# Analysis of pair trading with financial market data

Marcell P. Granát

2021 09 09

# Contents

**Abstract**

If one attends to the extremely large literature of demographic trends in the developed world, then the uncertainty about the effect of economic and human development factors on the fertility rate cannot be covered for a long time. Several empirical studies argue for the existence of the J-shaped effect of the development, but many papers come up with statements to the opposite. The goal of this paper is to contribute to the literature with an advanced panel econometric model based on regional observations. Beyond the human development factors (living standard, education and health) I extend my analysis by using youth unemployment and family benefit indicators as dependent variables. Important to note that statistics about unemployment are available only for a critically short period in the case of many regions. To manage this highly unbalanced nature of the dataset – while not rejecting the possibility to control for youth unemployment – I estimate the model with two different modeling frames: one without youth unemployment and another one with it. As a result, the paper confirms the empirical evidence that increasing human development in developed countries has a positive effect on total fertility rates, and income is the most important component. This finding is robust to the mentioned two frameworks. In contrast, the research come up only with week evidence for the significant effect of expenditure on family on total fertility rates on the long run.

*Keywords* — fertility rates, human development

# List of Tables

# List of Figures
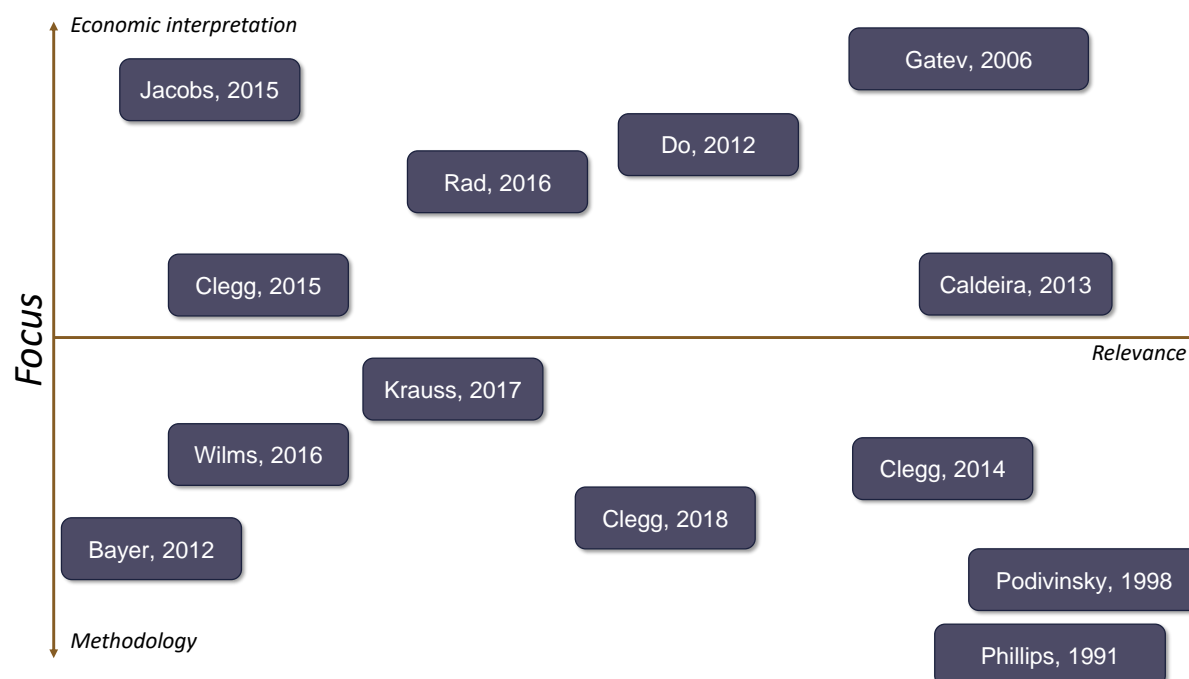
# Introduction

## Literature review



Figure 1: Classification of the core literature.

# Empirical usage of traditional econometric tools

## Explanatory data analysis

Engle-Granger method is a simple way to test cointegration in the bivariate case. Cointegration is diagnosed if the two tested series are integrated in the same order and a linear combination of them exist, which has an integration order of the original non-stationer series minus one [Kirchgässner and Wolters, 2007]. The most common is when the tested stock prices are I(1) and their linear combination is stationer.

The used stock prices are presented in figure 1. For a first glance, there is a high chance that some cointegrated pairs can be found in this set of series. To commit the tests the first step is to check the time-series integration order. For this purpose, I use ADF-test with a significance level of 5%. As a result, it is concluded that all the series are I(1) if any of their bivariate linear combinations is stationer, then cointegration is diagnosed. The first difference in the stock prices is shown in figure **??**.

## Engle-Granger method

The second step is to run OLS with all the possible pairs and check if there is a series of residuals stationer. Just as at the previous step the stationary test is augmented Dickey-Fuller test without constant or trend component in the auxiliary regression and $\alpha = 5\%$.

