

LEARNING FROM FORECAST ERRORS: A NEW APPROACH TO FORECAST COMBINATIONS

Tae-Hwy Lee¹ Ekaterina Seregina²

^{1,2}University of California, Riverside

40th International Symposium on Forecasting
October 27, 2020



OUTLINE

1. Motivation
2. Factor Graphical Model
3. Monte Carlo Simulation
4. Application
5. Conclusions

1. Motivation

- ▶ **Competing forecasts** of the univariate series y_t using p forecast models: $\hat{\mathbf{y}}_t = (\hat{y}_{1,t}, \dots, \hat{y}_{p,t})'$, $t = 1, \dots, T$.
- ▶ **Forecast errors**: $\mathbf{e}_t = (e_{1t}, \dots, e_{pt})' \sim \mathcal{N}(\mathbf{0}, \Sigma)$.
- ▶ Let $\Theta = \Sigma^{-1}$ be the *precision matrix*.
- ▶ **Goal**: Find the optimal forecast combination, $\hat{y}_t^c = \mathbf{w}'\hat{\mathbf{y}}_t$, that minimizes the MSFE of the combined forecast error:

$$\hat{e}_t^c = \mathbf{w}'\mathbf{e}_t$$

$$\begin{cases} \min_{\mathbf{w}} \text{MSFE} = \min_{\mathbf{w}} \mathbb{E} \left[\mathbf{w}'\mathbf{e}_t\mathbf{e}_t'\mathbf{w} \right] = \min_{\mathbf{w}} \mathbf{w}'\Sigma\mathbf{w} \\ \text{s.t. } \mathbf{w}'\boldsymbol{\iota}_p = 1, \end{cases}$$

$$\mathbf{w} = \frac{\Theta\boldsymbol{\iota}_p}{\boldsymbol{\iota}_p'\Theta\boldsymbol{\iota}_p}, \quad (1)$$

where $\boldsymbol{\iota}_p$ is a $p \times 1$ vector of ones.

NOTATION

Given a vector $\mathbf{u} \in \mathbb{R}^d$:

- ▶ $\|\mathbf{u}\|_1 = |u_1| + |u_2| + \dots + |u_d|$
- ▶ $\|\mathbf{u}\|_\infty = \max_{1 \leq i \leq d} |u_i|$

Given a matrix $\mathbf{U} \in \mathbb{R}^{p \times p}$:

- ▶ $\|\|\mathbf{U}\|_1 \equiv \max_{1 \leq j \leq p} \sum_{i=1}^p |\mathbf{U}_{i,j}|$ (maximum column sum)
- ▶ $\|\|\mathbf{U}\|_2^2 \equiv \Lambda_{\max}(\mathbf{U}\mathbf{U}')$ (the maximal singular value of \mathbf{U})

Abbreviations:

- ▶ EW – equal-weighted
- ▶ GLASSO and MB – graphical models that do not use factor structure
- ▶ Factor GLASSO and Factor MB – factor graphical models

SUCCESS OF EQUAL-WEIGHTED FORECASTS

- ▶ Let $\text{MSFE}(\mathbf{w}, \Sigma) = \mathbf{w}'\Sigma\mathbf{w}$

Estimation uncertainty in $\mathbf{w} \Rightarrow$ the “optimal” forecast combination is **not guaranteed** to outperform equal weights or improve the individual forecasts (Smith & Wallis, 2009; Claeskens et al., 2016)

$$\left| \text{MSFE}(\widehat{\mathbf{w}}, \widehat{\Sigma}) - \text{MSFE}(\mathbf{w}, \widehat{\Sigma}) \right| \leq \|\widehat{\mathbf{w}} - \mathbf{w}\|_1 \left\| \widehat{\Sigma}\mathbf{w} \right\|_\infty.$$

- ▶ Let $a = \boldsymbol{\iota}'_p \Theta \boldsymbol{\iota}_p / p$, and $\widehat{a} = \boldsymbol{\iota}'_p \widehat{\Theta} \boldsymbol{\iota}_p / p$ (Callot et al., 2019):

$$\|\widehat{\mathbf{w}} - \mathbf{w}\|_1 \leq \frac{a \frac{\|(\widehat{\Theta} - \Theta)\boldsymbol{\iota}_p\|_1}{p} + |a - \widehat{a}| \frac{\|\Theta\boldsymbol{\iota}_p\|_1}{p}}{|\widehat{a}|a},$$

- ▶ Consistent estimator of the precision matrix $\Theta \Rightarrow$ **control the estimation uncertainty in \mathbf{w}**

DO FORECASTERS MAKE COMMON MISTAKES?

