

# Big Data Econometrics

## Nowcasting and Early Estimates

### Task 10: “Big data handling tool developed as R package including supporting user documentation for deployment”

George Kapetanios\*    Massimiliano Marcellino<sup>†</sup>    Fotis Papailias<sup>‡</sup>  
Katerina Petrova<sup>§</sup>

December 8, 2017

#### Abstract

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

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

---

\*kapetaniosgeorge@gmail.com

<sup>†</sup>massimiliano.marcellino@unibocconi.it

<sup>‡</sup>fotis.papailias@kcl.ac.uk; fotis.papailias@quantf.com

<sup>§</sup>katerina.petrova@st-andrews.ac.uk

# Contents

<b>1</b>	<b>Task 2</b>	<b>2</b>
1.1	High Frequency Returns. File name: <a href="#">HF.R</a>	2
1.2	Mobile Phone Data. File name: <a href="#">MobPhone.R</a>	2
1.3	Web Prices. File name: <a href="#">Prices.R</a>	4
1.4	Google Trends. File name: <a href="#">Google.R</a>	5
1.5	Twitter Data. File name: <a href="#">Twitter.R</a>	5
1.6	Reuters written in Python	5
1.6.1	Index construction	8
<b>2</b>	<b>Task 3. File name: <a href="#">Outliers.R</a></b>	<b>10</b>
<b>3</b>	<b>Task 4</b>	<b>11</b>
3.1	Penalised Regression Models. File name: <a href="#">Ridge.R</a> , <a href="#">ElasticNet.R</a>	11
3.2	Spike and Slab Regression. File name: <a href="#">SpikeSlab.R</a>	13
3.3	Regression Trees and Forests. File name: <a href="#">Tree.R</a> , <a href="#">Forest.R</a>	14
3.4	Bayesian VAR models. File name: <a href="#">BVAR.R</a>	17
<b>4</b>	<b>Tasks 5 – 8</b>	<b>18</b>
4.1	Weekly Google Trends. File name: <a href="#">Weekly-Google.R</a>	19
4.2	Transformation	19
4.3	Averaging. File name: <a href="#">averaging.R</a>	20
4.4	Naive Estimate. File name: <a href="#">naive.R</a>	20
4.5	ARIMA. File name: <a href="#">ar.R</a>	21
4.6	Dynamic Factor Analysis. File name: <a href="#">dfa.R</a>	21
4.7	Factor Linear Regression. File name: <a href="#">Flinreg.R</a>	22
4.8	Partial Least Squares. File name: <a href="#">pls.R</a>	22
4.9	Sparse Principal Components. File name: <a href="#">spc.R</a>	23
4.10	Sparse Regression. File name: <a href="#">sparse.R</a>	23
4.11	Spike and Slab. File name: <a href="#">sparse.R</a>	24
4.12	Evaluation Statistics	25

