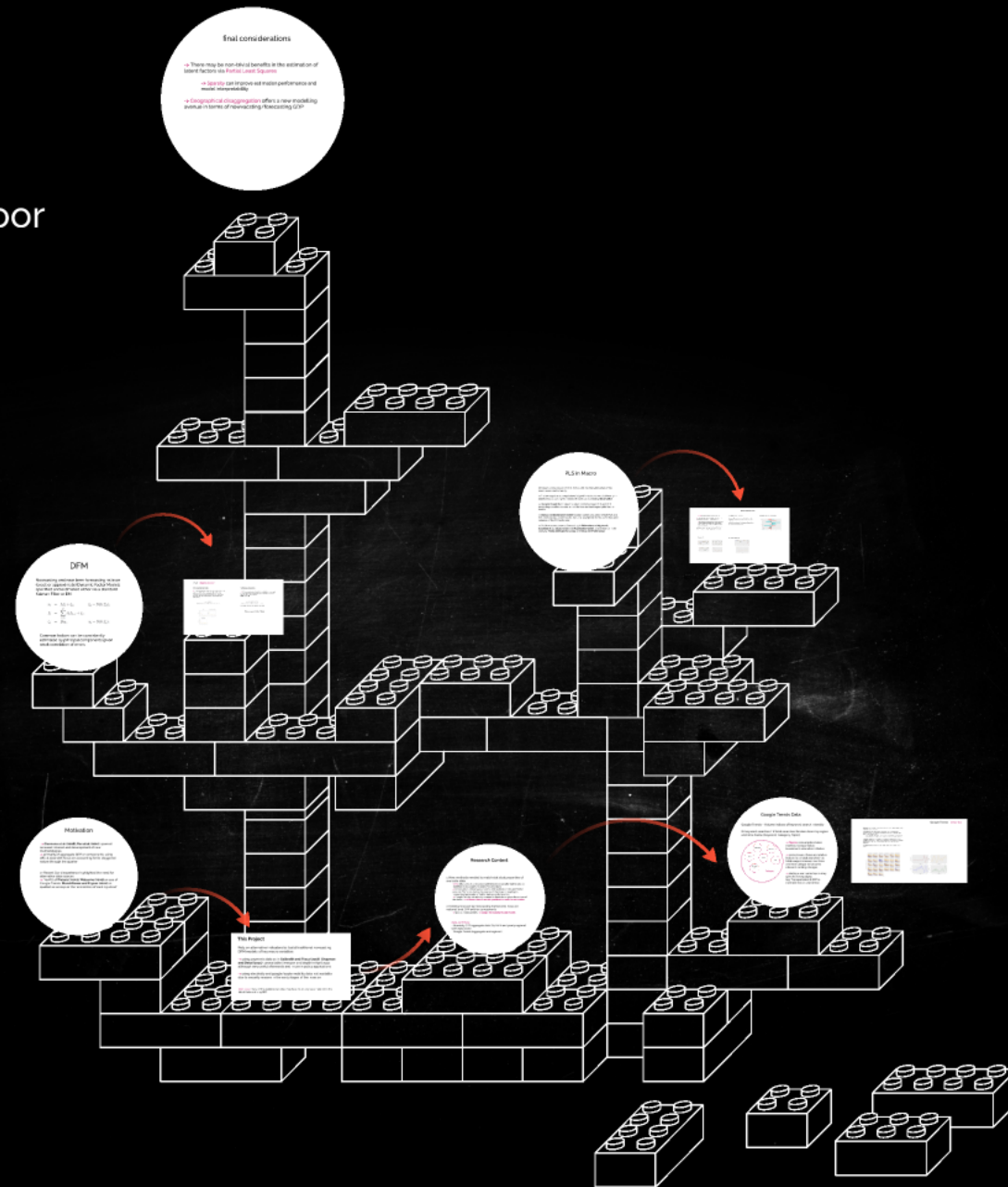


# Sparse Warcasting

## Forecasting in a data-rich but statistics-poor environment



# Motivation

- > **Giannone et al. (2008); Doz et al. (2011)** spurred renewed interest and development of new methodologies
  - > primarily of aggregate GDP or components, using official data with focus on accounting for its staggered nature through the quarter
- >> Recent Covid experience highlighted the need for alternative data sources
- > VoxEU of **Diebold (2020); Woloszko (2020)** on use of Google Trends; **Blanchflower and Bryson (2021)** on qualitative surveys or the "economics of walking about"

# Motivation

-> **Giannone et al. (2008); Doz et al. (2011)** spurred renewed interest and development of new methodologies

-> primarily of aggregate GDP or components, using official data with focus on accounting for its staggered nature through the quarter

>> Recent Covid experience highlighted the need for alternative data sources

-> VoxEU of **Diebold (2020); Woloszko (2020)** on use of Google Trends; **Blanchflower and Bryson (2021)** on qualitative surveys or the "economics of walking about"

>> The Feb 24th russian invasion led to a freeze of all official data gathering by local and national statistical agencies

--> only left with alternative data sources

# This Project

Rely on alternative indicators to build traditional nowcasting DFM models of key macro variables

-> using payments data as in **Galbraith and Tkacz (2018), Chapman and Desai (2021)** : unavailable timespan and depth in April 2022 although very useful afterwards and in use in policy applications

-> using electricity and google/apple mobility data: not available due to security reasons in the early stages of the invasion;

**Main issue:** many of the available variables may have no, or very weak, relation to the latent factors driving GDP

## Research Context

### 1. New methods needed to match statistical properties of available data

- > No official statistics released until mid 2022; use alternative data as identified in development economics literature
- > Lit considers primarily peace-time GDP estimation: NL and Twitter possibly "flip" signs during the course of the invasion leading to inconsistent parameters; Twitter data recently for a fee
- > Google Trends not entirely immune to these issue given the nature of the shock ->> **different factor model specification and/or estimation**

### 2. Existing forecasting/nowcasting frameworks focus on national level GDP and/or components

- > Spatial heterogeneity ->> **scope for regional factor model**

#### **data-summary:**

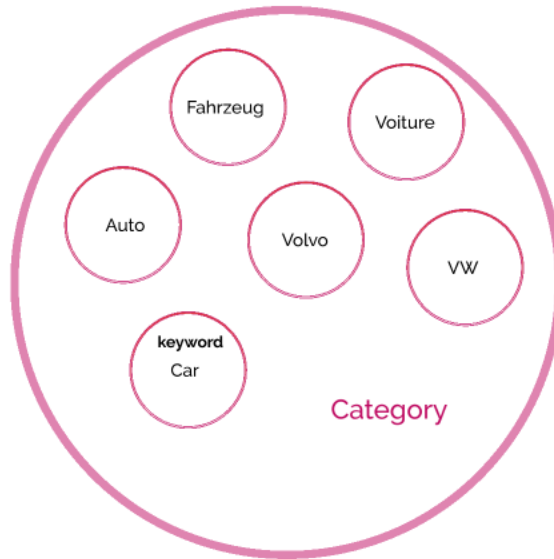
- Quarterly GDP aggregate data (Q4 2021) and yearly regional GDP data (2020)
- Google Trends (aggregate and regional)



# Google Trends Data

Google Trends = Volume indices of keyword search intensity

(# keyword searches / # total searches) broken down by region and time frame [Keyword, Category, Topic]



-> **Topics**: consumption, labor markets, transportation, investment, education, inflation

-> some issues: these are relative indices (w.r.t. total searches); as total usage increases over time and new categories become relevant, ranking changes.

-> during a war, some topics may give the wrong signal (eg. Transportation & GDP in normal times vs. war times)

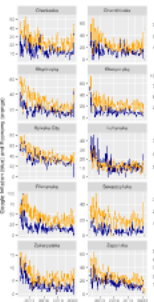
> Ettredge et al. (2005) is one of the first to use Google Trends data to forecast US unemployment.

> Askits Zimmermann (2009) "Google Trends and Forecasting" use Google searches to forecast figures several months ahead.

> This is particularly relevant as in 2008, Google Trends data was usually delayed several months as compared to official statistics.

> Choi and Varian (2010) "Predicting the future" strong point for the use of Google Trends data for variables such as car sales, unemployment, and confidence.

> Wu and Brynjolfsson (2010) leverage Google Trends data to predict prices



# Google Trends - empirics

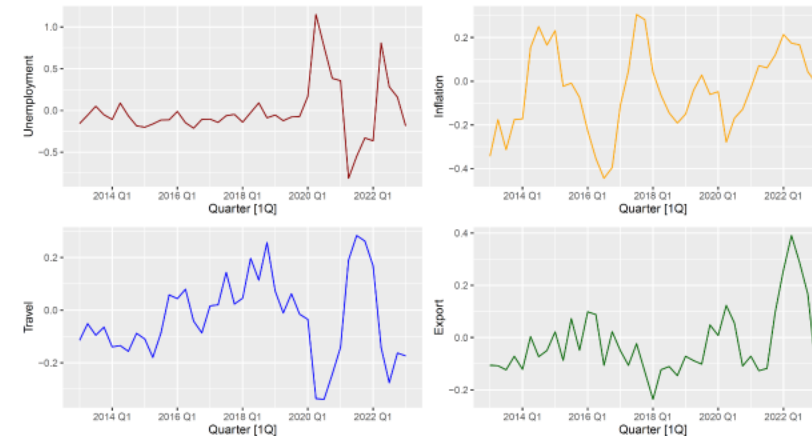
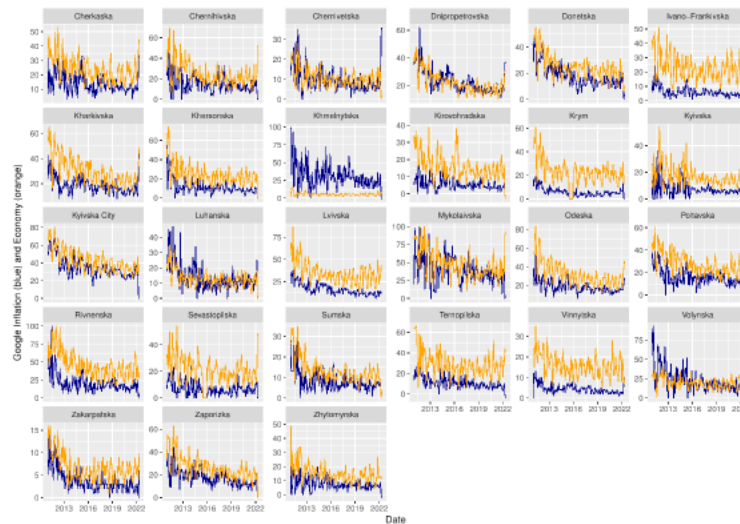
> **Ettredge et al. (2005)** is one of the earlier references using Google Search activity to forecast US unemployment.