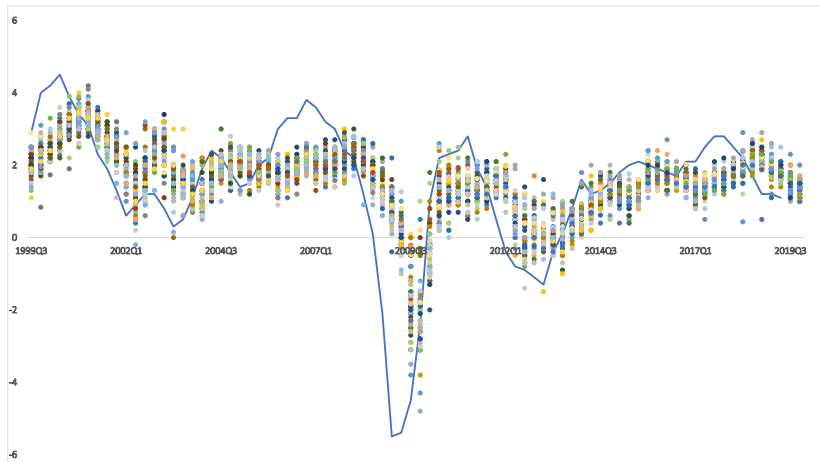


Figure 1: The ECB's SPF: Circles denote quarterly 1-year-ahead forecasts of the Euro-area real GDP growth, YOY percentage change. Blue line - actual series.

- ▶ Forecast errors follow an approximate q -factor model:

$$\underbrace{\mathbf{e}_t}_{p \times 1} = \mathbf{B} \underbrace{\mathbf{f}_t}_{q \times 1} + \varepsilon_t, \quad t = 1, \dots, T$$

- ▶ $\mathbf{f}_t = (f_{1t}, \dots, f_{qt})'$ - factors of the forecast errors.
 - ▶ \mathbf{B} - matrix of factor loadings.
 - ▶ ε_t - idiosyncratic component. Assume $\mathbb{E}[\varepsilon_t | \mathbf{f}_t] = 0$.
- ▶ Notation:

$$\mathbb{E}[\varepsilon_t \varepsilon_t'] = \Sigma_\varepsilon$$

$$\mathbb{E}[\mathbf{f}_t \mathbf{f}_t'] = \Sigma_f$$

$$\mathbb{E}[\mathbf{e}_t \mathbf{e}_t'] = \Sigma = \mathbf{B} \Sigma_f \mathbf{B}' + \Sigma_\varepsilon$$

$$\Theta_\varepsilon = \Sigma_\varepsilon^{-1}, \quad \Theta_f = \Sigma_f^{-1}$$

- ▶ **Challenge:** When forecast errors are driven by common factors, cannot assume sparse Θ .
- ▶ **Question:** How to estimate HD Θ under the factor structure?

EXISTING LITERATURE VS THIS PAPER

Existing Literature:

1. Graphical Models: estimate **precision matrix** directly (Nodewise-Regression by Meinshausen & Bühlmann (MB), 2006; Graphical Lasso (GLASSO) by Friedman et al., 2008).
 - ▶ Assumption: sparse precision matrix.
2. Factor Models:

$$\underbrace{\mathbf{e}_t}_{p \times 1} = \mathbf{B} \underbrace{\mathbf{f}_t}_{q \times 1} + \varepsilon_t, \quad t = 1, \dots, T \quad (2)$$

Idea: estimate **covariance matrix** using Eq (2), invert it.

This Paper: how to use graphical models under the factor structure to estimate Θ for the estimation of the optimal forecast combination weights, \mathbf{w} .

