

Big Data Econometrics Nowcasting and Early Estimates

Big data handling tool – Source code description

George Kapetanios* Massimiliano Marcellino[†] Fotis Papailias[‡]
Katerina Petrova[§]

Abstract

Functions and usability of the R code are discussed in the paper. Most of the code is originally written by the authors. Whenever a R package is used, it is cited appropriately.

Keywords: Big Data, Nowcasting, Density Forecasting, Density Evaluation, R Software.

*kapetaniosgeorge@gmail.com

[†]massimiliano.marcellino@unibocconi.it

[‡]fotis.papailias@kcl.ac.uk; fotis.papailias@quantf.com

[§]katerina.petrova@st-andrews.ac.uk

Contents

| | | |
|----------|---|-----------|
| 1 | Methods for feature extraction of Big Data sources to usable time-series for econometric modelling | 2 |
| 1.1 | High frequency returns – HF.R | 2 |
| 1.2 | Mobile phone data – Mobile.R | 2 |
| 1.3 | Web prices – Prices.R | 4 |
| 1.4 | Google Trends – Google.R | 5 |
| 1.5 | Twitter Data – Twitter.R | 5 |
| 1.6 | Reuters data scraping – Reuters.py | 5 |
| 2 | Filtering techniques for high frequency data | 10 |
| 2.1 | Outliers detection – Outliers.R | 10 |
| 3 | Modelling techniques for Big Data | 11 |
| 3.1 | Penalised Regression Models – Ridge.R , ElasticNet.R | 11 |
| 3.2 | Spike and Slab regression – SpikeSlab.R | 13 |
| 3.3 | Regression Trees and Forests – Tree.R , Forest.R | 13 |
| 3.4 | Bayesian VAR models – BVAR.R | 15 |
| 4 | Modelling strategies for nowcasting/early estimates purposes | 17 |
| 4.1 | Weekly Google Trends – Weekly-Google.R | 17 |
| 4.2 | Transformation | 18 |
| 4.3 | Averaging – averaging.R | 18 |
| 4.4 | Naive Estimate – naive.R | 19 |
| 4.5 | ARIMA – ar.R | 19 |
| 4.6 | Dynamic Factor Analysis – dfa.R | 20 |
| 4.7 | Factor Linear Regression – Flinreg.R | 20 |
| 4.8 | Partial Least Squares – pls.R | 21 |
| 4.9 | Sparse Principal Components – spc.R | 21 |
| 4.10 | Sparse Regression – sparse.R | 22 |
| 4.11 | Spike and Slab – spike.R | 23 |
| 4.12 | Evaluation Statistics | 24 |

1 Methods for feature extraction of Big Data sources to usable time-series for econometric modelling

1.1 High frequency returns – [HF.R](#)

Here we use the highfrequency R package, load edit and plot the necessary data as follows. We use various comments in order to explain the code and the inputs used in each function.

Realised volatility based on 5-min cleaned returns.

```
1 library("highfrequency")
2
3 # Load data
4 data(sample_returns_5min)
5
6 # Store the 5 min returns data
7 r <- sample_returns_5min
8
9 # Calculate the realised volatility based on 5 min returns
10 # Input. r: 5mins return data
11 RV <- rowSums(r^2)
12
13 # Creat a plot
14 # Input. RV: realised volatility based on 5mins return data as calculated above
15 plot(RV, type="l", main="EUR/USD 5-min RV", xlab="Time", ylab="RV")
```

1.2 Mobile phone data – [Mobile.R](#)

Here we load some .csv files which are also provided. The purpose of this code is to load the mobile phone activity data and aggregate it in terms of calls, SMS and internet activity. We use various comments in order to explain the code and the inputs used in each function.

Below we only demonstrate the part of the code which corresponds to Internet activity. The same procedure is replicated for call and SMS activity.

Mobile phone data aggregation and plots.

```
1 # Load the necessary data.
2 # Each file corresponds to a single day.
3 days <- c("sms-call-internet-mi-2013-11-01.csv",
4           "sms-call-internet-mi-2013-11-02.csv",
5           "sms-call-internet-mi-2013-11-03.csv",
6           "sms-call-internet-mi-2013-11-04.csv",
7           "sms-call-internet-mi-2013-11-05.csv",
8           "sms-call-internet-mi-2013-11-06.csv",
```

```

9         "sms-call-internet-mi-2013-11-07.csv")
10
11 # Create some labels to be used in the plots later
12 days.l <- c("2013-11-01, Friday", "2013-11-02, Saturday", "2013-11-03, Sunday",
13            "2013-11-04, Monday", "2013-11-04, Tuesday", "2013-11-05, Wednesday",
14            "2013-11-07, Thursday")
15 days.l2 <- c("2013-11-01", "2013-11-02", "2013-11-03",
16            "2013-11-04", "2013-11-04", "2013-11-05",
17            "2013-11-07")
18
19 # Create a sparse matrix to store the results
20 result1 <- array(NA, dim=c(24,5,NROW(days)))
21
22 # We start with a loop.
23 # j runs for each different day. In the above we have 7 days (NROW(days)).
24 for(j in 1:NROW(days)) {
25     # Load the data for day j.
26     x <- read.csv(days[j], header=TRUE)
27
28     # Identify unique users
29     ud <- unique(x[,1])
30
31     # Extract Total activity during the day across users
32     # Create a sparse matrix to store the results
33     totact <- matrix(NA, NROW(ud), 5)
34     rownames(totact) <- as.character(ud)
35     colnames(totact) <- c("smsin", "smsout", "callin", "callout", "internet")
36
37     for(i in 1:NROW(ud)) {
38         it <- ud[i]
39         choose <- which(x[,1]==it)
40         # Input. x which has the data and here we are looping across users.
41         totact[i,1] <- sum(x[choose, 4], na.rm=TRUE)
42         totact[i,2] <- sum(x[choose, 5], na.rm=TRUE)
43         totact[i,3] <- sum(x[choose, 6], na.rm=TRUE)
44         totact[i,4] <- sum(x[choose, 7], na.rm=TRUE)
45         totact[i,5] <- sum(x[choose, 8], na.rm=TRUE)
46     }
47
48     result1[,j] <- totact
49     cat("day ", j, " just done - still left ", NROW(days), "\n")
50 }
51
52 # We are now ready to create some plots
53 # First, generate the colours.
54 cols <- rainbow(NROW(days), alpha = 1)
55
56 # Calls Internet. Using the 5th column of array across days, "i", we can aggregate
57 # (sum) the internet activity.
58 i <- 1
59 plot(result1[,5,i], type="l", xlab="Hours", ylab="Call", col=cols[i], main="Total
    Internet Activity")
60
61 # Add the remaining days looping across "i"
62 for(i in 2:NROW(days)) {
63     lines(result1[,5,i], col=cols[i])
64 }
65
66 legend("topleft", legend=days.l, col=cols, lty=rep(1, NROW(days)), cex=0.6)