With the described parameters[1] the tests confirm only one cointegrated pair (see Figure 4), and that result holds only if the stock price of Bank of America is in regressor role, but it does not, when that is used as

---

[1]In my previously mentioned GitHub repository, you may find that I wrote an R function to commit the whole Engle-Granger method with specified parameters. It would be reasonable to see the results with a different stationary test or with a different
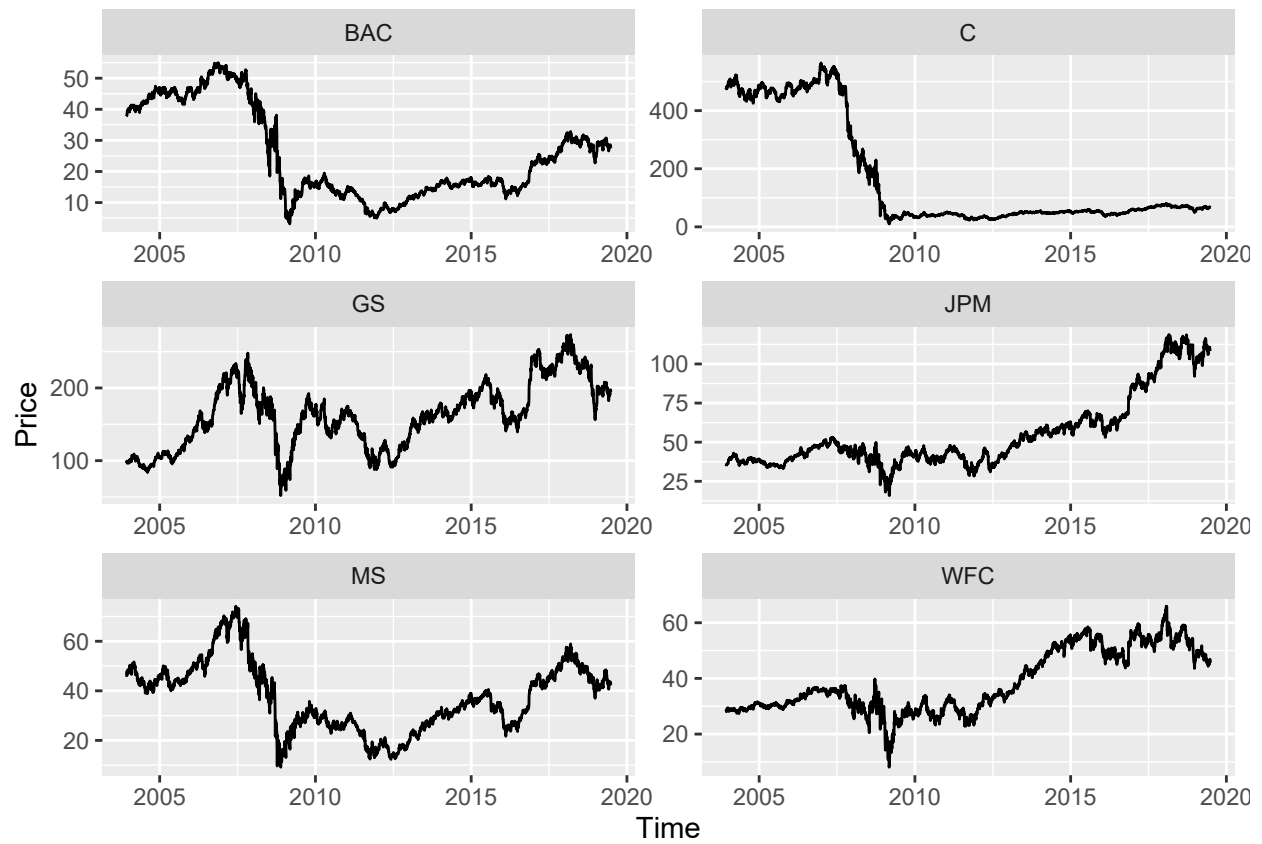
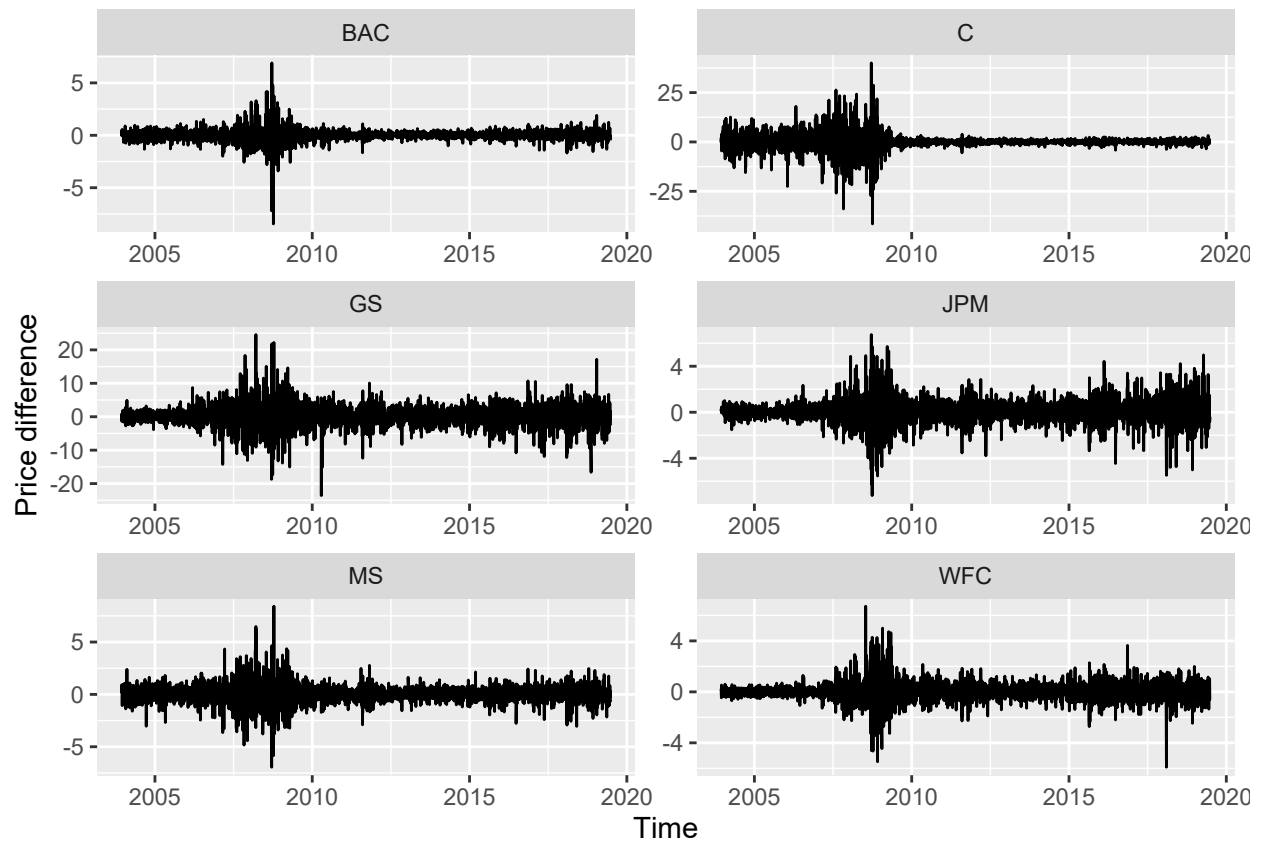Figure 2: Time-series used in this study
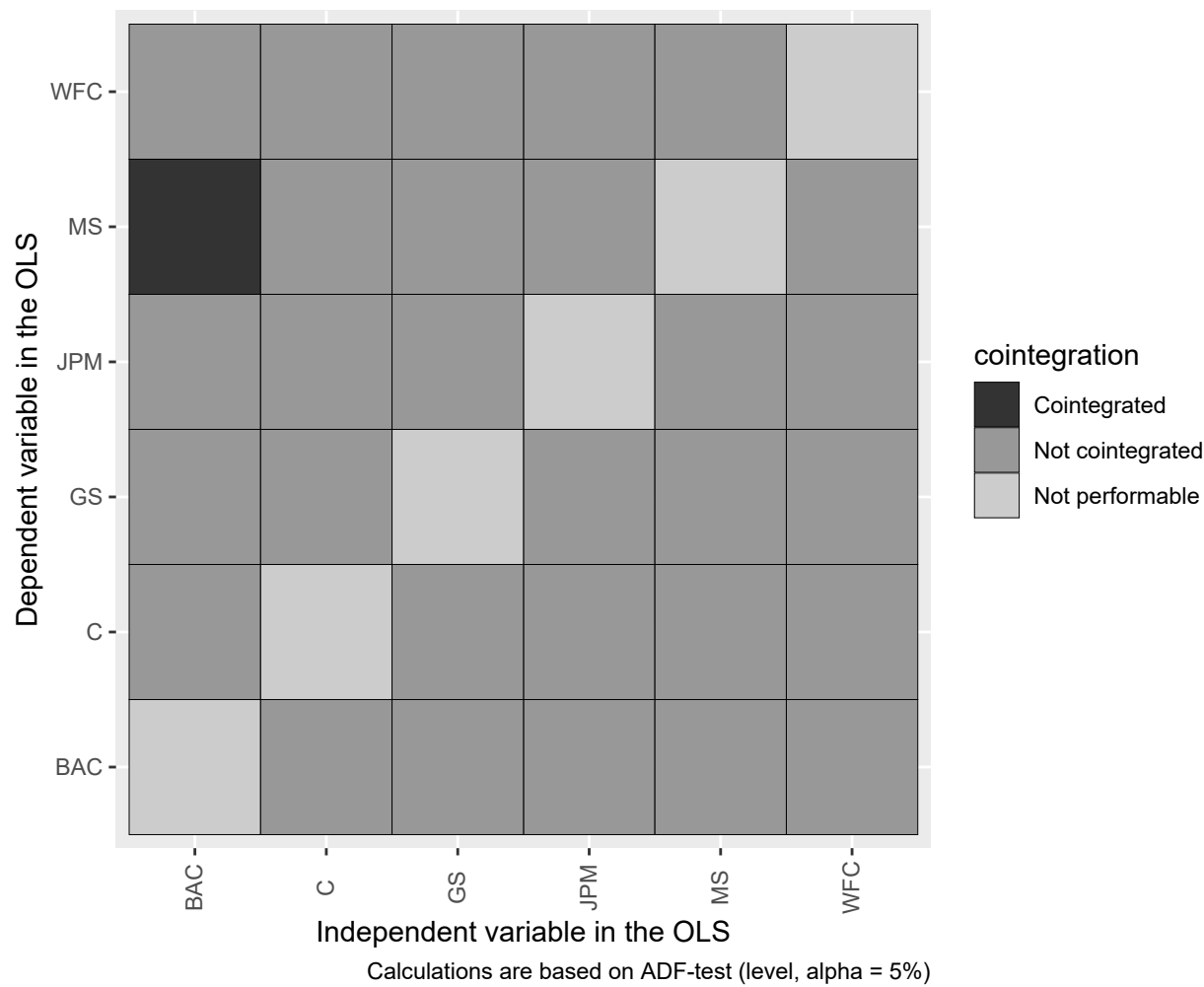
Figure 3: First difference of the time-series

