

# Author Contributions Checklist Form

This form documents the artifacts associated with the article (i.e., the data and code supporting the computational findings) and describes how to reproduce the findings.

## Part 1: Data

☐ This paper **does not** involve analysis of external data (i.e., no data are used or the only data are generated by the authors via simulation in their code).

☒ I certify that the author(s) of the manuscript have legitimate access to and permission to use the data used in this manuscript.

## Abstract

FRED-MD is a large macroeconomic database designed for the empirical analysis of “big data.” The datasets of monthly and quarterly observations mimic the coverage of datasets already used in the literature, but they add three appealing features. They are updated in real-time through the FRED database. They are publicly accessible, facilitating the replication of empirical work.

## Availability

- ☒ Data **are** publicly available  
☐ Data **cannot be made** publicly available

If the data are publicly available, see the *Publicly available data* section. Otherwise, see the *Non-publicly available data* section, below.

## Publicly available data

- ☒ Data are available online at: <https://research.stlouisfed.org/econ/mccracken/fred-databases/>
- ☐ Data are available as part of the paper’s supplementary material.
- ☐ Data are publicly available by request, following the process described here:

- ☐ Data are or will be made available through some other mechanism, described here:

## Non-publicly available data

Discussion of lack of publicly available data:

## Description

### File format(s)

- ☒ CSV or other plain text: current.csv
- ☐ Software-specific binary format (.Rda, Python pickle, etc.):
- ☐ Standardized binary format (e.g., netCDF, HDF5, etc.):
- ☐ Other (described here):

### Data dictionary

- ☒ Provided by the authors in the following file(s): Code\_Reproduction/Real\_Data/Pre-Process/current.csv
- ☐ Data file(s) is (are) self-describing (e.g., netCDF files)
- ☐ Available at the following URL:

### Additional information (optional)

## Part 2: Code

### Abstract

We provide the code for reproduce our simulation and real data application. See workflow and readme.md for more details.

### Description

#### Code format(s)

- ☒ Script files
  - ☒ R   ☐ Python   ☐ Matlab
  - ☐ Other:
- ☐ Package
  - ☐ R   ☐ Python   ☐ MATLAB toolbox
  - ☐ Other:
- ☐ Reproducible report
  - ☐ R Markdown   ☐ Jupyter notebook
  - ☐ Other:
- ☐ Shell script
- ☐ Other (described here):

### Supporting software requirements

#### Version of primary software used

R 4.2.2, RStudio 2022.07.2+576 "Spotted Wakerobin" Release  
(e7373ef832b49b2a9b88162cfe7eac5f22c40b34, 2022-09-06) for macOS  
Mozilla/5.0 (Macintosh; Intel Mac OS X 12\_5\_0) AppleWebKit/537.36 (KHTML, like Gecko)  
QtWebEngine/5.12.10 Chrome/69.0.3497.128 Safari/537.36

#### Libraries and dependencies used by the code

Glmnetc (4.1-6), SIS (0.8-8), MASS (7.3-58.1), Elastic-Net (1.3), randomForest (4.7-1.1),

### Supporting system/hardware requirements (optional)

### Parallelization used

- ☒ No parallel code used
- ☐ Multi-core parallelization on a single machine/node  
Number of cores used:
- ☐ Multi-machine/multi-node parallelization  
Number of nodes and cores used:

### License

- ☒ MIT License (default)
- ☐ BSD
- ☐ GPL v3.0
- ☐ Creative Commons
- ☐ Other (described here):

### Additional information (optional)

## Part 3: Reproducibility workflow

### Scope

The provided workflow reproduces:

- ☐ Any numbers provided in text in the paper
- ☒ The computational method(s) presented in the paper (i.e., code is provided that implements the method(s))
- ☐ All tables and figures in the paper
- ☐ Selected tables and figures in the paper, as explained and justified here:

### Workflow details

#### Location

The workflow is available:

