

# Forecasting With Many predictors

The Econometrics of Predictability

*This version: June 4, 2014*

June 3, 2014

- Dynamic Factor Models
- The 3-Pass Regression Filter
- Regularized Reduced Rank Regression
- Time permitting
  - Bagging
  - Filters and decompositions

## How Many is Many?

- Many here means 25 or more
- Often many more, 100s of series

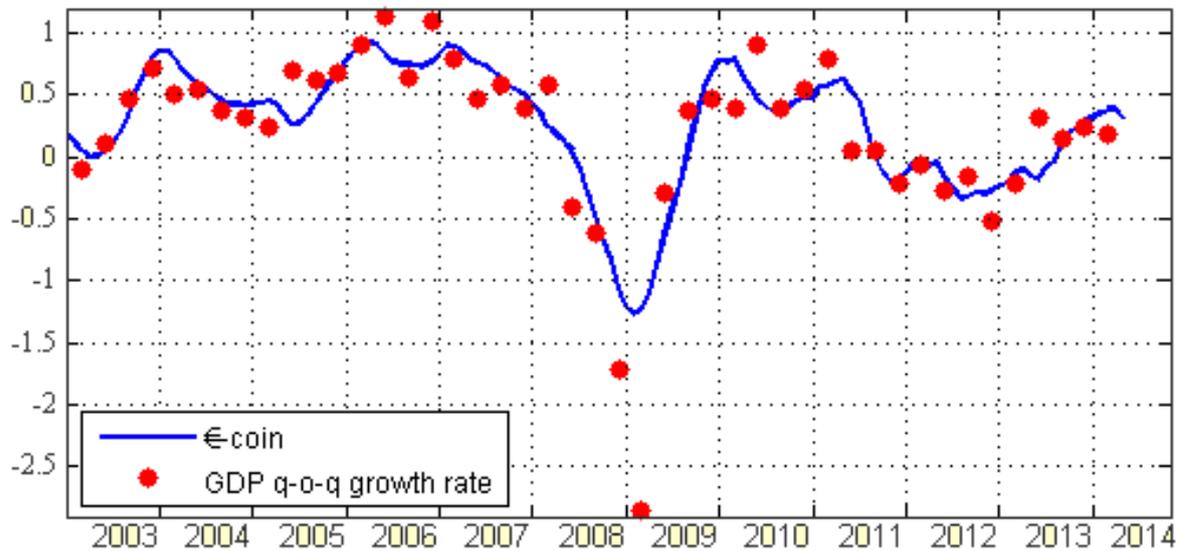
## Why factor models

- Are parsimonious while effectively including many regressors
- Can remove measurement error or other useless information from predictors
- Factor may be of interest
  - Leading indicators:
    - ▷ €-coin
    - ▷ Chicago Fed National Activity Index
    - ▷ Aruoba-Diebold-Scotti Business Conditions Index
  - Real and Nominal factors
  - Global and Local factors

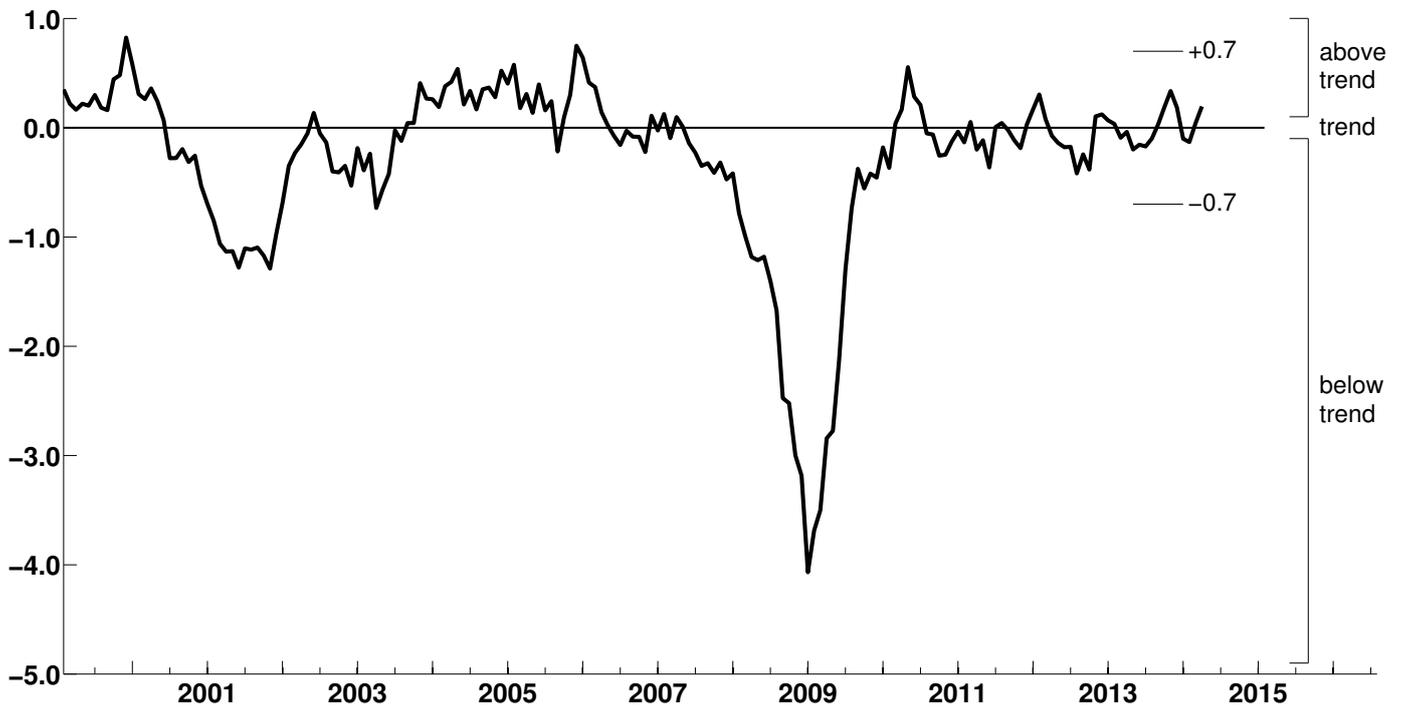
- European Coincident Indicator
- First factor in a Europe-wide model

## €-coin: the Euro Area Economy in One Figure – May 2014

€-coin and euro-area GDP



- Factor extracted from 85 series
- Based on research in forecasting inflation



- Based on factor model in Aruoba, Diebold & Scotti
- Extracts common factor in:
  - weekly initial jobless claims
  - monthly payroll employment
  - industrial production
  - personal income less transfer payments, manufacturing and trade sales
  - quarterly real GDP

## The Model

- Scalar *latent* factor

$$x_t = \sum_{i=1}^q \rho_i x_{t-i} + \eta_i$$

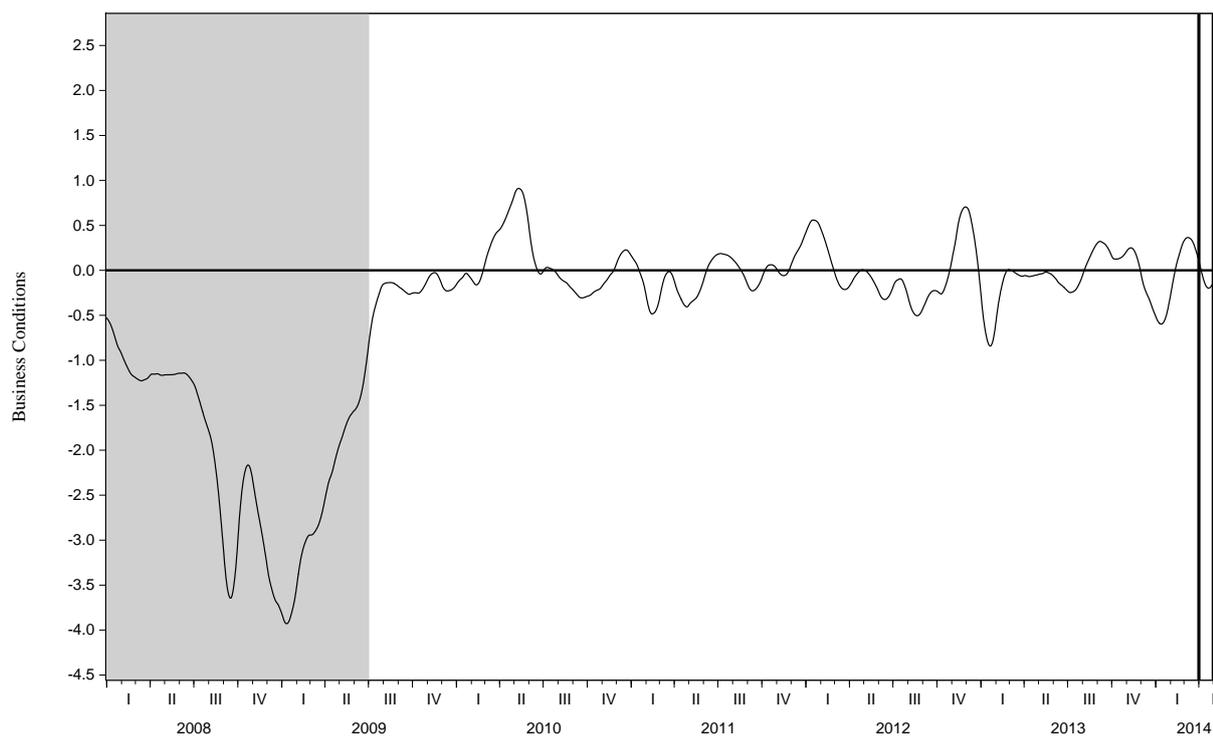
- Indicators

$$y_{it} = c_i + \beta_i x_t + \sum_{j=1}^{p_i} \gamma_j y_{it-\Delta_j} + \epsilon_i$$

- $\Delta_j$  allows series to have different observational frequencies



## Aruoba-Diebold-Scotti Business Conditions Index ( 12/31/2007- 05/24/2014)



- $T$  number of time series observations
- $k$  number of series available to forecast
- $\mathbf{y}_t$  series to be forecast,  $m$  by 1
  - $m$  will often be 1
- $\mathbf{x}_t$  series used to forecast,  $k$  by 1
  - Usually assume  $E[\mathbf{x}_t] = \mathbf{0}$  and  $\text{Cov}[\mathbf{x}_t] = \mathbf{I}_k$
  - Demeaned and standardized
  - Suppose  $\mathbf{x}_t = \Sigma_{\mathbf{x}}^{-1/2} (\tilde{\mathbf{x}}_t - \boldsymbol{\mu}_{\mathbf{x}})$
- $\mathbf{f}_t$  factors,  $r$  by 1
- $\mathbf{x}_t$  *may be*  $\mathbf{y}_t$ , but not necessarily
  - $\mathbf{y}_t$  could be subset of  $\mathbf{x}_t$  (common)
  - $\mathbf{y}_t$  could be excluded from factor estimation (uncommon)

- Factor models help avoid issues with large, kitchen-sink models
- Consider issue of parameter estimation error when forecasting
- Suppose correct model is linear

$$y_{t+1} = \boldsymbol{\beta} \mathbf{x}_t + \epsilon_t$$

- Forecast using OLS estimates is then

$$\begin{aligned} \hat{y}_{t+1|t} &= \hat{\boldsymbol{\beta}} \mathbf{x}_t \\ &= (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} + \boldsymbol{\beta}) \mathbf{x}_t \\ &= \underbrace{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \mathbf{x}_t}_{\text{estimation error}} + \underbrace{\boldsymbol{\beta} \mathbf{x}_t}_{\text{correct forecast}} \end{aligned}$$

- Suppose  $\epsilon_t, \mathbf{x}_t$  are independent and jointly normally distributed

$$\text{Cov} \begin{bmatrix} \epsilon_t \\ \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \sigma_\epsilon^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k \end{bmatrix}$$

- Standard assumptions have  $k$  fixed, so as  $T \rightarrow \infty$ ,  $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \xrightarrow{p} \mathbf{0}$

$$\hat{y}_{t+1|t} \sim N(\boldsymbol{\beta} \mathbf{x}_t, 0)$$

- Degenerate normal - no error since  $\boldsymbol{\beta}$  is effectively *known*
- What about the case when  $k$  is large
- Use *diagonal asymptotics*,  $k/T \rightarrow c$ ,  $0 < \underline{\kappa} < c < \bar{\kappa} < \infty$
- In this case

$$\hat{y}_{t+1|t} \sim N(\boldsymbol{\beta} \mathbf{x}_t, k/T \times \sigma_\epsilon^2)$$

- Is still random, even when  $T \rightarrow \infty$
- True even if all  $\boldsymbol{\beta} = \mathbf{0}$ !

- When the number of parameters is large, then almost all coefficients must be 0

$$y_t = \sum_{i=1}^k \beta_i x_{t,i} + \epsilon_i$$

- Variance of the LHS is the same as the RHS

$$V[y_t] = \sum_{i=1}^k \beta_i^2 + \sigma_\epsilon^2$$

- If  $k \rightarrow \infty$ ,  $\inf_i |\beta_i| > \underline{\kappa} > 0$ , then  $V[y_t] \rightarrow \infty$
- Even when  $T$  is very large, it will not usually make sense to have  $k$  extremely large
- Factor models will effectively have small  $\beta_i$  coefficient, only using two steps
  1. Construct average-like estimators of factors from  $\mathbf{x}_t$  – coefficients are  $O(1/k)$
  2. Weight these using a small number of relatively large coefficients

- Consider the cross-section of asset returns
- Model uses factors as RHS variables

$$x_{it} = \sum_{j=1}^r \lambda_{ij} f_{jt} + \epsilon_{it}$$

- $\lambda_{ij}$  are the factor loadings for series  $i$ , factor  $j$
- $\epsilon_{it}$  is the idiosyncratic error for series  $i$
- In vector notation,

$$\mathbf{x}_t = \mathbf{\Lambda} \mathbf{f}_t + \boldsymbol{\epsilon}_t$$

$k \times 1 \quad k \times r \quad r \times 1 \quad r \times 1$

- $\mathbf{\Lambda}$  is  $k$  by  $r$
- $\mathbf{f}_t$  is  $r$  by 1

- In matrix notation,

$$\mathbf{X} = \mathbf{F} \mathbf{\Lambda}' + \boldsymbol{\epsilon}$$

$T \times k \quad T \times r \quad r \times k \quad T \times k$

- $\mathbf{X}$  is  $T$  by  $k$
- $\mathbf{F}$  is  $T$  by  $r$
- $\boldsymbol{\epsilon}$  is  $k$  by 1
- When model is a strict (as opposed to approximate),  $E[\boldsymbol{\epsilon}_t] = \mathbf{0}$  and  $E[\boldsymbol{\epsilon}_t \boldsymbol{\epsilon}_t'] = \boldsymbol{\Sigma}_\epsilon = \text{diag}(\sigma_1^2, \dots, \sigma_m^2)$
- Covariance of  $\mathbf{x}_t$  is then

$$\boldsymbol{\Lambda} \boldsymbol{\Omega} \boldsymbol{\Lambda}' + \boldsymbol{\Sigma}_\epsilon$$

- $\boldsymbol{\Omega} = \text{Cov}[\mathbf{f}_t]$ ,  $r$  by  $r$
- Covariance will play a crucial role in estimation of factors

- Principal components can be used to estimate factors
- Formally, problem is

$$\min_{\beta, \mathbf{f}_1, \dots, \mathbf{f}_T} \sum_{t=1}^T (\mathbf{x}_t - \beta \mathbf{f}_t)' (\mathbf{x}_t - \beta \mathbf{f}_t) \text{ subject to } \beta' \beta = \mathbf{I}_r$$

- $\beta$  is  $k$  by  $r$ 
  - $\beta$  is related to but different from  $\Lambda$
  - $\Lambda$  is the DGP parameter
  - $\beta$  is a normalized and *rotated* version of  $\Lambda$

## Definition (Rotation)

A square matrix  $\mathbf{B}$  is said to be a rotation of a square matrix  $\mathbf{A}$  if  $\mathbf{B} = \mathbf{Q}\mathbf{A}$  and  $\mathbf{Q}\mathbf{Q}' = \mathbf{Q}'\mathbf{Q} = \mathbf{I}$ .

