooled regression and  $\hat{df} = |\hat{\beta}|_0 + N$  for the fixed effects regression, where  $|\cdot|_0$  is the  $\ell_0$ -norm defined as a number of non-zero coefficients; see [Zou, Hastie, and Tibshirani \(2007\)](#) for more details. The AIC is computed as

$$\text{AIC} = \frac{\|\mathbf{y} - \hat{\alpha}\nu - \mathbf{X}\hat{\beta}\|_{NT}^2}{\hat{\sigma}^2} + \frac{2}{NT} \times \hat{df}.$$

Lastly, the corrected Akaike information criteria is

$$\text{AICc} = \frac{\|\mathbf{y} - \hat{\alpha}\nu - \mathbf{X}\hat{\beta}\|_{NT}^2}{\hat{\sigma}^2} + \frac{2\hat{df}}{NT - \hat{df} - 1}.$$

The AICc might be a better choice when  $p$  is large compared to the sample size. For the fixed effects regressions, we replace  $\hat{\alpha}\nu$  with  $B\hat{\alpha}$  everywhere and adjust the degrees of freedom as described above. We report results for each of these four choices of the tuning parameters.

### 3 Oracle inequalities for sg-LASSO panel data regressions

In this section, we provide the theoretical analysis of the predictive performance of pooled and fixed effects panel data regressions with the sg-LASSO regularization, including the standard LASSO and the group LASSO regularizations. It is worth stressing that our analysis is not tied to the mixed-frequency data setting and applies to generic high-dimensional panel data regularized with the sg-LASSO penalty function. Importantly, we focus on panels consisting of  $\tau$ -mixing time series with polynomial (Pareto-type) tails.

### 3.1 Pooled regression

The pooled linear projection model is

$$y_{it} = \alpha + x_{it}^\top \beta + u_{it}, \quad \mathbb{E}[u_{it} z_{it}] = 0, \quad i \in [N], t \in [T],$$

where  $\alpha \in \mathbf{R}$  and  $\beta \in \mathbf{R}^p$  are unknown projection coefficients,  $z_{it} = (1, x_{it}^\top)^\top$ , and we use  $[J]$  to denote the set  $\{1, 2, \dots, J\}$  for arbitrary positive integer  $J$ . The vector of covariates  $x_{it} \in \mathbf{R}^p$  may include the time-varying covariates common for all entities (macroeconomic factors) as well as lags of  $y_{it}$  and lags of some baseline covariates. It is worth stressing that pooled regressions can also potentially accommodate heterogeneity provided that the data are clustered in a relatively small number of clusters of similar entities.

Put  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})^\top$ ,  $\mathbf{x}_i = (x_{i1}, \dots, x_{iT})^\top$ ,  $\mathbf{u}_i = (u_{i1}, \dots, u_{iT})^\top$ , and let  $\iota \in \mathbf{R}^T$  be a vector of ones. Then the regression equation after stacking time series observations is

$$\mathbf{y}_i = \alpha \iota + \mathbf{x}_i \beta + \mathbf{u}_i, \quad i \in [N].$$

Define further  $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_N^\top)^\top$ ,  $\mathbf{X} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_N^\top)^\top$ , and  $\mathbf{u} = (\mathbf{u}_1^\top, \dots, \mathbf{u}_N^\top)^\top$ . Then the regression equation after stacking all cross-sectional observations is

$$\mathbf{y} = \alpha \iota + \mathbf{X} \beta + \mathbf{u},$$

where  $\iota \in \mathbf{R}^{NT}$  is a vector of ones.

The pooled sg-LASSO estimator  $\hat{\beta}$  solves

$$\min_{(a, b) \in \mathbf{R}^{1+p}} \|\mathbf{y} - a\iota - \mathbf{X}b\|_{NT}^2 + 2\lambda\Omega(b), \quad (1)$$

where  $\|z\|_{NT}^2 = z^\top z / NT$  for  $z \in \mathbf{R}^{NT}$ , and

$$\Omega(b) = \gamma|b|_1 + (1 - \gamma)\|b\|_{2,1}$$

is the sg-LASSO penalty function. The penalty function  $\Omega$  interpolates between the LASSO penalty  $|b|_1 = \sum_{j=1}^p |b_j|$  and the group LASSO penalty  $\|b\|_{2,1} = \sum_{G \in \mathcal{G}} |b_G|_2$ , where  $\mathcal{G}$  is a partition of  $[p] = \{1, 2, \dots, p\}$  and  $|b_G|_2 = (\sum_{j \in G} |b_j|^2)^{1/2}$  is the  $\ell_2$  norm. The parameter  $\gamma \in [0, 1]$  determines the relative weights of the LASSO and the group LASSO penalization, while the amount of regularization is controlled by the regularization parameter  $\lambda > 0$ . Note that the group structure  $\mathcal{G}$  has to be specified by the econometrician, which in our setting is defined by the high-frequency lags of



different covariates. Throughout the paper we assume that groups have fixed size, which is well-justified in our empirical applications.<sup>5</sup>

For a random variable  $\xi$ , let  $\|\xi\|_q = (\mathbb{E}|\xi|^q)^{1/q}$  denote its  $L_q$ ,  $q \geq 1$  norm. Following Babii, Ghysels, and Striaukas (2020a) and Babii, Ghysels, and Striaukas (2020b), we consider  $\tau$ -mixing processes, measuring the temporal dependence with  $\tau$ -mixing coefficients. The  $\tau$ -mixing processes can be placed somewhere between the  $\alpha$ -mixing processes and mixingales – they are less restrictive than  $\alpha$ -mixing, yet at the same time are amenable to coupling similarly to  $\alpha$ -mixing processes, which is not the case for the mixingales; see Dedecker and Doukhan (2003), Dedecker and Prieur (2004), and Dedecker and Prieur (2005) for more details. This allows us to obtain sharp concentration inequalities in Section 4.

For a  $\sigma$ -algebra  $\mathcal{M}$  and a random vector  $\xi \in \mathbf{R}^l$ , the coupling  $\tau$  coefficient is defined as

$$\tau(\mathcal{M}, \xi) = \sup_{f \in \text{Lip}_1} \int_{\mathbf{R}} \|F_{f(\xi)|\mathcal{M}}(x) - F_{f(\xi)}(x)\|_1 dx,$$

where  $\text{Lip}_1$  is a set of 1-Lipschitz functions from  $\mathbf{R}^l$  to  $\mathbf{R}$ ,  $F_\zeta$  is the CDF of  $\zeta = f(\xi)$  and  $F_{\zeta|\mathcal{M}}$  is the CDF of  $\zeta$  conditionally on  $\mathcal{M}$ .<sup>6</sup> For a stochastic process  $(\xi_t)_{t \in \mathbf{Z}}$  with a natural filtration generated by its past  $\mathcal{M}_t = \sigma(\xi_t, \xi_{t-1}, \dots)$ , the  $\tau$ -mixing coefficients are defined as

$$\tau_k = \sup_{j \geq 1} \max_{l \in [j]} \frac{1}{l} \sup_{t+k \leq t_1 < \dots < t_l} \tau(\mathcal{M}_t, (\xi_{t_1}, \dots, \xi_{t_l})), \quad k \geq 0.$$

We say that the process is  $\tau$ -mixing if its  $\tau$ -mixing coefficients are decreasing to zero.

The following assumption imposes mild restrictions on the data.

**Assumption 3.1** (Data). *(i) for each  $t \in \mathbf{Z}$ ,  $\{u_{it}z_{it} : i \geq 1\}$  are i.i.d. and for each  $i \geq 1$ ,  $\{u_{it}z_{it} : t \in \mathbf{Z}\}$  is a stationary process; (ii)  $\max_{j \in [p]} \|u_{it}z_{it,j}\|_q = O(1)$  for some  $q > 2$ ; (iii) for every  $j \in [p]$ ,  $\tau$ -mixing coefficients of  $\{u_{it}z_{it,j} : t \in \mathbf{Z}\}$  satisfy  $\tau_k \leq ck^{-a}$ ,  $\forall k \geq 1$  with some universal constants  $c > 0$  and  $a > (q-1)/(q-2)$ ; (iii)  $\max_{j,k \in [p]} \|z_{it,j}z_{it,k}\|_{\tilde{q}} = O(1)$  for some  $\tilde{q} > 2$ ; (iv) for every  $j, k \in [p]$ ,  $\tau$ -mixing coefficients of  $\{z_{it,j}z_{it,k} : t \in \mathbf{Z}\}$  satisfy  $\tilde{\tau}_k \leq \tilde{c}k^{-\tilde{a}}$ ,  $\forall k \geq 1$  for some universal constants  $\tilde{c} > 0$  and  $\tilde{a} > (\tilde{q}-1)/(\tilde{q}-2)$ .*

It is worth mentioning that the stationarity hypothesis can be relaxed at costs of introducing heavier notation. We require that only  $2 + \epsilon$  moments exist with  $\epsilon > 0$ ,

<sup>5</sup>See Babii (2020) for a continuous-time mixed-frequency regression where the group size is allowed to increase with the sample size.

<sup>6</sup>See Dedecker and Prieur (2004) and Dedecker and Prieur (2005) for equivalent definitions.

which is a realistic assumption in our empirical application and more generally for datasets encountered in time series and financial econometrics applications. Note also that the temporal dependence is assumed to fade away relatively slow – at a polynomial rate as measured by the  $\tau$ -mixing coefficients.

Next, we assume that the matrix  $\Sigma = \mathbb{E}[z_{it}z_{it}^\top]$  is non-singular uniformly over  $p$ .

**Assumption 3.2** (Covariance matrix). *The smallest eigenvalue of  $\Sigma$  is uniformly bounded away from zero by some universal constant.*

Assumption 3.2 can also be relaxed to the restricted eigenvalue condition imposed on the population covariance matrix  $\Sigma$ ; see also Babii, Ghysels, and Striaukas (2020b).

Lastly, we assume that the regularization parameter  $\lambda$  scales appropriately with the number of covariates  $p$ , the length of the panel  $T$ , and the size of the cross-section  $N$ . The precise order of the regularization parameter is described by the Fuk-Nagaev inequality for long panels appearing in Theorem 4.1 of the next section. In what follows, we say that  $a \sim b$  if and only if  $\underline{c}b \leq a \leq \bar{c}b$  for some appropriately defined constants  $\underline{c}, \bar{c} > 0$ .

**Assumption 3.3** (Regularization). *The regularization parameter satisfies*

$$\lambda \sim \left( \frac{p}{\delta(NT)^{\kappa-1}} \right)^{1/\kappa} \vee \sqrt{\frac{\log(p/\delta)}{NT}}$$

for some  $\delta \in (0, 1)$  and  $\kappa = \frac{(a+1)q-1}{a+q-1}$ , where  $a, q$  are as in Assumptions 3.1.

Assumption 3.3 describes the theoretically optimal level of the regularization parameter.<sup>7</sup>

Our first result is the oracle inequality for the pooled sg-LASSO estimator described in Eq. 1. The result allows for misspecified regressions with a non-trivial approximation error in the sense that we consider more generally

$$\mathbf{y} = \mathbf{m} + \mathbf{u},$$

where  $\mathbf{m} \in \mathbf{R}^{NT}$  is approximated with  $\mathbf{Z}\rho$ ,  $\mathbf{Z} = (\iota, \mathbf{X})$ , and  $\rho = (\alpha, \beta^\top)^\top$ . The approximation error  $\mathbf{m} - \mathbf{Z}\rho$  might come from the fact the MIDAS weight function may not have the exact expansion in terms of the specified dictionary or from the

---

<sup>7</sup>An interesting challenging question is how the data-driven choices of the tuning parameter affect the performance of the LASSO-type estimators; see Chetverikov, Liao, and Chernozhukov (2020) for an example of such analysis in the case of cross-validation with i.i.d. sub-Gaussian data.

fact that some of the relevant covariates are not included in the regression equation. To state the result, let  $S_0 = \{j \in [p] : \beta_j \neq 0\}$  be the support of  $\beta$  and let  $\mathcal{G}_0 = \{G \in \mathcal{G} : \beta_G \neq 0\}$  be the group support of  $\beta$ . Consider the *effective sparsity* of the sparse-group structure, defined as  $s = (\gamma\sqrt{|S_0|} + (1-\gamma)\sqrt{|\mathcal{G}_0|})^2$ . Note that  $s$  simplifies to the sparsity of  $\beta$ ,  $|S_0|$ , when  $\gamma = 1$  and to the group sparsity  $|\mathcal{G}_0|$  when  $\gamma = 0$ .