- ☐ As part of the paper's supplementary material
- ☒ In this Git repository: [https://github.com/qooyqpqy123/FARM\\_Reproduction](https://github.com/qooyqpqy123/FARM_Reproduction)
- ☐ Other:

#### Format(s)

- ☐ Single master code file
- ☐ Wrapper (shell) script(s)
- ☐ Self-contained R Markdown file, Jupyter notebook, or other literate programming approach
- ☐ Text file (e.g., a readme-style file) that documents workflow
- ☐ Makefile
- ☒ Other (more detail in 'Instructions' below)

#### Instructions

```
##Same Content with Readme.md in the github path ##
```

```
# Reproduction for the methodology of "Are Factor Regression and Sparse Regression Adequate?"
```

```
# 1. **Real Data Analysis** <br />
```

## Data <br />

### 1.1. Data Description <br />

FRED-MD is a large macroeconomic database designed for the empirical analysis of "big data." The datasets of monthly and quarterly observations mimic the coverage of datasets already used in the literature, but they add three appealing features. They are updated in real-time through the FRED database. They are publicly accessible, facilitating the replication of empirical work. <br />

### 1.2. Data Availability <br />

The original dataset can be found in <https://research.stlouisfed.org/econ/mccracken/fred-databases/>. <br />

In our .zip file, the original Dataset is provided in folder ``Code\_Reproduction/Real\_Data/Pre\_process/current.csv." <br />

The description of the data is provided in "Code\_Reproduction/Real\_Data/Pre\_process/Description\_of\_Variables.pdf." <br /> This description describes the meaning of all columns of this dataset.

## Codes Description and Implementation Details:

There are in total four folders under the folder: ``Code\_Reproduction/Real\_Data." <br /> 1. Pre\_process, <br /> 2. Prediction\_Table\_3 <br /> 3. PCR\_Adequate\_Table4, <br /> 4. Sparse\_Adequate\_Table4. <br />

### Folder `Code\_Reproduction/Real\_Data/Pre\_process`:

It contains the Codes for real data pre-processing. The document `current.csv` in that folder is our downloaded original dataset from the aforementioned website. <br />

The code `**realdata_read.R**` contains the codes for pre-processing the data. The line 4 in "realdata\_read.R", we have a variable called: "choose\_time". If one lets choose\_time=1 and run the codes, the output is the processed dataset `**realdata_115_126.csv**`. If one lets choose\_time=2, and run the codes, the output is the processed dataset `**realdata_189_126.csv**`. <br />

`**Run time**` of `**realdata_read.R**` is less than 3 min.

Finally, for document "Description\_of\_Variables.pdf" inside that document, it contains the detailed description of all variables in the original dataset "current.csv". <br />

### Folder `Code\_Reproduction/Real\_Data/Prediction\_Table3`:

This folder contains the codes for reproducing the prediction results of Table 3 in our paper. <br />

\*For code file `**realdata_prediction_115.R**`, it first reads the data "realdata\_115\_126.csv" in that folder. On the line 54 of this file, there is a variable "j". If we set j=49, and run the codes, we reproduce the line "HOUSTNE" line during "08.2010-02.2020" in Table 3. If we set j=81 and run the codes, we reproduce the line "GS5" line during "08.2010-02.2020" in Table 3. <br />

\*For code file **realdata\_prediction\_189.R**, it first reads the data "realdata\_189\_126.csv" in that folder. On the line 54 of this file, there is a variable "j". If we set j=49, and run the codes, we reproduce the line "HOUSTNE" line during "02.1992-10.2007" in Table 3. If we set j=81 and run the codes, we reproduce the line "GS5" line during "02.1992-10.2007" in Table 3. <br />

For both .R files, the output are printed in line 149-161 and 189 for different methods :Lasso (sparse linear regression), PCR (latent factor regression), Ridge (Ridge regression), El-Net (Elastic Net), RF (Random Forest), and FarmSelect (Factor adjusted Lasso), respectively.