- $\mathbf{f}_t$  is  $r$  by 1
- $\beta' \beta = \mathbf{I}_r$  is a *normalization*, and is required
  - $\beta \mathbf{f}_t = ((\beta/2)(2\mathbf{f}_t))$
  - Generally, for full rank  $\mathbf{Q}$ ,  $(\beta \mathbf{Q})(\mathbf{Q}^{-1} \mathbf{f}_t) = \tilde{\beta} \tilde{\mathbf{f}}_t$

- If  $\beta$  was observable, solution would be OLS

$$\hat{\mathbf{f}}_t = (\beta' \beta)^{-1} \beta' \mathbf{x}_t$$

This can be substituted into the objective function

$$\sum_{t=1}^T (\mathbf{x}_t - \beta (\beta' \beta)^{-1} \beta' \mathbf{y}_t)' (\mathbf{x}_t - \beta (\beta' \beta)^{-1} \beta' \mathbf{x}_t) = \sum_{t=1}^T \mathbf{x}_t' (\mathbf{I} - \beta (\beta' \beta)^{-1} \beta') \mathbf{x}_t$$

- This works since  $\mathbf{I} - \beta (\beta' \beta)^{-1} \beta'$  is *idempotent*
  - $\mathbf{A}\mathbf{A} = \mathbf{A}$
- Some additional manipulation using the trace operator on a scalar leads to two equivalent expressions

$$\begin{aligned} \min_{\beta} \sum_{t=1}^T \mathbf{x}_t' (\mathbf{I} - \beta (\beta' \beta)^{-1} \beta') \mathbf{x}_t &= \max_{\beta} \text{tr} \left( (\beta' \beta)^{-1/2} \beta' \Sigma_{\mathbf{x}} \beta (\beta' \beta)^{-1/2} \right) \\ &= \max_{\beta} \beta' \Sigma_{\mathbf{x}} \beta \end{aligned}$$

- All subject to  $\beta' \beta = \mathbf{I}_r$
- Solution to last problem sets  $\beta$  to the *eigenvectors* of  $\Sigma_{\mathbf{x}}$

## Definition (Eigenvalue)

The eigenvalues of a real, symmetric matrix  $k$  by  $k$  matrix  $\mathbf{A}$  are the  $k$  solutions to

$$|\lambda \mathbf{I}_k - \mathbf{A}| = 0$$

where  $|\cdot|$  is the determinant.

### ■ Properties of eigenvalues

- ▶  $\det \mathbf{A} = \prod_{i=1}^r \lambda_i$
- ▶  $\text{tr} \mathbf{A} = \sum_{i=1}^r \lambda_i$
- ▶ For positive (semi) definite  $\mathbf{A}$ ,  $\lambda_i > 0$ ,  $i = 1, \dots, r$  ( $\lambda_i \geq 0$ )
- ▶ Rank
  - ▷ Full-rank  $\mathbf{A}$  implies  $\lambda_i \neq 0$ ,  $i = 1, \dots, r$
  - ▷ Rank  $q < r$  matrix  $\mathbf{A}$  implies  $\lambda_i \neq 0$ ,  $i = 1, \dots, q$  and  $\lambda_j = 0$ ,  $j = q + 1, \dots, r$

## Definition (Eigenvector)

An a  $k$  by 1 vector  $\mathbf{u}$  is an eigenvector corresponding to an eigenvalue  $\lambda$  of a real, symmetric matrix  $k$  by  $k$  matrix  $\mathbf{A}$  if

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$$

- Properties of eigenvectors
  - If  $\mathbf{A}$  is positive definite, then

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}'$$

where  $\mathbf{\Lambda}$  is diagonal and  $\mathbf{V}\mathbf{V}' = \mathbf{V}'\mathbf{V} = \mathbf{I}$

## Definition (Orthonormal Matrix)

A  $k$ -dimensional orthonormal matrix  $\mathbf{U}$  satisfies  $\mathbf{U}'\mathbf{U} = \mathbf{I}_k$ , and so  $\mathbf{U}' = \mathbf{U}^{-1}$ .

- Implication is

$$\mathbf{V}'\mathbf{A}\mathbf{V} = \mathbf{V}'\mathbf{V}\mathbf{\Lambda}\mathbf{V}'\mathbf{V} = \mathbf{\Lambda}$$

- $\mathbf{X}$  is  $T$  by  $k$  (assume demeaned)
- $\mathbf{X}'\mathbf{X}$  is real and symmetric with eigenvalues  $\Lambda = \text{diag}(\lambda_i)_{i=1,\dots,k}$
- Factors are estimated

$$\mathbf{X}'\mathbf{X} = \mathbf{V}\Lambda\mathbf{V}'$$

$$\mathbf{V}'\mathbf{X}'\mathbf{X}\mathbf{V} = \mathbf{V}'\mathbf{V}\Lambda\mathbf{V}'\mathbf{V}$$

$$(\mathbf{XV})'(\mathbf{XV}) = \Lambda \text{ since } \mathbf{V}' = \mathbf{V}^{-1}$$

$$\mathbf{F}'\mathbf{F} = \Lambda.$$

- $\mathbf{F} = \mathbf{XV}$  is the  $T$  by  $k$  matrix of factors
- $\boldsymbol{\beta} = \mathbf{V}'$  is the  $k$  by  $k$  matrix of factor loadings.
- All factors exactly reconstruct  $\mathbf{Y}$

$$\mathbf{F}\boldsymbol{\beta} = \mathbf{FV}' = \mathbf{YV}\mathbf{V}' = \mathbf{Y}$$

▸ Assumes  $k$  is large

- Note that both factors *and* loadings are orthogonal since

$$\mathbf{F}'\mathbf{F} = \Lambda \text{ and } \boldsymbol{\beta}'\boldsymbol{\beta} = \mathbf{I}$$

- Only loadings are normalized

- Consider simple example where

$$x_{it} = 1 \times f_t + \epsilon_{it}$$

- $f_t$  and  $\epsilon_{it}$  are all independent, standard normal
- Covariance of  $\mathbf{x}$  is  $\Sigma_{\mathbf{x}} = 1 + I_k$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

- First eigenvector is

$$(k^{-1/2}, k^{-1/2}, \dots, k^{-1/2})$$

- Form is due to normalization

$$\sum_{i=1}^k v_{ij}^2 = 1, \quad \sum_{i=1}^k v_{ij}v_{in} = 0$$

- $\sum_{i=1}^k (k^{-1/2})^2 = \sum_{i=1}^k k^{-1} = k k^{-1} = 1$

- Estimated factor is then

$$\hat{f}_t = \sum_{i=1}^k k^{-1/2} x_{it} = k^{1/2} \left( \frac{1}{k} \sum x_{it} \right) = k^{1/2} \bar{x} = \sum_{i=1}^k w_i x_i$$

- What about  $\bar{x}$

$$\begin{aligned} \bar{x} &= k^{-1} \left( \sum_{i=1}^k f_t + \epsilon_{it} \right) \\ &= \hat{f}_t + \bar{\epsilon}_t \\ &\approx \hat{f}_t \end{aligned}$$

- Normalization means factor is  $O_p(k^{1/2})$ 
  - Can always re-normalize factor to be  $O_p(1)$  using  $\hat{f}_t/k^{1/2}$
- Key assumption is that  $\bar{\epsilon}_t$  follows some form of LLN *in*  $k$
- In strict factor model, no correlation so simple

- Strict factor models require strong assumptions

$$\text{Cov}(\epsilon_{it}, \epsilon_{js}) = 0 \quad i \neq j, s \neq t$$

- These are easily rejectable in practice
- **Approximate Factor Models** relax these assumptions and allow:
  - (*Weak*) Serial correlation in  $\epsilon_t$

$$\sum_{s=0}^{\infty} |\gamma_s| < \infty$$

- (*Weak*) Cross-sectional correlation between  $\epsilon_{it}$  and  $\epsilon_{jt}$

$$\lim_{k \rightarrow \infty} \sum_{i \neq j}^k \mathbb{E} |\epsilon_{it} \epsilon_{jt}| < \infty$$

- Heteroskedasticity in  $\epsilon$
- Requires pervasive factors

$$\begin{aligned} \mathbf{x}_t &= \mathbf{\Lambda} \mathbf{f}_t + \boldsymbol{\epsilon}_t \\ \lim_{k \rightarrow \infty} \text{rank}(k^{-1} \mathbf{\Lambda}' \mathbf{\Lambda}) &= r \end{aligned}$$

- Key input for factor estimation is  $\Sigma_{\mathbf{x}}$ 
  - In most theoretical discussions of PCA, this is the covariance

$$\Sigma_{\mathbf{x}} = T^{-1} \sum_{t=1}^T (\mathbf{x}_t - \hat{\boldsymbol{\mu}})(\mathbf{x}_t - \hat{\boldsymbol{\mu}})$$

- Two other simple versions are used
  - Outer-product

$$T^{-1} \mathbf{X}'\mathbf{X} = T^{-1} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t'$$

- Similar to fitting OLS *without* a constant

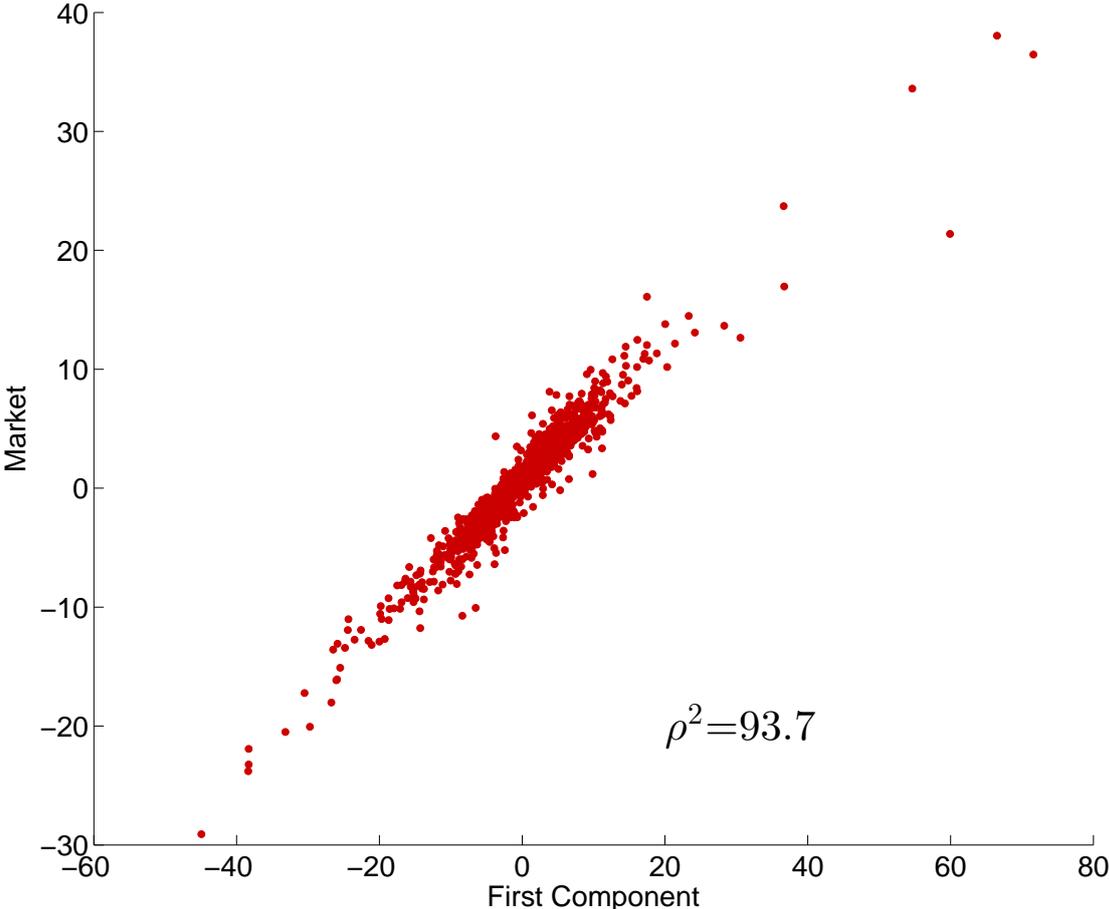
- Correlation matrix

$$\mathbf{R}_{\mathbf{x}} = T^{-1} \sum_{t=1}^T \mathbf{z}_t \mathbf{z}_t'$$

- $\mathbf{z}_t = (\mathbf{x}_t - \hat{\boldsymbol{\mu}}) \oslash \hat{\sigma}$  are the original data series, only **studentized**
  - **Important** since scale is not well defined for many economic data (e.g. indices)

- Initial exploration based on Fama-French data
  - 100 portfolios
    - Sorted on size and boot-to-market
  - 49 portfolios
    - Sorted on industry
- Equities are known to follow a strong factor model
  - Series missing more than 24 missing observations were dropped
    - 73 for 10 by 10 sort remaining
    - 41 of 49 industry portfolios
  - First 24 data points dropped for all series
  - July 1928 – December 2013
- $T = 1,026$
- $k = 114$
- Two versions, studentized and *raw*