## 1 Task 2

### 1.1 High Frequency Returns. File name: [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. File name: [MobPhone.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 {
26   # Load the data for day j.
27   x <- read.csv(days[j], header=TRUE)
28
29   # Identify unique users
30   ud <- unique(x[,1])
31
32   # Extract Total activity during the day across users
33   # Create a sparse matrix to store the results
34   totact <- matrix(NA, NROW(ud), 5)
35   rownames(totact) <- as.character(ud)
36   colnames(totact) <- c("smsin", "smsout", "callin", "callout", "internet")
37
38   for(i in 1:NROW(ud))
39   {
40     it <- ud[i]
41     choose <- which(x[,1]==it)
42     # Input. x which has the data and here we are looping across users.
43     totact[i,1] <- sum(x[choose, 4], na.rm=TRUE)
44     totact[i,2] <- sum(x[choose, 5], na.rm=TRUE)
45     totact[i,3] <- sum(x[choose, 6], na.rm=TRUE)
46     totact[i,4] <- sum(x[choose, 7], na.rm=TRUE)
47     totact[i,5] <- sum(x[choose, 8], na.rm=TRUE)
48   }
49   result1[,,j] <- totact
50   cat("day ", j, " just done - still left ", NROW(days), "\n")
51 }
52
53 # We are now ready to create some plots
54 # First, generate the colours.
55 cols <- rainbow(NROW(days), alpha = 1)
56
57 # Calls Internet
58 # Using the 5th column of array across days, "i", we can aggregate (sum) the
   internet activity.
59 i <- 1
60 plot(result1[,5,i], type="l", xlab="Hours", ylab="Call", col=cols[i],
61       main="Total Internet Activity")
62 # Add the remaining days looping across "i"
63 for(i in 2:NROW(days))
64 {
65   lines(result1[,5,i], col=cols[i])
66 }
67 legend("topleft", legend=days.l, col=cols, lty=rep(1, NROW(days)),
68       cex=0.6)

```

### 1.3 Web Prices. File name: [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 {
20   xcat <- which(x[,3]==uc[i])
21   xcat <- x[xcat,]
22
23   tmat <- matrix(NA, NROW(ud), NROW(uid))
24   colnames(tmat) <- c(uid)
25
26   # Create a doouble loop: (j) across id's and (jj) across dates
27   for(j in 1:NROW(uid))
28   {
29     xcat2 <- which(xcat[,2]==uid[j])
30     xcat2 <- xcat[xcat2,]
31
32     for(jj in 1:NROW(ud))
33     {
34       cw <- which(xcat2[,1]==ud[jj])
35
36       if(NROW(cw)==0){
37         tmat[jj,j] <- NA
38       }else{
39         tmat[jj,j] <- xcat2[cw,4]
40       }
41     }
42   }
43   # sure that the result is numeric
44   tmat <- apply(tmat, 2, as.numeric)
45
46   # Apply the methodology of Cavallo
47   R <- tmat[2:NROW(tmat),]/tmat[1:(NROW(tmat)-1),]
48   R <- apply(R, 1, prod, na.rm=TRUE)
49   R <- R^(1/NCOL(tmat))
50
51   cats[,i] <- R

```

```

52 }
53
54 # Cumulate across time
55 I <- apply(cats, 2, cumprod)
56 w <- as.matrix(rep(1/NCOL(I), NCOL(I)))
57
58 # use the appropriate weights as in Cavallo
59 S <- I %*% w
60
61 # Produce the final CPI estimate plot.
62 plot(as.Date(rownames(S)), S, type="l", main="Online Prices CPI", xlab="Time", ylab=
      "Index")

```

---

## 1.4 Google Trends. File name: [Google.R](#)

We use the “*gtrendsR*” R 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. File name: [Twitter.R](#)

We use the “*twitteR*” R library to download Twitter data in R.

### Twitter Data

---

```

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

```

---

## 1.6 Reuters written in Python

Lines:

- 1-5: import packages: BeautifulSoup (HTML parsing), Scrapy (web crawling), Logging (recording errors), DateTime (parsing dates);
- 7: open a log file to record 404 errors (page not found);
- 11: create a name for the web scraper within the Scrapy CrawlSpider class;
- 14-19: spider settings. Note that we disabled the “AUTOTHROTTLE” setting and defined parameters manually to limit speed (one page download every 0.2 seconds) and contemporaneous requests (4).
- 21-26: 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 below.
- 29-31: writes to a log file 404 errors;
- 33: passes the HTML code of the page to the parsing library lxml;
- 34-35: finds and isolates the permanent URL of the article;
- 36-37: finds and isolates the alternate URL of the article;
- 38-39: finds and isolates the article tags;
- 40-42: finds and isolates the publication date, headline and text of the article;
- 43-51: removes the remaining HTML formatting for clean reading;
- 52-53: parses the date and publication time to obtain a short-form date;
- 55-57: writes the collected data to a CSV file.