**Theorem 3.1.** *Suppose that Assumptions 3.1, 3.2, and 3.3 are satisfied. Then with probability at least  $1 - \delta - O\left(\frac{s^{\tilde{\kappa}} p^2}{(NT)^{\tilde{\kappa}-1} + p^2 e^{-cNT/s^2}}\right)$*

$$\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 \lesssim s\lambda^2 + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT}^2$$

and

$$|\hat{\alpha} - \alpha| + |\hat{\beta} - \beta|_1 \lesssim s\lambda + \lambda^{-1} \|\mathbf{m} - \mathbf{Z}\beta\|_{NT}^2 + s^{1/2} \|\mathbf{m} - \mathbf{Z}\beta\|_{NT},$$

for some  $c > 0$  and  $\tilde{\kappa} = \frac{(\bar{a}+1)\bar{q}-1}{\bar{a}+\bar{q}-1}$ .

The proof of this result can be found in the Appendix. Theorem 3.1 applies to panel data unlike the result of Babii, Ghysels, and Striaukas (2020b). It provides the oracle inequalities describing the prediction and the estimation accuracy in the environment where the number of regressors  $p$  is allowed to scale with the effective sample size  $NT$ . Importantly, the result is stated under the weak tail and mixing conditions in Assumption 3.1. Parameters  $\kappa$  and  $\tilde{\kappa}$  are the mixing-tails exponents for stochastic processes driving the regression score and the covariance matrix respectively.

To describe convergence rates, the following condition considers a simplified setting, where the effective sparsity  $s$  is constant, the approximation error vanishes sufficiently fast, and the total number of regressors scales appropriately with the effective sample size  $NT$ .

**Assumption 3.4.** *Suppose that (i)  $s = O(1)$ ; (ii)  $\|\mathbf{m} - \mathbf{Z}\beta\|_{NT}^2 = O_P(\lambda^2)$ ; (iii)  $p^2 = o((NT)^{\tilde{\kappa}-1})$ .*

In particular, Assumption 3.4 allows for 1)  $N \rightarrow \infty$  while  $T$  is fixed; 2)  $T \rightarrow \infty$  while  $N$  is fixed; and 3) both  $N \rightarrow \infty$  and  $T \rightarrow \infty$  without restricting the relative growth of the two. The following result describes the prediction and the estimation convergence rates in the asymptotic environment outlined in Assumption 3.4 and is an immediate consequence of Theorem 3.1.

**Corollary 3.1.** *Suppose that Assumptions 3.1, 3.2, 3.3, and 3.4 are satisfied. Then*

$$\|\mathbf{Z}(\hat{\beta} - \beta)\|_{NT}^2 = O_P\left(\frac{p^{2/\kappa}}{(NT)^{2-2/\kappa}} \vee \frac{\log p}{NT}\right)$$

and

$$|\hat{\beta} - \beta|_1 = O_P \left( \frac{p^{1/\kappa}}{(NT)^{1-1/\kappa}} \vee \sqrt{\frac{\log p}{NT}} \right).$$

Note that for large  $a$ , the mixing-tails exponent is  $\kappa \approx q$ . Therefore, for the data that are close to independent, the prediction accuracy is approximately of order  $O_P \left( \frac{p^{2/q}}{(NT)^{2-2/q}} \vee \frac{\log p}{NT} \right)$ , which is the rate one would obtain for the i.i.d. data applying directly Fuk and Nagaev (1971), Corollary 4,<sup>8</sup> so in this sense our result is sharp. If the data are sub-Gaussian, then moments of all order  $q \geq 2$  exist and for any particular sample size  $NT$ , the first term can be made arbitrarily small relative to the second taking large enough  $q$ . In this case we recover the  $O_P \left( \frac{\log p}{NT} \right)$  rate typically obtained for sub-Gaussian data. Therefore, the Fuk-Nagaev inequality provides a more accurate description of the performance of the LASSO-type estimators.

If the polynomial tail dominates, then we need  $p = o((NT)^{\kappa-1})$  for the prediction and the estimation consistency provided that  $\tilde{\kappa} \geq 2\kappa - 1$ . The pooled sg-LASSO estimator is expected to work well whenever the number of regressors  $p$  is small relative to  $(NT)^{\kappa-1}$ . This is a *significantly weaker* requirement compared to  $p = o(T^{\kappa-1})$  needed for time series regressions in Babii, Ghysels, and Striaukas (2020b). In particular since  $\kappa > 2$ ,  $p = o((NT)^{\kappa-1})$  can be significantly weaker than  $p = o(NT)$  condition needed in the QMLE/GMM framework without regularization. How much the sg-LASSO improves upon the (unregularized) QMLE depends on the heaviness of tails and persistence of the underlying stochastic processes as measured by the mixing-tails exponent  $\kappa$ . In particular, for light tails and weakly persistent series, the mixing-tails exponent  $\kappa$  is large, offsetting the dependence on  $p$ .

Lastly, it is worth mentioning that the oracle inequality is driven by the time series with the heaviest tail and that it might be possible to obtain sharper results allowing for heterogeneous tails at the costs of introducing a heavier notation.

## 3.2 Fixed effects

Pooled regressions are attractive since the effective sample size  $NT$  can be huge, yet the heterogeneity of individual time series may be lost. If the underlying series have substantial heterogeneity over  $i \in [N]$ , then taking this into account might reduce the projection error and improve the predictive accuracy. At a very extreme side, the cross-sectional structure can be completely ignored and individual time-series regressions can be used for prediction. The fixed effects panel data regressions

---

<sup>8</sup>Recall that the Fuk-Nagaev inequality provides sharper description of concentration compared to the simple Markov's bound in conjunction with the Rosenthal's moment inequality.

strike a good balance between the two extremes controlling for heterogeneity with entity-specific intercepts. The linear projection model with fixed effects is

$$y_{it} = \alpha_i + x_{it}^\top \beta + u_{it}, \quad \mathbb{E}[u_{it} z_{it}] = 0, \quad i \in [N], t \in [T],$$

where  $z_{it} = (1, x_{it}^\top)^\top$ . Note that the entity-specific intercepts  $\alpha_i$  are deterministic constants and the projection model is always well-defined. The fixed effects have to be estimated to construct the best linear predictor  $\alpha_i + x_{it}^\top \beta$ .

The fixed effects sg-LASSO estimators  $\hat{\alpha}$  and  $\hat{\beta}$  solve

$$\min_{(a,b) \in \mathbf{R}^{N+p}} \|\mathbf{y} - Ba - \mathbf{X}b\|_{NT}^2 + 2\lambda\Omega(b),$$

where  $B = I_N \otimes \iota$ ,  $I_N$  is  $N \times N$  identity matrix,  $\iota \in \mathbf{R}^T$  is the vector with all coordinates equal to one, and  $\Omega$  is the sg-LASSO penalty. It is worth stressing that the design matrix  $\mathbf{X}$  does not include the intercept and that we do not penalize the fixed effects. This is done because the sparsity over the fixed effects does not hold even in the special case where all intercepts are equal. By Fermat's rule, the first-order conditions are

$$\begin{aligned} \hat{\alpha} &= (B^\top B)^{-1} B^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \\ 0 &= \mathbf{X}^\top M_B (\mathbf{X}\hat{\beta} - \mathbf{y}) / NT + \lambda z^* \end{aligned}$$

for some  $z^* \in \partial\Omega(\hat{\beta})$ , where  $b \mapsto \partial\Omega(b)$  is the subdifferential of  $\Omega$  and  $M_B = I - B(B^\top B)^{-1} B^\top$  is the orthogonal projection matrix. It is easy to see from the first-order conditions that the estimator of  $\hat{\beta}$  is equivalent to: 1) penalized GLS estimator for the first-differenced regression; 2) penalized OLS estimator for the regression written in the deviation from time means; and 3) penalized OLS estimator where the fixed effects are partialled-out. Thus, the equivalence between the three approaches is not affected by the penalization, cf., [Arellano \(2003\)](#) for low-dimensional panels.

For the fixed effects regression, we define

$$\hat{\Sigma} = \begin{pmatrix} \frac{1}{T} B^\top B & \frac{1}{\sqrt{NT}} B^\top \mathbf{X} \\ \frac{1}{\sqrt{NT}} \mathbf{X}^\top B & \frac{1}{NT} \mathbf{X}^\top \mathbf{X} \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} I_N & \frac{1}{\sqrt{NT}} \mathbb{E}[B^\top \mathbf{X}] \\ \frac{1}{\sqrt{NT}} \mathbb{E}[\mathbf{X}^\top B] & \mathbb{E}[x_{it} x_{it}^\top] \end{pmatrix}. \quad (2)$$

We will assume that the smallest eigenvalue of  $\Sigma$  is uniformly bounded away from zero by some constant. Note that if  $x_{it} \sim N(0, I_p)$ , then  $\Sigma$  is approximately equal to the identity matrix for large  $N$ .

The order of the regularization parameter is governed by the Fuk-Nagaev inequality for long panels in [Theorem 4.1](#), with the only difference that it has to take into account the fact that the fixed effects parameters are estimated.

**Assumption 3.5** (Regularization). *The regularization parameter satisfies*

$$\lambda \sim \left( \frac{p \vee N^{\kappa/2}}{\delta(NT)^{\kappa-1}} \right)^{1/\kappa} \vee \sqrt{\frac{\log(p \vee N/\delta)}{NT}}$$

for some  $\delta \in (0, 1)$  and  $\kappa = \frac{(a+1)q-1}{a+q-1}$ , where  $a, q$  are as in Assumptions 3.1.

Similarly to the pooled regressions, we state the oracle inequality allowing for the approximation error. For fixed effects regressions we redefine  $\mathbf{Z} = (B, \mathbf{X})$ ,  $\rho = (\alpha, \beta^\top)^\top$ . Put also  $r_{N,T,p} = p(s \vee N)^{\tilde{\kappa}} T^{1-\tilde{\kappa}} (N^{1-\tilde{\kappa}/2} + pN^{1-\tilde{\kappa}}) + p(p \vee N)e^{-cNT/(s \vee N)^2}$  with  $\tilde{\kappa} = \frac{(\tilde{a}+1)\tilde{q}-1}{\tilde{a}+\tilde{q}-1}$  and some  $c > 0$ . Recall also that  $\Sigma$  in Assumption 3.2 is redefined according to Eq. 2, so that  $\Sigma$  is non-singular uniformly over  $p, N, T$ .

**Theorem 3.2.** *Suppose that Assumptions 3.1, 3.2, and 3.5 are satisfied. Then with probability at least  $1 - \delta - O(r_{N,T,p})$*

$$\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 \lesssim (s \vee N)\lambda^2 + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT}^2.$$

Theorem 3.2 states the oracle inequalities for the prediction error in the fixed effects panel data regressions estimated with the sg-LASSO. To see clearly, how the prediction accuracy scales with the sample size, we make the following assumption.

**Assumption 3.6.** *Suppose that (i)  $s = O(1)$ ; (ii)  $\|\mathbf{m} - \mathbf{Z}\beta\|_{NT}^2 = O_P(N\lambda^2)$ ; (iii)  $(p + N^{\tilde{\kappa}/2})pN/T^{\tilde{\kappa}-1} = o(1)$  and  $p(p \vee N)e^{-cT/N} = o(1)$ .*

The following corollary is an immediate consequence of Theorem 3.2.

**Corollary 3.2.** *Suppose that Assumptions 3.1, 3.2, 3.5, and 3.6 are satisfied. Then*

$$\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 = O_P \left( \frac{p^{2/\kappa} \vee N}{N^{1-2/\kappa} T^{2-2/\kappa}} \vee \frac{\log(p \vee N)}{T} \right).$$

Note that this result allows for  $p, N, T \rightarrow \infty$  at appropriate rates and that we pay additional price for estimating  $N$  fixed effects which plays a similar role to the effective dimension of covariates. Therefore, in order to achieve accurate prediction, the panel has to be sufficiently long to offset the estimation error of the individual fixed effects.