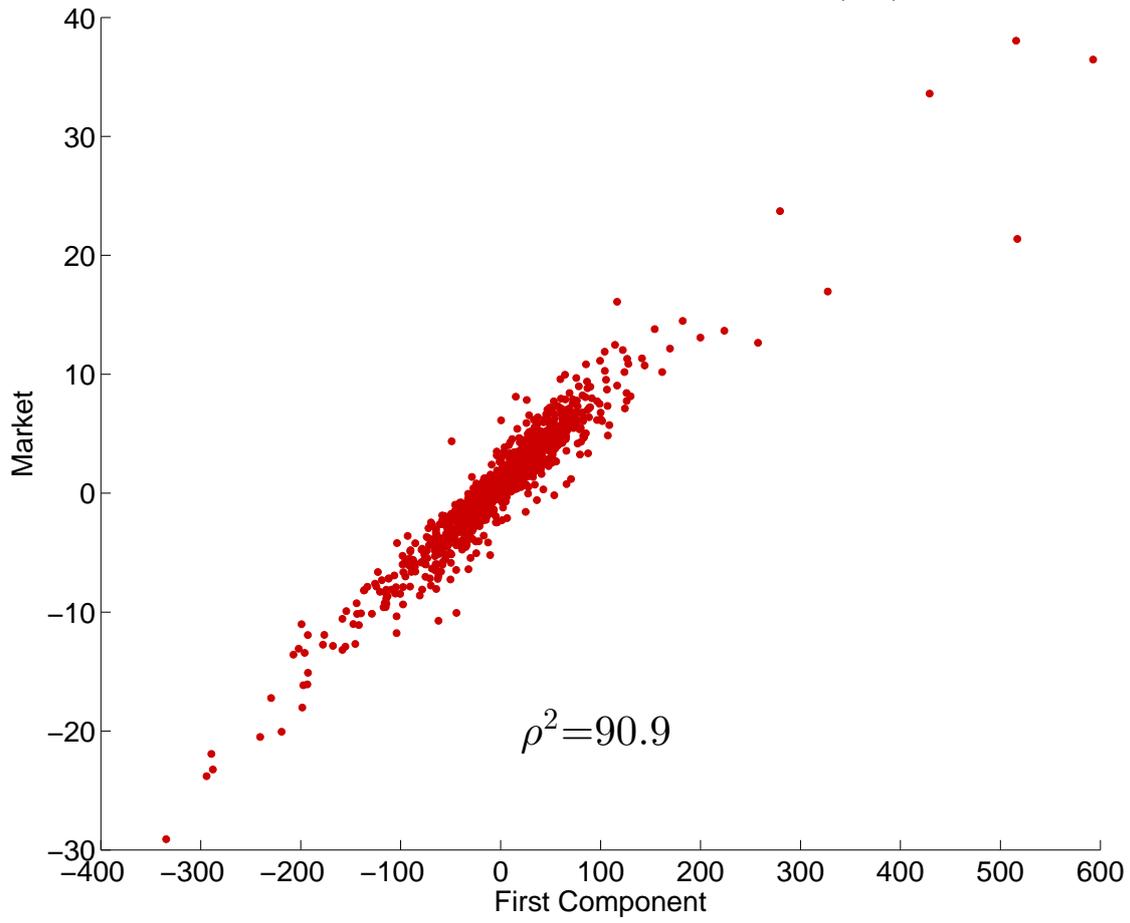
Scatter Plot of Excess Market and 1st PC



# First Factor from FF Data (Raw)



Scatter Plot of Excess Market and 1st PC (raw)



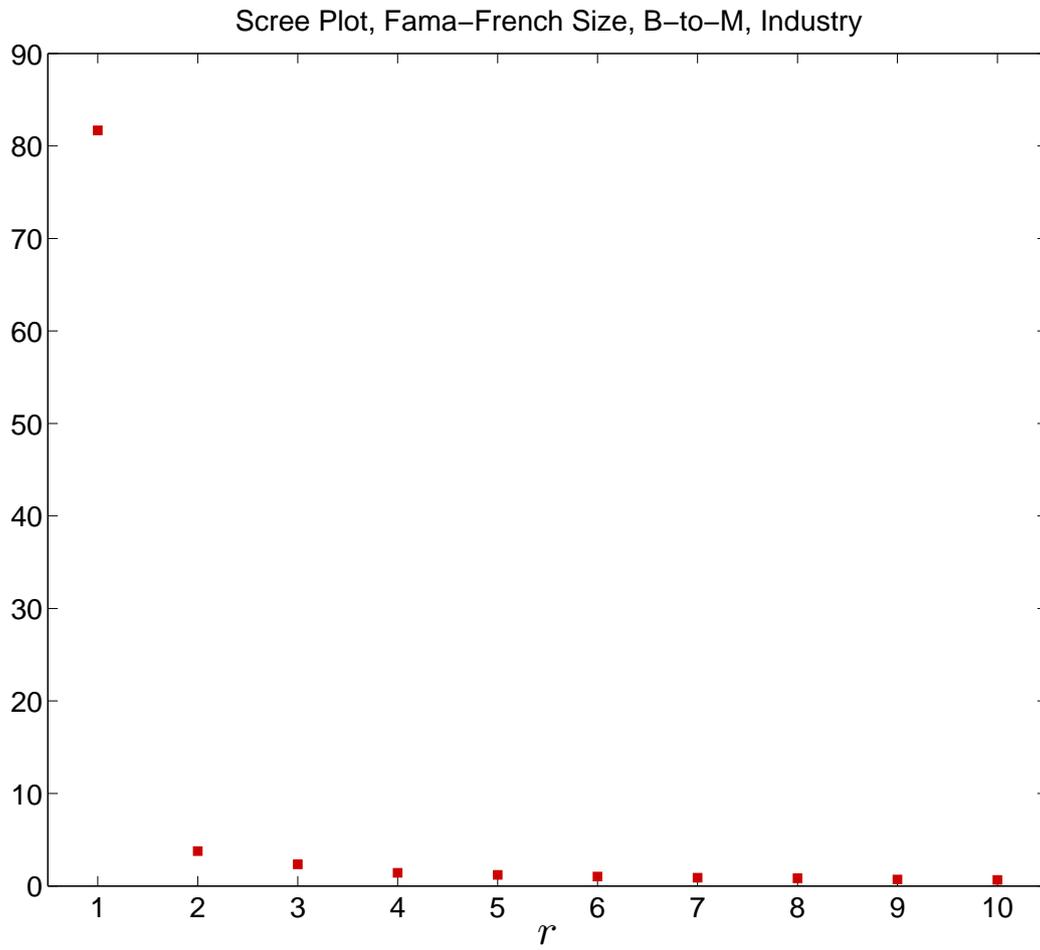
- So far have assumed  $r$  is known
- In practice  $r$  has to be estimated
- Two methods
  - Graphical using **Scree plots**
    - Plot of ordered eigenvalues, usually standardized by sum of all
    - Interpret this as the  $R^2$  of including  $r$  factors
    - Recall  $\sum_{i=1}^l \lambda_i = k$  for correlation matrix (Why?)
    - Closely related to system  $R^2$ ,

$$R^2(r) = \frac{\sum_{i=1}^r \lambda_i}{\sum_{j=1}^k \lambda_j}$$

- Information criteria-based
  - Similar to AIC/BIC, only need to account for both  $k$  and  $T$

## Stylized Fact(ors)

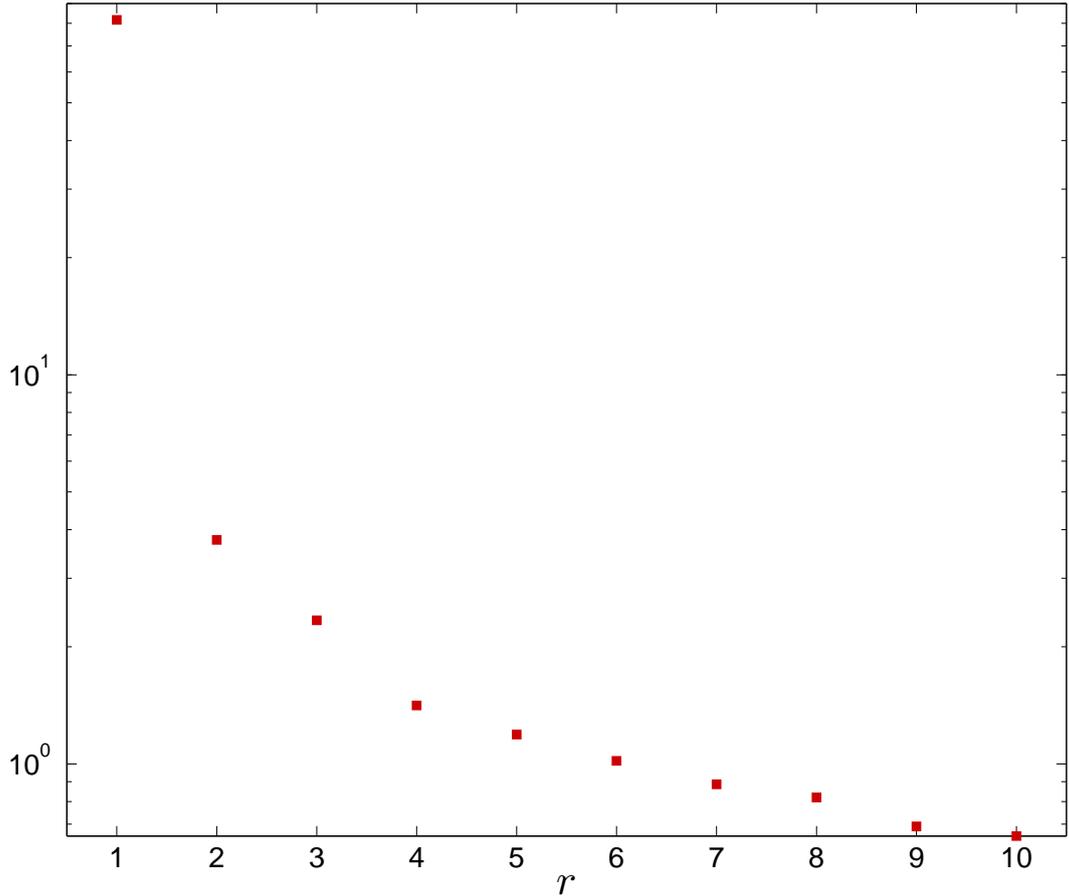
If in doubt, all known economic panels have between 1 and 6 factors



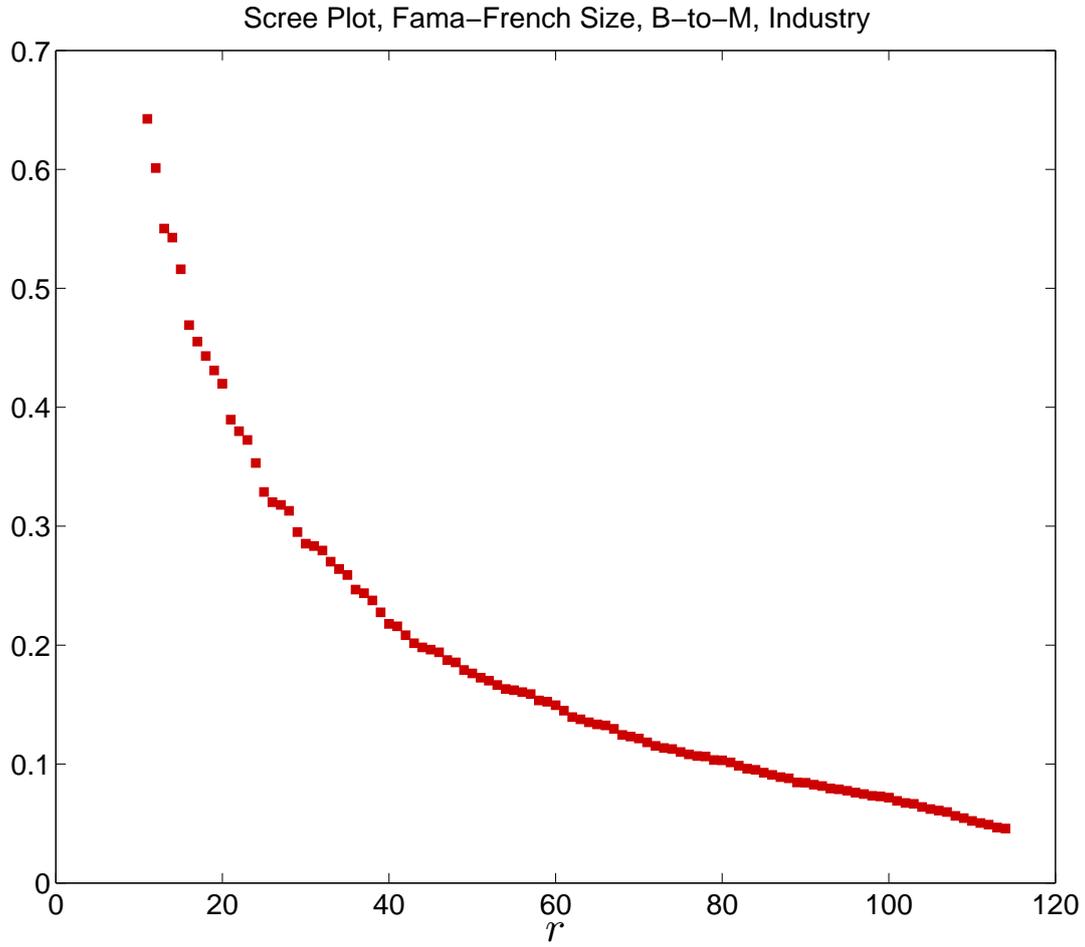
# Scree Plot: Fama-French



Scree Plot, Fama-French Size, B-to-M, Industry (Log)



# Scree Plot: Fama-French (Non-Factors)



- Bai & Ng (2002) studied the problem of selecting the correct number of factors in an approximate factor model
- Proposed a number of information criteria with the form

$$\ln \widehat{V}(r) + r \times g(k, T)$$

$$\widehat{V}(r) = \sum_{t=1}^T (\mathbf{x}_t - \widehat{\boldsymbol{\beta}}(r) \mathbf{f}_t(r))' (\mathbf{x}_t - \widehat{\boldsymbol{\beta}}(r) \mathbf{f}_t(r))$$