Figure 4: Results of Engle-Granger method

dependent variable[2].

## Johansen-test

Johansen test is adequate cointegration test when there are more than two tested series at the same time. This test is performed to estimate the number of cointegrated vectors (r) in the system. If there is any cointegration in the model then $0 < r < k$, where k is the number of tested time-series. The system decomposition is not unique, so we can only estimate the cointegration rank r [Kirchgässner and Wolters, 2007]. The method can be performed with several tests, in this paper I chose the Lmax test. It gives a vector of the test statistics as a result and that may be compared to critical values. The null hypothesis is that $r \leq x$, where $x = 0, 1, 2, ..., k - 1$. The number of cointegrated vectors is the smallest $x$, under which the null hypothesis is not rejected. The empirical analysis in this study shows that the r in this system is 1 on the full time-interval[3], which confirms the identical result like the one found with the Engel-Granger method.

## Engle-Granger method with rolling window

In this section, I expound the results of the previously presented Engle-Granger method performed with a rolling window. The size of the windows is 250 days. Important to note, it is not sure that a stock price has the same integration order in each window. It can happen that a cointegration test is not performable, because in that period the integration orders do not match. Since this calculation is heavily time-consuming, only three of the six stock will be tested in this paper. This means that the maximum number of cointegrated pairings is 6 ($3 \times 3 - 3$). The test parameters are the same as described before, results are shown in figure 6.

In figure 6 it can be seen that the number of cointegrated pairings reaches the maximum number at the end of 2008, 2012 and in the middle of 2008, 2016. In 2008 there is also a long period when there are 4 cointegrated pairings. This result suggests a pattern that in recession cointegration may be more frequent.

## Johansen test with rolling window

Performing the Johansen test with a rolling window is a similar extension as the one presented in the previous chapter. The calculations were performed with the same 250 window size and $r$ is examined at the significance level of 1%, 5% and 10%. The result can be seen in figure **??**.

In figure **??** the period of recession is also visualized. It looks like the $r = 1$ result at that time is more frequent than most of the case when there is no recession, similarly the $r = 2$ result. One deviation from this pattern is at 2018, where $r = 1$ result is extremely frequent.

Looking at the distribution of the results controlling for the period of recession also confirms this hypothesis. During a recession, the proportion of $r = 2$ result (2.19%) is twice as much as the proportion when there is not recession (1.08%) with 10% significance level. Similarly $r = 1$ is the result of 15.31% of the total tests performed with $\alpha = 10\%$ in periods of recession, while 7.34% is when there is expansion. With different significance level, identical results can be concluded.

---

significance level (especially if calculating its profitability is also in focus). With the written function, it is possible to modify the test parameters and see how the results change.

[2]The matrix of the results is not a symmetrical.