---

```
1
2 from bs4 import BeautifulSoup
3 import scrapy
4 import logging
5 from scrapy.spiders import CrawlSpider
6 from datetime import datetime
7
```

```

8 logging.basicConfig(filename='log_50.log', level=logging.ERROR)
9
10
11 class SpiderReuters (CrawlSpider):
12     name = 'crawler_final_v21_tags_2'
13     handle_httpstatus_list = [404]
14
15     custom_settings = {
16         'AUTOTHROTTLER_ENABLED': False,
17         'CONCURRENT_REQUESTS_PER_DOMAIN': '4',
18         'DOWNLOAD_DELAY': '0.2',
19         'USER_AGENT': "Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:38.0) Gecko/
20         20100101 Firefox/38.0"
21     }
22
23     def start_requests(self):
24         f = open(r"/home/PATH-to-file/list_of_URLs.txt", 'r')
25         start_urls = [url.strip() for url in f.readlines()]
26         f.close()
27         for urls in start_urls:
28             yield scrapy.Request(url=urls, callback=self.parse)
29
30     def parse(self, response):
31         if response.status == 404:
32             with open('log_404.txt', 'a') as l:
33                 l.write(str(response) + '\n')
34         else:
35             soup = BeautifulSoup(response.body, 'lxml')
36             extractor_a = soup.find('link', rel='canonical')
37             a_return = extractor_a['href']
38             extractor_alt = soup.find('link', rel='alternate')
39             alt_return = extractor_alt['href']
40             extractor_tags = soup.find('meta', property="og:article:tag")
41             tags_return = extractor_tags['content']
42             extractor_1 = soup.find_all('span', class_='timestamp')
43             extractor_2 = soup.find_all('h1', class_='article-headline')
44             extractor = soup.find_all('span', id='article-text')
45             composer = [extractor_1, a_return, alt_return, extractor_2, extractor,
46             tags_return]
47             composer_2 = []
48             for element in composer:
49                 element_2 = ' '.join(str(element).strip('[]').splitlines())
50                 composer_2.append(element_2)
51             composer_text = "|".join([str(item).replace('"', "") for item in
52             composer_2])
53             clean_page = BeautifulSoup(composer_text, 'lxml')
54             all_text = ' '.join(clean_page.findAll(text=True))
55             all_text_uni = all_text.decode('unicode_escape').encode('utf-8', 'strict
56             ')
57             date = all_text_uni[4:17].replace(', ', ' ').strip(' ')
58             date_object = datetime.strptime(date, '%b %d %Y')
59             page = str(response).replace("/", "").replace(".", "").replace(":", "").
60             replace("<", "").replace(">", "")
61             with open(r'/home/PATH/output_%s.csv' % page, 'a') as a:
62                 a.write(date_object.strftime('%b %d %Y') + '|')
63                 a.write(' '.join(all_text_uni.splitlines()).replace(" ", ""))

```



### 1.6.1 Index construction

- Lines 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);
- 7: starts timer;
- 9: initializes article counter;
- 11: dictionary of search terms (regular expressions);
- 13-17: creates a CSV file for articles with column headers 'date', 'url', 'headline' and 'uncertainty';
- 20-21: iterates through directories and obtains a list of files in each folder;
- 23-26: opens and counts each article;
- 27-30: excludes articles containing the word "SPORT" among the tags;
- 31-40: 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;
- 41-49: otherwise, date, url, tags and a "0" binary indicator are written to the same csv file;
- 50-53: handles csvError exceptions (i.e. ignores files in folder that are not csv format);
- 57-64: 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.

---

```
1 import time
2 import os
3 import csv
4 import re
5 import pandas as pd
6
7 start_time = time.time()
8
```

```

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

```

```

55 print('end list article list')
56
57 df = pd.read_csv('art_list_RISK.csv')
58 print('dataset in memory')
59 df["date"] = pd.to_datetime(df["date"])
60 frequency_table = pd.crosstab(index=df["date"], columns=df['uncertainty'], margins=
    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 Task 3. File name: **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*” file.

### Using *scores()* from the “*outliers*” package.