```

1.3 Web prices – Prices.R

Here we load a .csv sample file with the appropriate data. We use various comments in order to explain the code and the inputs used in each function.

Web prices aggregation and plots.

```

1  # Load the data
2  x <- read.csv("arg-sample2.csv", header=TRUE)
3  x <- as.matrix(x)
4  x <- x[,c(5, 1, 7, 6)] # Extract the necessary columns
5  colnames(x) <- c("date", "id", "cat", "price")
6
7  d <- as.Date(x[,1])      # create the dates sequence
8  ud <- sort(unique(d))    # extract unique dates for time series aggregation
9  uid <- unique(x[,2])     # extract unique id's
10 uc <- unique(x[,3])      # extract unique categories
11
12 # Create sparse matrix to store the results
13 cats <- matrix(NA, NROW(ud)-1, NROW(uc))
14 colnames(cats) <- uc
15 rownames(cats) <- as.character(ud[2:NROW(ud)])
16
17 # Start the loop across categories
18 for(i in 1:NROW(uc)) {
19   xcat <- which(x[,3]==uc[i])
20   xcat <- x[xcat,]
21
22   tmat <- matrix(NA, NROW(ud), NROW(uid))
23   colnames(tmat) <- c(uid)
24
25   # Create a doouble loop: (j) across id's and (jj) across dates
26   for(j in 1:NROW(uid)) {
27     xcat2 <- which(xcat[,2]==uid[j])
28     xcat2 <- xcat[xcat2,]
29
30     for(jj in 1:NROW(ud)) {
31       cw <- which(xcat2[,1]==ud[jj])
32
33       if(NROW(cw)==0){
34         tmat[jj,j] <- NA
35       } else {
36         tmat[jj,j] <- xcat2[cw,4]
37       }
38     }
39   }
40   # ensure that the result is numeric
41   tmat <- apply(tmat, 2, as.numeric)
42
43   # Apply the methodology of Cavallo
44   R <- tmat[2:NROW(tmat),]/tmat[1:(NROW(tmat)-1),]
45   R <- apply(R, 1, prod, na.rm=TRUE)
46   R <- R^(1/NCOL(tmat))
47   cats[,i] <- R
48 }
49
50 # Cumulate across time
51 I <- apply(cats, 2, cumprod)

```

```
52 w <- as.matrix(rep(1/NCOL(I), NCOL(I)))
53
54 # use the appropriate weights as in Cavallo
55 S <- I %*% w
56
57 # Produce the final CPI estimate plot.
58 plot(as.Date(rownames(S)), S, type="l", main="Online Prices CPI", xlab="Time", ylab="Index")
```

1.4 Google Trends – [Google.R](#)

We use the `gtrendsR` library to download *Google Trends* in R. Currently, *Google Trends* are offered in monthly frequency.

Google Trends Download.

```
1 # User Input
2 kwd <- "gbp"           # Keyword(s)
3 reg <- "GB"            # Region
4 tsd <- "all"           # Time Frame: "all" (since 2004),
5                        # "today+XXX-y" (last XXX years)
6 stp <- "web"           # "web", "news", "images", "froogle", "youtube"
7 cct <- 0               # category
8
9 # Call the function and download data
10 gt <- gtrends(kwd, reg, tsd, stp, cct)
11
12 # Make a plot
13 plot(gt)
```

1.5 Twitter Data – [Twitter.R](#)

We use the `twitterR` library to download Twitter data in R.

Twitter data fetching.

```
1 # Inputs:
2 # user defined keyword, below we use #gbpusd feed
3 # n: the number of nodes to download
4 rdmTweets <- searchTwitter('#gbpusd', n=6500)
```

1.6 Reuters data scraping – [Reuters.py](#)

We provide below with a stepwise description of the Python code used for scraping Reuters data:

- L. 1-5: import packages: BeautifulSoup (HTML parsing), Scrapy (web crawling), Logging (recording errors), DateTime (parsing dates);
- L. 7: open a log file to record 404 errors (page not found);
- L. 10: create a name for the web scraper within the Scrapy CrawlSpider class;
- L. 13-18: spider settings. Note that we disabled the AUTOTHROTTLER setting and defined parameters manually to limit speed (one page download every 0.2 seconds) and contemporaneous requests (4).
- L. 20-25: define the input URLs for the web scraper. In a typical Scrapy crawler, one would only define a single starting page for the spider to move from, but we have already obtained the full list of pages to scrape in a text file. Therefore, we copy all the URLs in the file to a list and we pass each item to the parse function.
- L. 28-30: writes to a log file 404 errors;
- L. 32: passes the HTML code of the page to the parsing library lxml;
- L. 33-34: finds and isolates the permanent URL of the article;
- L. 35-36: finds and isolates the alternate URL of the article;
- L. 37-38: finds and isolates the article tags;
- L. 39-41: finds and isolates the publication date, headline and text of the article;
- L. 42-50: removes the remaining HTML formatting for clean reading;
- L. 51-52: parses the date and publication time to obtain a short-form date;
- L. 54-56: writes the collected data to a CSV file.