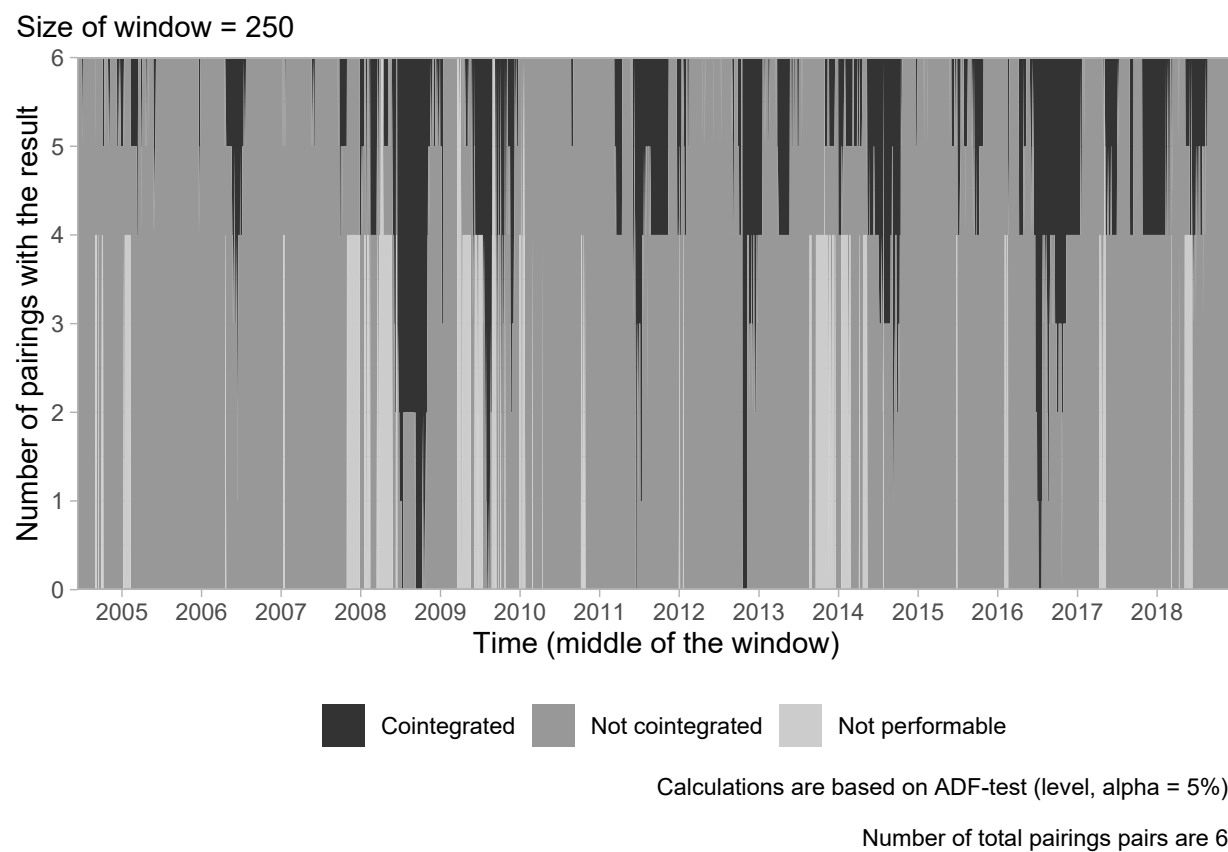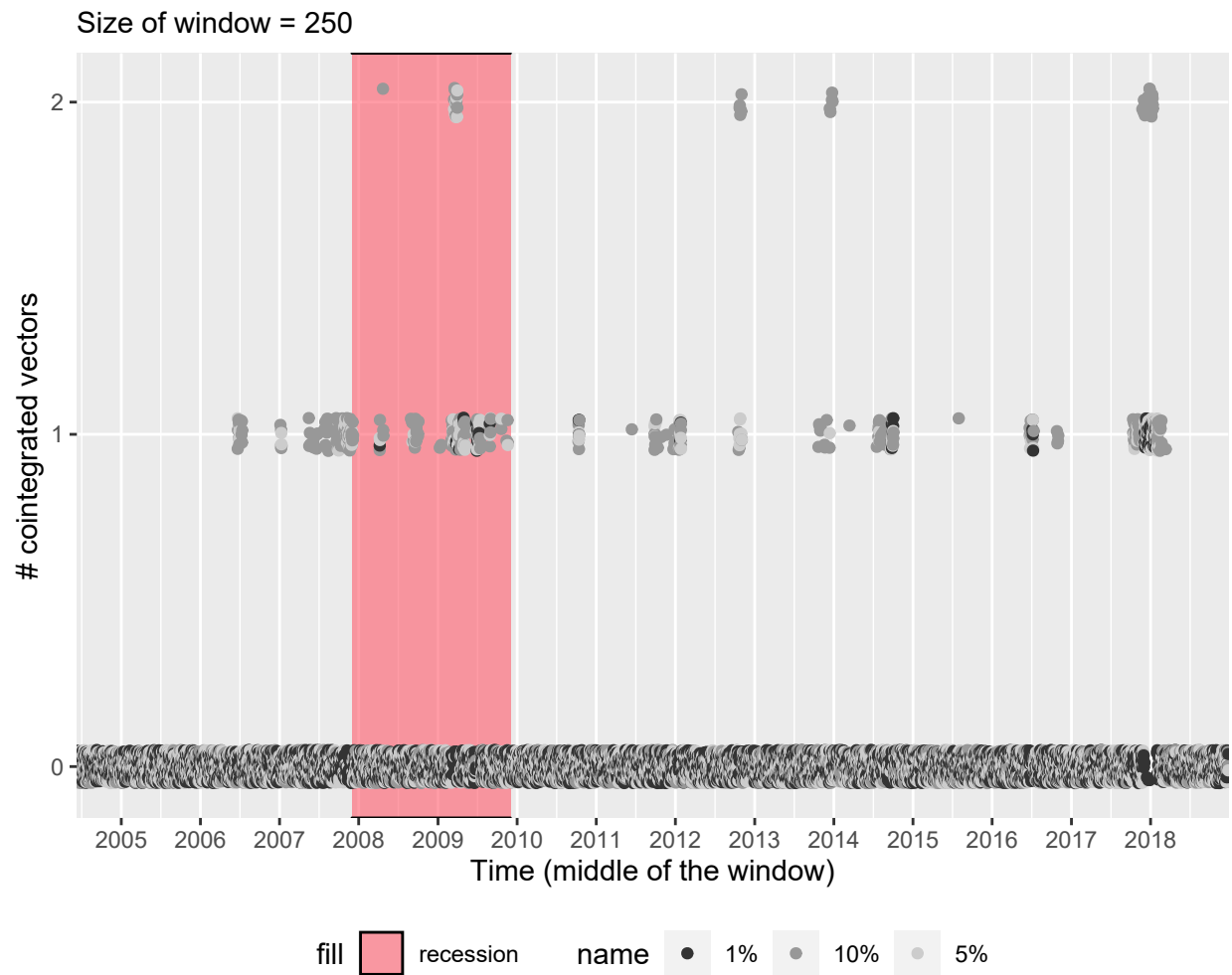[3]Same result is stated on 1%, 5% and 10% significance level.