> **Askits Zimmermann (2009)** "Google Econometrics and Unemployment Forecasting" use Google searches related to unemployment to forecast official figures several months ahead.

> This is particularly relevant as in 2008-2009, data releases on key macrovariables are usually delayed several months as compared to observed macro and financial shocks

> **Choi and Varian (2010)** "Predicting the Present with Google Trends" makes a strong point for the use of Google trends data to nowcast a multitude of economic variables such as car sales, unemployment claims, travel, and consumer confidence.

> **Wu and Brynjolfsson (2010)** leverage Google Search Data to forecast house prices



> **Ettredge et al. (2005)** is one of the earlier references using Google Search activity to forecast US unemployment.

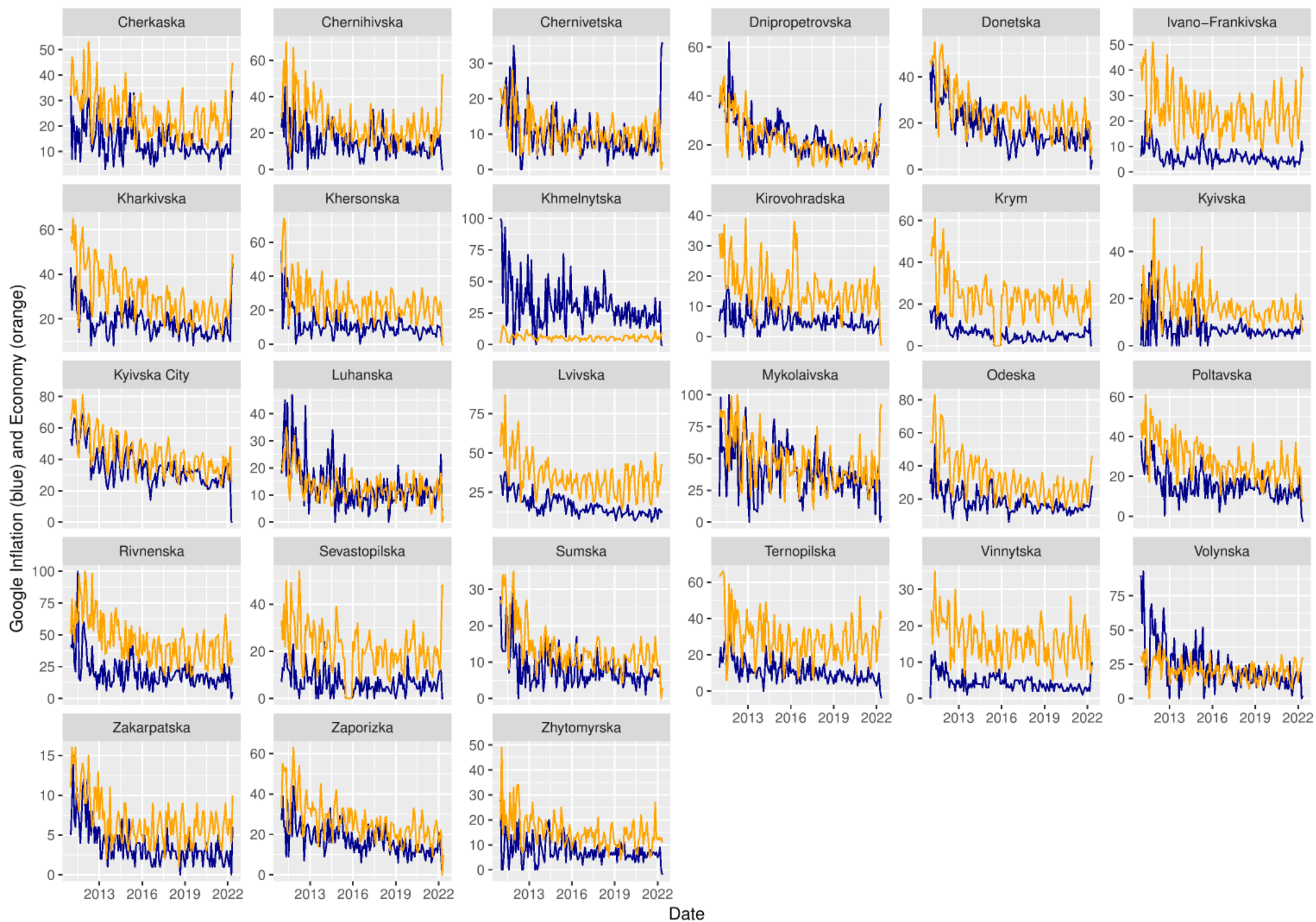
> **Askits Zimmermann (2009)** "Google Econometrics and Unemployment Forecasting" use Google searches related to unemployment to forecast official figures several months ahead.