Reuters data scraping.

```
1 from bs4 import BeautifulSoup
2 import logging
3 import scrapy
4 from scrapy.spiders import CrawlSpider
5 from datetime import datetime
6
7 logging.basicConfig(filename='log_50.log', level=logging.ERROR)
```

```

8
9 class SpiderReuters (CrawlSpider):
10     name = 'crawler_final_v21_tags_2'
11     handle_httpstatus_list = [404]
12
13     custom_settings = {
14         'AUTOTHROTTLE_ENABLED': False,
15         'CONCURRENT_REQUESTS_PER_DOMAIN': '4',
16         'DOWNLOAD_DELAY': '0.2',
17         'USER_AGENT': "Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:38.0) Gecko/
20100101 Firefox/38.0"
18     }
19
20     def start_requests(self):
21         f = open(r"/home/PATH-to-file/list_of_URLs.txt", 'r')
22         start_urls = [url.strip() for url in f.readlines()]
23         f.close()
24         for urls in start_urls:
25             yield scrapy.Request(url=urls, callback=self.parse)
26
27     def parse(self, response):
28         if response.status == 404:
29             with open('log_404.txt', 'a') as l:
30                 l.write(str(response) + '\n')
31         else:
32             soup = BeautifulSoup(response.body, 'lxml')
33             extractor_a = soup.find('link', rel='canonical')
34             a_return = extractor_a['href']
35             extractor_alt = soup.find('link', rel='alternate')
36             alt_return = extractor_alt['href']
37             extractor_tags = soup.find('meta', property="og:article:tag")
38             tags_return = extractor_tags['content']
39             extractor_1 = soup.find_all('span', class_='timestamp')
40             extractor_2 = soup.find_all('h1', class_='article-headline')
41             extractor = soup.find_all('span', id='article-text')
42             composer = [extractor_1, a_return, alt_return, extractor_2, extractor,
43             tags_return]
44             composer_2 = []
45             for element in composer:
46                 element_2 = ' '.join(str(element).strip('[]').splitlines())
47                 composer_2.append(element_2)
48             composer_text = "|".join([str(item).replace('"', "") for item in
49             composer_2])
50             clean_page = BeautifulSoup(composer_text, 'lxml')
51             all_text = ' '.join(clean_page.findAll(text=True))
52             all_text_uni = all_text.decode('unicode_escape').encode('utf-8', 'strict
53             ')
54             date = all_text_uni[4:17].replace(',', ' ').strip(' ')
55             date_object = datetime.strptime(date, '%b %d %Y')
56             page = str(response).replace("/", "").replace(".", "").replace(":", "").
57             replace("<", "").replace(">", "")
58             with open(r'/home/PATH/output_%s.csv' % page, 'a') as a:
59                 a.write(date_object.strftime('%b %d %Y') + '|')
60                 a.write(' '.join(all_text_uni.splitlines()).replace(" ", " "))

```

We also provide below with a stepwise description of the code used for the index construction:

- L. 1-5: imports packages: time (algorithm timer), os (interfacing with operating system), re (regular expressions), csv (comma separated values file reading and writing), pandas (time series statistical package);
- L. 7: starts timer;
- L. 9: initializes article counter;
- L. 11: dictionary of search terms (regular expressions);
- L. 13-15: creates a CSV file for articles with column headers 'date', 'url', 'headline' and 'uncertainty';
- L. 18-19: iterates through directories and obtains a list of files in each folder;
- L. 21-23: opens and counts each article;
- L. 24-27: excludes articles containing the word "SPORT" among the tags;
- L. 28-34: if the article body contains one of the terms in line 11, the date, url, tags and a "1" binary indicator are written to a csv file;
- L. 35-40: otherwise, date, url, tags and a "0" binary indicator are written to the same csv file;
- L. 41-44: handles csvError exceptions (i.e. ignores files in folder that are not csv format);
- 48-55: the tagged article list just obtained is loaded in a Pandas dataframe, dates are converted to a machine-readable format and daily frequencies are obtained and written to a csv file; the actual index is computed at line 62.