## 4 Fuk-Nagaev inequality for panel data

In this section we obtain new Fuk-Nagaev concentration inequality for panel data reflecting the concentration jointly over  $N$  and  $T$ . It is worth stressing that the inequality does not follow directly from the Fuk-Nagaev inequality of [Babii, Ghysels, and Striaukas \(2020a\)](#) and is of independent interest for the high-dimensional panel data.<sup>9</sup>

**Theorem 4.1.** *Let  $\{\xi_{it} : i \geq 1, t \in \mathbf{Z}\}$  be an array of centered random vectors in  $\mathbf{R}^p$  such that  $\{\xi_{i1}, \dots, \xi_{iT} : i \geq 1\}$  are i.i.d. and for each  $i \geq 1$ ,  $(\xi_{it})_{t \in \mathbf{Z}}$  is a stationary stochastic process such that (i)  $\max_{j \in [p]} \|\xi_{it,j}\|_q = O(1)$  for some  $q > 2$ ; (ii) for every  $j \in [p]$ ,  $\tau$ -mixing coefficients of  $(\xi_{it,j})_{t \in \mathbf{Z}}$  satisfy  $\tau_k^{(j)} \leq ck^{-a}, \forall k \geq 1$  for some universal constants  $c > 0$  and  $a > \frac{q-1}{q-2}$ . Then for every  $u > 0$*

$$\Pr \left( \left| \sum_{i=1}^N \sum_{t=1}^T \xi_{it} \right|_{\infty} > u \right) \leq c_1 p N T u^{-\kappa} + 4 p e^{-c_2 u^2 / N T}$$

for some  $c_1, c_2 > 0$  and  $\kappa = \frac{(a+1)q-1}{a+q-1}$ .

It follows from Theorem 4.1 that there exists  $C > 0$  such that for every  $\delta \in (0, 1)$

$$\Pr \left( \left| \frac{1}{N T} \sum_{t=1}^T \sum_{i=1}^N \xi_{it} \right|_{\infty} \leq C \left( \frac{p}{\delta (N T)^{\kappa-1}} \right)^{1/\kappa} \vee \sqrt{\frac{\log(8p/\delta)}{N T}} \right) \geq 1 - \delta.$$

Note that the inequality reflects the concentration jointly over  $N$  and  $T$  and that tails and persistence play an important role through the mixing-tails exponent  $\kappa$ . The inequality is a key technical tool that allows us to handle panel data with heavier than Gaussian tails and non-negligible  $T$  and  $N$ . The proof of this result can be found in the Appendix and is based on the blocking technique, cf., [Bosq \(1993\)](#) combined with the  $\tau$ -coupling lemma of [Dedecker and Prieur \(2004\)](#).

For short panels with small  $T$ , the following inequality might be a better choice.

**Theorem 4.2.** *Let  $\{\xi_{it} : i \geq 1, t \in \mathbf{Z}\}$  be an array of centered random vectors in  $\mathbf{R}^p$  such that  $\{\xi_{i,1}, \dots, \xi_{it} : i \geq 1\}$  are i.i.d. and for each  $i \geq 1$ ,  $(\xi_{it})_{t \in \mathbf{Z}}$  is a stationary stochastic process such that (i)  $\max_{j \in [p]} \|\xi_{it,j}\|_q = O(1)$  for some  $q > 2$ ; (ii) for*

---

<sup>9</sup>The direct application of the time series Fuk-Nagaev inequality of [Babii, Ghysels, and Striaukas \(2020a\)](#) leads to inferior concentration results for panel data.

every  $j \in [p]$ ,  $\tau$ -mixing coefficients of  $(\xi_{it,j})_{t \in \mathbf{Z}}$  satisfy  $\tau_k^{(j)} \leq ck^{-a}, \forall k \geq 1$  for some universal constants  $c > 0$  and  $a > \frac{q-1}{q-2}$ . Then for every  $u > 0$

$$\Pr \left( \left| \sum_{i=1}^N \sum_{t=1}^T \xi_{it} \right|_{\infty} > u \right) \leq c_1 p N u^{-q} + 4 p e^{-c_2 u^2 / NT}$$

for some  $c_1, c_2 > 0$ .

It follows from Theorem 4.2 that there exists  $C > 0$  such that for every  $\delta \in (0, 1)$

$$\Pr \left( \left| \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N \xi_{it} \right|_{\infty} \leq C \left( \frac{p}{\delta N^{q-1}} \right)^{1/q} \vee \sqrt{\frac{\log(8p/\delta)}{NT}} \right) \geq 1 - \delta.$$

The proof of this result can be found in the Appendix and is a straightforward application of the Fuk-Nagaev inequality for i.i.d. data and the Rosenthal's moment inequality, in contrast to Theorem 4.1. This inequality does not capture the concentration over  $T$  and may be a suboptimal choice for long panels which is the case in our empirical application.

## 5 Empirical Application

In our empirical application, we consider nowcasting the P/E ratios of 210 US firms using a set of predictors that are sampled at mixed frequencies. We use 24 predictors, including traditional macro and financial series as well as non-standard series generated by the textual analysis. We apply pooled and fixed effects sg-LASSO-MIDAS panel data models and compare them with several benchmarks such as random walk (RW), analysts consensus forecasts, and unstructured elastic net. We also compute predictions using individual-firm high-dimensional time series regressions and provide results for several choices of the tuning parameter. Lastly, we provide results for low-dimensional single-firm MIDAS regressions using forecast combination techniques used by [Andreou, Ghysels, and Kourtellis \(2013\)](#) and [Ball and Ghysels \(2018\)](#). The latter is particularly relevant regarding the analysis in the current paper as it also deals with nowcasting price earnings ratios. The forecast combination methods consist of estimating ADL-MIDAS regressions with each of the high-frequency covariates separately. In our case this leads to 24 predictions, corresponding to the number of predictors. Then a combination scheme, typically discounted mean squared error type, produces a single nowcast. One could call this a pre-machine learning large dimensional approach. It will, therefore, be interesting to assess how



this approach compares to the regularized MIDAS panel regression machine learning approach introduced in the current paper.

We start with a short review of the data, with more detailed descriptions and tables available appearing in Appendix Section D, followed by a summary of the methods used and the empirical results obtained.

## 5.1 Data description

The full sample consists of observations between 1<sup>st</sup> of January, 2000 and 30<sup>th</sup> of June, 2017. Due to the lagged dependent variables in the models, our effective sample starts the third fiscal quarter of 2000. We use the first 25 observations for the initial sample, and use the remaining 42 observations for evaluating the out-of-sample forecasts, which we obtain by using an expanding window forecasting scheme. We collect the data from CRSP and I/B/E/S to compute quarterly P/E ratios and firm-specific financial covariates; RavenPack is used to compute daily firm-level textual-analysis-based data; real-time monthly macroeconomic series are obtained from FRED-MD dataset, see McCracken and Ng (2016) for more details; FRED is used to compute daily financial markets data and, lastly, monthly news attention series extracted from the *Wall Street Journal* articles is retrieved from Bybee, Kelly, Manela, and Xiu (2019).<sup>10</sup> Appendix Section D provides a detailed description of the data sources. In particular, firm-level variables, including P/E ratios, are described in Appendix Table A.3, and the other predictor variables in Appendix Table A.4. The list of all firms we consider in our analysis appears in Appendix Table A.5.

**P/E ratio and analysts' forecasts sample construction.** Our target variable is the P/E ratio for each individual firm. To compute it, we use CRSP stock price data and I/B/E/S earnings data. Earnings data are subject to release delays of 1 to 2 months depending on the firm and quarter. Therefore, to reflect the real-time information flow, we separately compute the dependent variable, analysts' consensus forecasts, and the target variable using stock prices that were available in real-time. We also take into account that different firms have different fiscal quarters, which also affects the real-time information flow.

For example, suppose for a particular firm the fiscal quarters are at the end of the third month in a quarter, i.e. end of March, June, September, and December. Our dependent variable used in regression models is computed by taking the end of quarter prices and dividing them by the respective earnings value. The consensus

---

<sup>10</sup>The dataset is publicly available at <http://www.structureofnews.com/>.

forecast of the P/E ratio is computed using the same end of quarter price data which is divided by the earnings consensus forecast value. The consensus is computed by taking all individual prediction values up to the end of the quarter and aggregating those values by taking either the mean or the median. To compute the target variable which we use to measure the prediction performance, we adjust for publication lags and use prices of the publication date instead of the end of fiscal quarter prices. More precisely, suppose we predict the P/E ratio for the first quarter. Earnings are typically published with 1 to 2 months delay; say for example for a particular firm the data is published on the 25th of April. In this case, we record the stock price for this particular firm on 25th of April, and divide it by the realized earnings value.

## 5.2 Models and main results

To compute forecasts, we estimate several regression models. First, we estimate firm-specific sg-LASSO-MIDAS regressions, which in Table 1 we refer to as *Individual*. The model is written as

$$\mathbf{y}_i = \iota\alpha_i + \mathbf{x}_i\beta_i + \mathbf{u}_i, \quad i = 1, \dots, N,$$

and the firm-specific predictions are computed as  $\hat{y}_{i,T+1} = \hat{\alpha}_i + x_{i,T+1}^\top \hat{\beta}_i$ . As noted in Section 2,  $\mathbf{x}_i$  contains lags of the low-frequency target variable and MIDAS weights for each of the high-frequency covariate. We then estimate the following pooled and fixed effects sg-LASSO-MIDAS panel data models

$$\begin{aligned} \mathbf{y} &= \alpha\iota + \mathbf{X}\beta + \mathbf{u} && \text{Pooled} \\ \mathbf{y} &= B\alpha + \mathbf{X}\beta + \mathbf{u} && \text{Fixed Effects} \end{aligned}$$

and compute predictions as

$$\begin{aligned} \hat{y}_{i,T+1} &= \hat{\alpha} + x_{i,T+1}^\top \hat{\beta} && \text{Pooled} \\ \hat{y}_{i,T+1} &= \hat{\alpha}_i + x_{i,T+1}^\top \hat{\beta} && \text{Fixed Effects.} \end{aligned}$$

We benchmark firm-specific and panel data regression-based nowcasts against two simple alternatives. First, we compute forecasts for the RW model as

$$\hat{y}_{i,T+1} = y_{i,T}.$$

Second, we consider predictions of P/E implied by analysts earnings nowcasts using the information up to time  $T + 1$ , i.e.

$$\hat{y}_{i,T+1} = \bar{y}_{i,T+1},$$

where  $\bar{y}$  indicates that the forecasted P/E ratio is based on consensus earnings forecasts made at the end of  $T + 1$  quarter, and the stock price is also taken at the end of  $T + 1$ .

To measure the forecasting performance, we compute the mean squared forecast errors (MSE) for each method. Let  $\bar{\mathbf{y}}_i = (y_{iT_{is}+1}, \dots, y_{iT_{os}})^\top$  represent the out-of-sample realized P/E ratio values, where  $T_{is}$  and  $T_{os}$  denote the last initial in-sample observation and the last out-of-sample observation respectively, and let  $\hat{\mathbf{y}}_i = (\hat{y}_{iT_{is}+1}, \dots, \hat{y}_{iT_{os}})$  collect the out-of-sample forecasts from a specific method. Then, the mean squared forecast errors are computed as

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N \frac{1}{T - T_{is} + 1} (\bar{\mathbf{y}}_i - \hat{\mathbf{y}}_i)^\top (\bar{\mathbf{y}}_i - \hat{\mathbf{y}}_i).$$

RW	MSE An.-mean	MSE An.-median	sg-LASSO						
2.331	2.339	2.088	$\gamma =$	0	0.2	0.4	0.6	0.8	1
<u>Panel A. Cross-validation</u>									
	Individual	1.545	1.551	1.567	1.594	1.614	1.606		
	Pooled	1.459	1.456	<b>1.455</b>	1.456	1.455	1.459		
	Fixed Effects	1.500	1.489	1.487	1.501	1.480	1.489		
<u>Panel B. BIC</u>									
	Individual	1.657	1.634	1.609	1.543	1.561	1.610		
	Pooled	1.482	1.498	1.491	1.495	1.493	1.483		
	Fixed Effects	1.515	1.496	<b>1.472</b>	1.512	1.483	1.476		
<u>Panel C. AIC</u>									
	Individual	1.622	1.589	1.560	1.603	1.674	1.688		
	Pooled	1.494	1.492	1.488	1.487	1.490	1.492		
	Fixed Effects	1.504	1.487	1.486	1.504	<b>1.479</b>	1.489		
<u>Panel D. AICc</u>									
	Individual	2.025	2.122	2.272	2.490	2.923	3.255		
	Pooled	1.494	1.484	1.488	1.487	1.490	1.492		
	Fixed Effects	1.491	1.488	1.486	1.504	<b>1.479</b>	1.489		

Table 1: Prediction results – The table reports average over firms MSEs of out-of-sample predictions. The nowcasting horizon is the current month, i.e. we predict the P/E ratio using information up to the end of current fiscal quarter. Each Panel A-D block represents different ways of calculating the tuning parameter  $\lambda$ . Bold entries are the best results in a block.