- $\widehat{V}(r)$  is the value of the objective function with  $r$  factors

- Three versions

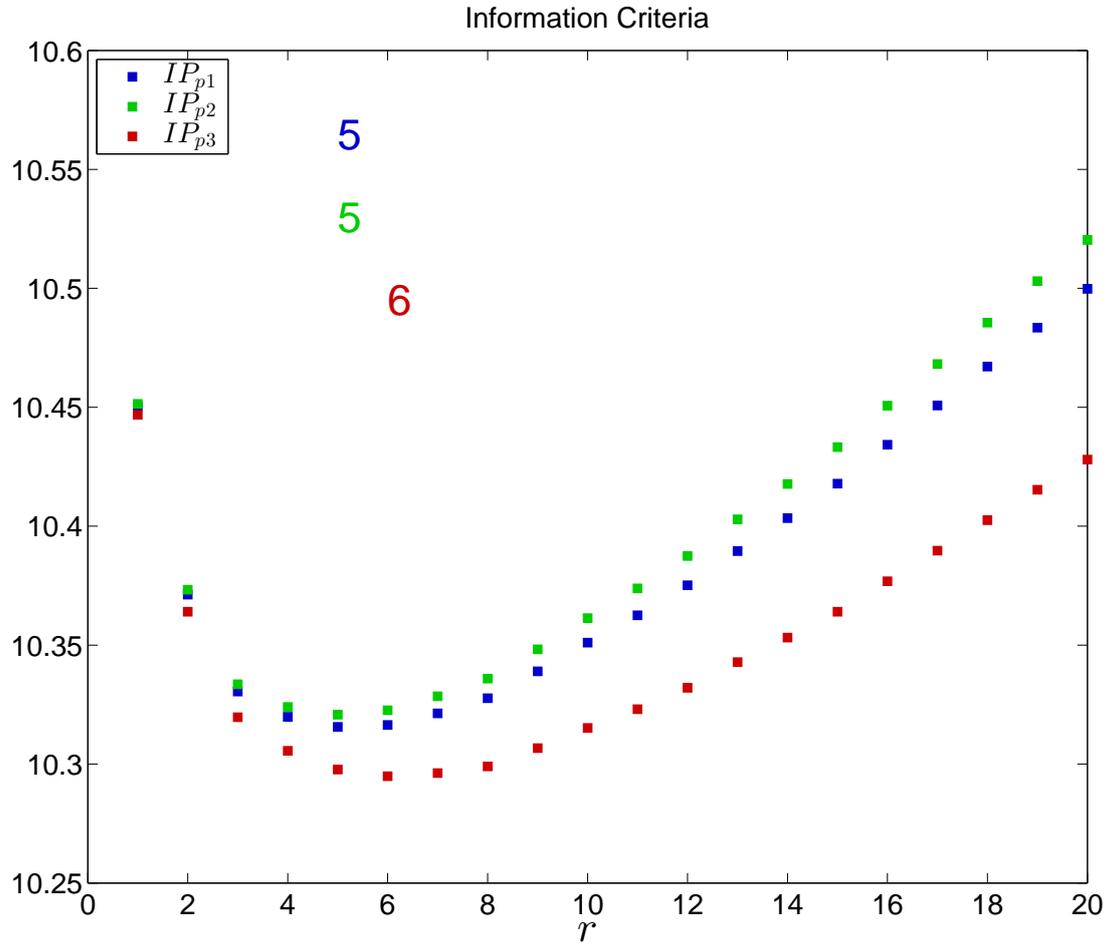
$$IC_{p_1} = \ln \widehat{V}(r) + r \left( \frac{k+T}{kT} \right) \ln \left( \frac{kT}{k+T} \right)$$

$$IC_{p_2} = \ln \widehat{V}(r) + r \left( \frac{k+T}{kT} \right) \ln (\min(k, T))$$

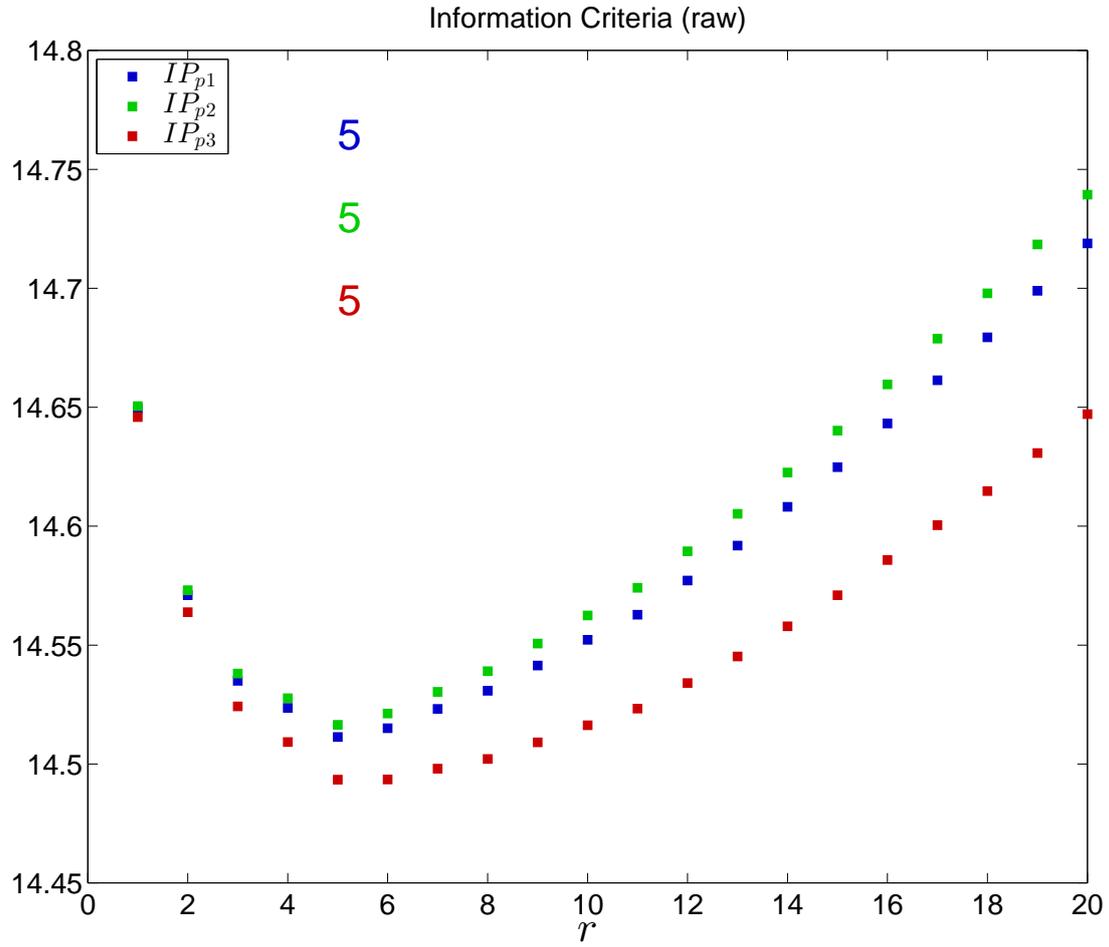
$$IC_{p_3} = \ln \widehat{V}(r) + r \left( \frac{\ln (\min(k, T))}{\min(k, T)} \right)$$

- Suppose  $k \approx T$ ,  $IC_{p_2}$  is BIC-like

$$IC_{p_2} = \ln \widehat{V}(r) + 2r \left( \frac{\ln T}{T} \right)$$



# Information Criteria: Fama-French (Raw)



- Fit can be assessed both globally and for individual series
- Least squares objective leads to natural  $R^2$  measurement of fit
- Global fit

$$R_{\text{global}}^2(r) = 1 - \frac{\text{tr}(\mathbf{X} - \hat{\boldsymbol{\beta}}(r)\mathbf{F}(r))'(\mathbf{X} - \hat{\boldsymbol{\beta}}(r)\mathbf{F}(r))}{\text{tr}(\mathbf{X}'\mathbf{X})}$$

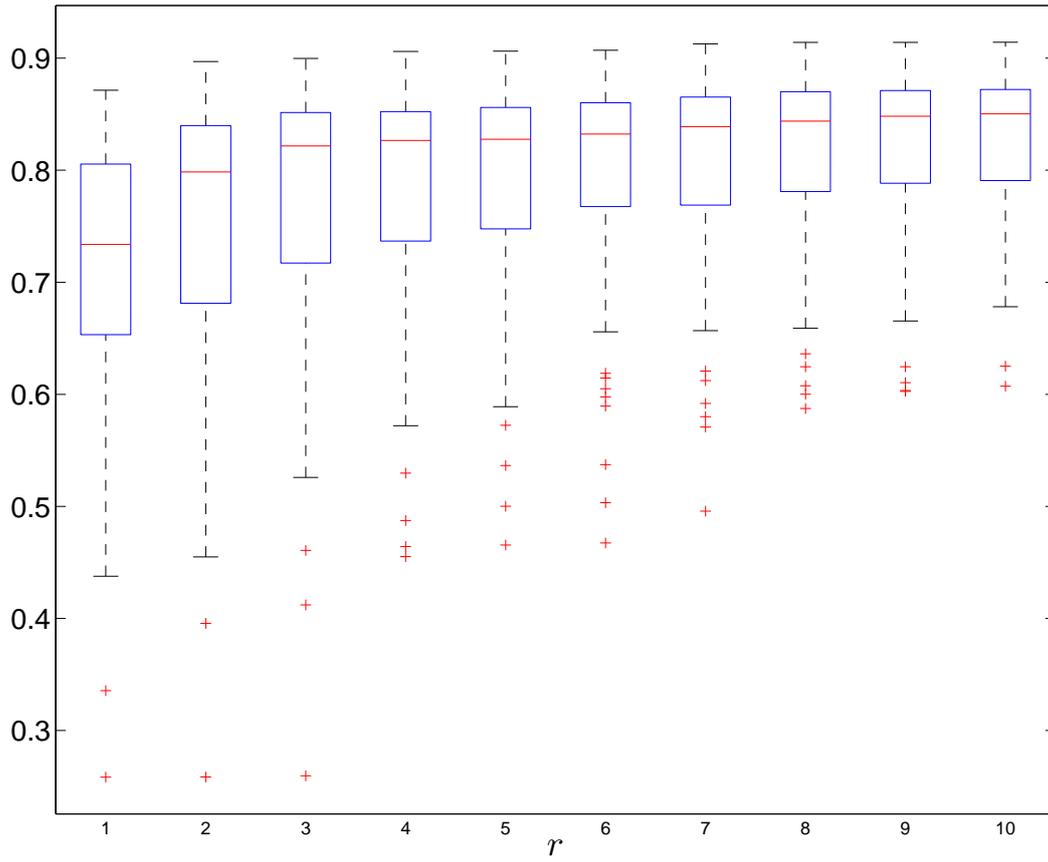
$$= \frac{\sum_{i=1}^r \lambda_i}{\sum_{j=1}^k \lambda_j}$$

- Numerator is just  $\widehat{V}(r) = \sum_{i=1}^k \sum_{t=1}^T \left(x_{it} - \sum_{j=1}^r \hat{\beta}_{ij}f_{jt}\right)^2$
- When  $\mathbf{x}$  has been studentized,  $\text{tr}(\mathbf{X}'\mathbf{X}) = \sum_{j=1}^k \lambda_j = Tk$
- Individual fit

$$R_i^2(r) = 1 - \frac{\sum_{t=1}^T \left(x_{it} - \sum_{j=1}^r \hat{\beta}_{ij}f_{jt}\right)^2}{\sum_{t=1}^T x_{it}^2}$$

- Useful for assessing series not well described by factor model

Individual  $R^2$  using  $r$  factors





- Dynamic factors model specify dynamics in the factors
- Basic DFM is

$$\mathbf{x}_t = \sum_{i=0}^s \Phi_i \mathbf{f}_t + \boldsymbol{\epsilon}_t$$
$$\mathbf{f}_t = \sum_{j=1}^q \Psi \mathbf{f}_{t-j} + \boldsymbol{\eta}_t$$

- Observed data depend on contemporaneous and lagged factors
- Factors have VAR-like dynamics
- Assumed that  $\mathbf{f}_t$  and  $\boldsymbol{\epsilon}_t$  are stationary, so  $\mathbf{x}_t$  is also stationary
  - **Important:** must transform series appropriately when applying to data
- $\boldsymbol{\epsilon}_t$  can have weak dependence in both the cross-section and time-series
- $E[\boldsymbol{\epsilon}_t, \boldsymbol{\eta}_s] = \mathbf{0}$  for all  $t, s$

$$\mathbf{x}_t = \sum_{i=0}^s \Phi_i \mathbf{f}_{t-i} + \epsilon_t, \quad \mathbf{f}_t = \sum_{j=1}^q \Psi_j \mathbf{f}_{t-j} + \eta_t$$

- Optimal forecast can be derived

$$\begin{aligned} E [x_{it+1} | \mathbf{x}_t, \mathbf{f}_t, \mathbf{x}_{t-1}, \mathbf{f}_{t-1}, \dots] &= E \left[ \sum_{i=0}^s \phi_i \mathbf{f}_{t+1-i} + \epsilon_{it+1} | \mathbf{x}_t, \mathbf{f}_t, \mathbf{x}_{t-1}, \mathbf{f}_{t-1}, \dots \right] \\ &= E_t \left[ \sum_{i=0}^s \phi_i \mathbf{f}_{t+1-i} \right] + E_t [\epsilon_{it+1}] \\ &= \sum_{i=1}^{s'} \mathbf{A}_i \mathbf{f}_{t-i+1} + \sum_{j=1}^n \mathbf{B}_j x_{it-j+1} \end{aligned}$$

- Predictability in both components
  - Lagged factors predict factors
  - Lagged  $x_{it}$  predict  $\epsilon_{it}$

- DFM is really factors plus moving average
- Moving average processes can be replaced with AR processes when invertible

$$\begin{aligned}y_t &= \epsilon_t + \theta \epsilon_{t-1} \\y_t - \theta y_{t-1} &= \epsilon_t + \theta \epsilon_{t-1} - \theta (\theta \epsilon_{t-2} + \epsilon_{t-1}) \\&= \epsilon_t - \theta^2 \epsilon_{t-2} \\y_t - \theta y_{t-1} + \theta^2 y_{t-2} &= \epsilon_t - \theta^2 \epsilon_{t-2} + \theta^2 (\theta \epsilon_{t-3} + \epsilon_{t-2}) \\&= \epsilon_t + \theta^2 (\theta \epsilon_{t-3} + \epsilon_{t-2}) \\ \sum_{i=0}^{\infty} (-\theta)^i y_{t-i} &= \epsilon_t \\y_t &= \sum_{i=1}^{\infty} -(-\theta)^i y_{t-i} + \epsilon_t\end{aligned}$$

- Can approximate finite MA with finite AR
- Quality will depend on the persistence of the MA component

- Superficially dynamic factor models appear to be more complicated than static factor models
- Dynamic Factor models can be directly estimated using Kalman Filter or spectral estimators that account for serial correlation in factors
  - Latter are not useful for forecasting since 2-sided
- (Big) However, DFM can be converted to Static model by relabeling
- In DFM, factors are

$$[\mathbf{f}_t, \mathbf{f}_{t-1}, \dots, \mathbf{f}_{t-s}]$$

- Total of  $r(s + 1)$  factors in model
- Equivalent to static model with *at most*  $r(s + 1)$  factors
  - Redundant factors will not appear in static version

- Consider basic DFM

$$\begin{aligned}x_{it} &= \phi_{i1}f_t + \phi_{i2}f_{t-1} + \epsilon_{it} \\f_t &= \psi f_{t-1} + \eta_t\end{aligned}$$

- Model can be expressed as

$$\begin{aligned}x_{it} &= \phi_{i1}(\psi f_{t-1} + \eta_t) + \phi_{i2}f_{t-1} + \epsilon_{it} \\&= \phi_{i1}\eta_t + \phi_{i2}(1 + (\phi_{i1}/\phi_{i2})\psi)f_{t-1} + \epsilon_{it}\end{aligned}$$

- One version of static factors are  $\eta_t$  and  $f_{t-1}$ 
  - In this particular version,  $\eta_t$  is not “dynamic” since it is WN
  - $f_{t-1}$  follows an AR(1) process
- Other *rotations* will have different dynamics



- Basic simulation

$$\begin{aligned}x_{it} &= \phi_{i1}f_t + \phi_{i2}f_{t-1} + \epsilon_{it} \\f_t &= \psi f_{t-1} + \eta_t\end{aligned}$$