Index construction.

```
1 import os, re, time
2 import csv
3 import pandas as pd
4
5 start_time = time.time()
6
7 article_count = 0
8
9 combined_semantic = "\\risk\\b\\b\\brisks\\b\\b\\brisks\\b"
```

```

10
11 with open('art_list_RISK.csv', 'w', encoding='utf-8') as el:
12     spamwriter = csv.DictWriter(el, delimiter=',', fieldnames=['date', 'url', 'alt_
13     url', 'tags', 'filename', 'uncertainty'], lineterminator='\n')
14     spamwriter.writeheader()
15
16 # search iteration
17 for root, dirs, filename in os.walk(r'/home/mattia/Desktop/All_articles'):
18     for sub_file in filename:
19         try:
20             with open(os.path.join(root, sub_file), 'r', encoding="utf8", errors='
21             ignore') as f:
22                 spamreader = csv.DictReader(f, delimiter='|', fieldnames=['date', '
23                 longdate', 'url', 'alt_url', 'header', 'article', 'tags'])
24                 article_count += 1
25                 for row in spamreader:
26                     if 'SPORT' in [x.strip() for x in row['tags'].split(',')]:
27                         print('sport', row['url'])
28                         break
29                     match = re.search(combined_semantic, row['article'], re.IGNORECASE)
30                     if match:
31                         print('progress', '%.3f' % ((article_count / 2803677) * 100), '%'
32                         ,)
33                         line_with_dummy = dict(url=row['url'], alt_url=row['alt_url'],
34                         date=row['date'], tags=row['tags'], filename=sub_file, uncertainty='1')
35                         with open('art_list_RISK.csv', 'a', encoding='UTF-8') as out:
36                             spamwriter_2 = csv.DictWriter(out, delimiter=',', fieldnames
37                             =['date', 'url', 'alt_url', 'tags', 'filename', 'uncertainty'], lineterminator='
38                             \n')
39                             spamwriter_2.writerow(line_with_dummy)
40                     else:
41                         row.update({'uncertainty': '0'})
42                         line_with_zero = dict(url=row['url'], alt_url=row['alt_url'],
43                         date=row['date'], tags=row['tags'], filename=sub_file, uncertainty='0')
44                         with open('art_list_RISK.csv', 'a', encoding='UTF-8') as out:
45                             spamwriter_3 = csv.DictWriter(out, delimiter=',', fieldnames
46                             =['date', 'url', 'alt_url', 'tags', 'filename', 'uncertainty'], lineterminator='
47                             \n')
48                             spamwriter_3.writerow(line_with_zero)
49                     except csv.Error:
50                         with open('csverror.txt', 'a', encoding='utf-8') as csverr:
51                             csverr.write(root + sub_file)
52                             pass
53
54 print('end list article list')
55
56 df = pd.read_csv('art_list_RISK.csv')
57 print('dataset in memory')
58 df["date"] = pd.to_datetime(df["date"])
59 frequency_table = pd.crosstab(index=df["date"], columns=df['uncertainty'], margins=
60     True)
61 df2 = pd.DataFrame(frequency_table)
62 df2['ratio'] = df2[1]/df2['All']
63 df2.columns = ['no_uncertainty', 'uncertainty', 'all', 'ratio']
64 df2.to_csv('Marcellino_daily_RISK_index.csv', encoding='utf-8')
65
66 elapsed_time = time.time() - start_time
67
68 print(elapsed_time/60, 'minutes')

```

2 Filtering techniques for high frequency data

2.1 Outliers detection – [Outliers.R](#)

In this task we were mostly concerned with outliers detection, seasonalities and data cleaning. Below, we present the functions we use in the Outliers.R script.

Using scores from the outliers package.

```

1 # Inputs
2 # x is a vector of data, unstructured or aggregated depending on the theme.
3 # type: "z" calculates normal scores (differences between each value and the mean
  divided
4 # by sd)
5 # prob: the corresponding p-values are returned level choice depends on the user
6 od <- scores(x, type="z", prob=0.99)
7
8 # identify the position of outliers
9 which(od==TRUE)

```

Now, we use the built-in stl to identify and extract the trend and seasonal component. This will lead to the cleaned series.

Seasonal decomposition of time series by Loess.

```

1 # Input: aggx is a numeric vector of aggregated time series in weekly frequency
2 # First, we transform the numeric vector in a time series (ts) object correctly
3 # specifying the frequency.
4 tsaggx <- ts(aggx, frequency=7)
5
6 # Then, we use the ts object, tsaggx, as the input in the LOESS function.
7 # s.window: can be a string "periodic" or "per" which reads the frequency from the
8 # ts transformation, otherwise it can be a user choice.
9 ss <- stl(tsaggx, s.window="per")
10
11 # Plot the output
12 plot(ss, main="Daily Aggregation, Weekly Pattern")
13
14 # Extract the seasonal component (xs) and the trend component (xt)
15 xs <- ss$time.series[,1]
16 xt <- ss$time.series[,2]
17
18 # Calculate the cleaned series (xc)
19 xc <- tsaggx-xs-xt
20
21 # And create a plot
22 plot(xc, type="l", xlab="", ylab="", main="Detrended and Deseasonalised")

```

3 Modelling techniques for Big Data

In this section, we present the codes for Ridge, Lasso, Elastic Net penalised regressions, regression trees and random forests. For the penalised regression models, we use several functions from the `glmnet` package.

3.1 Penalised Regression Models – [Ridge.R](#), [ElasticNet.R](#)

For the estimation of Ridge and Elastic Net regression, we use the `glmnet` package. The main function is `glmnet(x, y, family=".", alpha=.)` and it implements penalised regression of an $N \times p$ matrix of explanatory variables x on a N -dimensional vector y .

The package also performs k-fold cross validation, with the function `cv.glmnet`. Finally, to generate predictions the following command can be used: `predict(fit.info, newdata=x.test)`, where `fit.info` is the output from the `glmnet` (e.g. coefficients/ confidence intervals) etc.

Ridge penalised regression.

```
1 # Example Code Ridge
2 rm(list=ls())
3 install.packages("glmnet") # download and install package
4
5 library (glmnet)
6
7 #generate some artificial data
8 set.seed(1)
9 n <- 200 # Number of observations
10 p <- 300 # Number of predictors included in model
11
12 beta<- c(1/(1:p)^2)
13 # approximately sparse model, slope coefficients small but not zero
14 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
15 y <- x%*%beta + rnorm(n)
16
17 #generate the dependant variable y
18 fit.ridge <- glmnet(x, y, family="gaussian", alpha=0)
19
20 #select ridge penalty
21 nforecast=5
22 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
23
24 predict(fit.ridge, newdata=xnew) # predictions based on estimated coefficients
```

Lasso regression.

```
1 # Example Code Lasso
2 rm(list=ls())
3 install.packages("glmnet") # download and install package
```

```
4
5 library (glmnet)
6
7 #generate some artificial data
8 set.seed(1)
9 n <- 200 # Number of observations
10 p <- 300 # Number of predictors included in model
11
12 beta<- matrix(c(rep(1,p/2),rep(0,p/2)))
13 # sparse model, some slope coefficients are zero
14 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
15
16 y <- x%%beta + rnorm(n)
17 # generate the dependant variable y
18
19 fit.lasso <- glmnet(x, y, family="gaussian", alpha=1) #select Lasso penalty
20
21 nforecast=5
22 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
23
24 predict(fit.lasso, newdata=xnew) # predictions based on estimated coefficients
```