The main results are reported in Table 1, while additional results for unstructured LASSO estimators and the forecast combination approach appear in Appendix Tables A.1-A.2. First, we document that analysts-based predictions have much larger mean squared forecast errors (MSEs) compared to model-based predictions. The sharp increase in quality of model- versus analyst-based predictions indicates the usefulness of machine learning methods to nowcast P/E ratios, see Tables 1 and A.1. A better

performance is achieved for almost all machine learning methods - single firm or panel data regressions - and all tuning parameter choices. Unstructured panel data methods and forecast combination approach also yield more accurate forecasts, see Appendix Table A.1-A.2. The latter confirms the findings of [Ball and Ghysels \(2018\)](#).

Turning to the comparison of model-based predictions, we see from the results in Table 1 that sg-LASSO-MIDAS panel data models improve the quality of predictions over individual sg-LASSO-MIDAS models irrespective of the  $\gamma$  weight or the tuning parameter choice. This indicates that panel data structures are relevant for nowcasting P/E ratios. We also report similar findings for unstructured estimators. Within the panel data framework, we observe that fixed effects improve over pooled regressions in most cases except when cross-validation is used; compare Table 1-A.2 Panel A with Table 1-A.2 Panel B-D. The pooled model tuned by cross-validation seems to yield the best overall performance. In general, one can expect that cross-validation improves prediction performance over different tuning methods as it is directly linked to empirical risk minimization. In the case of fixed effects, however, we may lose the predictive gain due to smaller samples with each fold used in estimating the model. Lastly, the best results per tuning parameter block seem to be achieved when  $\gamma \notin \{0, 1\}$ , indicating that both sparsity within the group and at the group level matters for prediction performance.

In Appendix Figure A.1, we plot the sparsity pattern of the selected covariates for the two best-performing methods: a) pooled sg-LASSO regressions, tuned using cross-validation with  $\gamma = 0.4$ , and b) fixed effects sg-LASSO model with BIC tuning parameter and the same  $\gamma$  parameter. We also plot the forecast combination weights which are averaged over firms. The plots in Figure A.1 reveal that the fixed effects estimator yields sparser models compared to pooled regressions, and the sparsity pattern is clearer. In the fixed-effects case, the Revenue growth and the first lag of the dependent variable are selected throughout the out-of-sample period. BAA less AAA bond yield spread, firm-level volatility, and Aggregate Event Sentiment index are also selected very frequently. Similarly, these variables are selected in the pooled regression, but the pattern is less apparent. The forecast combination weights seem to yield similar, yet a more blurred pattern.<sup>11</sup> In this case, Revenue growth and firm-level stock returns covariates obtain relatively larger weights compared to the rest of covariates, particularly for the first part of the out-of-sample period. Therefore, the gain of machine learning methods - both single-firm and panel data - can be

---

<sup>11</sup>Note that forecast combination weights start in 2009 Q1 due to the first eight quarters being used as a pre-sample to estimate weights, see [Ball and Ghysels \(2018\)](#) for further details. Also, the forecast combination weights figure does not contain autoregressive lags; all four lags are always included in all forecasting regressions.

associated with sparsity imposed on the regression coefficient vector.

It is also worth noting that the textual news data analytics also appear in the models according the results appear in Figure A.1. These are the ESS, AES, AEV, CSS and NEP regressors described in detail in Appendix Section D. Among them, as already noted, AES - the Aggregate Event Sentiment index - features most prominently in the sg-LASSO models. It is worth emphasizing that the time series of news data is sparse since for many days are without firms-specific news. For such days, we impute zero values. The nice property of our mixed frequency data treatment with dictionaries, imputing zeros also implies that non-zero entries get weights with a decaying pattern for distant past values in comparison to the most recent daily news data.

### 5.3 Significance test

To test for the superior forecast performance, we use the Diebold and Mariano (1995) test for the pool of P/E ratio nowcasts. We compare the mean and median consensus forecasts versus panel data machine learning regressions with the smallest forecast error per tuning parameter block in Table 1. We report the forecast accuracy test results in Table 2.

When testing the full sample of pooled nowcasts, the gain in prediction accuracy is not significant even though the MSEs are much lower for the panel data sg-LASSO regressions relative to the consensus forecasts. The result may not be surprising, however, as some firms have a large number of outlier observations and the Diebold and Mariano (1995) test statistic is affected by the inevitably heavy-tailed forecast errors for such firms. However, when we equally split the pooled sample of nowcasts into firms with high versus low variance P/E ratios, the gain in forecast accuracy is (not) significant for all panel data machine learning regressions for (high) low variance P/E firms.

## 6 Conclusions

This paper introduces a new class of high-dimensional panel data regression models with dictionaries and sparse-group LASSO regularization. This type of regularization is an especially attractive choice for the predictive panel data regressions, where the low- and/or the high-frequency lags define a clear group structure, and dictionaries are used to aggregate time series lags. The estimator nests the LASSO and the group LASSO estimators as special cases, as discussed in our theoretical analysis.

	Full sample	Large variance	Low variance
Pooled (Cross-validation) vs An.-mean	0.852	0.567	2.300
Pooled (Cross-validation) vs An.-median	0.694	0.386	2.190
Fixed-effects (BIC) vs An.-mean	0.793	0.508	2.312
Fixed-effects (BIC) vs An.-median	0.628	0.319	2.202
Fixed-effects (AIC) vs An.-mean	0.825	0.540	2.312
Fixed-effects (AIC) vs An.-median	0.663	0.355	2.202
Fixed-effects (AICc) vs An.-mean	0.825	0.540	2.312
Fixed-effects (AICc) vs An.-median	0.663	0.355	2.202

Table 2: Forecasting performance significance – The table reports the [Diebold and Mariano \(1995\)](#) test statistic for pooled nowcasts comparing machine learning panel data regressions with analysts’ implied consensus forecasts, where An.-mean and An.-median denote mean and median consensus forecasts respectively. We compare panel models that have the smallest forecast error per tuning parameter block in Table 1.

Our theoretical treatment allows for the heavy-tailed data frequently encountered in time series and financial econometrics. To that end, we obtain a new panel data concentration inequality of the Fuk-Nagaev type for  $\tau$ -mixing processes.

Our empirical analysis sheds light on the advantage of the regularized panel data regressions for nowcasting corporate earnings. We focus on nowcasting the P/E ratio of 210 US firms and find that the regularized panel data regressions outperform several benchmarks, including the analysts’ predictions. Furthermore, we find that the regularized machine learning regressions outperform the forecast combinations and that the panel data approach improves upon the predictive time series regressions for individual firms.

While nowcasting earnings is a leading example of applying panel data MIDAS machine learning regressions, one can think of many other applications of interest in finance. Beyond earnings, analysts are also interested in sales, dividends, etc. Our analysis can also be useful for other areas of interest, such as regional and international panel data settings.

## References

- ANDREOU, E., E. GHYSELS, AND A. KOURTELLOS (2013): “Should macroeconomic forecasters use daily financial data and how?,” *Journal of Business and Economic Statistics*, 31(2), 240–251.
- ARELLANO, M. (2003): *Panel data econometrics*. Oxford University Press.
- BABII, A. (2020): “High-dimensional mixed-frequency IV regression,” *arXiv preprint arXiv:2003.13478*.
- BABII, A., E. GHYSELS, AND J. STRIAUKAS (2020a): “Inference for high-dimensional regressions with heteroskedasticity and autocorrelation,” *arXiv preprint arXiv:1912.06307*.
- (2020b): “Machine learning time series regressions with an application to nowcasting,” *arXiv preprint arXiv:2005.14057*.
- BALL, R. T., AND E. GHYSELS (2018): “Automated earnings forecasts: beat analysts or combine and conquer?,” *Management Science*, 64(10), 4936–4952.
- BELLONI, A., V. CHERNOZHUKOV, C. HANSEN, AND D. KOZBUR (2016): “Inference in high-dimensional panel models with an application to gun control,” *Journal of Business and Economic Statistics*, 34(4), 590–605.
- BOSQ, D. (1993): “Bernstein-type large deviations inequalities for partial sums of strong mixing processes,” *Statistics*, 24(1), 59–70.
- BYBEE, L., B. T. KELLY, A. MANELA, AND D. XIU (2019): “The structure of economic news,” Available at SSRN 3446225.
- CHERNOZHUKOV, V., J. A. HAUSMAN, AND W. K. NEWEY (2019): “Demand analysis with many prices,” National Bureau of Economic Research Discussion paper 26424.
- CHETVERIKOV, D., Z. LIAO, AND V. CHERNOZHUKOV (2020): “On cross-validated Lasso,” *Annals of Statistics (forthcoming)*.
- CHIANG, H. D., J. RODRIGUE, AND Y. SASAKI (2019): “Post-selection inference in three-dimensional panel data,” *arXiv preprint arXiv:1904.00211*.
- DEDECKER, J., AND P. DOUKHAN (2003): “A new covariance inequality and applications,” *Stochastic Processes and their Applications*, 106(1), 63–80.

- DEDECKER, J., AND C. PRIEUR (2004): “Coupling for  $\tau$ -dependent sequences and applications,” *Journal of Theoretical Probability*, 17(4), 861–885.
- (2005): “New dependence coefficients. Examples and applications to statistics,” *Probability Theory and Related Fields*, 132(2), 203–236.
- DIEBOLD, F. X., AND R. S. MARIANO (1995): “Comparing predictive accuracy,” *Journal of Business and Economic Statistics*, 13(3), 253–263.
- FARRELL, M. H. (2015): “Robust inference on average treatment effects with possibly more covariates than observations,” *Journal of Econometrics*, 189(1), 1–23.
- FERNÁNDEZ-VAL, I., AND M. WEIDNER (2018): “Fixed effects estimation of large-T panel data models,” *Annual Review of Economics*, 10, 109–138.
- FOSTEN, J., AND R. GREENAWAY-MCGREY (2019): “Panel data nowcasting,” *Available at SSRN 3435691*.
- FUK, D. K., AND S. V. NAGAEV (1971): “Probability inequalities for sums of independent random variables,” *Theory of Probability and Its Applications*, 16(4), 643–660.
- GHYSELS, E., P. SANTA-CLARA, AND R. VALKANOV (2006): “Predicting volatility: getting the most out of return data sampled at different frequencies,” *Journal of Econometrics*, 131(1–2), 59–95.
- GHYSELS, E., A. SINKO, AND R. VALKANOV (2006): “MIDAS regressions: Further results and new directions,” *Econometric Reviews*, 26(1), 53–90.
- HARDING, M., AND C. LAMARCHE (2019): “A panel quantile approach to attrition bias in Big Data: Evidence from a randomized experiment,” *Journal of Econometrics*, 211(1), 61–82.
- KHALAF, L., M. KICHIAN, C. J. SAUNDERS, AND M. VOIA (2020): “Dynamic panels with MIDAS covariates: Nonlinearity, estimation and fit,” *Journal of Econometrics* (forthcoming).
- KOCK, A. B. (2013): “Oracle efficient variable selection in random and fixed effects panel data models,” *Econometric Theory*, 29(1), 115–152.
- (2016): “Oracle inequalities, variable selection and uniform inference in high-dimensional correlated random effects panel data models,” *Journal of Econometrics*, 195(1), 71–85.