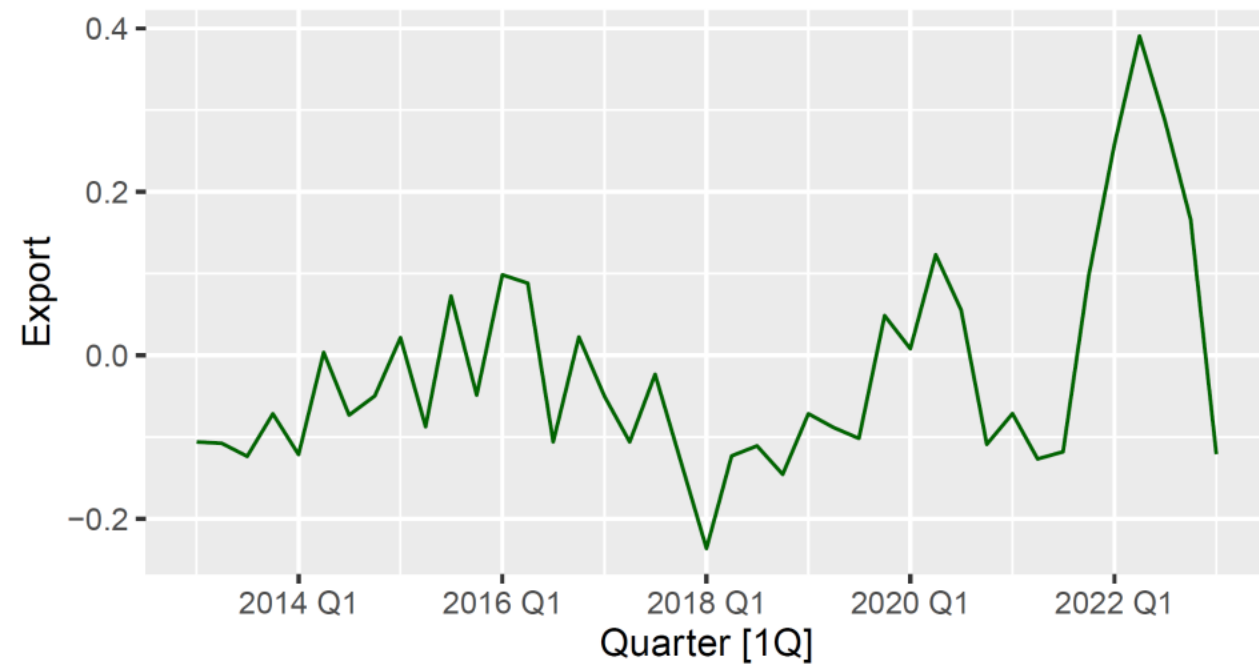
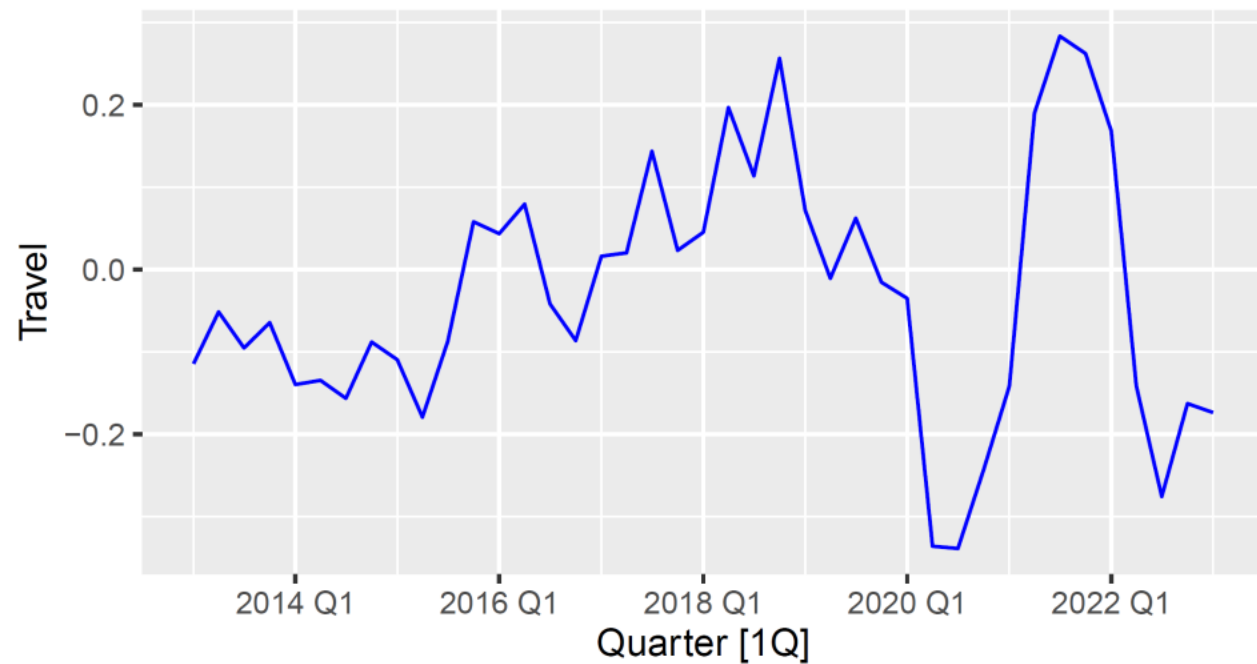
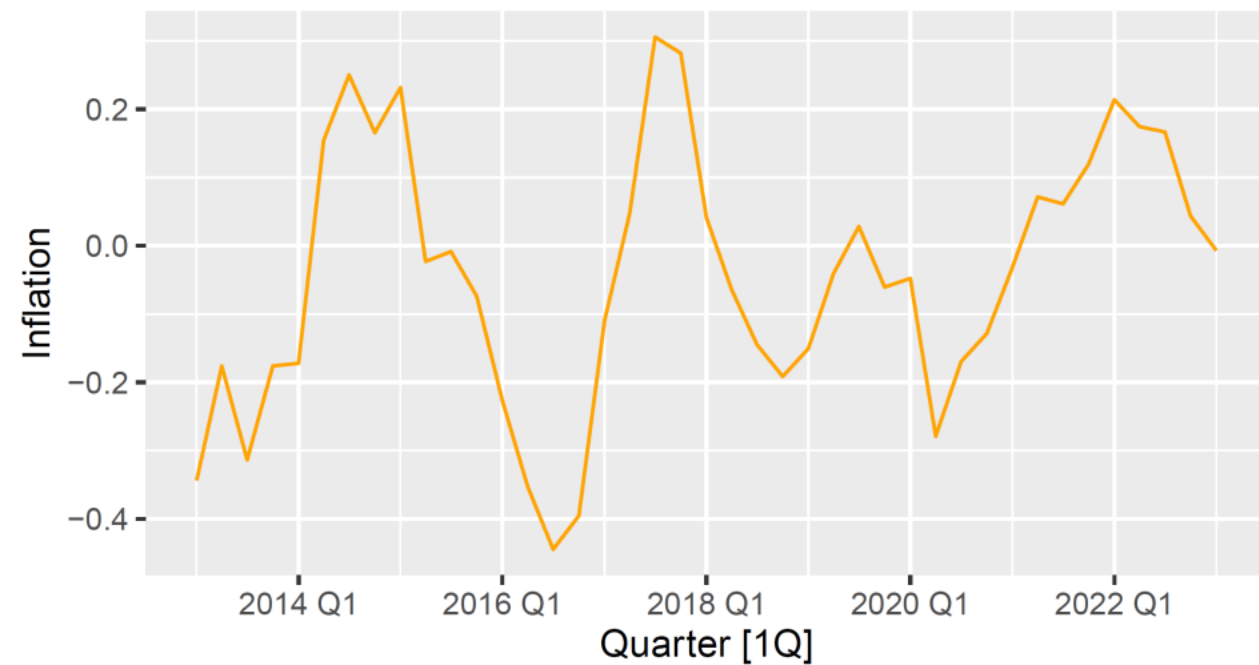
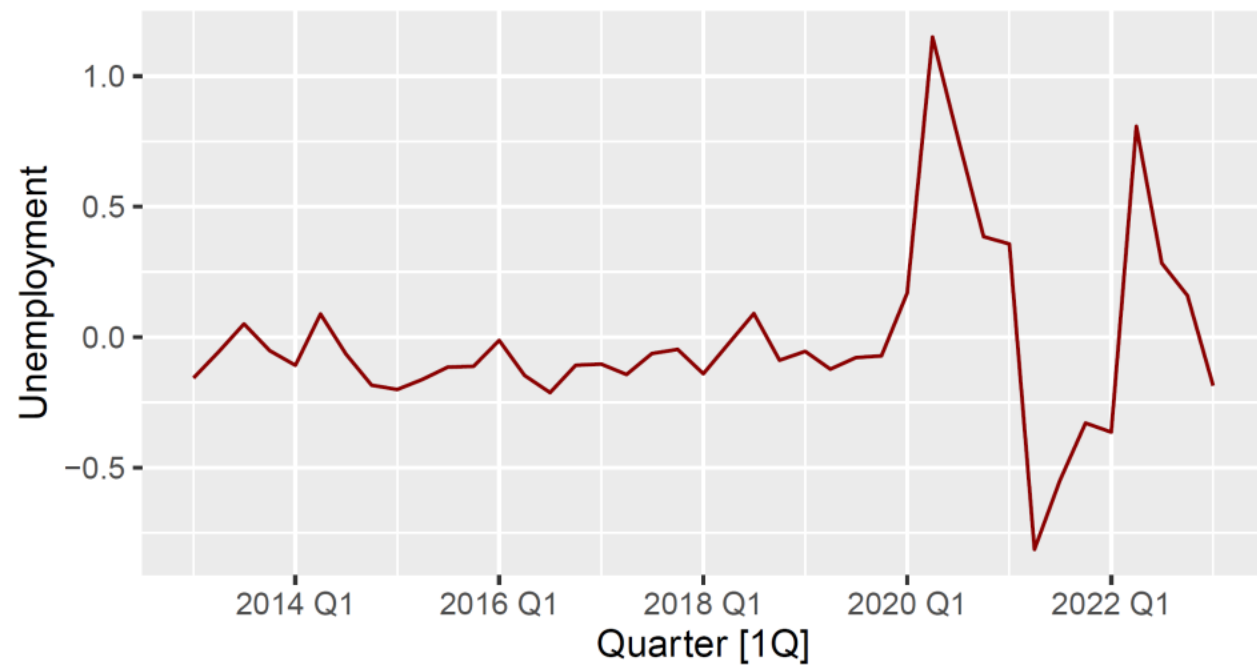
> This is particularly relevant as in 2008-2009, data releases on key macrovariables are usually delayed several months as compared to observed macro and financial shocks