Figure 5: Results of Engle-Granger method with rolling window

Figure 6: Results of Johansen-test with rolling window across time

# References

[Bayer and Hanck, 2013] Bayer, C. and Hanck, C. (2013). Combining non-cointegration tests. *Journal of Time Series Analysis*, 34(1):83–95.

[Caldeira and Moura, 2013a] Caldeira, J. and Moura, G. V. (2013a). Selection of a portfolio of pairs based on cointegration: A statistical arbitrage strategy.

[Caldeira and Moura, 2013b] Caldeira, J. and Moura, G. V. (2013b). Selection of a portfolio of pairs based on cointegration: A statistical arbitrage strategy. *SSRN Electronic Journal*.

[Clegg, 2014] Clegg, M. (2014). On the persistence of cointegration in pairs trading. *SSRN Electronic Journal*.

[Clegg and Krauss, 2018] Clegg, M. and Krauss, C. (2018). Pairs trading with partial cointegration. *Quantitative Finance*, 18(1):121–138.

[Do and Faff, 2012] Do, B. and Faff, R. (2012). Are pairs trading profits robust to trading costs? *Journal of Financial Research*, 35(2):261–287.

[Gonzalo and Lee, 1998] Gonzalo, J. and Lee, T.-H. (1998). Pitfalls in testing for long run relationships. *Journal of Econometrics*, 86(1):129–154.

[Huck, 2015] Huck, N. (2015). Pairs trading: does volatility timing matter? *Applied Economics*, 47(57):6239–6256.

[Huck and Afawubo, 2014] Huck, N. and Afawubo, K. (2014). Pairs trading and selection methods: is cointegration superior? *Applied Economics*, 47(6):599–613.

[Jacobs and Weber, 2015] Jacobs, H. and Weber, M. (2015). On the determinants of pairs trading profitability. *Journal of Financial Markets*, 23:75–97.

[Kirchgässner and Wolters, 2007] Kirchgässner, G. and Wolters, J. (2007). *Introduction to modern time series analysis.* Springer, Berlin and New York.

[Krauss and Herrmann, 2017] Krauss, C. and Herrmann, K. (2017). On the power and size properties of cointegration tests in the light of high-frequency stylized facts. *Journal of Risk and Financial Management*, 10(1):7.

[Neely et al., 2014] Neely, C. J., Rapach, D. E., Tu, J., and Zhou, G. (2014). Forecasting the equity risk premium: The role of technical indicators. *Management Science*, 60(7):1772–1791.

[Rad et al., 2016] Rad, H., Low, R. K. Y., and Faff, R. (2016). The profitability of pairs trading strategies: distance, cointegration and copula methods. *Quantitative Finance*, 16(10):1541–1558.