- $\phi_{i1} \sim N(1, 1), \phi_{i2} \sim N(.2, 1)$ 
  - Smaller signal makes it harder to find second factor
- $\psi = 0.5$ 
  - Higher persistence makes it harder since  $\text{Corr}[f_t, f_{t-1}]$  is larger
- Everything else standard normal
- $k = 100, T = 100$ 
  - Also  $k = 200$  and  $T = 200$  (separately)
- All estimation using PCA on correlation

## Number of Factors for Forecasting

Better to have  $r$  above  $r^*$  than below

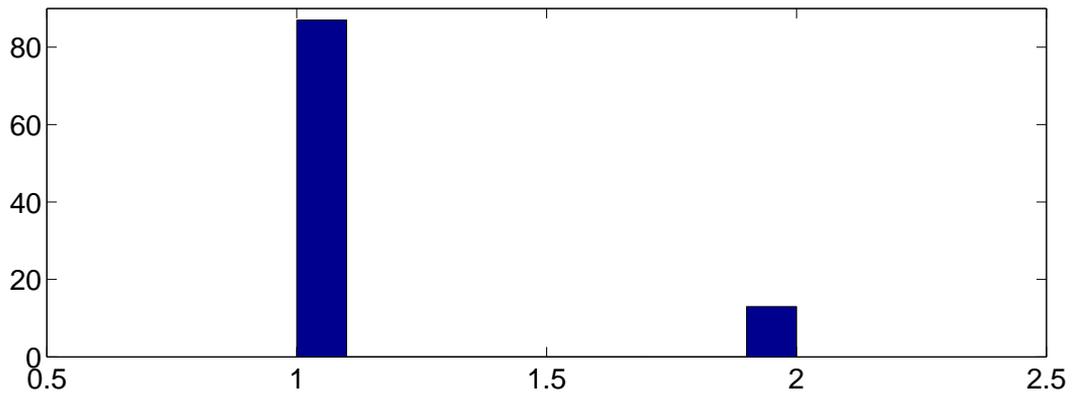
- Factors are not point identified
  - Can use an arbitrary rotation and model is equivalent
- Natural measure of similarity between original (GDP) factors and estimated factors is global  $R^2$

$$\hat{\mathbf{f}}_t = \mathbf{A}\mathbf{f}_t + \boldsymbol{\eta}_t$$
$$R^2 = 1 - \frac{\sum_{t=1}^T \hat{\boldsymbol{\eta}}_t' \hat{\boldsymbol{\eta}}_t}{\sum_{t=1}^T \mathbf{f}_t' \mathbf{f}_t}$$

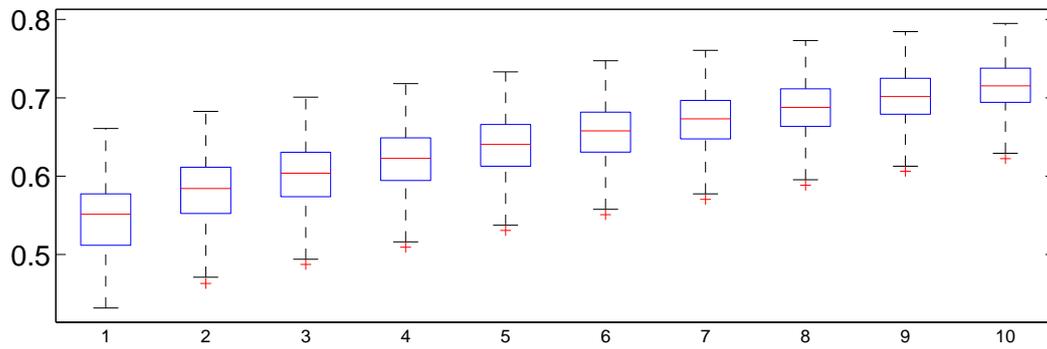
- Note that  $\mathbf{A}$  is a 2 by 2 matrix of regression coefficients



$IC_{p2}$  Selected  $r$ ,  $T=100$ ,  $k=100$

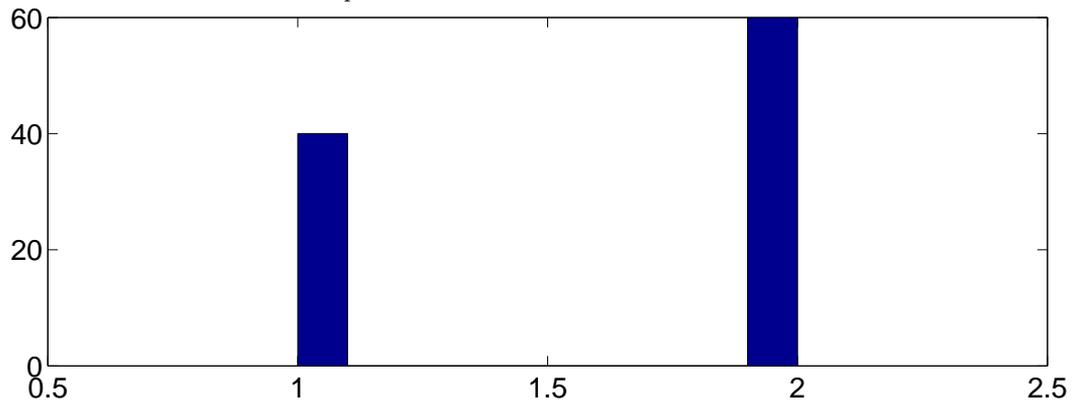


$R^2$  as a function of  $r$

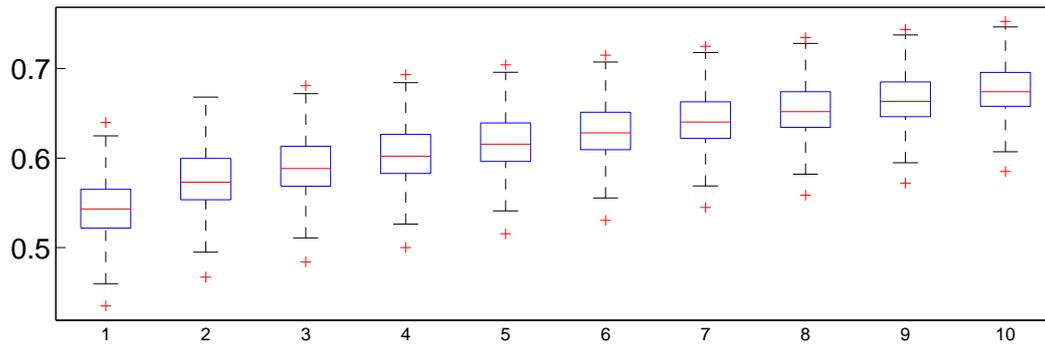




$IC_{p2}$  Selected  $r$ ,  $T=100$ ,  $k=200$

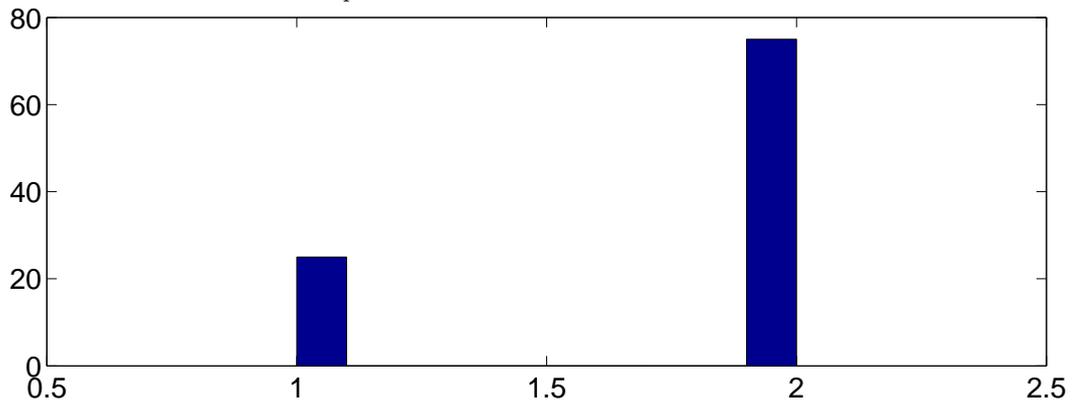


$R^2$  as a function of  $r$

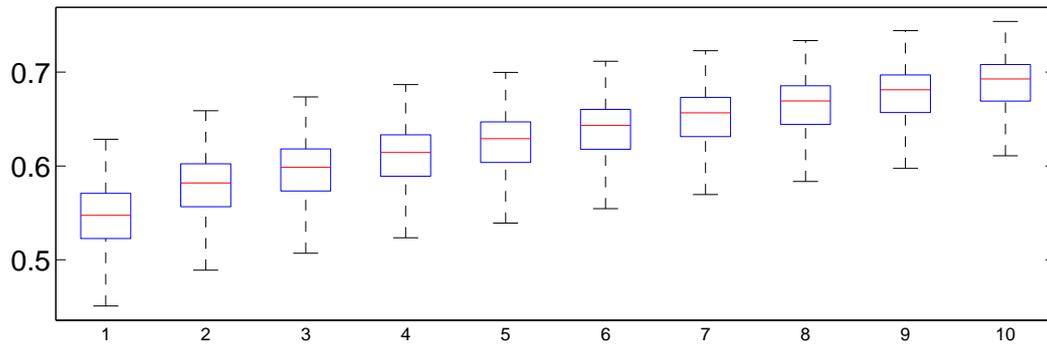


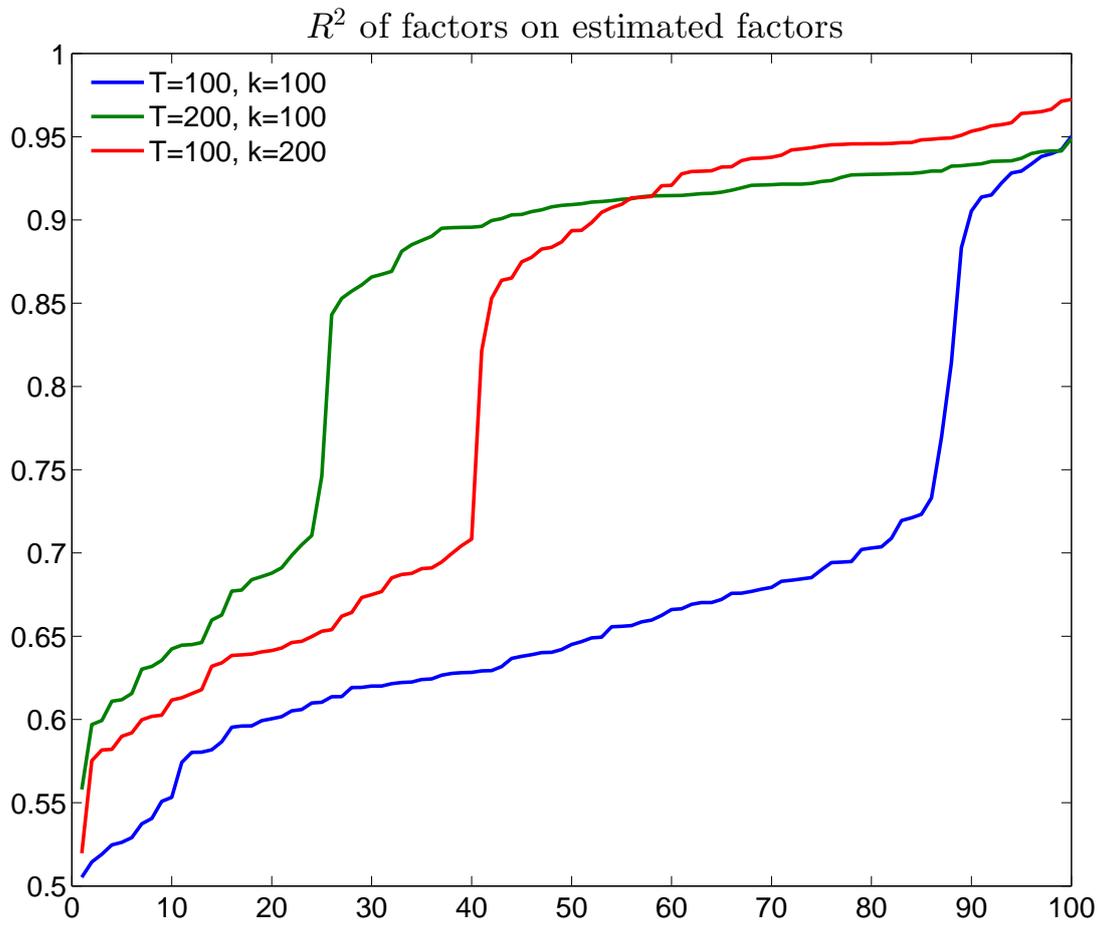


$IC_{p2}$  Selected  $r$ ,  $T=200$ ,  $k=100$



$R^2$  as a function of  $r$







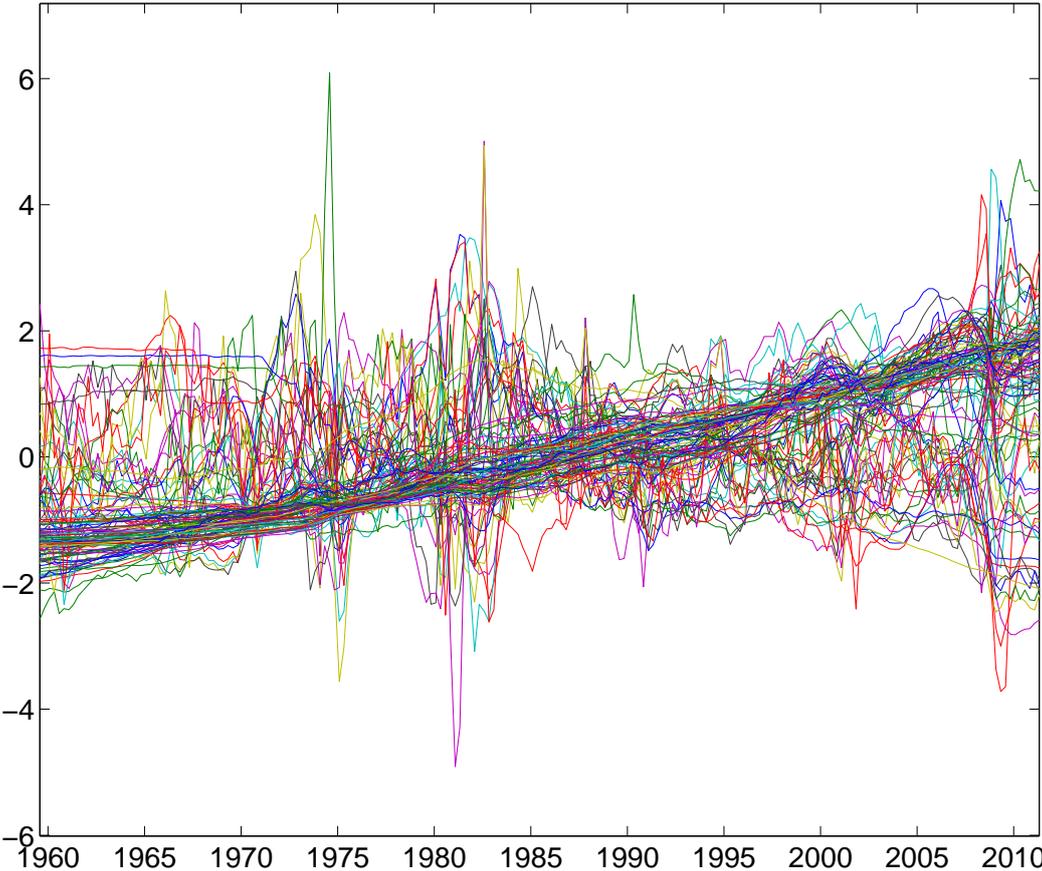
- Stock & Watson have been at the forefront of factor model development
- Data is from 2012 paper “Disentangling the Channels of the 2007-2009 Recession”
- Dataset consists of 137 monthly and 74 quarterly series
  - Not all used for factor estimation
  - Aggregates not used if disaggregated series available
- Monthly series are aggregated to quarterly, which is frequency of data
- Series with missing observations are dropped for simplicity
  - Before dropping those with missing values data set has 132 series
  - After 107 series remain