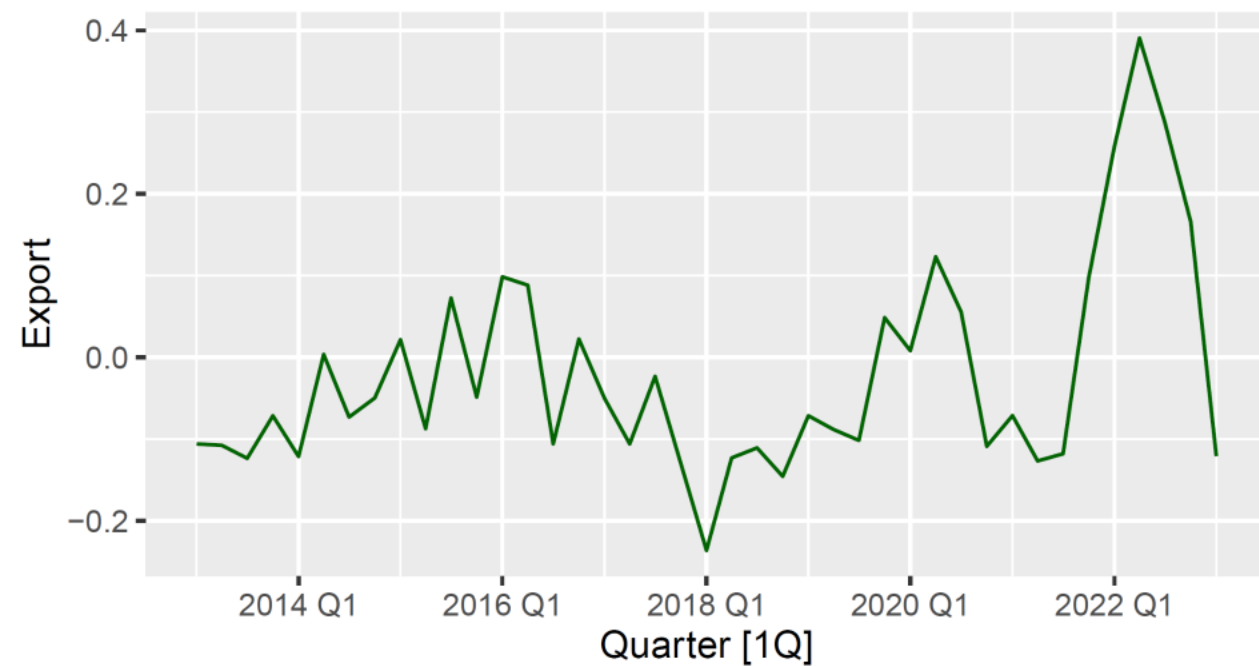
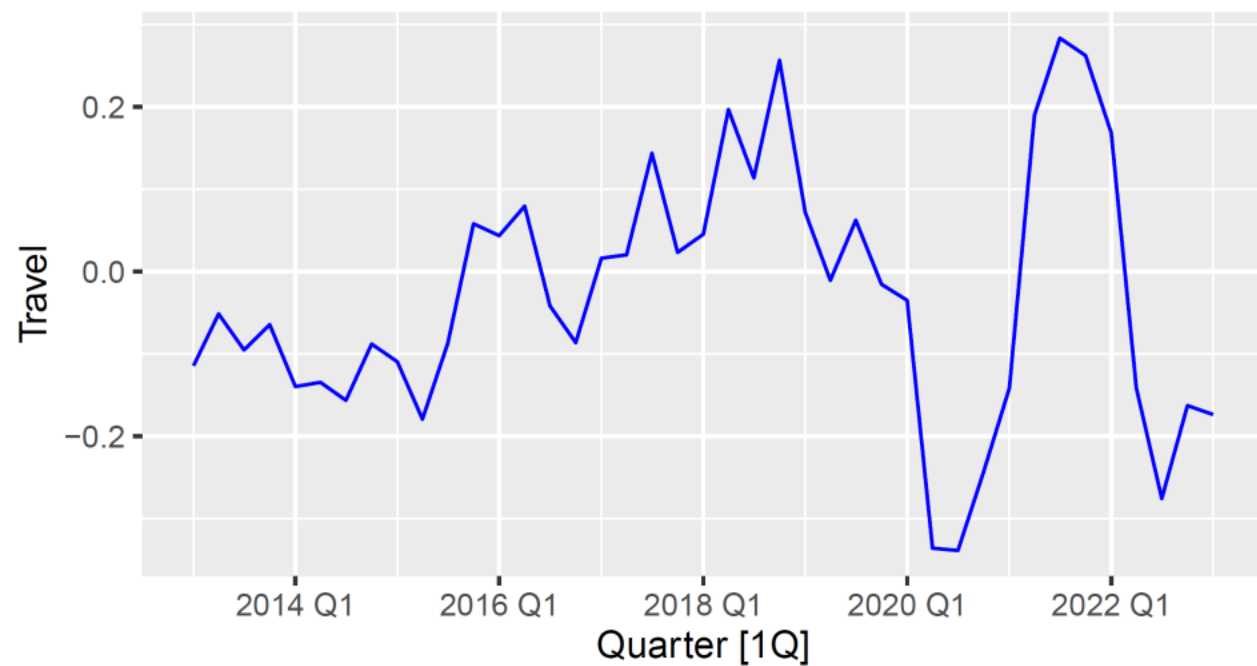
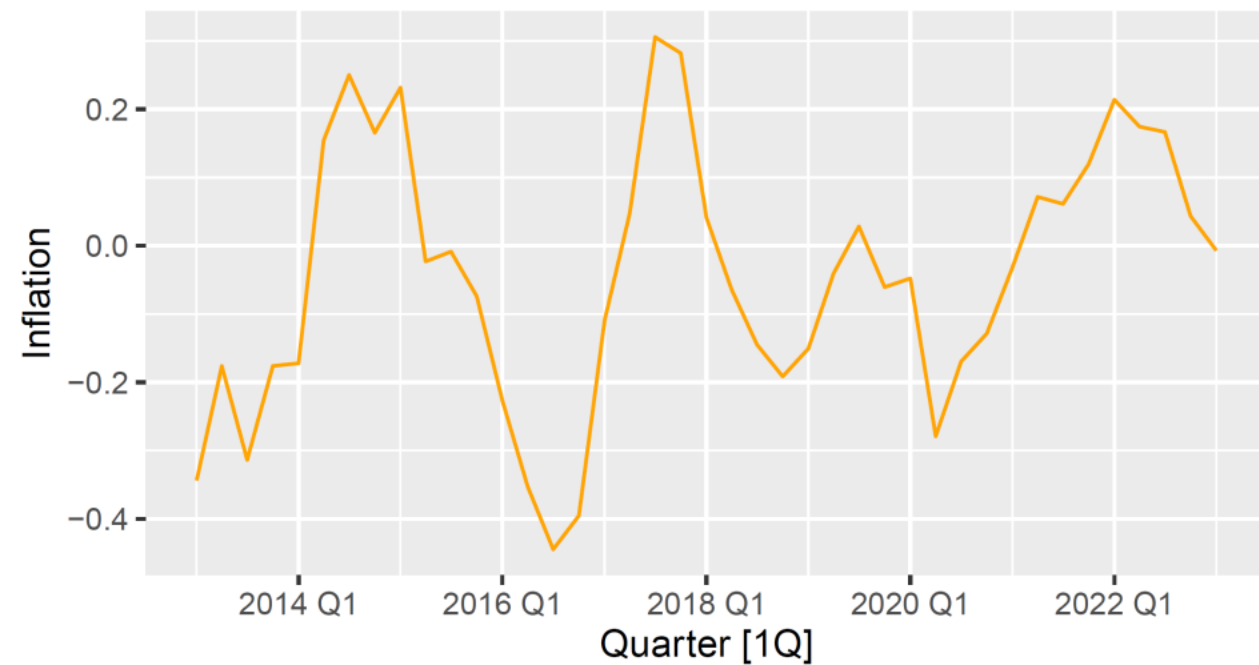
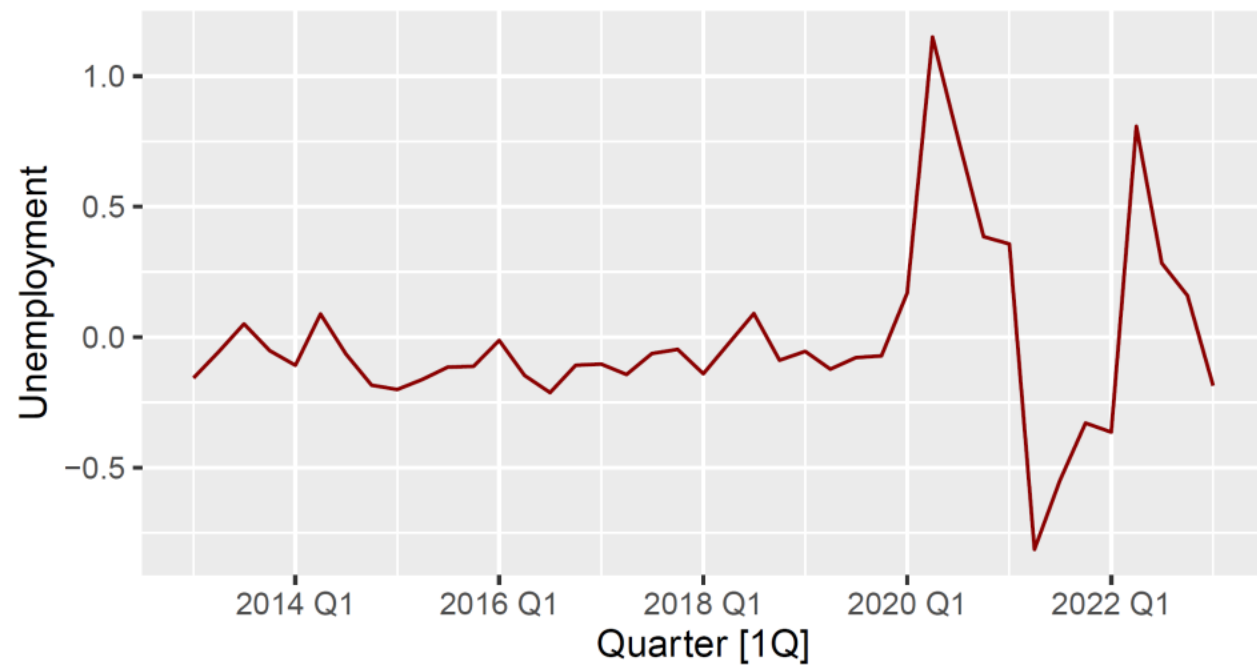
> **Choi and Varian (2010)** "Predicting the Present with Google Trends" makes a strong point for the use of Google trends data to nowcast a multitude of economic variables such as car sales, unemployment claims, travel, and consumer confidence.

> **Wu and Brynjolfsson (2010)** leverage Google Search Data to forecast house prices









# Google Trends - empirics

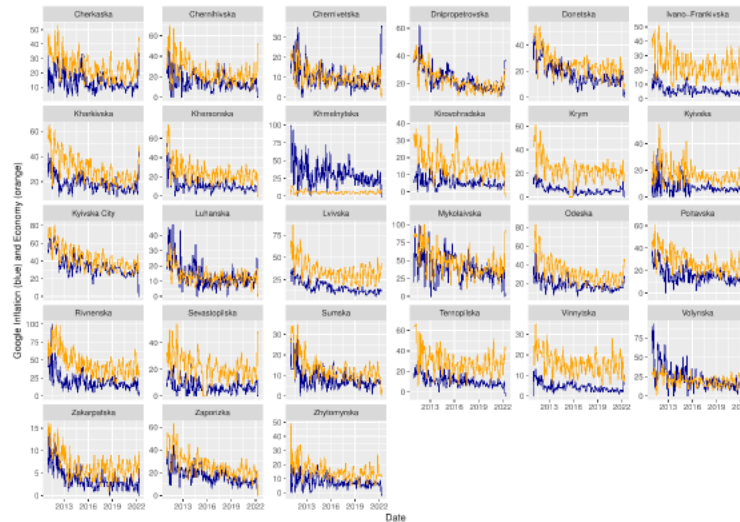
> **Ettredge et al. (2005)** is one of the earlier references using Google Search activity to forecast US unemployment.

> **Askits Zimmermann (2009)** "Google Econometrics and Unemployment Forecasting" use Google searches related to unemployment to forecast official figures several months ahead.

> This is particularly relevant as in 2008-2009, data releases on key macrovariables are usually delayed several months as compared to observed macro and financial shocks

> **Choi and Varian (2010)** "Predicting the Present with Google Trends" makes a strong point for the use of Google trends data to nowcast a multitude of economic variables such as car sales, unemployment claims, travel, and consumer confidence.

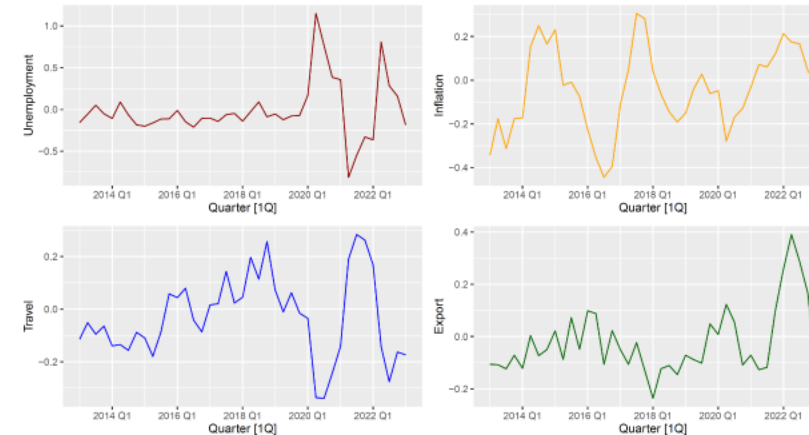
> **Wu and Brynjolfsson (2010)** leverage Google Search Data to forecast house prices



-> About 35 trends are used

-> Monthly time series contain plenty of variation (too much?) >> Q

-> Number implies some shrinkage/ dimensionality reduction is needed



# DFM

Nowcasting and near term forecasting rely on (exact or approximate) Dynamic Factor Models specified and estimated either via a standard Kalman Filter or EM

$$\begin{aligned}x_t &= \Lambda f_t + \xi_t, & \xi_t &\sim \mathbb{N}(0, \Sigma_\xi), \\f_t &= \sum_{i=1}^p A_i f_{t-i} + \zeta_t, \\ \zeta_t &= B\eta_t, & \eta_t &\sim \mathbb{N}(0, I_q).\end{aligned}$$

Common factors can be consistently estimated by principal components given weak correlation of errors