Elastic Net regression.

```
1 # Example Code Elastic Net
2 rm(list=ls())
3 install.packages("glmnet") # download and install package
4
5 library (glmnet)
6
7 #generate some artificial data
8 set.seed(1)
9 n <- 200 # Number of observations
10 p <- 300 # Number of predictors included in model
11
12 beta1<- c(1/(1:p/2)^2)
13 beta<- matrix(c(beta1,rep(0,p/2)))
14 #combination of sparse and approximately sparse coefficient vector
15 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
16
17 # generate the dependant variable y
18 y <- x%%beta + rnorm(n)
19
20 fit.elnet <- glmnet(x, y, family="gaussian", alpha=0.4) # 40% weight to Lasso
    penalty
21
22 nforecast=5
23 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
24
25 predict(fit.elnet, newdata=xnew) # generate predictions based on estimated
    coefficients
```

3.2 Spike and Slab regression – [SpikeSlab.R](#)

For the Spike and Slab regression model, we use the BoomSpikeSlab package. The main function is `lm.spike(y x, iter)` where x is an $N \times p$ matrix of explanatory variables, y is an N -dimensional vector and `iter` is the number of Metropolis draws from the posterior distribution of the parameters.

Slab and Spike regression.

```

1  # Example Code Slab and Spike
2  rm(list=ls())
3  install.packages("BoomSpikeSlab") # download and install package
4
5  library (BoomSpikeSlab)
6
7  #generate some artificial data
8  set.seed(1)
9  n = 200 #sample size
10 p = 300 # number of variables
11 nonzero = 3 # nubmer of variables with non-zero coefficients
12 niter <- 1000 # nubmer of MCMC draws
13 sigma <- .8
14
15 x <- cbind(1, matrix(rnorm(n * (p-1)), nrow=n))
16 beta <- c(rep(2,ngood),rep(0, p-nonzero ))
17 y <- rnorm(n, x %*% beta, sigma)
18 x <- x[,-1]
19
20 # estimate spike and Slab regression
21 model <- lm.spike(y ~ x, niter=niter)
22
23 # plots of coefficients
24 plot.ts(model$beta)
25 hist(model$sigma) ## should be near 8
26 plot(model)
27 summary(model)
28
29 # plot residuals
30 plot(model, "residuals")
31
32 Xnew = cbind( matrix(rnorm(n * (p-1)), nrow=n))
33
34 # if out-of-sample forecasts are required
35 yhat.slabs.new = predict.lm.spike(model, newdata=Xnew) #out-of-sample prediction

```

3.3 Regression Trees and Forests – [Tree.R](#), [Forest.R](#)

For the regression trees and forests, we make use of four R packages ISLS, randomForest, rpart and rpart.plot. To implement a standard regression tree with default settings: `fit.trees<- rpart(y x)` where x is an $N \times p$ matrix of explanatory variables and y is N -

dimensional vector.

To prune a tree, the following instructions can be used:

Tree pruning.

```
1 bestcp      <- trees$cptable[which.min(trees$cptable[, "xerror"]), "CP"]
2 fit.prunedtree <- prune(fit.trees, cp=bestcp)
3 prp(fit.prunedtree)
```

The main function to optimally estimate a forest is `RFfit<- tuneRF(x, y)`. Finally, the packages can be used to generate predictions. For regression trees, this can be achieved with the command:

```
1 yhat.pt<- predict(fit.prunedtree, newdata=as.data.frame(...))
```

while, for forests, this can be achieved with the command:

```
1 yhat.rf2<- predict(RFfit, newdata=...)
```

Below, we illustrate with an example.

Standard Regression Tree

```
1 rm(list=ls())
2
3 # Regression Tree
4
5 library (ISLR)
6 library(randomForest)
7 library(rpart)
8 library(rpart.plot)
9
10 #generate some artificial data
11 set.seed(1)
12 n <- 200 # Number of observations
13 p <- 300 # Number of predictors included in model
14
15 beta<- c(10/(1:p)^2)
16 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
17 y <- x%*%beta + rnorm(n)*4
18
19 # estimate Standard Regression tree on the atrificial data
20 fit.trees<- rpart(y~x)
21
22 prp(fit.trees)
23
24 # estimate pruned Regression tree on the atrificial data
25 bestcp      <- trees$cptable[which.min(trees$cptable[, "xerror"]), "CP"]
26 fit.prunedtree <- prune(fit.trees, cp=bestcp)
27
28 prp(fit.prunedtree)
```

Random Forest

```

1  rm(list=ls())
2
3  # Random Forest
4
5  library (ISLR)
6  library(randomForest)
7  library(rpart)
8  library(rpart.plot)
9
10 #generate some artificial data
11 set.seed(1)
12 n <- 200 # Number of observations
13 p <- 300 # Number of predictors included in model
14
15 beta<- c(10/(1:p)^2)
16 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
17 y <- x%*%beta + rnorm(n)*4
18
19 #Estimate a Random forest on the atrificial data
20 RFfit<- tuneRF(x, y, mtryStart=floor(sqrt(ncol(x))),stepFactor=1.5, improve=0.05,
21               nodesize=5, ntree=2000, doBest=TRUE)
22
23 #Find the best fir for the Random forest on the atrificial data
24 min <- RFfit$mtry
25 fit.rf2 <-randomForest(x, y, nodesize=5, mtry=min, ntree=2000)

```

3.4 Bayesian VAR models – BVAR.R

We make use the of the MSBVAR package for Bayesian vectorautoregressive models. The main function is `szbvar`, which takes the following inputs:

- Y is $T \times M$ matrix of time series;
- p is lag length;
- z is $T \times N$ matrix of exogenous vars, can be NULL;
- λ_0 is the overall tightness between 0 and 1;
- λ_1 is the standard deviation or tightness of the prior around the AR(1) parameters;
- λ_3 is the lag decay (> 0 , with 1=harmonic);
- λ_4 is the standard deviation or tightness around the intercept > 0 ;

- `lambda5` is the standard deviation or tightness around the exogenous variable coefficients;
- `mu5` is the sum of coefficients prior weight (larger values imply difference stationarity);
- `mu6` is dummy initial observations or drift prior (larger values allow for common trends);
- `nu` is the prior degrees of freedom, $m + 1$;
- `qm` is the frequency of the data for lag decay equivalence;
- prior can be of three values: 0 = Normal-Wishart prior, 1 = Normal-flat prior, 2 = flat-flat prior;
- `posterior.fit` is a logical, FALSE implies no estimation of log-posterior fit measures.

The package can be used to generate out-of-sample forecasts, using the function `forecast`, for example:

```
1 forecasts <- forecast(fit.bvar, nsteps)
```

with `nsteps` the numbers of horizons for the out-of-sample forecasts.

Bayesian VAR

```
1 rm(list=ls())
2
3 # Bayesian VAR
4 install.packages("MSBVAR") # download and install package
5
6 library (MSBVAR)
7
8 #generate some artificial data
9 set.seed(1)
10 n = 200 #sample size
11 p = 10 # number of variables
12 X0 = rep(0,p)
13 beta = rep(0.5,p)
14 B=diag(beta)
15 y=matrix(0,nrow=p,ncol=n)
16 for (i in 1:n) {
17   e = rnorm(p)
18   y[,i]=B%%X0+e
19   X0=y[,i]
20 }
21
22 # Reference prior model -- Normal-IW prior pdf
```

```
23 Bvar.Model <- szbvar(y, p=6, z=NULL, lambda0=0.6, lambda1=0.1, lambda3=2, lambda4
    =0.5, lambda5=0, mu5=0, mu6=0, nu=ncol(KEDS)+1, qm=4, prior=0, posterior.fit=F)
24
25 # Forecast -- this gives back the sample PLUS the forecasts.
26 forecasts <- forecast(Bvar.Model, nsteps=10, burnin=3000, gibbs=5000, exog=NULL)
27
28 # Conditional forecasts
29 conditional.forcs.ref <- hc.forecast(Bvar.Model, yhat, nsteps,
30 burnin=3000, gibbs=5000, exog=NULL)
```

4 Modelling strategies for nowcasting/early estimates purposes

The codes presented in this section address the following issues:

- (i) data manipulation,
- (ii) nowcasting, and
- (iii) post-processing of results.

The first part refers to all codes we used to manipulate the downloaded data. This part is “data-specific” and it applies to data downloaded from the same sources as in these tasks. If a researcher uses Bloomberg, Reuters, Macrobond or other data collection software, the original data is different and cannot be used as input in these functions. Therefore, we do not discuss them here. In the following section, we take as granted that the user has already downloaded and cleaned the data.

4.1 Weekly Google Trends – [Weekly-Google.R](#)

As mentioned in a previous section, *Google Trends* are downloaded in monthly frequency. We have developed a function which loads Google Trends over smaller time frames at weekly frequency, and then scales the data in order to obtain the weekly *Google Trends* for the overall period. The function uses the main function from *gtrendsR* package.

Weekly Google Trends.

```
1 # Set dates for all data
2 dfrom <- "2004-01-01" # starting date
3 dto <- "2017-09-01" # ending date
```

```
4
5 # (the news doesn't really have a lot of data, so stick to the web)
6 stp <- "web"           # web; news;
7 cct <- 0               # category
8
9 # General Indexes - web
10 reg <- ""             # Region, blank for all regions or use "GB" for the UK, etc.
11 kwd <- "uncertainty"
12
13 # Download the weekly trends and store them in c1 variable
14 # Input: kwd (keyword), reg (region)
15 #         cct (category), stp (domain)
16 #         dfrom (start), dto (end)
17 c1 <- weekly.GOOGLE(kwd, reg, cct, stp, dfrom, dto)
```

4.2 Transformation

We transform the final nowcast/forecast estimates to levels according to the nature of the series. In order to avoid repetition, the code is described here.

From change to levels.

```
1 # YTRANSF: if 3, we translate from growth to levels
2 #           if 2, we translate from 1-st diff to levels
3 #   zlast: is the last observed value for the dependent variable
4 #   zboot: are the bootstrap estimate for densities
5 #   zave: (or different name) is the final estimate in levels.
6 if(YTRANSF==3){
7   zlast <- YSAV[NROW(YSAV)]
8   zboot <- zlast*(1+zboot)
9   zave <- zlast*(1+zave)
10 }
11 if(YTRANSF==2){
12   zlast <- YSAV[NROW(YSAV)]
13   zboot <- zlast +zboot
14   zave <- zlast + zave
15 }
```

4.3 Averaging – [averaging.R](#)

We calculate the average value of the last observations and also use Bootstrap to calculate the corresponding density estimates.

Averaging and Bootstrap for density.

```
1 # Input:
2 # z is the dependent variable, vector of data
3 # zp: window length, e.g. zp=4, then we have the 4-period MA.
4 zave <- mean(z[(NROW(z)-zp+1):NROW(z),])
5
```

```
6 # Bootstrap for density estimation
7 # b: bootstrap window length
8 # B: number of bootstraps
9 b <- round(NROW(z)^(1/3))
10 boots <- matrix(NA, NROW(z), B)
11 for(j in 1:B){
12   boots[,j] <- MBB(z, b) # MBB: Moving Block Bootstrap
13 }
14
15 # Calculate the mean value for the same zp
16 zboot <- colMeans(boots[(NROW(boots)-zp+1):NROW(boots),])
```

4.4 Naive Estimate – [naive.R](#)

We calculate the naive forecast and also use bootstrap to calculate the corresponding density estimates.