- KOENKER, R. (2004): “Quantile regression for longitudinal data,” *Journal of Multivariate Analysis*, 91(1), 74–89.
- KOLANOVIC, M., AND R. KRISHNAMACHARI (2017): “Big data and AI strategies: Machine learning and alternative data approach to investing,” JP Morgan Global Quantitative & Derivatives Strategy Report.
- LAMARCHE, C. (2010): “Robust penalized quantile regression estimation for panel data,” *Journal of Econometrics*, 157(2), 396–408.
- LU, X., AND L. SU (2016): “Shrinkage estimation of dynamic panel data models with interactive fixed effects,” *Journal of Econometrics*, 190(1), 148–175.
- MCCRACKEN, M. W., AND S. NG (2016): “FRED-MD: A monthly database for macroeconomic research,” *Journal of Business and Economic Statistics*, 34(4), 574–589.
- SIMON, N., J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI (2013): “A sparse-group LASSO,” *Journal of Computational and Graphical Statistics*, 22(2), 231–245.
- SU, L., Z. SHI, AND P. C. PHILLIPS (2016): “Identifying latent structures in panel data,” *Econometrica*, 84(6), 2215–2264.
- TIBSHIRANI, R. (1996): “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society, Series B (Methodological)*, 58(1), 267–288.
- ZOU, H., T. HASTIE, AND R. TIBSHIRANI (2007): “On the “degrees of freedom” of the lasso,” *Annals of Statistics*, 35(5), 2173–2192.

# APPENDIX

## A Proofs of oracle inequalities

*Proof of Theorem 3.1.* The proof is similar to the proof of Babii, Ghysels, and Striaukas (2020b), Theorem 3.1 and is omitted. The main difference in the proof is that instead of applying the Fuk-Nagaev inequality from Babii, Ghysels, and Striaukas (2020a), Theorem 3.1, we apply the concentration inequality from Theorem 4.1 to

$$\left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T u_{it} z_{it} \right|_{\infty} \quad \text{and} \quad \max_{j,k \in [p]} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{it,j} z_{it,k} - \Sigma_{j,k} \right|$$

under Assumptions 3.1, 3.2, and 3.3. □

*Proof of Theorem 3.2.* Put  $r = (a^\top, b^\top)^\top$ . Then we solve

$$\min_{r \in \mathbf{R}^{N+p}} \|\mathbf{y} - \mathbf{Z}r\|_{NT}^2 + 2\lambda\Omega(b).$$

By Fermat's rule the solution to this problem satisfies

$$\mathbf{Z}^\top(\mathbf{Z}\hat{\rho} - \mathbf{y})/NT + \lambda z^* = \mathbf{0}_{N+p}$$

for some  $z^* = \begin{pmatrix} 0_N \\ z_b^* \end{pmatrix}$ , where  $0_N$  is  $N$ -dimensional vector of zeros,  $z_b^* \in \partial\Omega(\hat{\beta})$ ,  $\hat{\rho} = (\hat{\alpha}^\top, \hat{\beta}^\top)^\top$ , and  $\partial\Omega(\hat{\beta})$  is the sub-differential of  $b \mapsto \Omega(b)$  at  $\hat{\beta}$ . Taking the inner product with  $\rho - \hat{\rho}$

$$\begin{aligned} \langle \mathbf{Z}^\top(\mathbf{y} - \mathbf{Z}\hat{\rho}), \rho - \hat{\rho} \rangle_{NT} &= \lambda \langle z^*, \rho - \hat{\rho} \rangle \\ &= \lambda \langle z_b^*, \beta - \hat{\beta} \rangle \\ &\leq \lambda \left\{ \Omega(\beta) - \Omega(\hat{\beta}) \right\}, \end{aligned}$$

where the last line follows from the definition of the sub-differential. Rearranging

this inequality and using  $\mathbf{y} = \mathbf{m} + \mathbf{u}$

$$\begin{aligned}
\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 - \lambda \left\{ \Omega(\beta) - \Omega(\hat{\beta}) \right\} &\leq \langle \mathbf{Z}^\top \mathbf{u}, \hat{\rho} - \rho \rangle_{NT} + \langle \mathbf{Z}(\mathbf{m} - \mathbf{Z}\rho), \hat{\rho} - \rho \rangle_{NT} \\
&= \langle B^\top \mathbf{u}, \hat{\alpha} - \alpha \rangle_{NT} + \langle \mathbf{X}^\top \mathbf{u}, \hat{\beta} - \beta \rangle_{NT} \\
&\quad + \langle \mathbf{Z}(\mathbf{m} - \mathbf{Z}\rho), \hat{\rho} - \rho \rangle_{NT} \\
&\leq |B^\top \mathbf{u}/NT|_\infty |\hat{\alpha} - \alpha|_1 + \Omega^*(\mathbf{X}^\top \mathbf{u}/NT) \Omega(\hat{\beta} - \beta) \\
&\quad + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT} \|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT} \\
&\leq |B^\top \mathbf{u}/\sqrt{NT}|_\infty \vee \Omega^*(\mathbf{X}^\top \mathbf{u}/NT) \left\{ |\hat{\alpha} - \alpha|_1/\sqrt{N} + \Omega(\hat{\beta} - \beta) \right\} \\
&\quad + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT} \|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}, \tag{A.1}
\end{aligned}$$

where the second line follows by the dual norm inequality and the Cauchy-Schwartz inequality, and  $\Omega^*$  is the dual norm of  $\Omega$ . By [Babii, Ghysels, and Striaukas \(2020b\)](#), Lemma A.2.1.

$$\begin{aligned}
|B^\top \mathbf{u}/\sqrt{NT}|_\infty \vee \Omega^*(\mathbf{X}^\top \mathbf{u}/NT) &\leq C \max\{| \mathbf{X}^\top \mathbf{u}/NT |_\infty, |B^\top \mathbf{u}/\sqrt{NT}|_\infty\} \\
&= \max \left\{ \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T u_{it} x_{it} \right|_\infty, \max_{i \in [N]} \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T u_{it} \right| \right\}
\end{aligned}$$

for some  $C > 0$ , where the first inequality follows since  $\max_{G \in \mathcal{G}} |G| \lesssim 1$ . Under Assumption 3.1 by Theorem 4.1 and [Babii, Ghysels, and Striaukas \(2020a\)](#), Theorem 3.1 and Lemma A.1.1. for every  $u > 0$

$$\begin{aligned}
&\Pr(|B^\top \mathbf{u}/NT|_\infty \vee \Omega^*(\mathbf{X}^\top \mathbf{u}/NT) > u) \\
&\leq \Pr \left( \left| \frac{C}{NT} \sum_{i=1}^N \sum_{t=1}^T u_{it} x_{it} \right|_\infty > u \right) + \Pr \left( \max_{i \in [N]} \left| \frac{C}{\sqrt{NT}} \sum_{t=1}^T u_{it} \right| > u \right) \\
&\lesssim p(NT)^{1-\kappa} u^{-\kappa} + p e^{-c_1 u^2 NT} + N^{1-\kappa/2} T^{1-\kappa} u^{-\kappa} + 4N e^{-c_2 u^2 NT} \\
&\lesssim (pN^{1-\kappa} \vee N^{1-\kappa/2}) T^{1-\kappa} u^{-\kappa} + (p \vee N) e^{-c_2 u^2/NT}
\end{aligned}$$

for some  $c_1, c_2, C > 0$ . Therefore, under Assumption 3.5 with probability at least  $1 - \delta$

$$|B^\top \mathbf{u}/NT|_\infty \vee \Omega^*(\mathbf{X}^\top \mathbf{u}/NT) \lesssim \left( \frac{(pN^{1-\kappa}) \vee N^{1-\kappa/2}}{\delta T^{\kappa-1}} \right)^{1/\kappa} \vee \sqrt{\frac{\log(p \vee N/\delta)}{NT}} \lesssim \lambda.$$

In conjunction with the inequality in Eq. A.1, this gives

$$\begin{aligned}\|\mathbf{Z}\Delta\|_{NT}^2 &\leq c^{-1}\lambda\left\{|\hat{\alpha} - \alpha|_1/\sqrt{N} + \Omega(\hat{\beta} - \beta)\right\} + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT}\|\mathbf{Z}\Delta\|_{NT} + \lambda\left\{\Omega(\beta) - \Omega(\hat{\beta})\right\} \\ &\leq (c^{-1} + 1)\lambda\left\{|\hat{\alpha} - \alpha|_1/\sqrt{N} + \Omega(\hat{\beta} - \beta)\right\} + \|\mathbf{m} - \mathbf{Z}\rho\|_{NT}\|\mathbf{Z}\Delta\|_{NT}\end{aligned}\tag{A.2}$$

for some  $c > 1$  and  $\Delta = \hat{\rho} - \rho$ , where the second line follows by the triangle inequality. Note that the sg-LASSO penalty function can be decomposed as a sum of two semi-norms  $\Omega(b) = \Omega_0(b) + \Omega_1(b)$ ,  $\forall b \in \mathbf{R}^p$  and that we have  $\Omega_1(\beta) = 0$  and  $\Omega_1(\hat{\beta}) = \Omega_1(\hat{\beta} - \beta)$ . Then

$$\begin{aligned}\Omega(\beta) - \Omega(\hat{\beta}) &= \Omega_0(\beta) - \Omega_0(\hat{\beta}) - \Omega_1(\hat{\beta}) \\ &\leq \Omega_0(\hat{\beta} - \beta) - \Omega_1(\hat{\beta} - \beta).\end{aligned}\tag{A.3}$$

Suppose that  $\|\mathbf{m} - \mathbf{Z}\rho\|_{NT} \leq \frac{1}{2}\|\mathbf{Z}\Delta\|_{NT}$ . Then it follows from the first inequality in Eq. A.2 and Eq. A.3 that

$$\|\mathbf{Z}\Delta\|_{NT}^2 \leq 2c^{-1}\lambda\left\{|\hat{\alpha} - \alpha|_1/\sqrt{N} + \Omega(\hat{\beta} - \beta)\right\} + 2\lambda\left\{\Omega_0(\hat{\beta} - \beta) - \Omega_1(\hat{\beta} - \beta)\right\}.$$

Since the left side of this equation is  $\geq 0$ , this shows that

$$(1 - c^{-1})\Omega_1(\hat{\beta} - \beta) \leq (1 + c^{-1})\Omega_0(\hat{\beta} - \beta) + c^{-1}|\hat{\alpha} - \alpha|_1/\sqrt{N}$$

or equivalently

$$\Omega_1(\hat{\beta} - \beta) \leq \frac{c+1}{c-1}\Omega_0(\hat{\beta} - \beta) + (c-1)^{-1}|\hat{\alpha} - \alpha|_1/\sqrt{N}.\tag{A.4}$$

Put  $\Delta_N = ((\hat{\alpha} - \alpha)^\top / \sqrt{N}, (\hat{\beta} - \beta)^\top)^\top$ . Then under Assumption 3.2

$$\begin{aligned}
|\Delta_N|_1 &\lesssim \Omega(\hat{\beta} - \beta) + |\hat{\alpha} - \alpha|_1 / \sqrt{N} \\
&\leq \frac{2c}{c-1} \Omega_0(\hat{\beta} - \beta) + \frac{c}{c-1} |\hat{\alpha} - \alpha|_1 / \sqrt{N} \\
&\lesssim |\hat{\alpha} - \alpha|_2 + \sqrt{s} |\hat{\beta} - \beta|_2 \\
&\leq \sqrt{s \vee N} |\Delta_N|_2^2 \\
&\lesssim \sqrt{s \vee N} |\Sigma^{1/2} \Delta_N|_2^2 \\
&= \sqrt{s \vee N} \left\{ \|\mathbf{Z}\Delta\|_{NT}^2 + \Delta_N^\top (\hat{\Sigma} - \Sigma) \Delta_N \right\} \\
&\leq \sqrt{s \vee N} \left\{ \|\mathbf{Z}\Delta\|_{NT}^2 + |\Delta_N|_1^2 |\text{vech}(\hat{\Sigma} - \Sigma)|_\infty \right\} \\
&\lesssim \sqrt{s \vee N} \left\{ \lambda |\Delta_N|_1 + |\Delta_N|_1^2 |\text{vech}(\hat{\Sigma} - \Sigma)|_\infty \right\}.
\end{aligned}$$