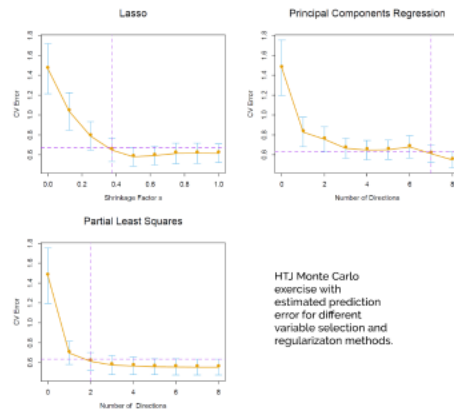
# PCR - digging deeper

## Principal Components

-> Given a data matrix  $\mathbf{X}$  (N obs x p variables). PCA will perform a SVD of the centered matrix  $\mathbf{X}^*$  to find directions in the column space of  $\mathbf{X}^*$  that have small variance (with direction vectors  $\mathbf{v}$  independent of each other)

$$\max_{\alpha} \text{Var}(\mathbf{X}\alpha)$$

subject to  $\|\alpha\| = 1, \alpha^T \mathbf{S} \mathbf{v}_{\ell} = 0, \ell = 1, \dots, m-1,$



## Partial Least Squares

-> PLS is a supervised method which identifies the components or factors ( $\phi$ ) to be independent of each other but also have high correlation with a target  $\mathbf{y}$   
**{Wold et al. 1984}**

$$\max_{\alpha} \text{Corr}^2(\mathbf{y}, \mathbf{X}\alpha) \text{Var}(\mathbf{X}\alpha)$$

subject to  $\|\alpha\| = 1, \alpha^T \mathbf{S} \hat{\phi}_{\ell} = 0, \ell = 1, \dots, m-1.$

Hastie, Tibshirani, Friedman 2nd ed. "The Elements of Statistical Learning"

## PLS as a Latent Factor Model

# Principal Components

-> Given a data matrix  $\mathbf{X}$  (N obs x p variables). PCA will perform a SVD of the centered matrix  $\mathbf{X}^*$  to find directions in the column space of  $\mathbf{X}^*$  that have small variance (with direction vectors  $v$  independent of each other)

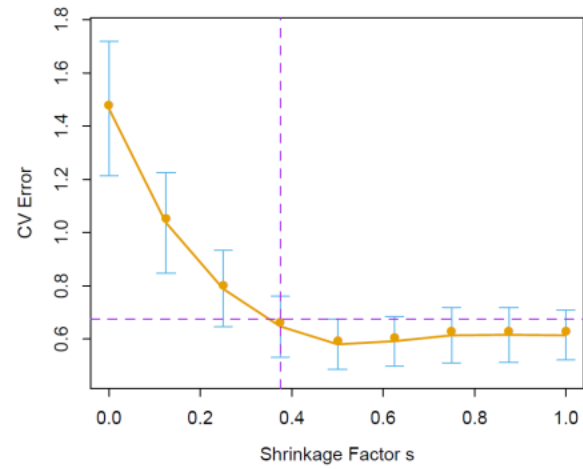
$$\begin{aligned} & \max_{\alpha} \text{Var}(\mathbf{X}\alpha) \\ & \text{subject to } \|\alpha\| = 1, \alpha^T \mathbf{S}v_{\ell} = 0, \ell = 1, \dots, m-1, \end{aligned}$$

# Partial Least Squares

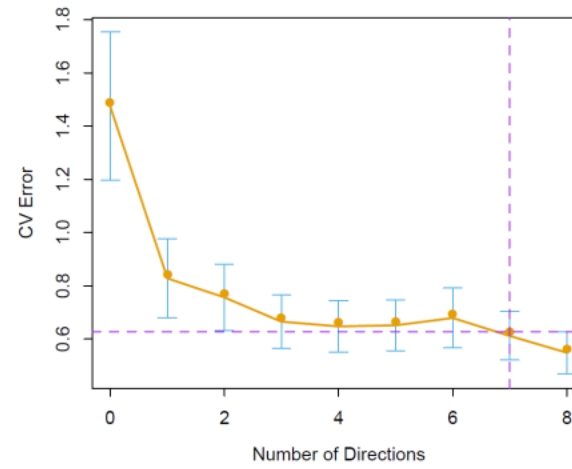
-> PLS is a supervised method which identifies the components or factors ( $\phi$ ) to be independent of each other but also have high correlation with a target  $\mathbf{y}$   
**{Wold et al. 1984}**

$$\begin{aligned} & \max_{\alpha} \text{Corr}^2(\mathbf{y}, \mathbf{X}\alpha) \text{Var}(\mathbf{X}\alpha) \\ & \text{subject to } \|\alpha\| = 1, \alpha^T \mathbf{S} \hat{\varphi}_{\ell} = 0, \ell = 1, \dots, m - 1. \end{aligned}$$

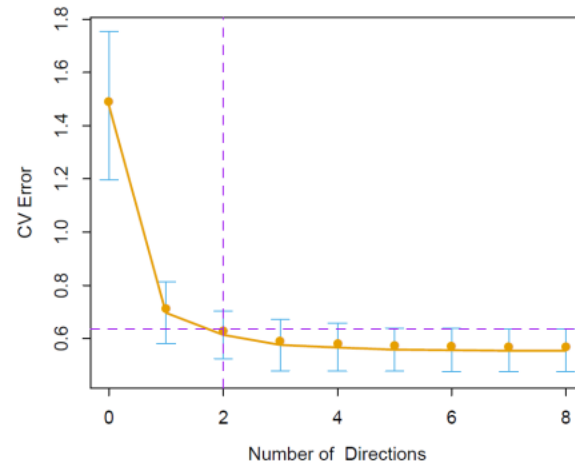
### Lasso



### Principal Components Regression



### Partial Least Squares



HTJ Monte Carlo  
exercise with  
estimated prediction  
error for different  
variable selection and  
regularization methods.

# PLS as a Latent Factor Model



# PLS as a Latent Factor Model

# PLS as a Latent Factor Model

$$T = X \times W^*, \quad X \in \mathbb{R}^{n \times p}, T \in \mathbb{R}^{n \times K}, K < p$$

$$X = T \times P' + \epsilon, \quad P \in \mathbb{R}^{p \times K}$$

$$y = T \times C' + \xi, \quad C \in \mathbb{R}^{1 \times K}$$

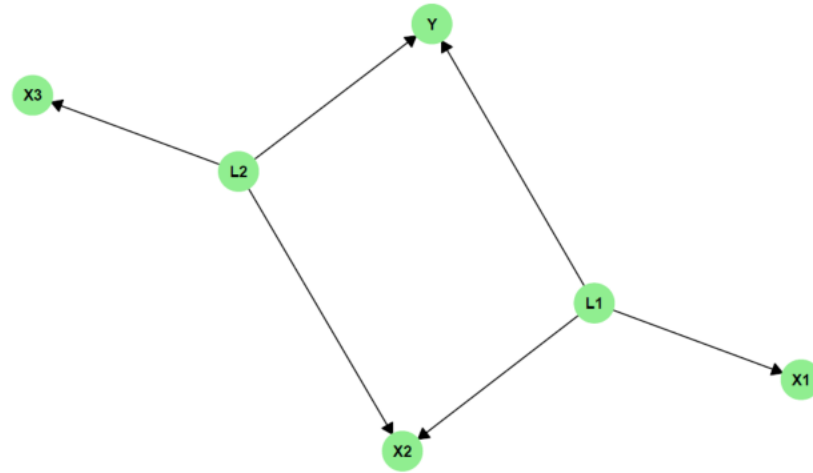


Figure 2 – Network Graph of a Simple Latent Factor Model

# PLS in Macro

Although a close cousin of PCR, PLS is a bit the Harry&Meghan of the macro econometric family

>> first emerged as a *computational algorithm* to tackle multicollinearity in bioinformatics lacking the needed theoretical scaffolding **Wold (1984)**

>> **Helland (1990)** illuminates the relationship between PLS and PCR presenting conditions under which the two methodologies yield similar results.

>> **Stoica and Söderström (1998)** explore conditions under which PCR and PLS produce equivalent results, deriving asymptotic formulas for bias and variance of the PLS estimator.

>> PLS in macroeconomic forecasting in **Eickmeier and Ng (2011)**, **Cubadda et al. (2013)**, **Groen and Kapetanios (2016)**, and in finance in the study by **Preda and Saporta (2005)** and **Kelly and Pruitt (2015)**

Although a close cousin of PCR, PLS is a bit the Harry&Meghan of the macro econometric family

>> first emerged as a *computational algorithm* to tackle multicollinearity in bioinformatics lacking the needed theoretical scaffolding **Wold (1984)**

>> **Helland (1990)** illuminates the relationship between PLS and PCR presenting conditions under which the two methodologies yield similar results.

>> **Stoica and Söderström (1998)** explore conditions under which PCR and PLS produce equivalent results, deriving asymptotic formulas for bias and variance of the PLS estimator.

>> PLS in macroeconomic forecasting in **Eickmeier and Ng (2011)**, **Cubadda et al. (2013)**, **Groen and Kapetanios (2016)**, and in finance in the study by **Preda and Saporta (2005)** and **Kelly and Pruitt (2015)**

# ML model overview

## Sparsity to tackle asymptotic inconsistency risk

-> Chun and Keles (2010) indicate challenges to asymptotic consistency of the PLS estimator in a "large p small n" context, with *fixed p1 relevant and increasing p - p1 irrelevant variables*.

-> The intuition for the lack of asymptotic consistency comes from the ridge-like nature of the PLS algorithm. Given that PLS latent factors load on all variables available in X, a larger fraction of irrelevant variables weaken the ability of the algorithm to identify the true factor directions.