Naive forecasting.

```
1 # Input:
2 # z is the dependent variable, vector of data
3 # here we extract the last observed value
4 zf <- z[NROW(z)]
5
6 # Bootstrap for density estimation
7 # b: bootstrap window length
8 # B: number of bootstraps
9 b <- round(NROW(z)^(1/3))
10 boots <- matrix(NA, NROW(z), B)
11 for(j in 1:B){
12   boots[,j] <- MBB(z, b)
13 }
14
15 # Extract the corresponding naives
16 zboot <- boots[NROW(boots),]
```

4.5 ARIMA – [ar.R](#)

We calculate various ARIMA model-based forecasts using the forecast package.

ARIMA forecasting.

```
1 # If an AR order is set to zero, we use AIC
2 if(arp==0){ arp <- NROW(ar.ols(z, aic=TRUE)$ar) }
3
4 # Estimate ARIMA using
5 # z: the vector of the observed dependent variable
6 # order: ARIMA order
7 # method: conditional sum of squares
8 fit <- Arima(z,order=c(arp,0,0), method=c("CSS"))
9 fout <- NULL
```

```
10
11 # Calculate percentiles
12 try(fout <- forecast(fit,h=1, bootstrap=TRUE, npaths=B, level=seq(51,99,1)), silent=
    TRUE)
13
14 # Put all percentiles together in the correct order
15 zout <- c(fout$mean, rev(as.numeric(fout$lower))[1], rev(fout$lower), fout$mean,
    fout$upper, as.numeric(fout$upper)[49])
```

4.6 Dynamic Factor Analysis – [dfa.R](#)

We calculate forecasts/nowcasts based on Dynamic Factor Analysis. Codes are adapted versions of original Giannone and Reichlin's Matlab programs.

Dynamic Factor Analysis forecasting.

```
1 # Extract factors using:
2 # XX: matrix of observed regressors
3 # q: dynamic rank
4 # r: static rank (r>=q)
5 # p: ar order of the state vector
6 fac.out <- FactorExtraction(XX,q=fq,r=fr,p=fp)
7 Fac <- fac.out$Fac
8 rownames(Fac) <- rownames(XX)
9
10 # Then use the linreg with the factors
11 z <- YY; f <- Fac; vlag <- 0; source("../linreg.R")
```

4.7 Factor Linear Regression – [Flinreg.R](#)

We calculate forecasts/nowcasts based on linear regression using factors which come from standardised matrices.

Factor Linear regression.

```
1 # jj: step-ahead
2 # z: target, vector
3 # f: factor which comes from standardised data
4 # ysd: the standard deviation of target
5 # ymu: the mean of target
6 jj <- 1
7 zreg <- as.matrix(z[(jj+1):NROW(z)],)
8 freg <- as.matrix(f[1:(NROW(f)-jj)],)
9 out <- lm(zreg~freg)
10 b <- out$coefficients;
11
12 # Now make sure to use ysd, ymu as we used factors from standardised input
13 outf <- ((f[NROW(f),]%*%b[2:NROW(b)]) + b[1])*ysd+ymu
```

4.8 Partial Least Squares – [pls.R](#)

We calculate forecasts/nowcasts based on partial least squares methodology. We are using the `pls` package.

Factor Linear regression.

```

1  # Lag the matrix of Regressors if necessary
2  # vlag: lag order
3  # XX : matrix, panel of regressors
4  # YY : vector, observed target
5  vlag <- 1
6  if(vlag>0){
7    XX <- cbind(lagf(YY, vlag)[,2:(vlag+1)], XX)
8    XX <- as.matrix(XX[(vlag+1):NROW(XX),])
9    YY <- as.matrix(YY[(vlag+1):NROW(YY),])
10 }
11
12 # Standardise
13 xxin <- xstd(XX)
14 ymu <- mean(YY)
15 ysd <- sd(YY)
16 yyin <- (YY-ymu)/ysd
17
18 # Extract factors
19 pp <- plsr(yyin~xxin, ncomp=qncomp, scale=FALSE)
20 f <- as.matrix(pp$scores)
21 z <- as.matrix(yyin)
22
23 # Use Factor linear regression
24 source("Flinreg.R")

```

4.9 Sparse Principal Components – [spc.R](#)

We calculate forecasts/nowcasts based on sparse principal components methodology. We are using the `nsprcomp` package.

Sparse Principal Components.

```

1  # jj: step-ahead
2  # z: target, vector
3  # f: factor which comes from standardised data
4  # ysd: the standard deviation of target
5  # ymu: the mean of target
6  vlag <- 1
7  if(vlag>0){
8    # XX <- lagmv(XX, vlag) # no because of small T dimension
9    XX <- cbind(lagf(YY, vlag)[,2:(vlag+1)], XX)
10
11    XX <- as.matrix(XX[(vlag+1):NROW(XX),])
12    YY <- as.matrix(YY[(vlag+1):NROW(YY),])
13 }
14

```

```

15 # Standardise
16 xxin <- xstd(XX)
17 ymu <- mean(YY)
18 ysd <- sd(YY)
19 yyin <- (YY-ymu)/ysd
20
21 # Extract factors
22 pc.out <- nsprcomp(x=xxin, retx=TRUE, ncomp=qncomp, nneg = FALSE, center=FALSE,
23                   scale.=FALSE)
24 f <- as.matrix(pc.out$x)
25 z <- as.matrix(yyin)
26
27 # Use Factor linear regression
28 source("../Flinreg.R")

```

4.10 Sparse Regression – [sparse.R](#)

We calculate forecasts/nowcasts based on sparse regression methodology. We are using the `glmnet` package.