2. Factor Graphical Model (FGM)

GRAPHICAL LASSO

- ▶ Given a sample $\{\mathbf{e}_t\}_{t=1}^T$, let $\mathbf{S} = (1/T) \sum_{t=1}^T (\mathbf{e}_t)(\mathbf{e}_t)'$ denote the sample covariance matrix.
- ▶ Let $\mathbf{W} = \mathbf{S} + \lambda \mathbf{I}$ and $\widehat{\mathbf{D}}^2 \equiv \text{diag}(\mathbf{W})$;
- ▶ Weighted penalized log-likelihood (Jankova & van de Geer, 2018):

$$\widehat{\Theta} = \arg \min_{\Theta = \Theta'} \text{trace}(\mathbf{W}\Theta) - \log \det(\Theta) + \lambda \sum_{i \neq j} \widehat{\mathbf{D}}_{ii} \widehat{\mathbf{D}}_{jj} |\Theta_{ij}|, \quad (3)$$

Idea of GL: Complete columns of Θ using the gradient of Eq (3)

NODEWISE REGRESSION

- ▶ Let \mathbf{e}_j be a $T \times 1$ vector of observations for the j -th regressor
- ▶ The remaining covariates are collected in a $T \times p$ matrix \mathbf{E}_{-j} .

For each $j = 1, \dots, p$ we run the following Lasso regressions:

$$\hat{\gamma}_j = \arg \min_{\gamma \in \mathbb{R}^{p-1}} \left(\|\mathbf{e}_j - \mathbf{E}_{-j}\gamma\|_2^2 / T + 2\lambda_j \|\gamma\|_1 \right), \quad (4)$$

where $\hat{\gamma}_j = \{\hat{\gamma}_{j,k}; j = 1, \dots, p, k \neq j\}$.

- ▶ For $j = 1, \dots, p$, define

$$\hat{\tau}_j^2 = \|\mathbf{e}_j - \mathbf{E}_{-j}\hat{\gamma}_j\|_2^2 / T + \lambda_j \|\hat{\gamma}_j\|_1 \quad (5)$$

NODEWISE REGRESSION

- Define

$$\hat{\mathbf{C}} = \begin{pmatrix} 1 & -\hat{\gamma}_{1,2} & \cdots & -\hat{\gamma}_{1,p} \\ -\hat{\gamma}_{2,1} & 1 & \cdots & -\hat{\gamma}_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{\gamma}_{p,1} & -\hat{\gamma}_{p,2} & \cdots & 1 \end{pmatrix}$$

and write

$$\hat{\mathbf{T}}^2 = \text{diag}(\hat{\tau}_1^2, \dots, \hat{\tau}_p^2)$$

- The approximate inverse is defined as

$$\hat{\Theta} = \hat{\mathbf{T}}^{-2} \hat{\mathbf{C}}. \quad (6)$$

FGM

- ▶ **Forecast errors:** $\mathbf{e}_t = (e_{1t}, \dots, e_{pt})' \sim \mathcal{N}(\mathbf{0}, \Sigma)$

$$\mathbf{e}_t = \mathbf{B}\mathbf{f}_t + \boldsymbol{\varepsilon}_t, \quad t = 1, \dots, T$$

$$\Sigma = \mathbf{B}\Sigma_f\mathbf{B}' + \Sigma_\varepsilon$$

$$\Theta = \Sigma^{-1}, \quad \Theta_\varepsilon = \Sigma_\varepsilon^{-1}, \quad \Theta_f = \Sigma_f^{-1}$$

- ▶ **Goal:** find the optimal forecast combination weights

$$\mathbf{w} = \frac{\Theta \boldsymbol{\iota}_p}{\boldsymbol{\iota}_p' \Theta \boldsymbol{\iota}_p}.$$

- ▶ **Challenge:** when factors are present, the precision matrix of forecast errors cannot be sparse.

FGM

$$\widehat{\Sigma}_\varepsilon = \frac{1}{T} \sum_{t=1}^T (\widehat{\varepsilon}_t - \bar{\varepsilon})(\widehat{\varepsilon}_t - \bar{\varepsilon})'; \quad \widehat{\Theta}_\varepsilon \leftarrow \text{Gr.Mdl: GLASSO or MB,}$$

$$\widehat{\Sigma}_f = \frac{1}{T} \sum_{t=1}^T (\widehat{\mathbf{f}}_t - \bar{\mathbf{f}})(\widehat{\mathbf{f}}_t - \bar{\mathbf{f}})'; \quad \widehat{\Theta}_f = \widehat{\Sigma}_f^{-1},$$