-> Sparsity is achieved via variable selection in a multitude of ways, depending on the joint specificities of data sample and machine learning model (Lasso like, FWD or BWD Variable Selection, GA)

## Variable Selection - overview

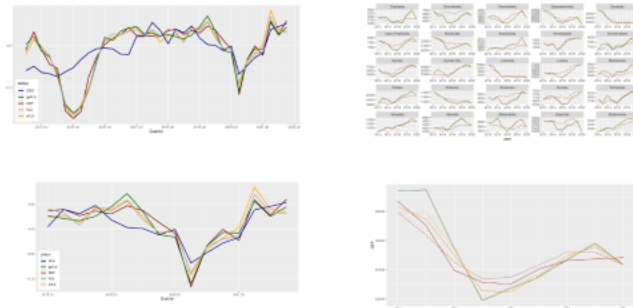
-> A wrapper (GA) and an embedded method (sPLS) are used to induce sparsity: as a results, variable selection leads to improved interpretability

-> sPLS of Chun and Keles (2010) introduces a LASSO penalty in the optimization problem and jointly selects the optimal number of latent factors and the amount of penalty

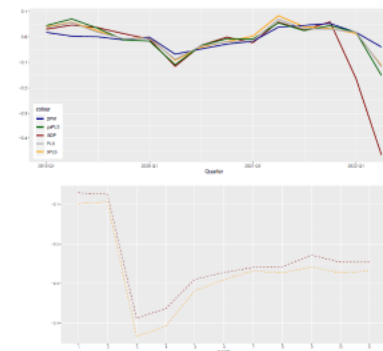
## Variable Selection - the GA algorithm



## Insample Fit



## Out of sample Forecast



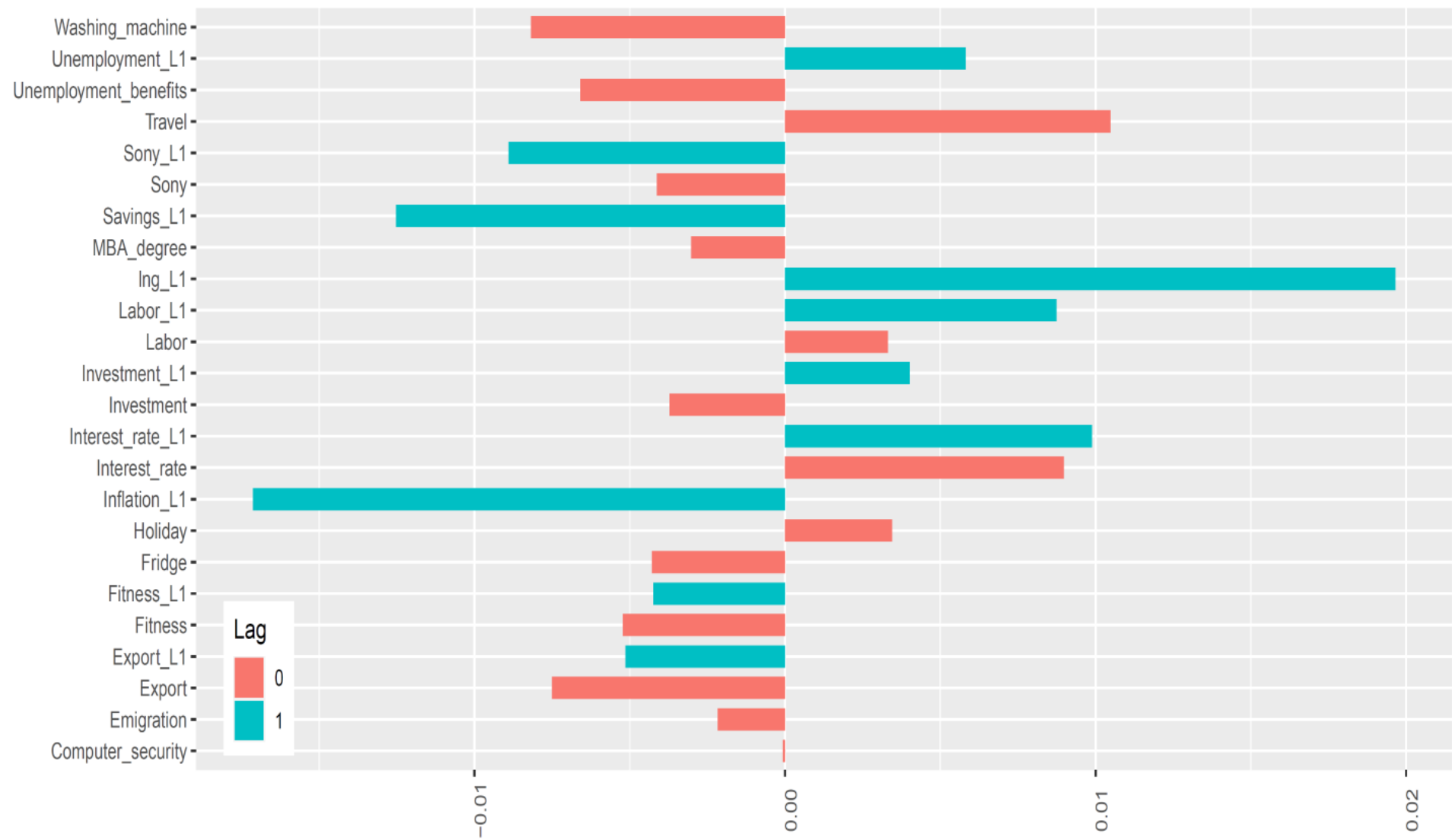


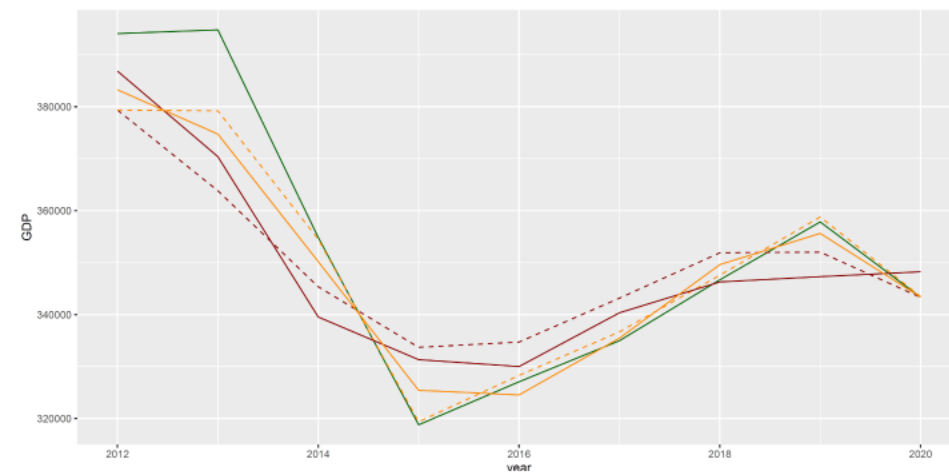
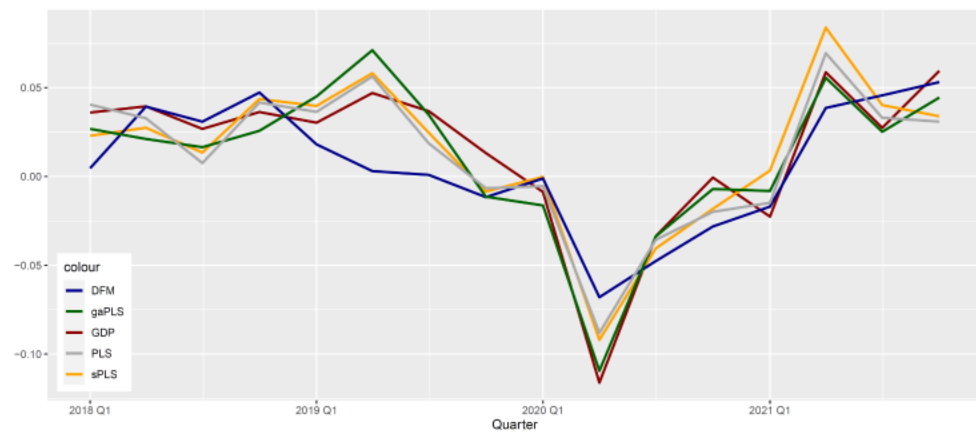
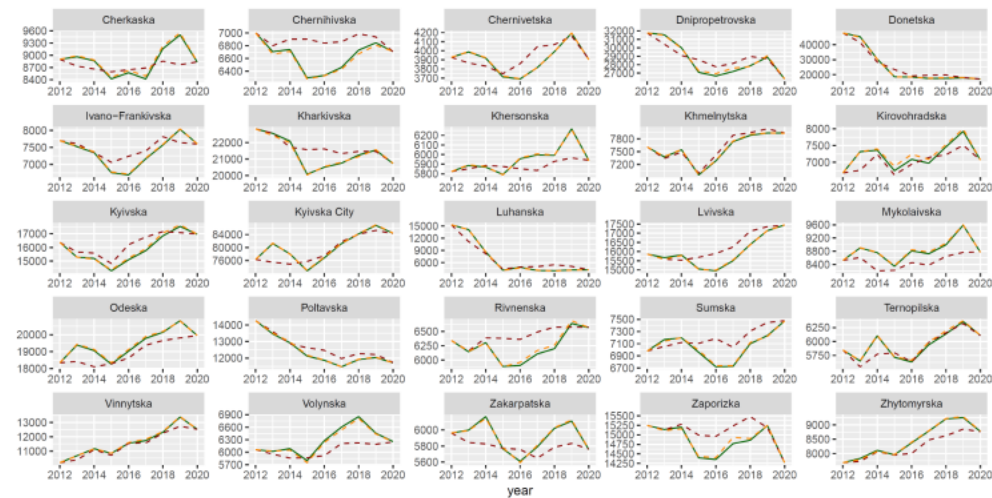
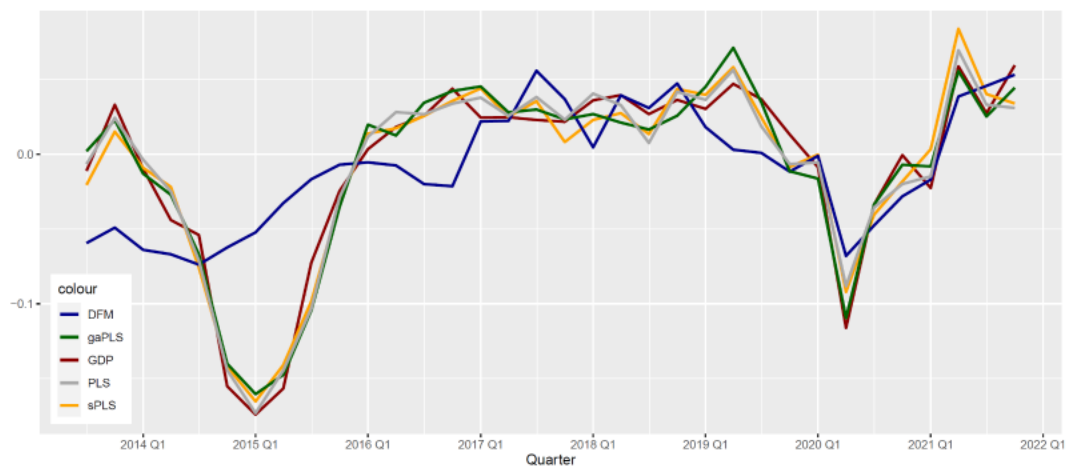
-> Chun and Keles (2010) indicate challenges to asymptotic consistency of the PLS estimator in a "large  $p$  small  $n$ " context, with *fixed  $p_1$  relevant and increasing  $p - p_1$  irrelevant* variables.