**\*\*Run time\*\*** of **realdata\_prediction\_115.R** and **realdata\_prediction\_189.R** are less than 10 min.

### Folder `Code\_Reproduction/Real\_Data/PCR\_Adequate\_Table4`:

This folder contains the codes for reproducing the hypothesis testing results of column "LA\_factor" of Table 4 in our paper. <br />

\*For code file **Real\_data\_PCR\_115.R**, it first reads the data "realdata\_115\_126.csv" in that folder. On the line 17 of this file, there is a variable "j". If we set j=49, and run the codes, we reproduce the line "HOUSTNE", column "LA\_factor," during "08.2010-02.2020" in Table 4. If we set j=81 and run the codes, we reproduce the line "GS5" line, column "LA\_factor", during "08.2010-02.2020" in Table 4. <br />

\*For code file **Real\_data\_PCR\_189.R**, it first reads the data "realdata\_189\_126.csv" in that folder. On the line 17 of this file, there is a variable "j". If we set j=49, and run the codes, we reproduce the line "HOUSTNE", column "LA\_factor," during "02.1992-10.2007" in Table 4. If we set j=81 and run the codes, we reproduce the line "GS5" line, column "LA\_factor", during "02.1992-10.2007" in Table 4. <br />

For both .R files, the output are printed p-values of the test on whether PCR is adequate.

**\*\*Run time\*\*** of **Real\_data\_PCR\_115.R** and **Real\_data\_PCR\_189.R** are less than 10 min.

### Folder `Code\_Reproduction/Real\_Data/Sparse\_Adequate\_Table4`:

This folder contains the codes for reproducing the hypothesis testing results of column "SP\_Linear" of Table 4 in our paper.

\*For code file **Real\_data\_Sparse\_115.R**, it first reads the data "realdata\_115\_126.csv" in that folder. On the line 94 of this file, there is a variable "j". If we set j=49, and run the codes, we reproduce the line "HOUSTNE", column "SP\_Linear," during "08.2010-02.2020" in Table 4. If we set j=81 and run the codes, we reproduce the line "GS5" line, column "SP\_Linear", during "08.2010-02.2020" in Table 4. <br />

\*For code file **Real\_data\_Sparse\_189.R**, it first reads the data "realdata\_189\_126.csv" in that folder. On the line 94 of this file, there is a variable "j". If we set j=49, and run the codes, we

reproduce the line "HOUSTNE", column "SP\_Linear," during "02.1992-10.2007" in Table 4. If we set  $j=81$  and run the codes, we reproduce the line "GS5" line, column "SP\_Linear", during "02.1992-10.2007" in Table 4. <br />

For both .R files, the output are printed p-values of the test on whether Sparse Regression is adequate.

**\*\*Run time\*\*** of **\*\*Real\_data\_Sparse\_115.R\*\*** and **\*\*Real\_data\_Sparse\_189.R\*\*** are less than 10 min.

# 2. **\*\*Simulation Studies\*\*** <br />

## Codes Description and Implementation Details:

### Folder `Code\_Reproduction/simulation/Consistency\_Figure1/`:

There are three sub-folders inside this folder, namely, Gaussian, Uniform, Heavy-tail. <br />

For folder Gaussian and Uniform, they contain codes **\*\*Gauss\_consistency.R\*\*** and **\*\*Uniform\_consistency.R\*\***, respectively. By running these codes directly, we reproduce the methodology for generating results for Figure 1 (a) and (b). <br />

For folder Heavy-tail: it contains codes: **\*\*huber\_estimation.R\*\*** and **\*\*T-consistency.R\*\***. By running these codes directly, we reproduce the methodology for generating the red line and blue line of (c) of Figure 1.

**\*\*Run time\*\*** of **\*\*Gauss\_consistency.R\*\***, **\*\*Uniform\_consistency.R\*\***, **\*\*huber\_estimation.R\*\*** and **\*\*T-consistency.R\*\*** are around 5 hrs, respectively.