- **Solution:** use Sherman-Morrison-Woodbury (SMW) formula to estimate the precision of forecast errors:

$$\text{FGr. Mdl} \rightarrow \widehat{\Theta} = \underbrace{\widehat{\Theta}_\varepsilon}_{\text{Gr. Mdl}} - \underbrace{\widehat{\Theta}_\varepsilon \widehat{\mathbf{B}}}_{\text{F.Mdl}} \underbrace{[\widehat{\Theta}_f + \widehat{\mathbf{B}}' \widehat{\Theta}_\varepsilon \widehat{\mathbf{B}}]^{-1}}_{\text{F.Mdl}} \underbrace{\widehat{\mathbf{B}}'}_{\text{F.Mdl}} \widehat{\Theta}_\varepsilon.$$

$$\widehat{\mathbf{w}} = \frac{\widehat{\Theta} \boldsymbol{\nu}_p}{\boldsymbol{\nu}_p' \widehat{\Theta} \boldsymbol{\nu}_p},$$

- If Gr. Mdl \equiv GL \Rightarrow Factor GLASSO;
- If Gr. Mdl \equiv MB \Rightarrow Factor MB

4. Monte Carlo Simulation

THEORETICAL RESULTS: SUMMARY

Recall:

$$\|\widehat{\mathbf{w}} - \mathbf{w}\|_1 \leq \frac{a \frac{\|(\widehat{\Theta} - \Theta)\boldsymbol{\nu}_p\|_1}{p} + |a - \widehat{a}| \frac{\|\Theta\boldsymbol{\nu}_p\|_1}{p}}{|\widehat{a}|a},$$

- ▶ **Consistency of Factor GLASSO** (under certain sparsity restrictions on Θ_ε): $\left\| \widehat{\Theta} - \Theta \right\|_\eta = o_P(1)$, $\eta = 1, 2$ (Lee, Seregina, 2020);
- ▶ **Consistency of Factor MB** (under certain sparsity restrictions on $\Theta_{j,\varepsilon}$): $\max_{1 \leq j \leq p} \left\| \widehat{\Theta}_j - \Theta_j \right\|_\eta = o_P(1)$, $\eta = 1, 2$ (Seregina, 2020)

$$\Rightarrow \|\widehat{\mathbf{w}} - \mathbf{w}\|_1 = o_P(1)$$

DGP1 FOR ESTIMATION

$$\mathbf{e}_t = (e_{1t}, \dots, e_{pt})' \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

$$\mathbf{f}_t = \phi_f \mathbf{f}_{t-1} + \zeta_t$$

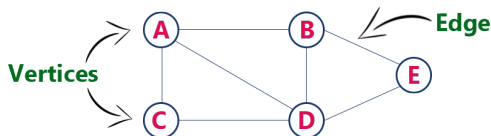
$$\underbrace{\mathbf{e}_t}_{p \times 1} = \mathbf{B} \underbrace{\mathbf{f}_t}_{q \times 1} + \varepsilon_t, \quad t = 1, \dots, T$$

- ▶ \mathbf{f}_t - $q \times 1$ vector of factors, $\phi_f = 0.2$.
- ▶ $\zeta_t \sim \mathcal{N}(0, 1)$, $\varepsilon_t \sim \mathcal{N}(0, \Sigma_\varepsilon)$, with sparse Θ_ε that has a **random graph structure** (next slide).
- ▶ \mathbf{B} : the first q columns of an upper triangular matrix from a Cholesky decomposition of the $p \times p$ Toeplitz matrix:

$$\mathbf{Q} = (\mathbf{Q})_{ij}, \text{ where } (\mathbf{Q})_{ij} = \rho^{|i-j|}, i, j \in 1, \dots, p; \rho = 0.2.$$

- ▶ Set $p = T^{0.85}$, $q = 2(\log(T))^{0.5}$, $T = \lceil 2^\kappa \rceil$, $\kappa = 7, 7.5, 8, \dots, 9.5$.