Sparse regression.

```

1 # jj: step-ahead
2 # z: target, vector
3 # f: factor which comes from standardised data
4 # ysd: the standard deviation of target
5 # ymu: the mean of target
6 vlag <- 1
7 if(vlag>0){
8   XX <- cbind(lagf(YY, vlag)[,2:(vlag+1)], XX)
9   XX <- as.matrix(XX[(vlag+1):NROW(XX),])
10  YY <- as.matrix(YY[(vlag+1):NROW(YY),])
11 }
12 z <- YY
13 f <- XX
14
15 # jj: step ahead, then correctly lead/lag the variables
16 jj <- 1
17 zreg <- as.matrix(z[(jj+1):NROW(z),])
18 freg <- as.matrix(f[1:(NROW(f)-jj),])
19
20 # Calculate beta which comes from sparse
21 # freg: regressors, matrix
22 # zreg: target, vector
23 # type.measure: "mse" for the calculation of lambda
24 # alpha: 1 for Lasso, 0.5 for LAR
25 fit.lasso.cv <- cv.glmnet(freg, zreg, type.measure="mse", alpha=salpha, family="
26   gaussian", standardize=TRUE)
27 s <- fit.lasso.cv$lambda.min
28 b <- as.numeric(coef(fit.lasso.cv, s))
29 outf <- ((f[NROW(f),]%*%b[2:NROW(b)]) + b[1])
30
31 # Calculate percentiles
32 sigmah <- sd(zreg-freg*%*%b[2:NROW(b)]-b[1])
33 spseq <- qnorm(seq(0.51, 0.99, 0.01))*sigmah

```

```
33 zout <- c(outf, outf-rev(spseq)[1], outf-rev(spseq), outf, outf+spseq, outf+spseq[
    NROW(spseq)])
```

4.11 Spike and Slab – [spike.R](#)

We calculate forecasts/nowcasts based on sparse regression methodology. We are using the BoomSpikeSlab package.

Spike and Slab Regression.

```
1  # jj: step-ahead
2  # z: target, vector
3  # f: factor which comes from standardised data
4  # ysd: the standard deviation of target
5  # ymu: the mean of target
6  lag <- 1
7  if(vlag>0){
8      # XX <- lagmv(XX, vlag) # no because of small T dimension
9      XX <- cbind(lagf(YY, vlag)[,2:(vlag+1)], XX)
10     XX <- as.matrix(XX[(vlag+1):NROW(XX),])
11     YY <- as.matrix(YY[(vlag+1):NROW(YY),])
12 }
13 # Standardise
14 xxin <- xstd(XX)
15 ymu <- mean(YY)
16 ysd <- sd(YY)
17 yyin <- (YY-ymu)/ysd
18
19 z <- yyin
20 f <- xxin
21
22 jj <- 1
23 zreg <- as.matrix(z[(jj+1):NROW(z),])
24 freg <- as.matrix(f[1:(NROW(f)-jj),])
25
26 # Spike and Slab
27 # niter: The number of MCMC iterations to run
28 # ping: output printing parameter
29 out <- lm.spike(zreg ~ freg, niter=niters, ping=B)
30
31 # keep the last B rounds
32 b <- out$beta
33 b <- b[(NROW(b)-B+1):NROW(b),]
34 bf <- colMeans(b)
35 outf <- ((f[NROW(f),]%*%bf[2:NROW(bf)]) + bf[1])*ysd+ymu
36
37 #Calculate percentiles
38 sigmah <- sd((zreg-freg%*%bf[2:NROW(bf)]-bf[1])*ysd+ymu)
39 spseq <- qnorm(seq(0.51, 0.99, 0.01))*sigmah
40 zout <- c(outf, outf-rev(spseq)[1], outf-rev(spseq), outf, outf+spseq, outf+spseq[
    NROW(spseq)])
```

4.12 Evaluation Statistics

We calculate various evaluation statistics. Check the main files for more details.

Mean Absolute Error.

```
1 # err: matrix with errors of various methods
2 mae <- as.matrix(colMeans(abs(err)))
3
4 # relative MAE
5 benchmark <- "AR(1)"
6 maeR <- mae/mae[benchmark,1]
```

Root Mean Squared Forecast Error.

```
1 # err: matrix with errors of various methods
2 rmsfe <- as.matrix(sqrt(colMeans(err^2)))
3
4 # relative RMSFE
5 benchmark <- "AR(1)"
6 rmsfeR <- rmsfe/rmsfe[benchmark,1]
```

Two Sided Diebold-Mariano.

```
1 # Using the package "forecast"
2 # err: matrix with errors of various methods
3 # i: method i stored in column i of err
4 # h: horizon
5 # power: power
6 benchmark <- "AR(1)"
7 dmp <- dm.test(err[,i], err[,benchmark], alternative=c("two.sided"), h=1, power=2)
8
9 # Extract the p-value
10 dmp <- as.numeric(dmp$p.value)
```

Sign Success Ratio.

```
1 # err: matrix with errors of various methods
2 # forc: signs of forecast direction compared to the last period for a given method
3 # sgnt: signs of direction of the target
4 ssr <- sum(sgnt==sign(forc))
5
6 # relative SSR
7 benchmark <- "AR(1)"
8 ssrR <- ssr/ssr[benchmark,1]
```

Berkowitz LR Test.

```
1 # xcloud: all percentile forecasts
2 xcloud <- allf
3 z <- pnorm(xcloud[,1], mean=apply(xcloud[,2:NCOL(xcloud)], 1, mean), sd=apply(xcloud
  [,2:NCOL(xcloud)], 1, sd))
4 Z <- qnorm(z)
5
```

```
6 # Extract Berkowitz P-value
7 berk <- BerkowitzTest(Z, lags=1, significance = 0.05)$LRp
```