```

1 # Input: x is a vector of data, unstructured or
2 #         aggregated depending on the theme.
3 #         type: "z" calculates normal scores
4 #             (differences between each value and
5 #             the mean divided by sd)
6 #         prob: the corresponding p-values are returned
7 #             level choice depends on the user
8 od <- scores(x, type="z", prob=0.99)
9
10 # identify the position of outliers
11 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
2 #         aggregated time series in weekly frequency
3 #         First, we transform the numeric vector in
4 #         a time series (ts) object correctly
5 #         specifying the frequency.
6 tsaggx <- ts(aggx, frequency=7)
7
8 # Then, we use the ts object, tsaggx, as
9 #         the input in the LOESS function.
10 # s.window: can be a string "periodic"
11 #            or "per" which reads the frequency
12 #            from the ts transformation, otherwise

```

```
13 #           it can be a user choice.
14 ss <- stl(tsaggx, s.window="per")
15
16 # Plot the output
17 plot(ss, main="Daily Aggregation, Weekly Pattern")
18
19 # Extract the seasonal component (xs)
20 #   and the trend component (xt)
21 xs <- ss$time.series[,1]
22 xt <- ss$time.series[,2]
23
24 # Calculate the cleaned series (xc)
25 xc <- tsaggx - xs - xt
26
27 # And create a plot
28 plot(xc, type="l", xlab="", ylab="", main="Detrended and Deseasonalised")
```

---

### 3 Task 4

In Section 4, 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. File name: **Ridge.R, ElasticNet.R**

For the estimation of Ridge and Elastic Net regression, we use the “*glmnet*” package. The lasso model is discussed later in Section 4.10. The main function is `glmnet(x, y, family=".", alpha=.)` and it implements penalised regression of an  $N$  by  $p$  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
8 #generate some artificial data
9
10 set.seed(1)
11
12 n <- 200 # Number of observations
13
14 p <- 300 # Number of predictors included in model
15
16 beta<- c(1/(1:p)^2)
17 #approximately sparse model, slope coefficients small but not zero
18 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
19
20 y <- x%*%beta + rnorm(n)
21 #generate the dependant variable y
22
23 fit.ridge <- glmnet(x, y, family="gaussian", alpha=0)
24 #select ridge penalty
25
26 nforecast=5
27 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
28
29 predict(fit.ridge, newdata=xnew) # generate predictions based on estimated
    coefficients

```

---

## Lasso Regression.

---

```

1
2 # Example Code Lasso
3 rm(list=ls())
4 install.packages("glmnet") # download and install package
5
6 library (glmnet)
7
8
9 #generate some artificial data
10
11 set.seed(1)
12
13 n <- 200 # Number of observations
14
15 p <- 300 # Number of predictors included in model
16
17 beta<- matrix(c(rep(1,p/2),rep(0,p/2)))
18 #sparse model, some slope coefficients are zero
19 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
20
21 y <- x%*%beta + rnorm(n)
22 # generate the dependant variable y
23
24 fit.lasso <- glmnet(x, y, family="gaussian", alpha=1) #select Lasso penalty
25
26
27 nforecast=5
28 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
29
30 predict(fit.lasso, newdata=xnew) # generate predictions based on estimated
    coefficients

```

## Elastic Net Regression.

```

1
2 # Example Code Elastic Net
3 rm(list=ls())
4 install.packages("glmnet") # download and install package
5
6 library (glmnet)
7
8
9 #generate some artificial data
10
11 set.seed(1)
12
13 n <- 200 # Number of observations
14
15 p <- 300 # Number of predictors included in model
16
17 beta1<- c(1/(1:p/2)^2)
18 beta<- matrix(c(beta1,rep(0,p/2)))
19 #combination between sparse and approximately sparse coefficient vector
20 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
21
22 y <- x%*%beta + rnorm(n)
23 # generate the dependant variable y
24
25 fit.elnet <- glmnet(x, y, family="gaussian", alpha=0.4) #gives 40% weight to Lasso
    penalty
26
27 nforecast=5
28 xnew <- matrix(rnorm(nforecast*p), nrow= nforecast, ncol=p)
29
30 predict(fit.elnet, newdata=xnew) # generate predictions based on estimated
    coefficients