DGP₁ FOR ESTIMATION



- ▶ **Random graph structure** (Erdős–Rényi model) for Θ_ε
Let \mathbf{A}_ε be a $p \times p$ adjacency matrix:

$$\mathbf{A}_\varepsilon^{ij} = \begin{cases} 1, & \text{for } i \neq j \text{ with probability } \pi, \\ 0, & \text{otherwise.} \end{cases}$$

edges in a graph $\equiv s_T = p(p-1)\pi/2$. To control sparsity, we set $\pi = 1/(pT^{0.8}) \Rightarrow s_T = \mathcal{O}(T^{0.05})$.

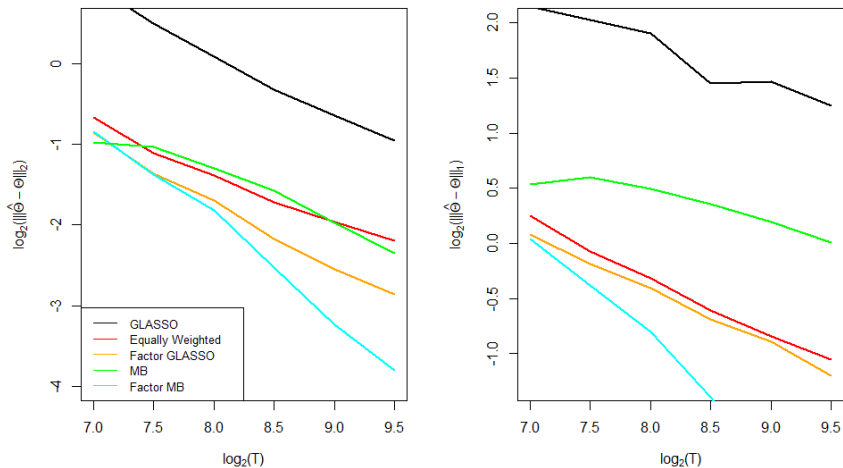


Figure 2: Averaged errors of the estimators of Θ on logarithmic scale (base 2): $p = T^{0.85}$, $q = 2(\log(T))^{0.5}$, $s_T = \mathcal{O}(T^{0.05})$.

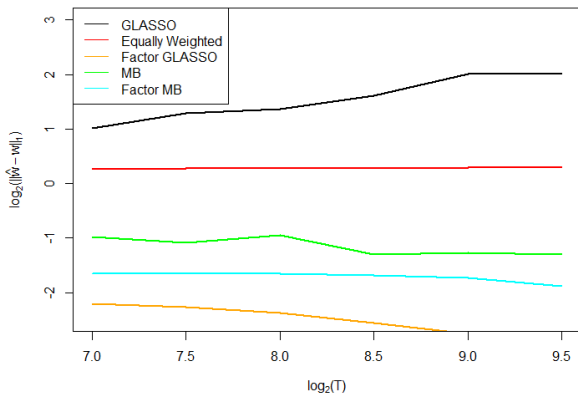


Figure 3: Averaged errors of the estimator of w (base 2) on logarithmic scale: $p = T^{0.85}$, $q = 2(\log(T))^{0.5}$, $s_T = \mathcal{O}(T^{0.05})$.

DGP2 FOR FORECASTING

$$\mathbf{x}_t = \Lambda \mathbf{g}_t + \mathbf{v}_t$$

$$\mathbf{g}_t = \phi \mathbf{g}_{t-1} + \boldsymbol{\xi}_t$$

$$y_{t+1} = \mathbf{g}'_t \boldsymbol{\alpha} + \sum_{s=1}^{\infty} \theta_s \epsilon_{t+1-s} + \epsilon_{t+1}$$

$$\theta_s = (1 + s)^{c_1} c_2^s, \quad c_1 \in \{0, 0.75\} \text{ and } c_2 \in \{0.6, 0.7, 0.8, 0.9\}$$

- ▶ \mathbf{x}_t - $N \times 1$ vector of predictors.
- ▶ \mathbf{g}_t - $r \times 1$ vector of factors.
- ▶ $\mathbf{v}_t \sim \mathcal{N}(0, \sigma_v^2)$, $\boldsymbol{\xi}_t \sim \mathcal{N}(0, \sigma_\xi^2)$, $\epsilon_{t+1} \sim \mathcal{N}(0, 1)$, $\boldsymbol{\alpha} \sim \mathcal{N}(1, 1)$.
- ▶ Λ : the first r rows of an upper triangular matrix from a Cholesky decomposition of the $N \times N$ Toeplitz matrix parameterized by ρ .