Consider the following event  $E = \{|\text{vech}(\hat{\Sigma} - \Sigma)|_\infty < 1/(2s \vee N)\}$ . Under Assumption 3.1 by Theorem 4.1 and Babii, Ghysels, and Striaukas (2020a), Theorem 3.1

$$\begin{aligned}
\Pr(E^c) &\leq \Pr \left( \max_{i \in [N], j \in [p]} \left| \frac{1}{\sqrt{NT}} \sum_{t=1}^T \{x_{it,j} - \mathbb{E}[x_{it,j}]\} \right| \geq \frac{1}{2s \vee N} \right) \\
&\quad + \Pr \left( \max_{1 \leq j \leq k \leq p} \left| \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T x_{it,j} x_{it,k} \right| \geq \frac{1}{2s \vee N} \right) \\
&\lesssim p(s \vee N)^{\tilde{\kappa}} T^{1-\tilde{\kappa}} (N^{1-\tilde{\kappa}/2} + pN^{1-\tilde{\kappa}}) + p(p \vee N) e^{-cNT/(s \vee N)^2}.
\end{aligned}$$

Therefore, on the event  $E$

$$|\hat{\alpha} - \alpha|_1 / \sqrt{N} + |\hat{\beta} - \beta|_1 = |\Delta_N|_1 \lesssim (s \vee N) \lambda,$$

and whence from Eq. A.2 we obtain

$$\begin{aligned}
\|\mathbf{Z}\Delta\|_{NT}^2 &\lesssim \lambda \left\{ |\hat{\alpha} - \alpha|_1 / \sqrt{N} + \Omega(\hat{\beta} - \beta) \right\} \\
&\lesssim \lambda |\Delta_N|_1 \\
&\leq (s \vee N) \lambda^2.
\end{aligned}$$

Suppose now that  $\|\mathbf{m} - \mathbf{Z}\rho\|_{NT} > \frac{1}{2} \|\mathbf{Z}\Delta\|_{NT}$ . Then, obviously,

$$\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 \leq 4 \|\mathbf{m} - \mathbf{Z}\rho\|_{NT}^2.$$

Therefore, on the event  $E$ , we always have

$$\|\mathbf{Z}(\hat{\rho} - \rho)\|_{NT}^2 \lesssim (s \vee N)\lambda^2 + 4\|\mathbf{m} - \mathbf{Z}\rho\|_{NT}^2,$$

which proves the statement of the theorem.  $\square$

## B Proofs of Fuk-Nagaev inequalities

*Proof of Theorem 4.1.* Suppose first that  $p = 1$ . For  $a \in \mathbf{R}$  with some abuse of notation, let  $[a]$  denote its integer part. For each  $i = 1, 2, \dots, N$ , split partial sums into blocks with at most  $J \in \mathbf{N}$  summands

$$\begin{aligned} V_{i,k} &= \xi_{i,(k-1)J+1} + \dots + \xi_{i,kJ}, & k &= 1, 2, \dots, [T/J] \\ V_{i,[T/J]+1} &= \xi_{i,[T/J]J+1} + \dots + \xi_{i,T}, \end{aligned}$$

where we set  $V_{i,[T/J]+1} = 0$  if  $[T/J]J = T$ . Let  $\{U_{i,t} : i, t \geq 1\}$  be i.i.d. random variables uniformly distributed on  $[0, 1]$  and independent of  $\{\xi_{i,t} : i, t \geq 1\}$ . Put  $\mathcal{M}_{i,t} = \sigma(V_{i,1}, \dots, V_{i,t-2})$  with  $t \geq 3$ . For  $t = 1, 2$ , set  $V_{i,t}^* = V_{i,t}$ , while for  $t \geq 3$ , by [Dedecker and Prieur \(2004\)](#), Lemma 5, there exist random variables  $V_{i,t}^* =_d V_{i,t}$  such that

1.  $V_{i,t}^*$  is  $\mathcal{M}_{i,t} \vee \sigma(V_{i,t}) \vee \sigma(U_{i,t})$ -measurable.
2.  $V_{i,t}^*$  is independent of  $\mathcal{M}_{i,t}$ .
3.  $\|V_{i,t} - V_{i,t}^*\|_1 = \tau(\mathcal{M}_{i,t}, V_{i,t})$ .

By 1. there exists a measurable function  $f_i$  such that

$$V_{i,t}^* = f_i(V_{i,t}, V_{i,t-2}, \dots, V_{i,1}, U_{i,t}).$$

Therefore, by 2.,  $(V_{i,2t}^*)_{t \geq 1}$  and  $(V_{i,2t-1}^*)_{t \geq 1}$  are sequences of independent random variables for every  $i = 1, \dots, N$ . Moreover,  $\{V_{i,2t}^* : i = 1, \dots, N, t \geq 1\}$  and  $\{V_{i,2t-1}^* : i = 1, \dots, N, t \geq 1\}$  are sequences of independent random variables since  $\{\xi_{i,t} : t = 1, \dots, T\}$  are independent over  $i = 1, \dots, N$ .

Decompose

$$\begin{aligned} \left| \sum_{i=1}^N \sum_{t=1}^T \xi_{i,t} \right| &\leq \left| \sum_{i=1}^N \sum_{t \geq 1} V_{i,2t}^* \right| + \left| \sum_{i=1}^N \sum_{t \geq 1} V_{i,2t-1}^* \right| + \sum_{i=1}^N \sum_{t=3}^{[T/J]+1} |V_{i,t} - V_{i,t}^*| \\ &\triangleq I + II + III. \end{aligned}$$

By [Fuk and Nagaev \(1971\)](#), Corollary 4 there exist constants  $c_1, c_2 > 0$  such that

$$\begin{aligned} \Pr(I > u/3) &\leq c_1 u^{-q} N \sum_{t \geq 1} \mathbb{E}|V_{i,2t}^*|^q + 2 \exp\left(-\frac{c_2 u^2}{N \sum_{t \geq 1} \text{Var}(V_{i,2t}^*)}\right) \\ &\leq c_1 u^{-q} N \sum_{t \geq 1} \mathbb{E}|V_{i,2t}|^q + 2 \exp\left(-\frac{c_2 u^2}{NT}\right), \end{aligned}$$

where we use  $V_{i,t}^* =_d V_{i,t}$  and

$$\sum_{t \geq 1} \text{Var}(V_{i,2t}) \leq \sum_{t \geq 1} \text{Var}(V_{i,t}) = O(T),$$

which follows by [Babii, Ghysels, and Striaukas \(2020a\)](#), Lemma A.1.2. Similarly,

$$\Pr(II > u/3) \leq c_3 u^{-q} N \sum_{t \geq 1} \mathbb{E}|V_{i,2t}|^q + 2 \exp\left(-\frac{c_4 u^2}{NT}\right)$$

for some constants  $c_3, c_4 > 0$ . Lastly, since  $\mathcal{M}_{i,t}$  and  $V_{i,t}$  are separated by  $J + 1$  lags of  $\xi_{i,t}$ , we have  $\tau(\mathcal{M}_{i,t}, V_{i,t}) \leq J\tau_J(J + 1)$ . By Markov's inequality and property 3., this gives

$$\begin{aligned} \Pr(III > u/3) &\leq \frac{3N}{u} \sum_{t=3}^{[T/J]+1} \|V_{i,t} - V_{i,t}^*\|_1 \\ &\leq \frac{3NT}{u} \tau_{J+1}. \end{aligned}$$

Combining all estimates together

$$\begin{aligned} \Pr\left(\left|\sum_{i=1}^N \sum_{t=1}^T \xi_{i,t}\right| > u\right) &\leq \Pr(I > u/3) + \Pr(II > u/3) + \Pr(III > u/3) \\ &\leq c_1 u^{-q} N \sum_{t \geq 1} \|V_{i,t}\|_q^q + 4e^{-c_2 u^2/NT} + \frac{3NT}{u} \tau_{J+1} \\ &\leq c_1 u^{-q} J^{q-1} NT \|\xi_{i,t}\|_q^q + \frac{3NT}{u} (J + 1)^{-a} + 4e^{-c_2 u^2/NT} \end{aligned}$$

for some constants  $c_1, c_2 > 0$ . To balance the first two terms, we shall choose the length of blocks  $J \sim u^{\frac{q-1}{q+a-1}}$ , in which case we get

$$\Pr\left(\left|\sum_{i=1}^N \sum_{t=1}^T \xi_{i,t}\right| > u\right) \leq c_1 NT u^{-\kappa} + 4e^{-c_2 u^2/NT}$$

for some  $c_1, c_2 > 0$ .

Finally, for  $p > 1$ , the result follows by the union bound.  $\square$

*Proof of Theorem 4.2.* Put

$$M_{q,N,T} \triangleq \max_{j \in [p]} \max_{i \in [N]} \mathbb{E} \left| \sum_{t=1}^T \xi_{i,t,j} \right|^q \quad \text{and} \quad B_{N,T}^2 \triangleq \max_{j \in [p]} \sum_{i=1}^N \text{Var} \left( \sum_{t=1}^T \xi_{i,t,j} \right).$$

By Jensen's inequality under the stationarity and the i.i.d. hypotheses

$$M_{q,N,T} \leq \max_{j \in [p]} T^q \mathbb{E} |\xi_{i,t,j}|^q \lesssim T^q,$$

where the last inequality follows under assumption (i). Similarly,

$$B_{N,T}^2 \leq N \max_{j \in [p]} \sum_{t=1}^T \sum_{k=1}^T |\text{Cov}(\xi_{1,t,j}, \xi_{1,k,j})| \lesssim NT,$$

where the last inequality follows from the computations in Theorem 4.1 under assumptions (i)-(ii).

Using these estimates, by the union bound and [Fuk and Nagaev \(1971\)](#), Corollary 4, for every  $u > 0$

$$\begin{aligned} \Pr \left( \left| \sum_{t=1}^T \sum_{i=1}^N \xi_{i,t} \right|_{\infty} > u \right) &\leq c_1 M_{q,N,T} p N u^{-q} + 2p \exp \left( -\frac{c_2 u^2}{B_{N,T}^2} \right) \\ &\leq c_3 p N T^q u^{-q} + 2e^{-c_4 u^2 / NT} \end{aligned}$$

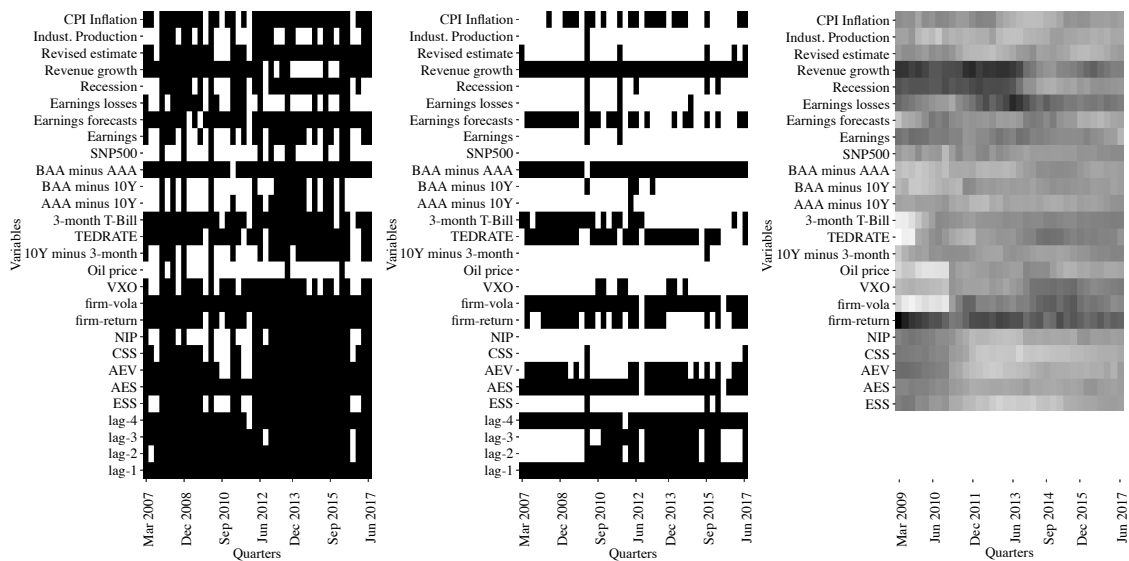
for some constants  $c_j > 0, j \in [4]$ . Therefore, there exists  $C > 0$  such that

$$\Pr \left( \left| \sum_{t=1}^T \sum_{i=1}^N \xi_{i,t} \right|_{\infty} \leq C (pNT^q/\delta)^{1/q} \vee \sqrt{NT \log(4p/\delta)} \right) \geq 1 - \delta.$$

$\square$

## C Additional empirical results





(a) Pooled sg-LASSO,  $\gamma = 0.4$ , cross-validation. (b) Fixed effects sg-LASSO,  $\gamma = 0.4$ , BIC. (c) Average forecast combination weights.