## Appendix: R codes

```r
# Setup -------------------------------------------------------------------

library(tidyverse)
library(urca)

WD <- getwd() %>% # root directory
  gsub(pattern = "PairsTrading.*", replacement = "PairsTrading")

load(str_c(WD, "/data.RData")) # financial assets data

theme_set(theme_light() + theme(
  legend.title = element_blank(),
  plot.title.position = "plot",
  plot.tag.position = "topright",
  plot.caption.position = "plot"
))

# EDA ---------------------------------------------------------------------

Bankdata %>%
  pivot_longer(-1) %>%
  ggplot(aes(x = Date, y = value)) +
  geom_line() +
  facet_wrap(vars(name), nrow = 3, scales = "free") +
  labs(
    x = "Time", y = "Price"
  )

Bankdata %>% select(-1) %>% cor() %>% data.frame() %>% rownames_to_column() %>%
  pivot_longer(-1) %>% mutate(
  value = ifelse(rowname == name, NA, value)
) %>%
  ggplot(aes(rowname, name, fill = value)) + geom_tile(color = "black") +
    scale_fill_gradient2(
    low = "#00A3AB", high = "#FF5B6B", space = "Lab", na.value = "grey50",
    guide = "legend", midpoint = 0, aesthetics = "fill", limits = c(-1,1)
  ) + labs(
    x = "", y = "", title = "Correlation-matrix", tag = "Not included"
  ) + theme(
    panel.border = element_blank()
)

# Engle-Granger method ----------------------------------------------------

Bankdata %>%
  select(-1) %>%
  apply(2, function(x) {
    # number of differences required for stationarity to each series
    forecast::ndiffs(x, test = "adf", alpha = 0.05, type = "level")
  })

Bankdata %>%
```

```
53    select(-1) %>%
54    apply(2, function(x) {
55      diff(x)
56    }) %>%
57    data.frame() %>%
58    mutate(
59      Date = tail(Bankdata$Date, -1)
60    ) %>%
61    pivot_longer(-Date) %>%
62    ggplot(aes(x = Date, y = value)) +
63    geom_line() +
64    facet_wrap(vars(name), nrow = 3, scales = "free") +
65    labs(
66      x = "Time", y = "Price difference"
67    )
68
69  cointegration_tests <- function(df, test, type, alpha) {
70    # test cointegrity for all combination in a df
71    ndiff_df <- df %>%
72      select(-1) %>%
73      apply(2, function(x) { # # of differences required for stationarity to each series
74        forecast::ndiffs(x, test = test, alpha = alpha, type = type)
75      })
76
77    v <- df %>% select(-1) %>% # remove year ---> IT MUST BE IN THE INPUT DF !
78      names(.)
79    df2 <- expand.grid(v, v) %>%
80      rename_all(funs(c("y", "x"))) %>%
81      mutate(
82        y = as.character(y),
83        x = as.character(x),
84        ndiff = ifelse(ndiff_df[y] == ndiff_df[x], ndiff_df[y], 0),
85        ndiff = ifelse(y == x, 0, ndiff) # if series are the same, put 0
86      )
87
88    v <- vector()
89    for (i in seq(nrow(df2))) {
90      if (df2[i, 3] != 0) {
91        if (lm(y ~ x, data = rename_all(data.frame(y = df[df2[i, 1]], x = df[df2[i, 2]]),
92                                        funs(c("y", "x")))) %>%
93          broom::augment() %>% .$.resid %>%
94          forecast::ndiffs(test = test, alpha = alpha, type = type) == df2[i, 3] - 1) {
95          v[i] <- 2 # 2 ---> series are cointegrated
96        } else {
97          v[i] <- 1 # 1 ---> not cointegrated, but test is commitable
98        }
99      } else {
100        v[i] <- 0 # 0 ---> test is not performable [I(0) OR not the same I() order OR
101        # series are the same]
102      }
103    }
104    df2 %>%
105      mutate(
```

```
106        cointegration = v
107      ) %>%
108      select(y, x, cointegration)
109 }
110
111 cointegration_tests_results <- cointegration_tests(df = Bankdata, test = "adf",
112                                                    type = "level", alpha = 0.05)
113
114 cointegration_tests_results %>%
115   mutate(
116     cointegration = case_when(
117       cointegration == 0 ~ "Not performable",
118       cointegration == 1 ~ "Not cointegrated",
119       cointegration == 2 ~ "Cointegrated"
120     ),
121     cointegration = factor(cointegration, levels = c("Cointegrated", "Not cointegrated",
122                                                      "Not performable"))
123   ) %>%
124   ggplot() +
125   geom_tile(aes(x = x, y = y, fill = cointegration), color = "black") +
126   scale_fill_grey() +
127   theme(
128     axis.text.x = element_text(angle = 90, vjust = 0.45),
129   ) +
130   labs(
131     y = "Dependent variable in the OLS",
132     x = "Independent variable in the OLS",
133     caption = "Calculations are based on ADF-test (level, alpha = 5%)"
134   ) + theme(
135   panel.border = element_blank()
136 )
137
138 # Johansen-test ---------------------------------------------------------------------
139
140 Bankdata %>%
141   select(-1) %>%
142   ca.jo(type = "eigen", K = 5, ecdet = "none", spec = "longrun") %>%
143   summary() # Number of cointegrated vectors = 1
144
145 ## Engle-Granger method with rolling window ----------------------------------------
146
147 for (i in 1:(nrow(Bankdata) - 249)) {
148   if (i == 1) {
149     cointegration_tests_rw <- mutate(
150       cointegration_tests(df = Bankdata[i:(i + 249), 1:4], test = "adf", type = "level",