```

## 3.2 Spike and Slab Regression. File name: [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 by 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
2 # Example Code Slab and Spike
3 rm(list=ls())
4 install.packages("BoomSpikeSlab") # download and install package
5
6 library (BoomSpikeSlab)
7
8 #generate some artificial data
9
10 set.seed(1)

```

```

11
12 n = 200 #sample size
13
14 p = 300 # number of variables
15
16 nonzerob = 3 # nubmer of variables with non-zero coefficients
17
18 niter <- 1000 # nubmer of MCMC draws
19
20 sigma <- .8
21
22 x <- cbind(1, matrix(rnorm(n * (p-1)), nrow=n))
23
24 beta <- c(rep(2,ngood),rep(0, p-nonzerob ))
25
26
27 y <- rnorm(n, x %*% beta, sigma)
28
29 x <- x[,-1]
30
31
32 #Estimate spike and Slab regression
33 model <- lm.spike(y ~ x, niter=niter)
34
35
36 #Plots of coefficients
37 plot.ts(model$beta)
38
39 hist(model$sigma) ## should be near 8
40
41 plot(model)
42
43 summary(model)
44
45
46 #Plot residuals
47 plot(model, "residuals")
48
49
50 Xnew = cbind( matrix(rnorm(n * (p-1)), nrow=n))
51
52
53 #if out-of-sample forecasts are required
54 yhat.slab.new = predict.lm.spike(model, newdata=Xnew) #out-of-sample prediction

```

### 3.3 Regression Trees and Forests. File name: **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  $N$  by  $p$  matrix of explanatory variables and  $y$  is  $N$ -dimensional vector.

To prune a tree, the following command 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:

### Main Command

---

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

---

For forests, this can be achieved with the command:

### Main 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 # download and install packages
6 install.packages("ISLR") # download and install package
7 install.packages("randomForest") # download and install package
8 install.packages("rpart") # download and install package
9 install.packages("rpart.plot") # download and install package
10
11 library (ISLR)
12
13 library(randomForest)
14
15 library(rpart)
16
17 library(rpart.plot)
18
19
20 #generate some artificial data
21
22 set.seed(1)
23
24
25 n <- 200 # Number of observations
26
27 p <- 300 # Number of predictors included in model
```



```

28
29 beta<- c(10/(1:p)^2)
30
31 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
32
33 y <- x%%beta + rnorm(n)*4
34
35
36 #Estimate Standard Regression tree
37 on the atrificial data
38 fit.trees<- rpart(y~x)
39
40 prp(fit.trees)
41
42
43 #Estimate pruned Regression tree on the atrificial data
44 bestcp <- trees$cptable[which.min(trees$cptable[, "xerror"]), "CP"]
45
46 fit.prunedtree <- prune(fit.trees, cp=bestcp)
47
48 prp(fit.prunedtree)

```

## Random Forest

```

1 rm(list=ls())
2
3 # Random Forest
4
5 # download and install packages
6 install.packages("ISLR") # download and install package
7 install.packages("randomForest") # download and install package
8 install.packages("rpart") # download and install package
9 install.packages("rpart.plot") # download and install package
10
11 library (ISLR)
12
13 library(randomForest)
14
15 library(rpart)
16
17 library(rpart.plot)
18
19
20 #generate some artificial data
21
22 set.seed(1)
23
24
25 n <- 200 # Number of observations
26
27 p <- 300 # Number of predictors included in model
28
29 beta<- c(10/(1:p)^2)
30
31 x <- matrix(rnorm(n*p), nrow=n, ncol=p)
32
33 y <- x%%beta + rnorm(n)*4
34

```

```

35
36 #Estimate a Random forest on the atrificial data
37
38 RFfit<- tuneRF(x, y, mtryStart=floor(sqrt(ncol(x))),stepFactor=1.5, improve=0.05,
    nodesize=5, ntree=2000, doBest=TRUE)
39
40
41 #Find the best fir for the Random forest on the atrificial data
42
43 min          <- RFfit$mtry
44
45 fit.rf2 <-randomForest(x, y, nodesize=5, mtry=min, ntree=2000)