### Folder `Code\_Reproduction/simulation/Simulation\_Table1/`:

There is one file **\*\*pcr\_adequate.R\*\*** under this folder. It reproduces the methodology of Table 1. <br />

There are several parameters in this file need to be pre-specified. They correspond to different simulation settings of the paper. <br />

1. First, on line 7, if we let the parameter "choose\_dim=1", it outputs the results with dimension  $p=200$ . If we let the parameter "choose\_dim=2", it outputs the results with dimension  $p=500$ . <br />

2. Second, on line 10, if we let the parameter "choose=1", it outputs the results under the setting where factors F and idiosyncratic components are generated via setting 1 in section 5.2 of this paper. If we let the parameter "choose=2", it outputs the results under the setting where factors F and idiosyncratic components are generated via setting 2 in section 5.2 of this paper . <br />

3. Third, on line 16, if we let the parameter "choose\_noise=1", it outputs the results when the noise is generated based on Gaussian distribution. If we let the parameter "choose\_noise=2", it outputs the results under the setting where the noise follows uniform distribution. <br />



**\*\*Run time\*\*** of **\*\*pcr\_adequate.R\*\*** for any given setting described above is around 12 hrs.

**###** Folder ``Code_Reproduction/simulation/Simulation_Table2/``:

There is one file **\*\*inference.R\*\*** under this folder. It reproduces the methodology of Table 2.

<br />

There are several parameters in this file need to be pre-specified. They correspond to different simulation settings of the paper. <br />

1. First, on line 93, if we let the parameter "choose\_p=1", it outputs the results with dimension p=250. If we let the parameter "choose\_p=2", it outputs the results with dimension p=600. <br />

2. Second, on line 94, if we let the parameter "choose\_noise=1", it outputs the results where the noise is generated based on Gaussian distribution. If we let the parameter "choose\_noise=2", it outputs the results under the setting where the noise follows uniform distribution. <br />

3. Third, on line 95, if we let the parameter "choose\_mix=1", it outputs the results under the setting where the covariates are generated i.i.d. If we let the parameter "choose\_mix=2", it outputs the results under the setting where the covariates are strong mixing. <br />

**\*\*Run time\*\*** of **\*\*inference.R\*\*** for any given setting described above is around 5 hrs.

**###** Folder ``Code_Reproduction/simulation/Prediction_Section_B.1/``:

There are four files inside this folder. **\*\*Prediction\_Table1.R\*\***, **\*\*Prediction\_Table2.R\*\***, **\*\*Prediction\_Table3.R\*\***, and **\*\*Prediction\_Table4.R\*\*** <br />

They reproduces the Table 1-4 in Appendix, respectively.

**\*\*Run time\*\*** of **\*\*Prediction\_Table1.R\*\***, **\*\*Prediction\_Table2.R\*\***, **\*\*Prediction\_Table3.R\*\***, and **\*\*Prediction\_Table4.R\*\*** for any given setting described above is less than 30 min.

**###** Folder ``Code_Reproduction/simulation/Figure7_appendix/``:

There is one file **\*\*qqplot\_figure7.R\*\*** in the folder. It reproduces the methodology to generate Figure 7 in Appendix. The implementation is to run codes in 1-160 lines to generate "sure.csv". After that, one needs to read in this "sure.csv" to plot the qq-plot (corresponds to lines 164-171).

**\*\*Run time\*\*** of **\*\*qqplot\_figure7.R\*\*** for any given setting described above is around 5 hrs.

## Expected run-time

Approximate time needed to reproduce the analyses on a standard desktop machine:

- ☐ <1 minute
- ☐ 1-10 minutes
- ☐ 10-60 minutes
- ☐ 1-8 hours
- ☒ >8 hours
- ☐ Not feasible to run on a desktop machine, as described here:

## Additional documentation (optional)

## Notes (optional)