Figure A.1: Sparsity patterns and forecast combination weights.

RW	MSE An.-mean	MSE An.-median	sg-LASSO	elnet-U	elnet
2.331	2.339	2.088			
<u>Panel A. Cross-validation</u>					
Individual	1.545	1.610	1.609		
Pooled	<b>1.455</b>	1.497	1.460		
Fixed Effects	1.480	1.517	1.514		
<u>Panel B. BIC</u>					
Individual	1.543	1.641	1.875		
Pooled	1.482	1.489	1.486		
Fixed Effects	<b>1.472</b>	1.493	1.494		
<u>Panel C. AIC</u>					
Individual	1.560	1.692	1.860		
Pooled	1.487	1.492	1.494		
Fixed Effects	<b>1.479</b>	1.493	1.499		
<u>Panel D. AICc</u>					
Individual	2.025	1.734	2.097		
Pooled	1.484	1.492	1.494		
Fixed Effects	<b>1.479</b>	1.494	1.499		

Table A.1: Prediction results – The table reports average over firms MSEs of out-of-sample predictions. The nowcasting horizon is the current month, i.e. we predict the P/E ratio using information up to the end of current fiscal quarter. Each Panel A-D block represents different ways of calculating the tuning parameter  $\lambda$ . Bold entries are the best results in a block. We report elastic net MSEs averaged over LASSO/ridge weight  $[0, 0.2, 0.4, 0.6, 0.8, 1]$ : elnet-U method is where high-frequency lags are unrestricted, elnet method is where we use only the first high-frequency lag for each covariate. We also report the best sg-LASSO specification for each tuning parameter method and each model specification, see Table 1.

## D Data description

### D.1 Firm-level data

The full list of firm-level data is provided in Table A.3. We also add two daily firm-specific stock market predictor variables: stock returns and a realized variance measure, which is defined as the rolling sample variance over the previous 60 days (i.e. 60-day historical volatility).

#### D.1.1 Firm sample selection

We select a sample of firms based on data availability. First, we remove all firms from I/B/E/S which have missing values in earnings time series. Next, we retain firms that we are able to match with CRSP dataset. Finally, we keep firms that we can match with the RavenPack dataset.

RW	MSE An.-mean	MSE An.-median	F.Comb	sg-LASSO						
2.794	2.836	2.539	2.405	$\gamma =$	0	0.2	0.4	0.6	0.8	1
							<u>Panel A. Cross-validation</u>			
	Individual	1.808	1.817	1.836	1.864	1.889	1.884			
	Pooled	1.692	1.689	<b>1.688</b>	<b>1.688</b>	<b>1.688</b>	1.689			
	Fixed Effects	1.743	1.726	1.725	1.743	1.712	1.726			
							<u>Panel B. BIC</u>			
	Individual	1.972	1.945	1.914	1.833	1.853	1.912			
	Pooled	1.723	1.741	1.733	1.738	1.736	1.724			
	Fixed Effects	1.760	1.734	<b>1.707</b>	1.756	1.717	1.710			
							<u>Panel C. AIC</u>			
	Individual	1.929	1.889	1.853	1.903	1.989	2.003			
	Pooled	1.737	1.735	1.729	1.728	1.732	1.734			
	Fixed Effects	1.747	1.724	1.724	1.747	<b>1.712</b>	1.726			
							<u>Panel D. AICc</u>			
	Individual	2.401	2.513	2.679	2.918	3.404	3.732			
	Pooled	1.737	1.725	1.729	1.728	1.732	1.734			
	Fixed Effects	1.732	1.725	1.724	1.747	<b>1.712</b>	1.726			

Table A.2: Prediction results – The table reports average over firms MSEs of out-of-sample predictions for the same models as in Table 1 - discarding the first 8 quarters to compute for forecast combination weights - with additional result of prediction errors using forecast combination approach of [Ball and Ghysels \(2018\)](#), denoted as *F.Comb*. Hence the out-of-sample quarters start at 2009 Q1. The nowcasting horizon is the current month, i.e. we predict the P/E ratio using information up to the end of current fiscal quarter. Each Panel A-D block represents different ways of calculating the tuning parameter  $\lambda$ . Bold entries are the best results in a block.

### D.1.2 Firm-specific text data

We create a link table of RavenPack ID and PERMNO identifiers which enables us to merge I/B/E/S and CRSP data with firm-specific textual analysis generated data from RavenPack. The latter is a rich dataset that contains intra-daily news information about firms. There are several editions of the dataset; in our analysis, we use the Dow Jones (DJ) and Press Release (PR) editions. The former contains relevant information from Dow Jones Newswires, regional editions of the Wall Street Journal, Barron's and MarketWatch. The PR edition contains news data, obtained from various press releases and regulatory disclosures, on a daily basis from a variety of newswires and press release distribution networks, including exclusive content from PRNewswire, Canadian News Wire, Regulatory News Service, and others. The DJ edition sample starts at 1<sup>st</sup> of January, 2000, and PR edition data starts at 17<sup>th</sup> of January, 2004.

We construct our news-based firm-level covariates by filtering only highly relevant news stories. More precisely, for each firm and each day, we filter out news that has the *Relevance Score* (REL) larger or equal to 75, as is suggested by the RavenPack News Analytics guide and used by practitioners, see for example [Kolanovic and Krishnamachari \(2017\)](#). REL is a score between 0 and 100 which indicates how strongly a news story is linked with a particular firm. A score of zero means that the entity is vaguely mentioned in the news story, while 100 means the opposite. A score of 75 is regarded as a significantly relevant news story. After applying the REL filter, we apply a novelty of the news filter by using the *Event Novelty Score* (ENS); we keep data entries that have a score of 100. Like REL, ENS is a score between 0 and 100. It indicates the novelty of a news story within a 24-hour time window. A score of 100 means that a news story was not already covered by earlier announced news, while subsequently published news story score on a related event is discounted, and therefore its scores are less than 100. Therefore, with this filter, we consider only novel news stories. We focus on *five sentiment indices* that are available in both DJ and PR editions. They are:

**Event Sentiment Score** (ESS), for a given firm, represents the strength of the news measured using surveys of financial expert ratings for firm-specific events. The score value ranges between 0 and 100 - values above (below) 50 classify the news as being positive (negative), 50 being neutral.

**Aggregate Event Sentiment** (AES) represents the ratio of positive events reported on a firm compared to the total count of events measured over a rolling

91-day window in a particular news edition (DJ or PR). An event with  $ESS > 50$  is counted as a positive entry while  $ESS < 50$  as negative. Neutral news ( $ESS = 50$ ) and news that does not receive an ESS score does not enter into the AES computation. As ESS, the score values are between 0 and 100.

**Aggregate Event Volume** (AEV) represents the count of events for a firm over the last 91 days within a certain edition. As in AES case, news that receives a non-neutral ESS score is counted and therefore accumulates positive and negative news.

**Composite Sentiment Score** (CSS) represents the news sentiment of a given news story by combining various sentiment analysis techniques. The direction of the score is determined by looking at emotionally charged words and phrases and by matching stories typically rated by experts as having short-term positive or negative share price impact. The strength of the scores is determined by intra-day price reactions modeled empirically using tick data from approximately 100 large-cap stocks. As for ESS and AES, the score takes values between 0 and 100, 50 being the neutral.

**News Impact Projections** (NEP) represents the degree of impact a news flash has on the market over the following two-hour period. The algorithm produces scores to accurately predict a relative volatility - defined as scaled volatility by the average of volatilities of large-cap firms used in the test set - of each stock price measured within two hours following the news. Tick data is used to train the algorithm and produce scores, which take values between 0 and 100, 50 representing zero impact news.

For each firm and each day with firm-specific news, we compute the average value of the specific sentiment score. In this way, we aggregate across editions and groups, where the later is defined as a collection of related news. We then map the indices that take values between 0 and 100 onto  $[-1, 1]$ . Specifically, let  $x_i \in \{ESS, AES, CSS, NIP\}$  be the average score value for a particular day and firm. We map  $x_i \mapsto \bar{x}_i \in [-1, 1]$  by computing  $\bar{x}_i = (x_i - 50)/50$ .

	id	Frequency	Source	T-code
Panel A.				
-	Price/Earnings ratio	quarterly	CRSP & I/B/E/S	1
-	Price/Earnings ratio consensus forecasts	quarterly	CRSP & I/B/E/S	1
Panel B.				
1	Stock returns	daily	CRSP	1
2	Realized variance measure	daily	CRSP/computations	1
Panel C.				
1	Event Sentiment Score (ESS)	daily	RavenPack	1
2	Aggregate Event Sentiment (AES)	daily	RavenPack	1
3	Aggregate Event Volume (AEV)	daily	RavenPack	1
4	Composite Sentiment Score (CSS)	daily	RavenPack	1
5	News Impact Projections (NIP)	daily	RavenPack	1

Table A.3: Firm-level data description table – The *id* column gives mnemonics according to data source, which is given in the second column *Source*. The column *frequency* states the sampling frequency of the variable. The column *T-code* denotes the data transformation applied to a time-series, which are: (1) not transformed, (2)  $\Delta x_t$ , (3)  $\Delta^2 x_t$ , (4)  $\log(x_t)$ , (5)  $\Delta \log(x_t)$ , (6)  $\Delta^2 \log(x_t)$ . Panel A. describes earnings data, panel B. describes quarterly firm-level accounting data, panel C. daily firm-level stock market data and panel D. daily firm-level sentiment data series.

	id	Frequency	Source	T-code
Panel A.				
1	Industrial Production Index	monthly	FRED-MD	5
2	CPI Inflation	monthly	FRED-MD	6
Panel B.				
1	Crude Oil Prices	daily	FRED	6
2	S&P 500	daily	CRSP	5
3	VXO Volatility Index	daily	FRED	1
4	Moodys Aaa - 10-Year Treasury	daily	FRED	1
5	Moodys Baa - 10-Year Treasury	daily	FRED	1
6	Moodys Baa - Aaa Corporate Bond	daily	FRED	1
7	10-Year Treasury - 3-Month Treasury	daily	FRED	1
8	3-Month Treasury - Effective Federal funds rate	daily	FRED	1
9	TED rate	daily	FRED	1
Panel C.				
1	Earnings	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1
2	Earnings forecasts	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1
3	Earnings losses	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1
4	Recession	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1
5	Revenue growth	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1
6	Revised estimate	monthly	Bybee, Kelly, Manela, and Xiu (2019)	1

Table A.4: Other predictor variables description table – The *id* column gives mnemonics according to data source, which is given in the second column *Source*. The column *frequency* states the sampling frequency of the variable. The column *T-code* denotes the data transformation applied to a time-series, which are: (1) not transformed, (2)  $\Delta x_t$ , (3)  $\Delta^2 x_t$ , (4)  $\log(x_t)$ , (5)  $\Delta \log(x_t)$ , (6)  $\Delta^2 \log(x_t)$ . Panel A. describes real-time monthly macro series, panel B. describes daily financial markets data and panel C. monthly news attention series.