```

### 3.4 Bayesian VAR models. File name: [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  $T$  by  $M$  matrix of time series;  $p$  is lag length;  $z$   $T$  by  $N$  matrix of exogenous vars, can be `NULL`;  $\lambda_0$  between 0 and 1, overall tightness;  $\lambda_1$  is the standard deviation or tightness of the prior around the  $AR(1)$  parameters;  $\lambda_3$  Lag decay ( $> 0$ , with  $1=\text{harmonic}$ );  $\lambda_4$  Standard deviation or tightness around the intercept  $> 0$ ;  $\lambda_5$  is the standard deviation or tightness around the exogeneous variable coefficients;  $\mu_5$  is the sum of coefficients prior weight (larger values imply difference stationarity);  $\mu_6$  is dummy initial observations or drift prior (larger values allow for common trends);  $\nu$  is the prior degrees of freedom,  $m + 1$ ;  $q_m$  is the frequency of the data for lag decay equivalence; the input 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: `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
9  #generate some artificial data

```

```

10
11  set.seed(1)
12
13  n = 200 #sample size
14
15  p = 10 # number of variables
16
17  X0 = rep(0,p)
18
19  beta = rep(0.5,p)
20
21  B=diag(beta)
22
23  y=matrix(0,nrow=p,ncol=n)
24
25  for (i in 1:n) {
26
27    e = rnorm(p)
28
29    y[,i]=B%*%X0+e
30
31    X0=y[,i]
32  }
33
34
35
36  # Reference prior model -- Normal-IW prior pdf
37
38  Bvar.Model <- szbvar(y, p=6, z=NULL, lambda0=0.6,
39
40    lambda1=0.1, lambda3=2, lambda4=0.5,
41
42    lambda5=0, mu5=0, mu6=0,
43
44    nu=ncol(KEDS)+1, qm=4, prior=0,
45
46    posterior.fit=F)
47
48
49  # Forecast -- this gives back the sample PLUS the forecasts.
50  forecasts <- forecast(Bvar.Model, nsteps=10, burnin=3000, gibbs=5000, exog=NULL)
51
52
53  # Conditional forecasts
54
55  conditional.forcs.ref <- hc.forecast(Bvar.Model, yhat, nsteps,
56
57    burnin=3000, gibbs=5000, exog=NULL)

```

## 4 Tasks 5 – 8

In Tasks 5, 6 and 7 we divide the general codes in three parts: (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. The following sections, we take it as granted that the user has already downloaded and cleaned the data.

## 4.1 Weekly Google Trends. File name: [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

In “averaging.R”, ..., we transform the final nowcast/forecast estimates to levels according to the nature of the series. In order to avoid repetition, we describe the code 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 #
4 #   zlast: is the last observed value for the dependent variable
5 #   zboot: are the bootstrap estimate for densities
6 #   zave: (or different name) is the final estimate in levels.
7 if(YTRANSF==3){
8   zlast <- YSAV[NROW(YSAV)]
9   zboot <- zlast*(1+zboot)
10  zave <- zlast*(1+zave)
11 }
12 if(YTRANSF==2){
13   zlast <- YSAV[NROW(YSAV)]
14   zboot <- zlast +zboot
15   zave <- zlast + zave
16 }

```

### 4.3 Averaging. File name: [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: "z" is the dependent variable, vector of data
2 #       zp: window length, e.g. zp=4, then we have the 4-period MA.
3 zave <- mean(z[(NROW(z)-zp+1):NROW(z),])
4
5 # Bootstrap for density estimation
6 # b: bootstrap window length
7 # B: number of bootstraps
8 b <- round(NROW(z)^(1/3))
9 boots <- matrix(NA, NROW(z), B)
10 for(j in 1:B){
11   boots[,j] <- MBB(z, b) # MBB: Moving Block Bootstrap
12 }
13 # Calculate the mean value for the same zp
14 zboot <- colMeans(boots[(NROW(boots)-zp+1):NROW(boots),])

```

### 4.4 Naive Estimate. File name: [naive.R](#)

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

#### Naive Forecasting.