National Income and Product Accounts (NIPA)	12
Industrial Production	9
Employment and Unemployment	30
Housing Starts	6
Inventories, Orders, and Sales	7
Prices	25
Earnings and Productivity	8
Interest Rates	10
Money and Credit	6
Stock Prices, Wealth, Household Balance Sheets	8
Housing Prices	3
Exchange Rates	6
Other	2

- Monthly series were aggregated to quarterly using
  - Average
  - End-of-quarter
- All series were transformed to be stationary using one of:
  - No transform
  - Difference
  - Double-difference
  - Log
  - Log-difference
  - Double-log-difference
- Most series checked for outliers relative to *IQR* (rare)
- Final series were Studentized in estimation of PC

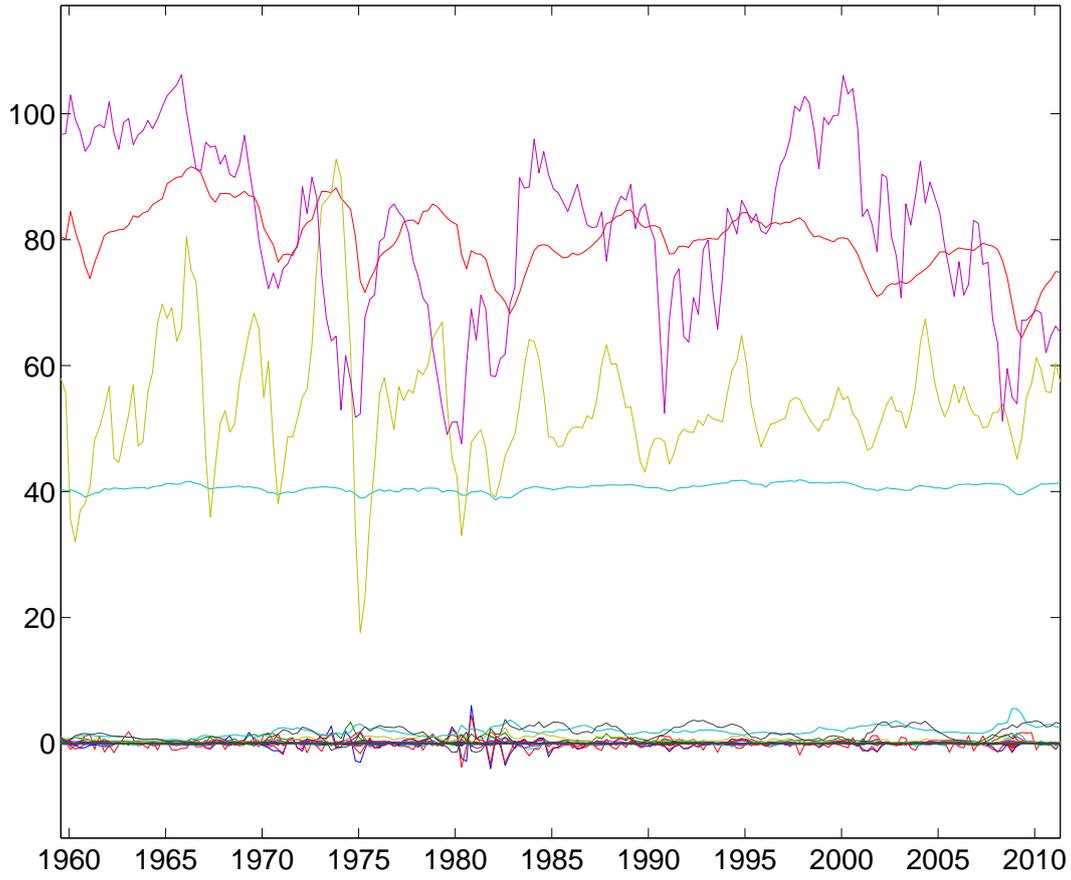


Untransformed SW Data (Studentized)



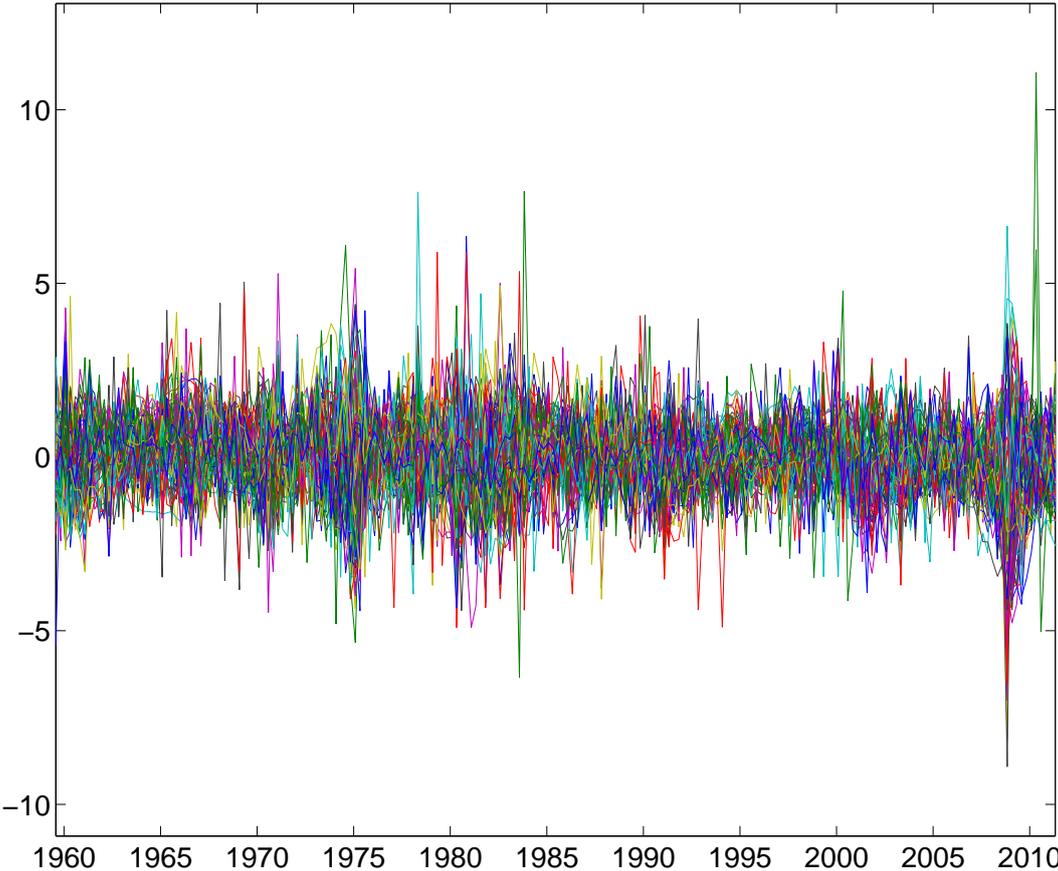


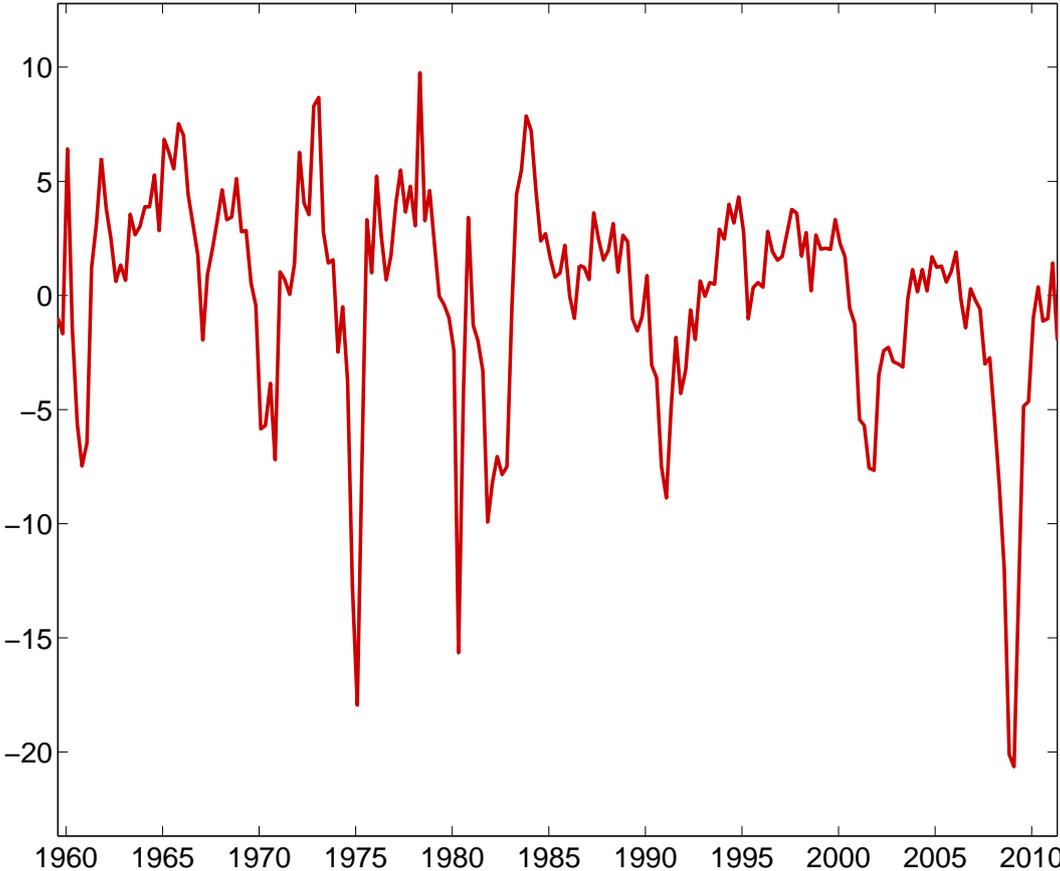
Transformed SW Data





Studentized SW Data

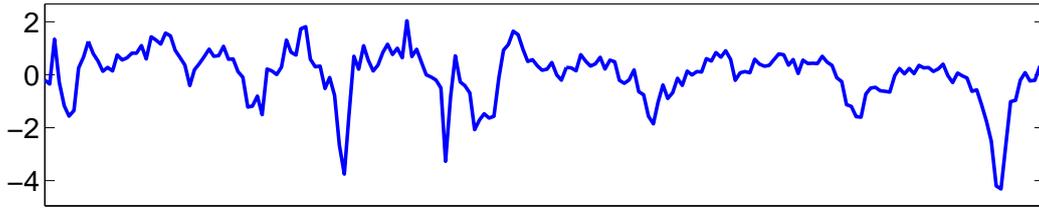




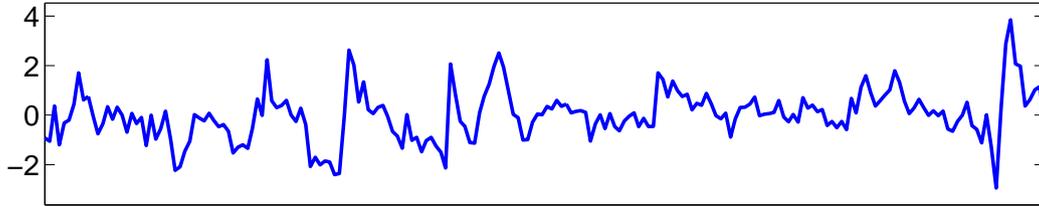
# First Three Components



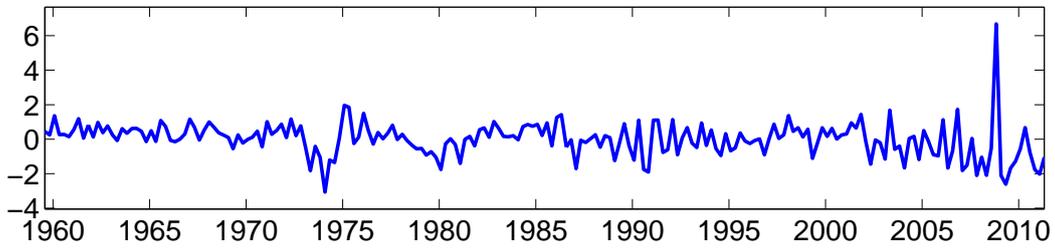
First Component (Standardized)



Second Component (Standardized)

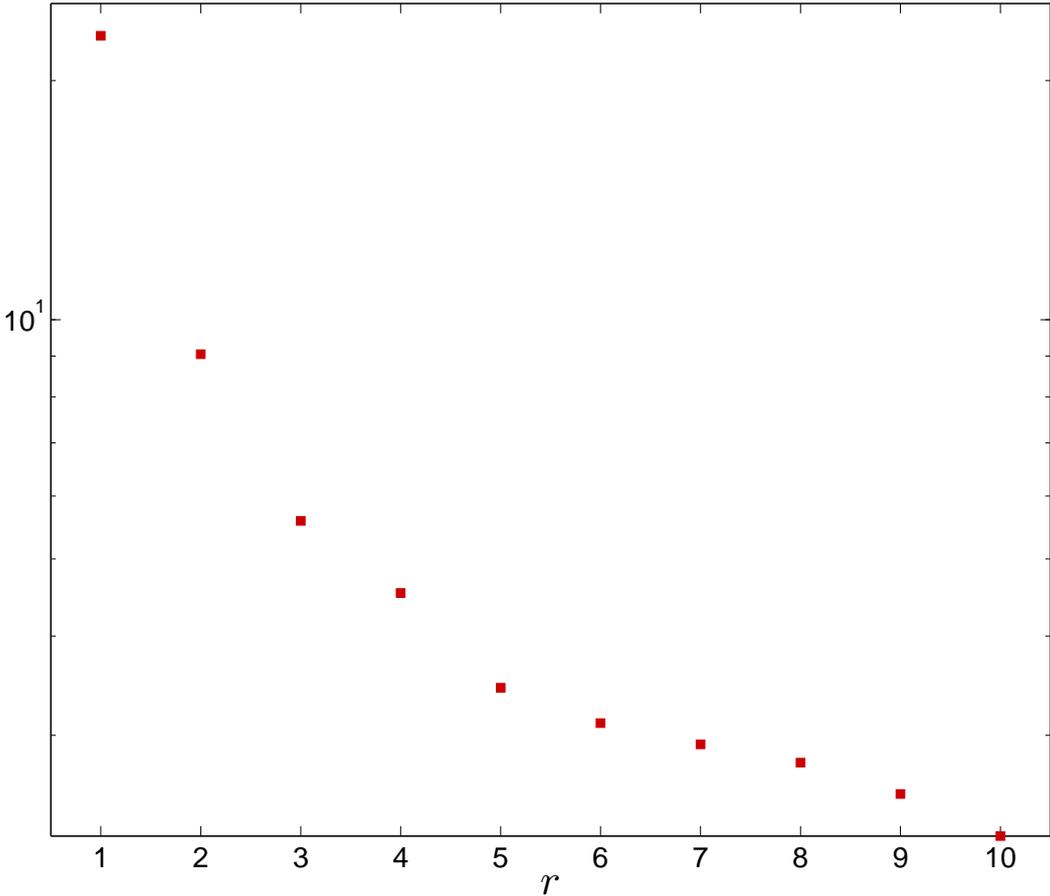


Third Component (Standardized)