	Ticker	Firm name	PERMNO	RavenPack ID
1	MMM	3M	22592	03B8CF
2	ABT	Abbott labs	20482	520632
3	AUD	Automatic data processing	44644	66ECFD
4	ADTN	Adtran	80791	9E98F2
5	AEIS	Advanced energy industries	82547	1D943E
6	AMG	Affiliated managers group	85593	30E01D
7	AKST	A K steel holding	80303	41588B
8	ATI	Allegheny technologies	43123	D1173F
9	AB	AllianceBernstein holding l.p.	75278	CB138D
10	ALL	Allstate corp.	79323	E1C16B
11	AMZN	Amazon.com	84788	0157B1
12	AMD	Advanced micro devices	61241	69345C
13	DOX	Amdocs ltd.	86144	45D153
14	AMKR	Amkor technology	86047	5C8D61
15	APH	Amphenol corp.	84769	BB07E4
16	AAPL	Apple	14593	D8442A
17	ADM	Archer daniels midland	10516	2B7A40
18	ARNC	Arconic	24643	EC821B
19	ATTA	AT&T	66093	251988
20	AVY	Avery dennison corp.	44601	662682
21	BHI	Baker hughes	75034	940C3D
22	BAC	Bank of america corp.	59408	990AD0
23	BAX	Baxter international inc.	27887	1FAF22
24	BBT	BB&T corp.	71563	1A3E1B
25	BDX	Becton dickinson & co.	39642	873DB9
26	BBBY	Bed bath & beyond inc.	77659	9B71A7
27	BHE	Benchmark electronics inc.	76224	6CF43C
28	BA	Boeing co.	19561	55438C
29	BK	Bank of new york mellon corp.	49656	EF5BED
30	BWA	BorgWarner inc.	79545	1791E7
31	BP	BP plc	29890	2D469F
32	EAT	Brinker international inc.	23297	732449
33	BMY	Bristol-Myers squibb co.	19393	94637C
34	BRKS	Brooks automation inc.	81241	FC01C0
35	CA	CA technologies inc.	25778	76DE40
36	COG	Cabot oil & gas corp.	76082	388E00
37	CDN	Cadence design systems inc.	11403	CC6FF5
38	COF	Capital one financial corp.	81055	055018
39	CRR	Carbo ceramics inc.	83366	8B66CE
40	CSL	Carlisle cos.	27334	9548BB
41	CCL	Carnival corporation & plc	75154	067779
42	CERN	Cerner corp.	10909	9743E5
43	CHRW	C.H. robinson worldwide inc.	85459	C659EB
44	SCHW	Charles schwab corp.	75186	D33D8C
45	CHKP	Check point software technologies ltd.	83639	531EF1
46	CHV	Chevron corp.	14541	D54E62
47	CI	CIGNA corp.	64186	86A1B9
48	CTAS	Cintas corp.	23660	BFAEB4
49	CLX	Clorox co.	46578	719477
50	KO	Coca-Cola co.	11308	EEA6B3
51	CGNX	Cognex corp.	75654	709AED
52	COLM	Columbia sportswear co.	85863	5D0337
53	CMA	Comerica inc.	25081	8CF6DD
54	CRK	Comstock resources inc.	11644	4D72C8
55	CAG	ConAgra foods inc.	56274	FA40E2
56	STZ	Constellation brands inc.	69796	1D1B07
57	CVG	Convergys corp.	86305	914819

58	COST	Costco wholesale corp.	87055	B8EF97
59	CCI	Crown castle international corp.	86339	275300
60	DHR	Danaher corp.	49680	E124EB
61	DRI	Darden restaurants inc.	81655	9BBFA5
62	DVA	DaVita inc.	82307	EFD406
63	DO	Diamond offshore drilling inc.	82298	331BD2
64	D	Dominion resources inc.	64936	977A1E
65	DOV	Dover corp.	25953	636639
66	DOW	Dow chemical co.	20626	523A06
67	DHI	D.R. horton inc.	77661	06EF42
68	EMN	Eastman chemical co.	80080	D4070C
69	EBAY	eBay inc.	86356	972356
70	EOG	EOG resources inc.	75825	A43906
71	EL	Estee lauder cos. inc.	82642	14ED2B
72	ETH	Ethan allen interiors inc.	79037	65CF8E
73	ETFC	E*TRADE financial corp.	83862	28DEFA
74	XOM	Exxon mobil corp.	11850	E70531
75	FII	Federated investors inc.	86102	73C9E2
76	FDX	FedEx corp.	60628	6844D2
77	FITB	Fifth third bancorp	34746	8377DB
78	FISV	Fiserv inc.	10696	190B91
79	FLEX	Flex ltd.	80329	B4E00D
80	F	Ford motor co.	25785	A6213D
81	FWRD	Forward air corp.	79841	10943B
82	BEN	Franklin resources inc.	37584	5B6C11
83	GE	General electric co.	12060	1921DD
84	GIS	General mills inc.	17144	9CA619
85	GNTX	Gentex corp.	38659	CC339B
86	HAL	Halliburton Co.	23819	2B49F4
87	HLIT	Harmonic inc.	81621	DD9E41
88	HIG	Hartford financial services group inc.	82775	766047
89	HAS	Hasbro inc.	52978	AA98ED
90	HLX	Helix energy solutions group inc.	85168	6DD6BA
91	HP	Helmerich & payne inc.	32707	1DE526
92	HSY	Hershey co.	16600	9F03CF
93	HES	Hess corp.	28484	D0909F
94	HON	Honeywell international inc.	10145	FF6644
95	JBHT	J.B. Hunt transport services Inc.	42877	72DF04
96	HBAN	Huntington bancshares inc.	42906	C9E107
97	IBM	IBM corp.	12490	8D4486
98	IEX	IDEX corp.	75591	E8B21D
99	IR	Ingersoll-Rand plc	12431	5A6336
100	IDTI	Integrated device technology inc.	44506	8A957F
101	INTC	Intel corp.	59328	17EDA5
102	IP	International paper co.	21573	8E0E32
103	IIN	ITT corp.	12570	726EEA
104	JAKK	Jakks pacific inc.	83520	5363A2
105	JNJ	Johnson & johnson	22111	A6828A
106	JPM	JPMorgan chase & co.	47896	619882
107	K	Kellogg co.	26825	9AF3DC
108	KMB	Kimberly-Clark corp.	17750	3DE4D1
109	KNGT	Knight transportation inc.	80987	ED9576
110	LSTR	Landstar system inc.	78981	FD4E8D
111	LSCC	Lattice semiconductor corp.	75854	8303CD
112	LLY	Eli lilly & co.	50876	F30508
113	LFUS	Littelfuse inc.	77918	D06755
114	LNC	Lincoln national corp.	49015	5C7601
115	LMT	Lockheed martin corp.	21178	96F126



116	MTB	M&T bank corp.	35554	D1AE3B
117	MANH	Manhattan associates inc.	85992	031025
118	MAN	ManpowerGroup inc.	75285	C0200F
119	MAR	Marriott international inc.	85913	385DD4
120	MMC	Marsh & McLennan cos.	45751	9B5968
121	MCD	McDonald's corp.	43449	954E30
122	MCK	McKesson corp.	81061	4A5C8D
123	MDU	MDU resources group inc.	23835	135B09
124	MRK	Merck & co. inc.	22752	1EBF8D
125	MTOR	Meritor inc	85349	00326E
126	MTG	MGIC investment corp.	76804	E28F22
127	MGM	MGM resorts international	11891	8E8E6E
128	MCHP	Microchip technology inc.	78987	CD FCC9
129	MU	Micron technology inc.	53613	49BBBC
130	MSFT	Microsoft corp.	10107	228D42
131	MOT	Motorola solutions inc.	22779	E49AA3
132	MSM	MSC industrial direct co.	82777	74E288
133	MUR	Murphy oil corp.	28345	949625
134	NBR	Nabors industries ltd.	29102	E4E3B7
135	NOI	National oilwell varco inc.	84032	5D02B7
136	NYT	New york times co.	47466	875F41
137	NFX	Newfield exploration co.	79915	9C1A1F
138	NEM	Newmont mining corp.	21207	911AB8
139	NKE	NIKE inc.	57665	D64C6D
140	NBL	Noble energy inc.	61815	704DAE
141	NOK	Nokia corp.	87128	C12ED9
142	NOC	Northrop grumman corp.	24766	FC1B7B
143	NTRS	Northern trust corp.	58246	3CCC90
144	NUE	NuCor corp.	34817	986AF6
145	ODEP	Office depot inc.	75573	B66928
146	ONB	Old national bancorp	12068	D8760C
147	OMC	Omnicom group inc.	30681	C8257F
148	OTEX	Open text corp.	82833	34E891
149	ORCL	Oracle corp.	10104	D6489C
150	ORBK	Orbotech ltd.	78527	290820
151	PCAR	Paccar inc.	60506	ACF77B
152	PRXL	Parexel international corp.	82607	EF8072
153	PH	Parker hannifin corp.	41355	6B5379
154	PTEN	Patterson-uti energy inc.	79857	57356F
155	PBCT	People's united financial inc.	12073	449A26
156	PEP	PepsiCo inc.	13856	013528
157	PFE	Pfizer inc.	21936	267718
158	PIR	Pier 1 imports inc.	51692	170A6F
159	PXD	Pioneer natural resources co.	75241	2920D5
160	PNCF	PNC financial services group inc.	60442	61B81B
161	POT	Potash corporation of saskatchewan inc.	75844	FFBF74
162	PPG	PPG industries inc.	22509	39FB23
163	PX	Praxair inc.	77768	285175
164	PG	Procter & gamble co.	18163	2E61CC
165	PTC	PTC inc.	75912	D437C3
166	PHM	PulteGroup inc.	54148	7D5FD6
167	QCOM	Qualcomm inc.	77178	CFF15D
168	DGX	Quest diagnostics inc.	84373	5F9CE3
169	RL	Ralph lauren corp.	85072	D69D42
170	RTN	Raytheon co.	24942	1981BF
171	RF	Regions financial corp.	35044	73C521
172	RCII	Rent-a-center inc.	81222	C4FBDC
173	RMD	ResMed inc.	81736	434F38

174	RHI	Robert half international inc.	52230	A4D173
175	RDC	Rowan cos. inc.	45495	3FFA00
176	RCL	Royal caribbean cruises ltd.	79145	751A74
177	RPM	RPM international inc.	65307	F5D059
178	RRD	RR R.R. donnelley & sons co.	38682	0BE0AE
179	SLB	Schlumberger ltd. n.v.	14277	164D72
180	SCTT	Scotts miracle-gro co.	77300	F3FCC3
181	SM	SM st. mary land & exploration co.	78170	6A3C35
182	SONC	Sonic corp.	76568	80D368
183	SO	Southern co.	18411	147C38
184	LUV	Southwest airlines co.	58683	E866D2
185	SWK	Stanley black & decker inc.	43350	CE1002
186	STT	State street corp.	72726	5BC2F4
187	TGNA	TEGNA inc.	47941	D6EAA3
188	TXN	Texas instruments inc.	15579	39BFF6
189	TMK	Torchmark corp.	62308	E90C84
190	TRV	The travelers companies inc.	59459	E206B0
191	TBI	TrueBlue inc.	83671	9D5D35
192	TUP	Tupperware brands corp.	83462	2B0AF4
193	TYC	Tyco international plc	45356	99333F
194	TSN	Tyson foods inc.	77730	AD1ACF
195	X	United states Steel corp.	76644	4E2D94
196	UNH	UnitedHealth group inc.	92655	205AD5
197	VIAV	Viavi solutions inc.	79879	E592F0
198	GWW	W.W. grainger inc.	52695	6EB9DA
199	WDR	Waddell & reed financial inc.	85931	2F24A5
200	WBA	Walgreens boots alliance inc.	19502	FACF19
201	DIS	Walt disney co.	26403	A18D3C
202	WAT	Waters corp.	82651	1F9D90
203	WBS	Webster financial corp.	10932	B5766D
204	WFC	Wells fargo & co.	38703	E8846E
205	WERN	Werner enterprises inc.	10397	D78BF1
206	WABC	Westamerica bancorp	82107	622037
207	WDC	Western digital corp.	66384	CE96E7
208	WHR	Whirlpool corp.	25419	BDD12C
209	WFM	Whole foods market inc.	77281	319E7D
210	XLNX	Xilinx inc.	76201	373E85

Table A.5: Final list of firms – The table contains the information about the full list of firms: tickers, firm names, CRSP PERMNO code and RavenPack ID. Tickers and firm names are taken as of June, 2017. PERMNO and RavenPack ID columns are used to match firms and firm news data.