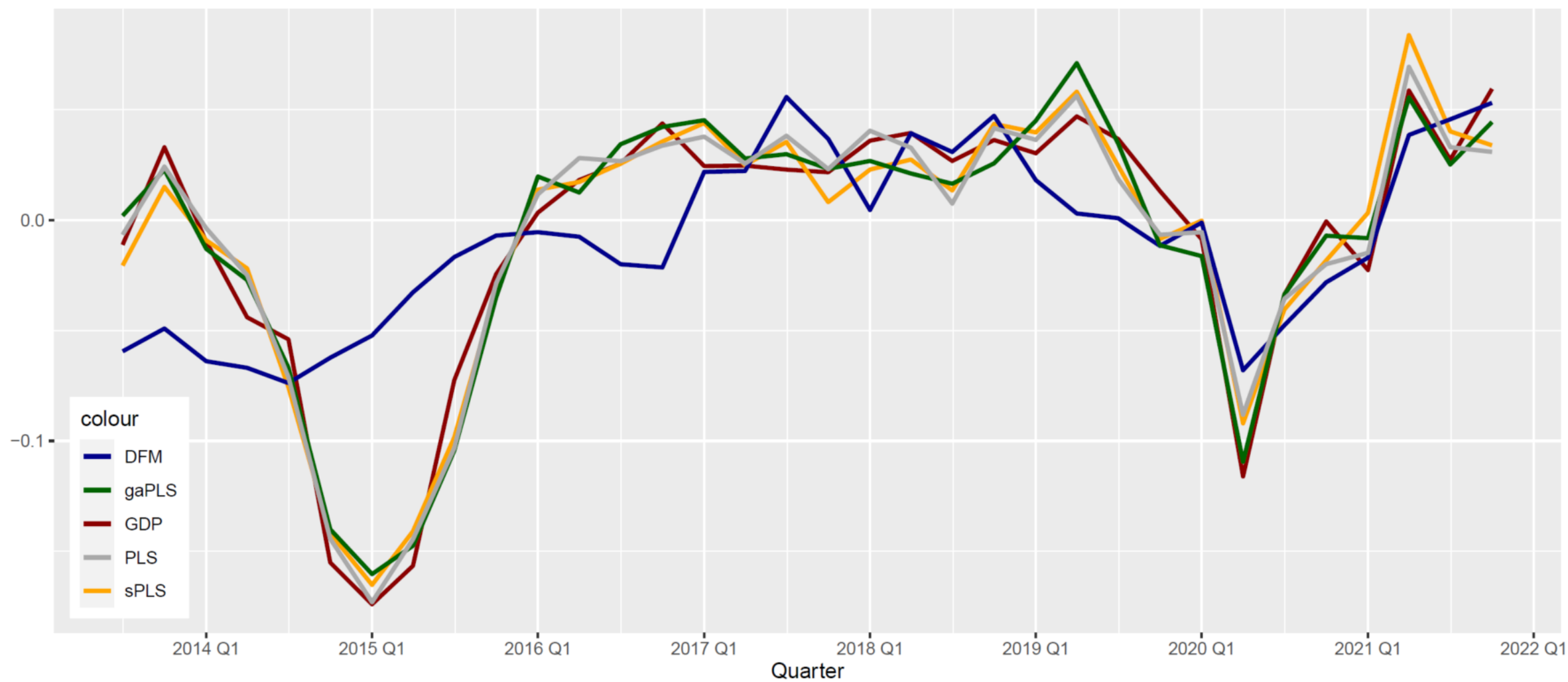
-> The intuition for the lack of asymptotic consistency comes from the ridge-like nature of the PLS algorithm. Given that PLS latent factors load on all variables available in  $X$ , a larger fraction of irrelevant variables weaken the ability of the algorithm to identify the true factor directions.

-> Sparsity is achieved via variable selection in a multitude of ways, depending on the joint specificities of data sample and machine learning model (Lasso like, FWD or BWD Variable Selection, GA)

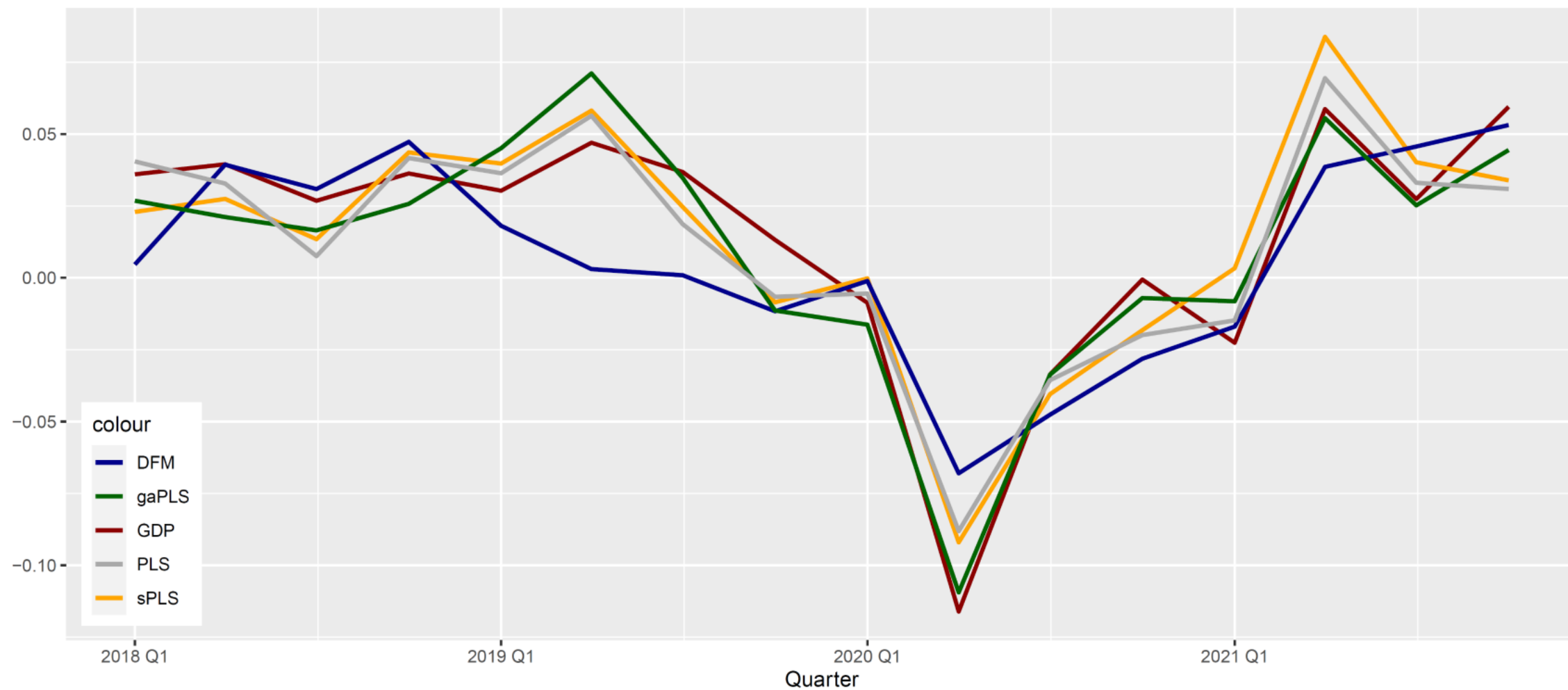
- > A wrapper (GA) and an embedded method (sPLS) are used to induce sparsity: as a results, variable selection leads to improved interpretability
- > sPLS of Chun and Keles (2010) introduces a LASSO penalty in the optimization problem and jointly selects the optimal number of latent factors and the amount of penaty