```

1 # Input: "z" is the dependent variable, vector of data
2 #       here we extract the last observed value
3 zf <- z[NROW(z)]
4
5 # Bootstrap for density estimation

```

```

6 # b: bootstrap window length
7 # B: number of bootstraps
8 b <- round(NROW(z)^(1/3))
9 boots <- matrix(NA, NROW(z), B)
10 for(j in 1:B){
11   boots[,j] <- MBB(z, b)
12 }
13 # Extract the corresponding naives
14 zboot <- boots[NROW(boots),]

```

---

## 4.5 ARIMA. File name: [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),
16   fout$mean, fout$upper, as.numeric(fout$upper)[49])

```

---

## 4.6 Dynamic Factor Analysis. File name: [dfa.R](#)

We calculate forecasts/nowcasts based on Dynamic Factor Analysis. Codes are adapted versions of Giannone and Reichlin original 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. File name: [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. File name: [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
9   XX <- as.matrix(XX[(vlag+1):NROW(XX),])
10  YY <- as.matrix(YY[(vlag+1):NROW(YY),])
11 }
12
13 # Standardise
14 xxin <- xstd(XX)
15 ymu <- mean(YY)
16 ysd <- sd(YY)
17 yyin <- (YY-ymu)/ysd
18
19 # Extract factors
20 pp <- pls(yyin~xxin, ncomp=qncomp, scale=FALSE)
21 f <- as.matrix(pp$scores)
22 z <- as.matrix(yyin)
23
24 # Use Factor linear regression
25 source("../Flinreg.R")

```

## 4.9 Sparse Principal Components. File name: [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. File name: [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
10     XX <- as.matrix(XX[(vlag+1):NROW(XX),])

```



```

11 YY <- as.matrix(YY[(vlag+1):NROW(YY),])
12 }
13 z <- YY
14 f <- XX
15
16 # jj: step ahead, then correctly lead/lag the variables
17 jj <- 1
18 zreg <- as.matrix(z[(jj+1):NROW(z),])
19 freg <- as.matrix(f[1:(NROW(f)-jj),])
20
21 # Calculate beta which comes from sparse
22 # freg: regressors, matrix
23 # zreg: target, vector
24 # type.measure: "mse" for the calculation of lambda
25 # alpha: 1 for Lasso, 0.5 for LAR
26 fit.lasso.cv <- cv.glmnet(freg, zreg, type.measure="mse", alpha=salpha, family="
    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
34 zout <- c(outf, outf-rev(spseq)[1], outf-rev(spseq), outf+spseq, outf+spseq[
    NROW(spseq)])

```

## 4.11 Spike and Slab. File name: [sparse.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
11   XX <- as.matrix(XX[(vlag+1):NROW(XX),])
12   YY <- as.matrix(YY[(vlag+1):NROW(YY),])
13 }
14 # Standardise
15 xxin <- xstd(XX)
16 ymu <- mean(YY)
17 ysd <- sd(YY)
18 yyin <- (YY-ymu)/ysd
19
20 z <- yyin
21 f <- xxin

```

```

22
23 jj <- 1
24 zreg <- as.matrix(z[(jj+1):NROW(z),])
25 freg <- as.matrix(f[1:(NROW(f)-jj),])
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 # keep the last B rounds
31 b <- out$beta
32 b <- b[(NROW(b)-B+1):NROW(b),]
33 bf <- colMeans(b)
34 outf <- ((f[NROW(f),]%*%bf[2:NROW(bf)]) + bf[1])*ysd+ymu
35 #Calculate percentiles
36 sigmah <- sd((zreg-freg)%*%bf[2:NROW(bf)]-bf[1])*ysd+ymu
37 spseq <- qnorm(seq(0.51, 0.99, 0.01))*sigmah
38 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
```

---