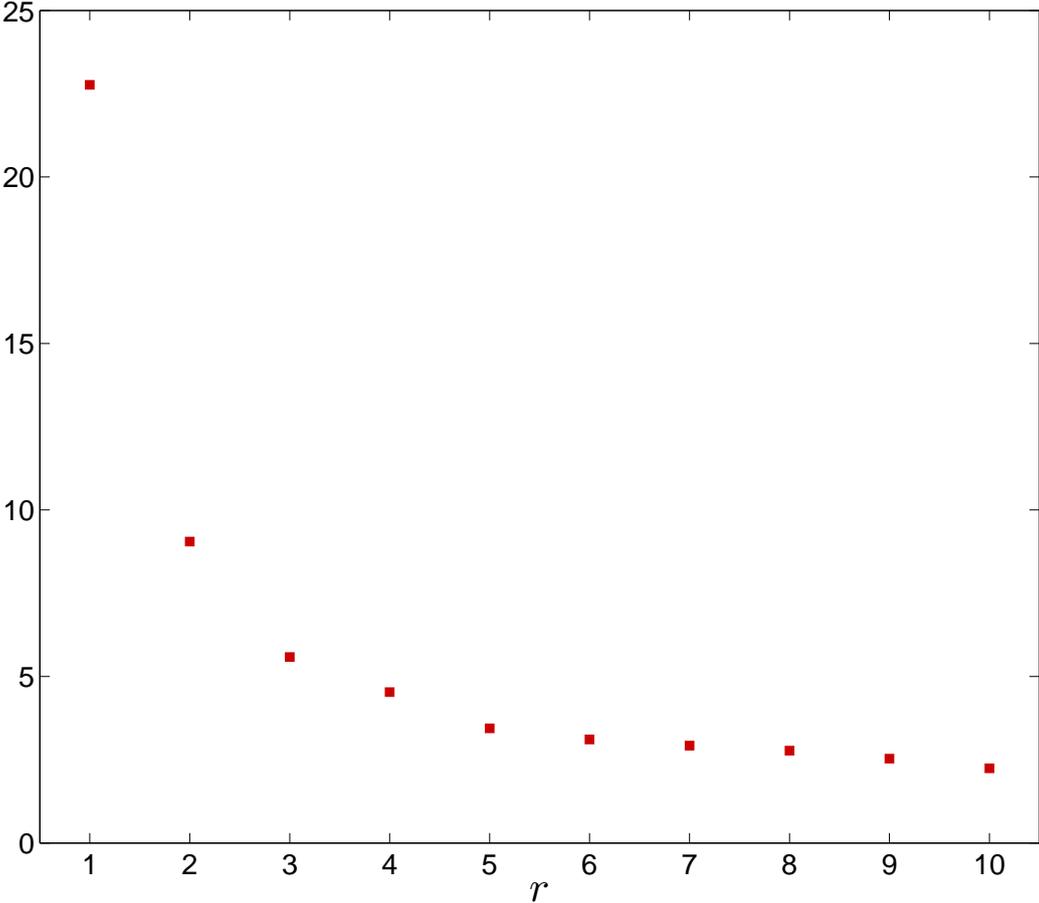


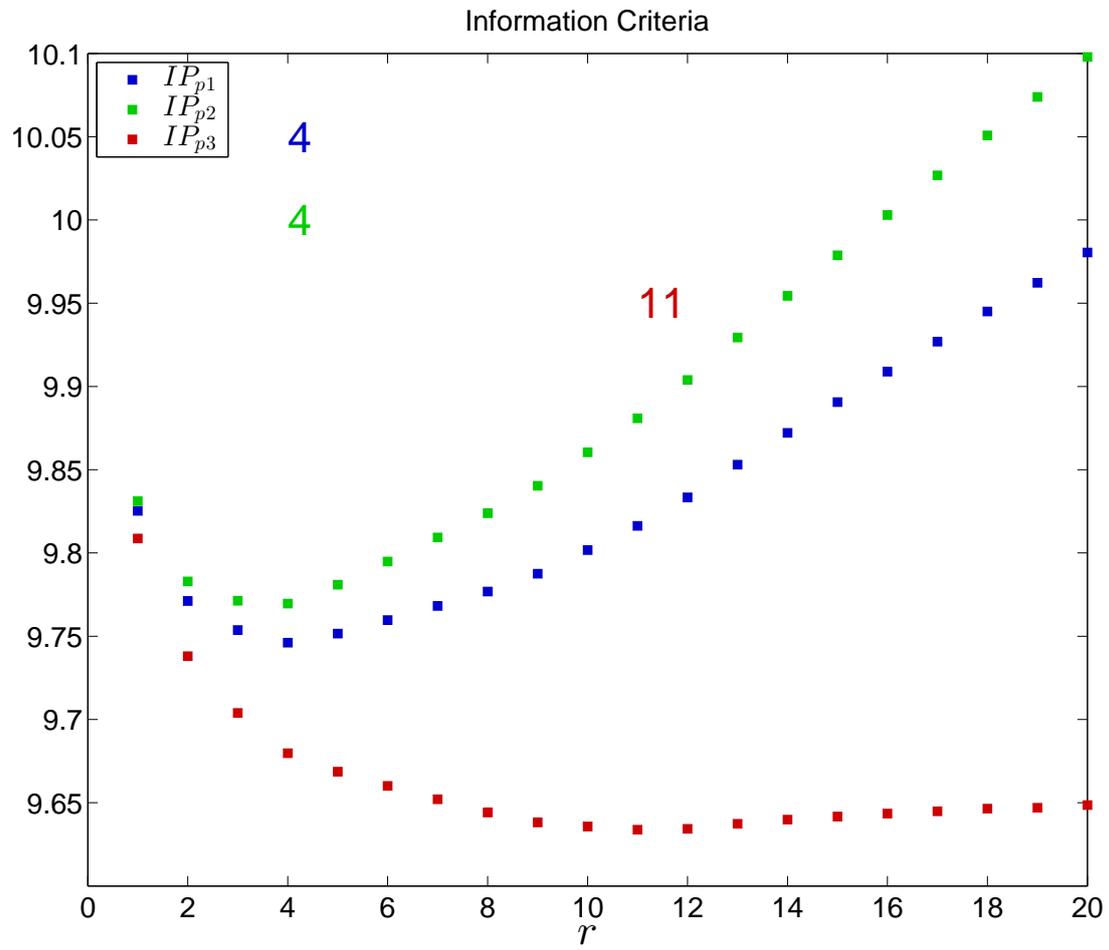
Scree Plot, Stock & Watson (Log)



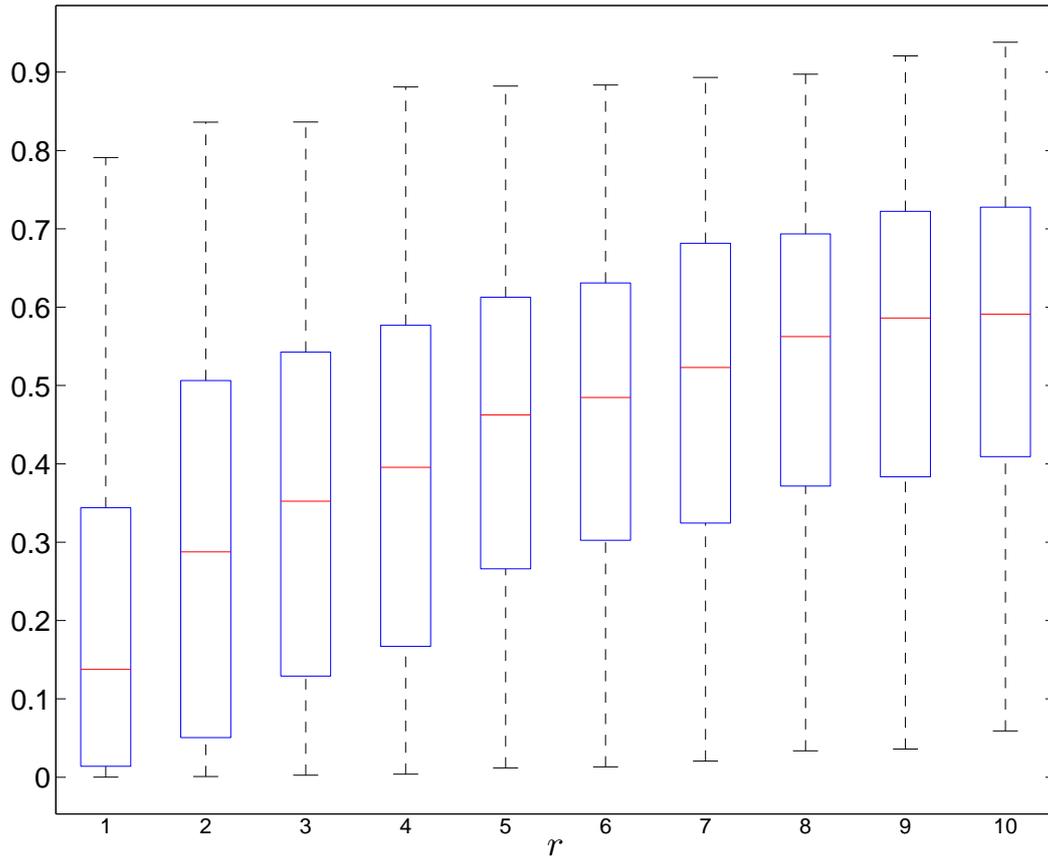


Scree Plot, Stock & Watson





Individual  $R^2$  using  $r$  factors



- Forecast problem is not meaningfully different from standard problem
- Interest is now in  $\mathbf{y}_t$ , which may or may not be in  $\mathbf{x}_t$ 
  - Note that stationary version of  $\mathbf{y}_t$  should be forecast, e.g.  $\Delta\mathbf{y}_t$  or  $\Delta^2\mathbf{y}_t$
- Two methods to forecast

## Unrestricted

$$y_{t+1} = \phi_0 + \sum_{i=1}^p \phi_i y_{t-i+1} + \boldsymbol{\theta}' \hat{\mathbf{f}}_t + \epsilon_{it}$$

- Treats factors as observed data, only makes sense if  $k$  is large
  - Uses an  $AR(P)$  to model residual dependence
  - Choice of number of factors to use, may be different from  $r$
  - Can also use lags of  $\mathbf{f}_t$  (uncommon)
  - Model selection is applicable as usual, e.g. BIC

## Restricted

- When  $\mathbf{y}_t$  is in  $\mathbf{x}_t$ ,  $\mathbf{y}_t = \boldsymbol{\beta} \hat{\mathbf{f}}_t + \epsilon_t$

$$\epsilon_t = \mathbf{y}_t - \boldsymbol{\beta} \hat{\mathbf{f}}_t$$

$$\begin{aligned} \hat{\mathbf{y}}_{t+1|t} &= \boldsymbol{\beta} \hat{\mathbf{f}}_{t+1|t} + \sum_{i=1}^p \phi_i \left( y_{t-i+1} - \boldsymbol{\beta} \hat{\mathbf{f}}_{t-i+1} \right) \\ &= \boldsymbol{\beta} \hat{\mathbf{f}}_{t+1|t} + \sum_{i=1}^p \phi_i \hat{\epsilon}_t \end{aligned}$$

- VAR to forecast  $\hat{\mathbf{f}}_{t+1}$  using lags of  $\hat{\mathbf{f}}_t$
- Univariate AR for  $\hat{\epsilon}_t$
- Usually found to be less successful than unrestricted
- Care is needed when using studentized data since forecasting recentered, rescaled version of  $y$

- When forecasting  $\Delta \mathbf{y}_t$ ,

$$\begin{aligned} E_t[\mathbf{y}_{t+1}] &= E_t[\mathbf{y}_{t+1} - \mathbf{y}_t + \mathbf{y}_t] \\ &= E_t[\Delta \mathbf{y}_{t+1}] + \mathbf{y}_t \end{aligned}$$

- At longer horizons,

$$E_t[\mathbf{y}_{t+h}] = \sum_{i=1}^h E_t[\Delta \mathbf{y}_{t+i}] + \mathbf{y}_t$$

- When forecasting  $\Delta^2 \mathbf{y}_t$

$$\begin{aligned} E_t[\mathbf{y}_{t+1}] &= E_t[\mathbf{y}_{t+1} - \mathbf{y}_t - \mathbf{y}_t + \mathbf{y}_{t-1} + 2\mathbf{y}_t - \mathbf{y}_{t-1}] \\ &= E_t[\Delta^2 \mathbf{y}_{t+1}] + 2\mathbf{y}_t - \mathbf{y}_{t-1} \end{aligned}$$

- ▶ In many cases interest is in  $\Delta \mathbf{y}_t$  when forecasting  $\Delta^2 \mathbf{y}_t$ 
  - ▷ For example CPI, inflation and change in inflation
  - ▷ Same as re-integrating  $\Delta \mathbf{y}_t$  to  $\mathbf{y}_t$

- Multistep can be constructed using either method
- Unrestricted requires additional VAR for  $\hat{\mathbf{f}}_t$
- Alternative use direct forecasting

$$y_{t+h|t} = \hat{\phi}_{(h)0} + \sum_{i=1}^{p^h} \hat{\phi}_{(h)i} y_{t-i+1} + \hat{\boldsymbol{\theta}}'_{(h)} \hat{\mathbf{f}}_t$$

- ▶  $(h)$  used to denote explicit parameter dependence on horizon
- ▶  $y_{t+h|t}$  can be either the period- $h$  value, or the  $h$ -period cumulative forecast (more common)
- Direct has been documented to be better than iterative in DFMs
  - ▶ [Problem dependent](#)