MODEL

- ▶ Factor-augmented autoregressive models of orders k, l , FAR(k, l):

$$\hat{y}_{t+1} = \hat{\mu} + \hat{\kappa}_1 \hat{g}_{1,t} + \cdots + \hat{\kappa}_k \hat{g}_{k,t} + \hat{\psi}_1 y_t + \cdots + \hat{\psi}_l y_{t+1-l},$$

where $k = 0, 1, \dots, K$ and $l = 0, 1, \dots, L$.

- ▶ The total number of forecasting models is:

$$p = (1 + K) \times (1 + L)$$

- ▶ Forecast errors:

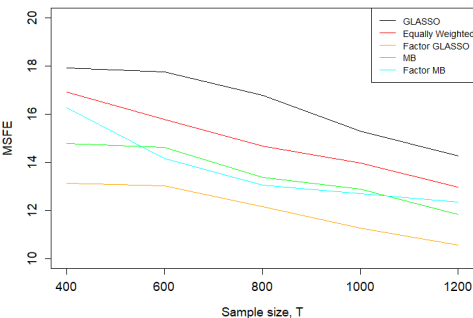
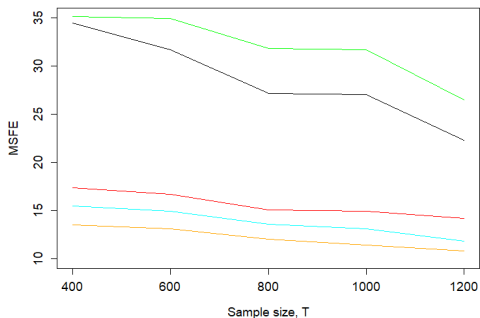
$$\underbrace{\mathbf{e}_t}_{p \times 1} = \mathbf{B} \underbrace{\mathbf{f}_t}_{q \times 1} + \boldsymbol{\varepsilon}_t, \quad t = 1, \dots, T$$

- ▶ Training sample: $m = T/2$. Test sample: $t = m, \dots, T - 1$.

- ▶ MSFE = $\frac{1}{T - m} \sum_{t=m}^{T-1} (y_{t+1} - \hat{\mathbf{w}}' \hat{\mathbf{y}}_t)^2$.

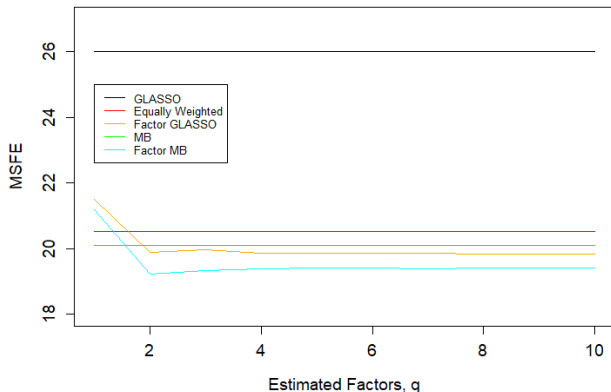
Plots of the MSFE over the sample size T

$c_1 \in \{0, 0.75\}$, $c_2 = 0.9$, $N = 100$, $r = 5$, $\sigma_\xi = 1$, $L = 7$,
 $K = 2$, $p = 24$, $q = 5$, $\rho = 0.9$, $\phi = 0.8$

(a) $c_1 = 0.75$ (b) $c_1 = 0$

Plots of the MSFE over the values of q

$c_1 = 0.75$, $c_2 = 0.9$, $T = 800$, $N = 100$, $r = 5$, $\sigma_\xi = 1$,
 $L = 12$, $K = 0$, $p = 13$, $q \in \{0, 1, \dots, 10\}$, $\rho = 0.9$, $\phi = 0.8$.