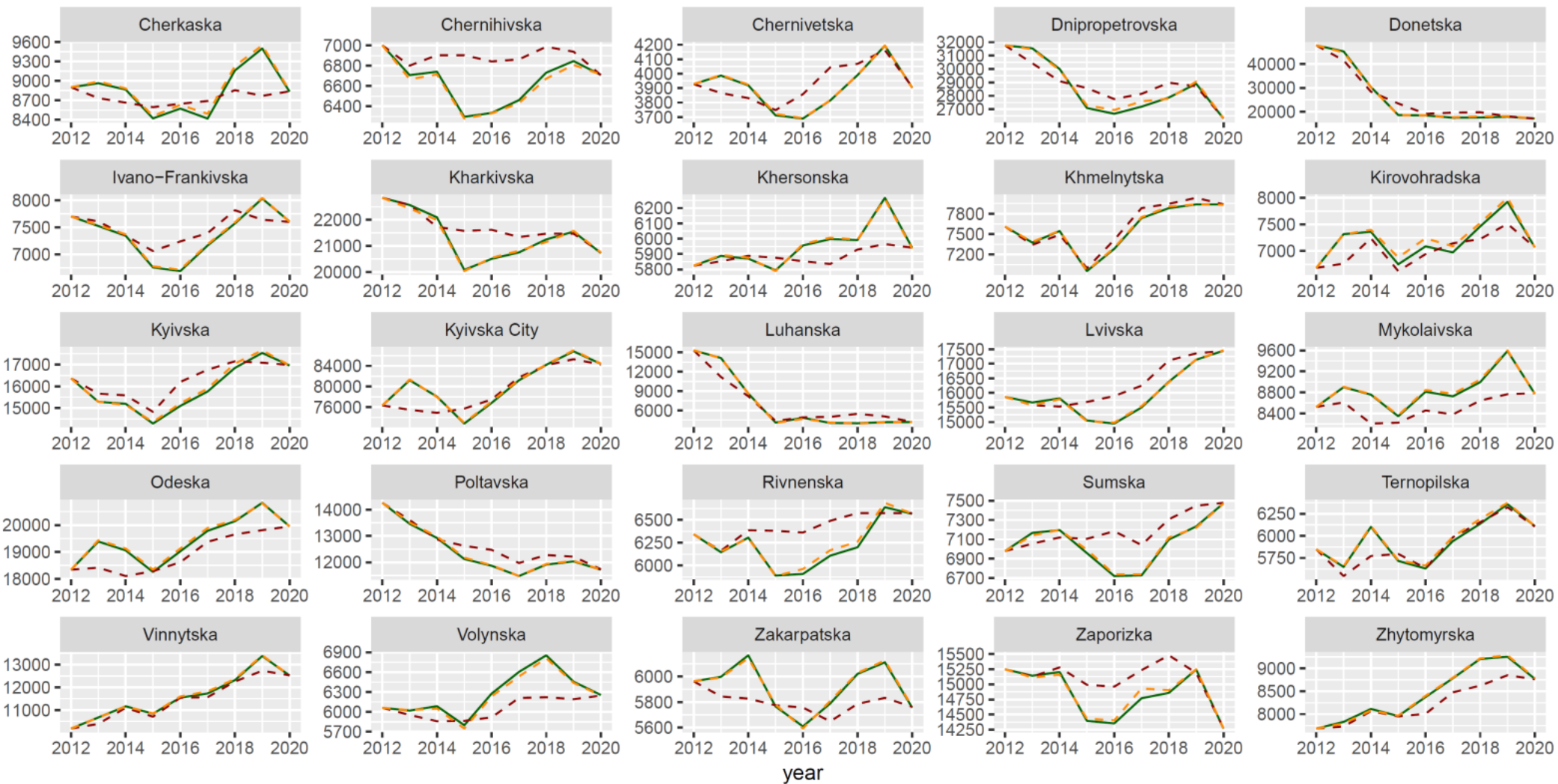


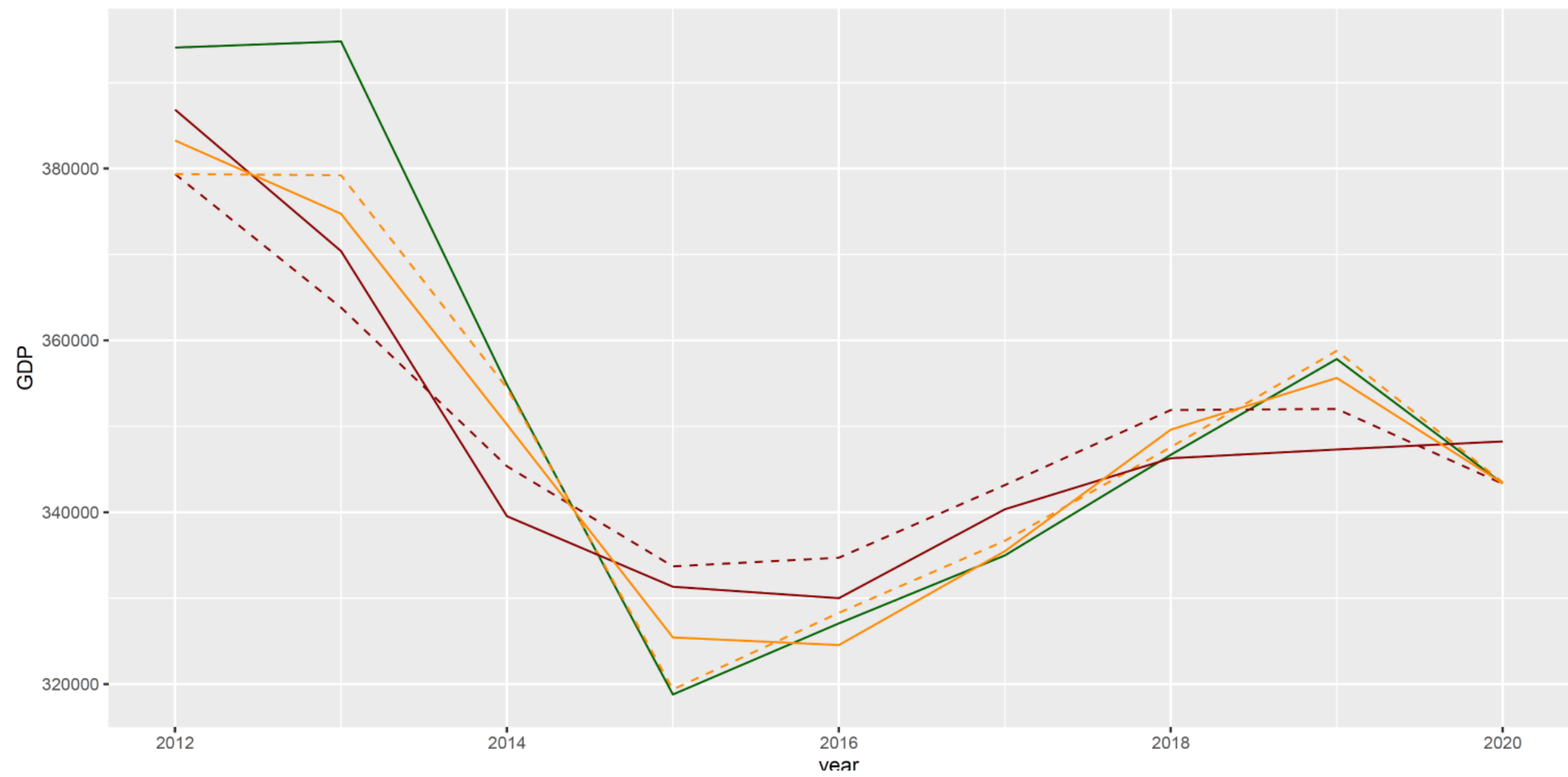


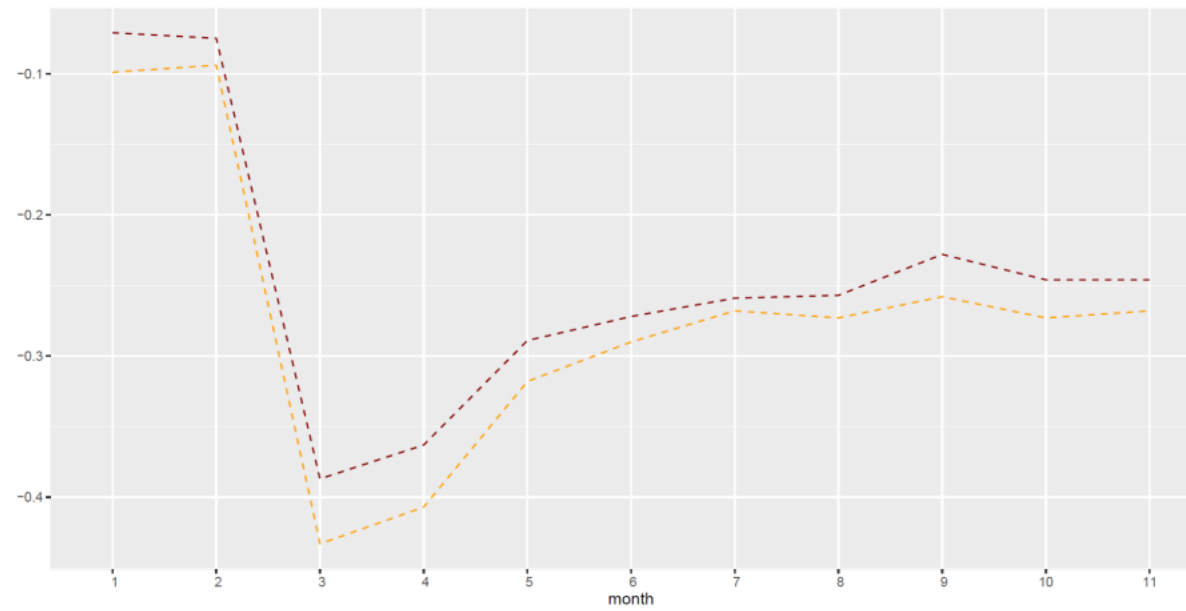
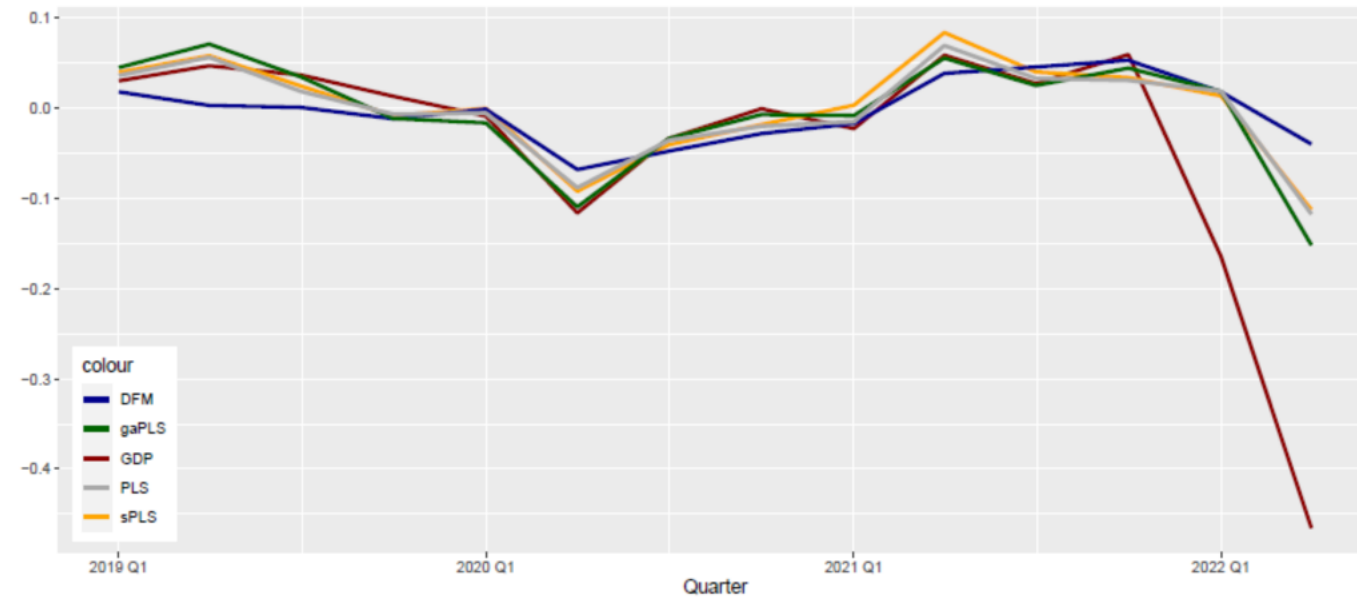












## final considerations

- > There may be non-trivial benefits in the estimation of latent factors via **Partial Least Squares**
- > **Sparsity** can improve estimation performance and model interpretability
- > **Geographical disaggregation** offers a new modelling avenue in terms of nowcasting/forecasting GDP

# Sparse Warcasting

Forecasting in a data-rich but statistics-poor environment