- Used BIC search across models
- 3 setups
  - GDP lags only (4), Components Only (6), Both

$$\sum_{j=1}^h \Delta g_{t+j} = \phi_0 + \sum_{s=1}^4 \gamma_s \Delta g_{t-s+1} + \sum_{n=1}^6 \psi_n f_{jt} + \epsilon_{ht}$$

	GDP Only		Components Only		Both		
	Lags	$R^2$	Lags	$R^2$	Lags	Lags	$R^2$
$h = 1$	1, 2, 4	.517	1, 2, 3, 4, 6	.662	1	1, 2, 3, 4, 6	.686
$h = 2$	1, 4	.597	1, 2, 3, 4, 6	.763	1	1, 2, 3, 4, 6	.771
$h = 3$	1, 4	.628	1, 2, 3, 4, 6	.785	1	1, 2, 3, 4, 6	.792
$h = 4$	1, 4	.661	1, 2, 3, 4, 6	.805	-	1, 2, 3, 4, 6	.805

- Basic PCA makes use of the covariance or more commonly correlation
- Correlation is technically a special case of *generalized PCA*

$$\min_{\beta, \mathbf{f}_t, \dots, \mathbf{f}_t} \sum_{t=1}^T (\mathbf{x}_t - \beta \mathbf{f}_t)' \Sigma_{\epsilon}^{-1} (\mathbf{x}_t - \beta \mathbf{f}_t) \text{ subject to } \beta' \beta = \mathbf{I}_r$$

- Clever choices of  $\Sigma_{\epsilon}$  lead to difference estimators
  - Using  $\text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_k^2)$  where  $\hat{\sigma}_j^2$  is variance of  $x_j$  leads to correlation
  - Tempting to use GLS version based on  $r$  principal components

## Algorithm (Principal Component Analysis using GLS)

1. Estimate  $\hat{\epsilon}_{it} = x_{it} - \hat{\beta}_i \hat{\mathbf{f}}_t$  using  $r$  factors
2. Estimate  $\hat{\sigma}_{\epsilon i}^2 = T^{-1} \sum \hat{\epsilon}_{it}^2$  and  $\mathbf{W} = \text{diag}(w_1, \dots, w_k)$  where

$$w_i = \frac{1/\hat{\sigma}_{\epsilon i}}{\sum_{j=1}^k 1/\hat{\sigma}_{\epsilon j}}$$

3. Compute PCA-GLS using  $\mathbf{W}\mathbf{X}$

- Absolute covariance weighting
  1. Compute complete residual covariance  $\hat{\Sigma}_\epsilon$  from residuals
  2. Replace  $\hat{\sigma}_{\epsilon i}^2$  in step 2 with  $\hat{\sigma}_{\epsilon i}^2 = \sum_{j=1}^k |\hat{\Sigma}_\epsilon(i, j)|$
- Down-weights series which have both large idiosyncratic variance *and* strong residual covariance
- Stock & Watson (2005) use more sophisticated method
  1. Estimate AR(P) on  $\hat{\epsilon}_{it}$  for all series

$$\hat{\epsilon}_{it} = \sum_{j=1}^{p_i} \phi_j \epsilon_{it-j} + \xi_{it}$$

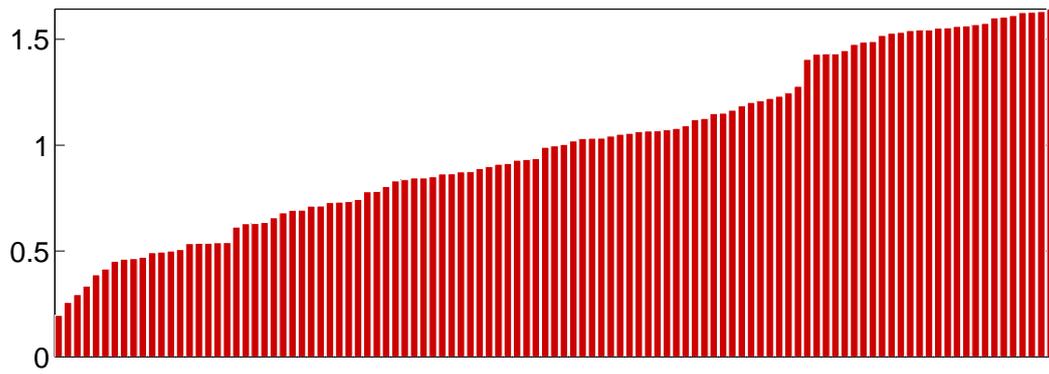
2. Construct quasi-differenced  $x_{it}$  using coefficients

$$\tilde{x}_{it} = x_{it} - \sum_{j=1}^{p_i} \hat{\phi}_j x_{it-j}$$

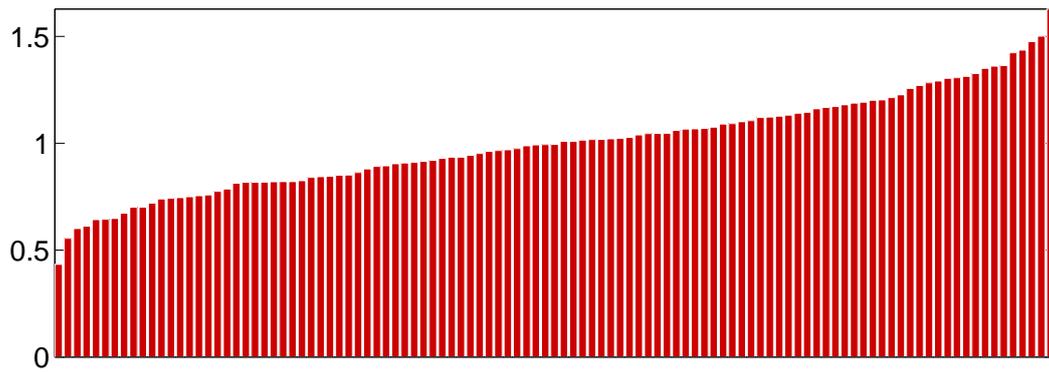
3. Estimate  $\hat{\sigma}_{\epsilon i}^2$  using  $\hat{\xi}_{it}$
4. Re-estimate factors using quasi-differenced data and weighting, iterate if needed

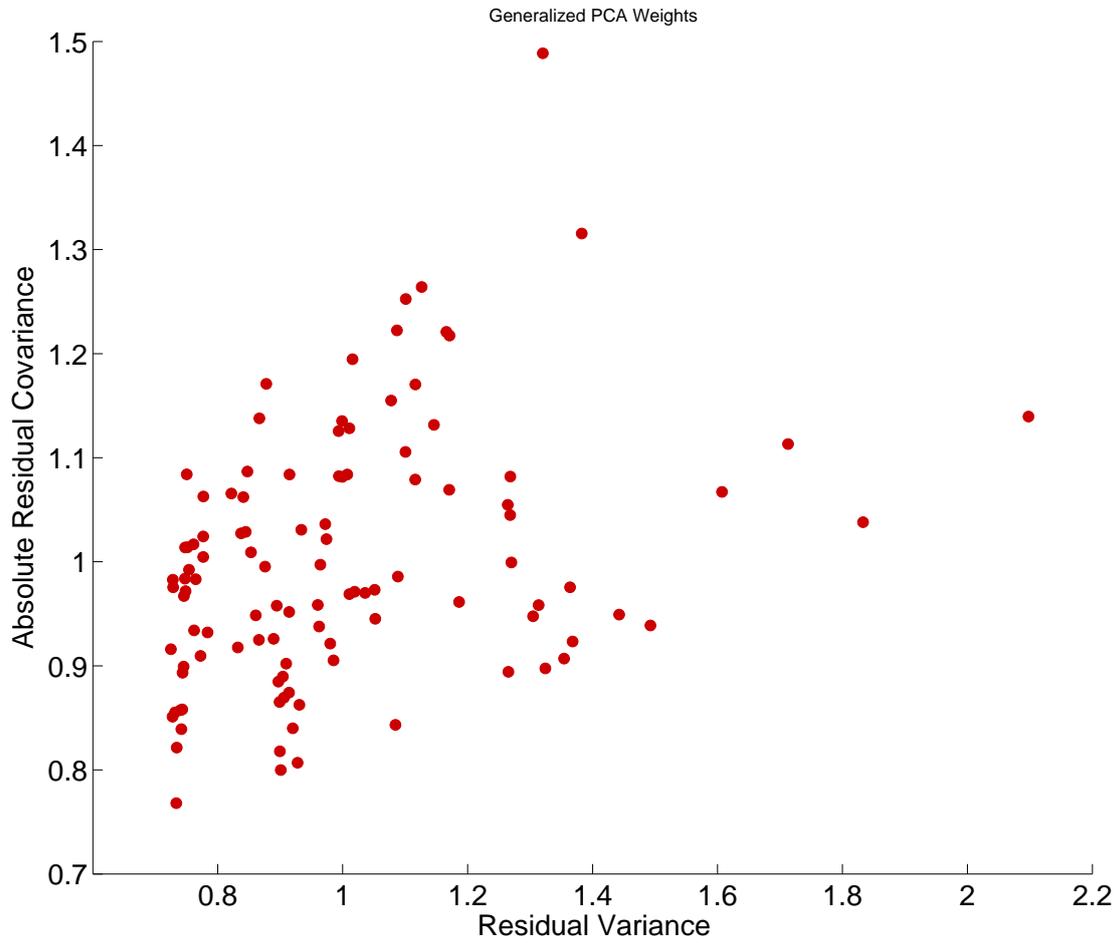


Normalized Residual Variance



Normalized Residual Absolute Covariance

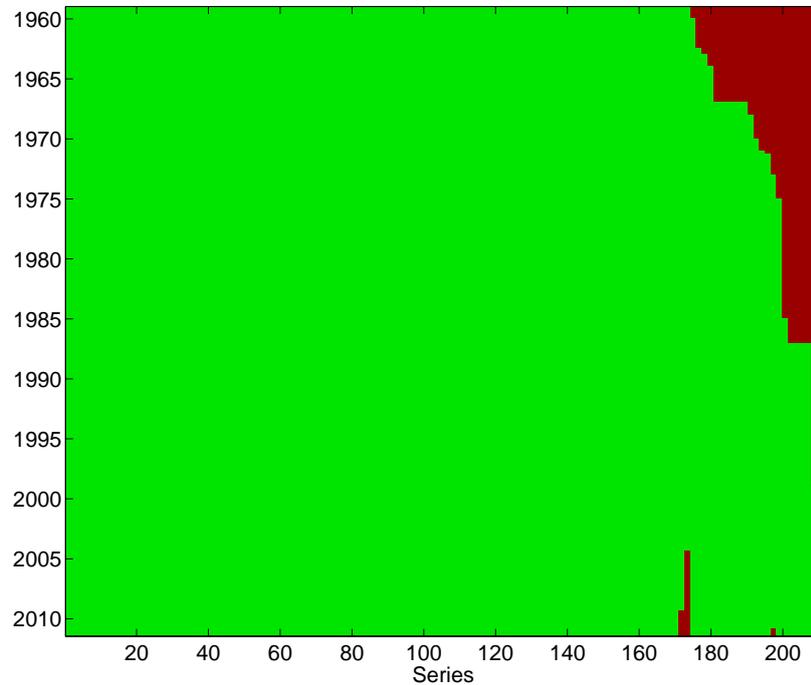




- Redundant factors can have adverse effects on common components
- Exactly redundant factors are identical to increasing the variance of a studentized data series
  - Including  $x_{it}$   $m$ -times is the same as using  $mx_{it}$
- Some evidence that excluding highly correlated factors is useful (Boivin & Ng 2006)
- Method
  1. For each series  $i$  find series with maximally correlated error, call index  $j_i$
  2. Drop series in  $\{j_i\}$  that are maximally correlated with more than 1 series
  3. For series which are each other's  $j_i$ , drop series with lower  $R^2$
- Can increase step 1 to two or even three series



- Two obvious solutions to missing data in PCA
  - Drop all series that have missing observations
  - Impute values for the missing values
- Missing data structure in SW 2012



- Some problem with unobserved states can be solved using the EM algorithm
- Consider problem of estimating means from an i.i.d. mixture

$$X_i = Y_i\mu_1 + (1 - Y_i)\mu_2 + Z_i$$

- ▶  $Y_i$  is i.i.d. Bernoulli( $p$ ),  $Z_i$  is standard normal
- ▶  $Y_i$  was observable, trivial problem (OLS)
- ▶ When  $Y_i$  is not observable, much harder
- ▶ EM algorithm will iterate across two steps:
  1. Construct “as-if”  $Y_i$  using expectations of  $Y_i$  given  $\mu_1$  and  $\mu_2$
  2. Compute

$$\hat{\mu}_1 = \frac{\sum \Pr(Y_i = 1)X_i}{\sum \Pr(Y_i = 1)} \quad \hat{\mu}_2 = \frac{\sum \Pr(Y_i = 0)X_i}{n - \sum \Pr(Y_i = 1)}$$

3. Return to 1, stopping if the means are not changing much
- ▶ Algorithm is initialized with “guesses” about  $\mu_1$  and  $\mu_2$ 
    - ▷ Example: Mean of data above median, mean of data below median
  - ▶ Consider case where  $\mu_1 = 10$ ,  $\mu_2 = -10$

- Ideally would like to solve PCA problem only for observed data
- Difficult in practice, no know closed form estimator
- Expectation-Maximization (EM) algorithm can be used to simply impute missing data
  - Replace missing with  $r$ -factor expectation (E)
  - Maximize the likelihood (M), or minimize sum of squares

## Algorithm (EM Algorithm for Imputing Missing Values in PCA)

1. Define  $w_{ij} = I [y_{ij} \text{ observed}]$  and set  $i = 0$
2. Construct  $\mathbf{X}^{(0)} = \mathbf{W} \odot \mathbf{X} + (1 - \mathbf{W}) \odot \mathbf{1}\bar{\mathbf{X}}$  where  $\mathbf{1}$  is a  $T$  by 1 vector of 1s
3. Until  $\left\| \mathbf{X}^{(i+1)} - \mathbf{X}^{(i)} \right\| < c$ :
  - a. Estimate  $r$  factors and factor loadings,  $\hat{\mathbf{F}}^{(i)}$  and  $\hat{\boldsymbol{\beta}}^{(i)}$  from  $\mathbf{X}^{(i)}$  using PCA
  - b. Construct  $\mathbf{X}^{(i+1)} = \mathbf{W} \odot \mathbf{X} + (1 - \mathbf{W}) \odot (\hat{\mathbf{F}}^{(i)} \hat{\boldsymbol{\beta}}^{(i)})$
  - c. Set  $i = i + 1$

- Can use partitioning to construct hierarchical factors
- Global and Local
  1. Extract 1 or more factors from all series
  2. For each regions or country  $j$ , regress series from country  $j$  on Global Factors, and extract 1 or more factors from residuals
    - Country factors uncorrelated with Global, but not local from other regions/countries
- Nominal and Real
  1. Extract 1 or more general factors
  2. For each group real/nominal series, regress on general factors and then extract factors from residuals