4. Application

Data

- ▶ McCracken and Ng (2016), FRED-MD, monthly, 1959:1-2020:07, $T = 726$
- ▶ $m = 120$, train sample, rolling windows
- ▶ $n \equiv T - m - h + 1$, test sample $t = m, \dots, T - h$
- ▶ Number of regressors in \mathbf{X} , $N = 128$

Models

- ▶ FAR(k, l) with $k = 0, 1, \dots, K = 9$, and $l = 0, 1, \dots, L = 11$
- ▶ Total number of forecasting models $p = 120$
- ▶ h -step-ahead forecasts ($h = 1, 2, 3, 4$)

Series for Forecasting

Let $\{Y_t\}_{t=1}^T$ be the series of interest for forecasting
(Coulombe et al. (2020))

- ▶ INDPROD and S&P500:

$$y_{t+h}^{(h)} = \frac{1}{h} \ln(Y_{t+h}/Y_t).$$

- ▶ UNRATE:

$$y_{t+h}^{(h)} = \frac{1}{h} (Y_{t+h}/Y_t).$$

- ▶ FEDFUNDS:

$$y_{t+h}^{(h)} = \ln(Y_{t+h}).$$

PREDICTION OF MONTHLY INDPROD AND S&P500

INDPROD					
h	EW	GLASSO	Factor GLASSO	MB	Factor MB
1	2.77E-04	1.51E-04	1.24E-04	2.23E-04	1.28E-04
2	3.26E-04	1.79E-04	5.59E-05	1.61E-04	1.38E-04
3	1.55E-04	9.77E-05	3.81E-05	1.17E-04	6.54E-05
4	1.18E-04	7.60E-05	2.38E-05	1.03E-04	2.65E-05
S&P500					
1	1.40E-03	1.39E-03	1.37E-03	1.34E-03	9.57E-03
2	1.71E-03	1.44E-03	8.95E-04	1.55E-03	1.01E-03
3	1.66E-03	1.34E-03	3.48E-04	1.43E-03	6.69E-04
4	1.27E-03	1.06E-03	3.95E-04	9.55E-04	7.91E-04

$$\text{MSFE} = \frac{1}{T - h - m + 1} \sum_{t=m}^{T-h} (y_{t+h}^h - \widehat{\mathbf{w}}' \widehat{\mathbf{y}}_t)^2$$

EW stands for the “Equal-Weighted” forecast, GLASSO and MB are the models that do not use the factor structure in the forecast errors. Factor GLASSO and Factor MB are our proposed Factor Graphical Models.

PREDICTION OF MONTHLY UNRATE AND FEDFUNDS

UNRATE					
h	EW	GLASSO	Factor GLASSO	MB	Factor MB
1	0.2531	0.0858	0.0109	0.0557	0.0107
2	0.3758	0.1334	0.0066	0.0448	0.0081
3	0.0743	0.0651	0.0066	0.0532	0.0051
4	2.1999	0.6871	0.1578	1.0973	0.2510
FEDFUNDS					
1	0.0609	0.1813	0.0205	0.0424	0.0448
2	0.1426	1.2230	0.0288	0.0675	0.0416
3	0.2354	1.2710	0.0508	0.1217	0.1038
4	0.3702	1.4672	0.0592	0.2470	0.1962

$$\text{MSFE} = \frac{1}{T - h - m + 1} \sum_{t=m}^{T-h} (y_{t+h}^h - \widehat{\mathbf{w}}' \widehat{\mathbf{y}}_t)^2$$

EW stands for the “Equal-Weighted” forecast, GLASSO and MB are the models that do not use the factor structure in the forecast errors. Factor GLASSO and Factor MB are our proposed Factor Graphical Models.

5. Conclusions

CONCLUSIONS

1. Learning from Forecast Errors:

- ▶ Different forecast models tend to make the same (common) mistakes.
- ▶ Forecast errors are driven by common factors.
- ▶ We cannot assume that the precision matrix of forecast errors is sparse.

2. A New Approach to Forecast Combinations:

- ▶ We decompose the forecast errors into the common and idiosyncratic errors.
- ▶ We assume the sparsity on the precision matrix of the idiosyncratic forecast errors.
- ▶ We develop the novel algorithm, Factor Graphical Models, for forecast combinations.

CONCLUSIONS

3. Simulation and Application:

- ▶ Factor GLASSO and Factor MB consistently estimate precision matrix of forecast errors and optimal combination weights.
- ▶ Factor GLASSO outperforms GLASSO, Factor MB outperforms MB.
- ▶ Both outperform EW.

Work in Progress: Time-Varying Factor Graphical Models, portfolio application.

thank you!

Questions? Please contact me at esere001@ucr.edu and I will be happy to address any questions.

More Info? Please visit my website at seregina.info.