151                           alpha = 0.05),
152       t = i
153     )
154   } else {
155     cointegration_tests_rw <- rbind(cointegration_tests_rw, mutate(
156       cointegration_tests(df = Bankdata[i:(i + 249), 1:4], test = "adf", type = "level",
157                           alpha = 0.05),
158       t = i
```

```
159        ))
160      }
161    }
162
163    cointegration_tests_rw %>%
164      filter(y != x) %>%
165      ggplot(aes(x = t, y = cointegration)) +
166      geom_point() +
167      facet_grid(cols = vars(x), rows = vars(y)) +
168      scale_y_continuous(breaks = c(0, 1, 2),
169                         labels = c("Not performable", "Not cointegrated", "Cointegrated")) +
170      labs(
171        subtitle = "Size of window = 250",
172        y = "Result of the test",
173        x = "# window",
174        caption = "Calculations are based on ADF-test (level, alpha = 5%)\n
175        Dependent variables (in the OLS) are placed horizontal, independents are vertical."
176      )
177
178    cointegration_tests_rw %>%
179      filter(cointegration == 2) %>%
180      mutate(cointegration = factor(cointegration)) %>%
181      group_by(y, x) %>%
182      tally() %>%
183      arrange(x) %>%
184      mutate(
185        n = n / max(cointegration_tests_rw$t),
186        n = scales::percent(n, accuracy = .01)
187      ) %>%
188      pivot_wider(id_cols = y, values_from = n, names_from = x, names_prefix = "x = ") %>%
189      arrange(y)
190
191    merge(expand.grid(1:(nrow(Bankdata) - 249), c(0, 1, 2)) %>%
192            rename_all(funs(c("t", "cointegration"))),
193      cointegration_tests_rw %>% filter(y != x) %>%
194        group_by(t, cointegration) %>%
195        summarise(n = n()),
196      all.x = T
197    ) %>%
198      mutate(
199        n = ifelse(is.na(n), 0, n),
200        cointegration = case_when(
201          cointegration == 0 ~ "Not performable",
202          cointegration == 1 ~ "Not cointegrated",
203          cointegration == 2 ~ "Cointegrated"
204        ),
205        cointegration = factor(cointegration, levels = c("Cointegrated", "Not cointegrated",
206                                                          "Not performable")),
207        t = as.Date(Bankdata$Date)[t + 125]
208      ) %>%
209      ggplot() +
210      geom_area(aes(x = t, y = n, fill = cointegration)) +
211      scale_y_continuous(expand = c(0, 0)) +
```

```r
212    scale_x_date(expand = c(0, 0), date_breaks = "1 year", date_labels = "%Y") +
213    theme(
214      legend.position = "bottom"
215    ) +
216    labs(
217      subtitle = "Size of window = 250",
218      y = "Number of pairings with the result",
219      x = "Time (middle of the window)",
220      caption = "Calculations are based on ADF-test (level, alpha = 5%).\n
221      Number of total pairings pairs are 6."
222    ) +
223    scale_fill_grey()
224
225  # Johansen test with rolling window -----------------------------------------------
226
227  johansen_tests_rw <- data.frame(t = 1:(nrow(Bankdata) - 249)) %>% mutate(
228    pct10 = NA, pct5 = NA, pct1 = NA
229  )
230
231  for (i in 1:(nrow(Bankdata) - 249)) {
232    if (i == 1) {
233      johansen_critical_values <- ca.jo(
234        x = Bankdata[i:(i + 249), 2:4], type = "eigen",
235        K = 5, ecdet = "none", spec = "longrun"
236      )@cval
237    }
238    johansen_tests_rw[i, 2] <- which.max(rev(ca.jo(
239      x = Bankdata[i:(i + 249), 2:4], type = "eigen",
240      K = 5, ecdet = "none", spec = "longrun"
241    )@teststat) < rev(johansen_critical_values[, 1])) - 1
242    johansen_tests_rw[i, 3] <- which.max(rev(ca.jo(
243      x = Bankdata[i:(i + 249), 2:4], type = "eigen",
244      K = 5, ecdet = "none", spec = "longrun"
245    )@teststat) < rev(johansen_critical_values[, 2])) - 1
246    johansen_tests_rw[i, 4] <- which.max(rev(ca.jo(
247      x = Bankdata[i:(i + 249), 2:4], type = "eigen",
248      K = 5, ecdet = "none", spec = "longrun"
249    )@teststat) < rev(johansen_critical_values[, 3])) - 1
250  }
251
252  ggplot() +
253    geom_ribbon(aes(
254      x = c(as.Date("2007-12-01"), as.Date("2009-12-01")),
255      ymin = -Inf,
256      ymax = Inf,
257      fill = "recession"), color = "black", alpha = .6) +
258    geom_jitter(data = johansen_tests_rw %>%
259                  pivot_longer(-1) %>%
260                  mutate(
261                    name = case_when(
262                      name == "pct1" ~ "1%",
263                      name == "pct5" ~ "5%",
264                      name == "pct10" ~ "10%"
```

```
265                    ),
266                    t = as.Date(Bankdata$Date)[t + 125]
267                  ),
268                aes(x = t, y = value, color = name),width = 0, height = 0.05) +
269     scale_color_grey() +
270     theme(
271       legend.position = "bottom"
272     ) +
273     scale_y_continuous(breaks = c(0, 1, 2)) +
274     scale_x_date(expand = c(0, 0), date_breaks = "1 year", date_labels = "%Y") +
275     labs(
276       subtitle = "Size of window = 250",
277       y = "# cointegrated vectors",
278       x = "Time (middle of the window)",
279       caption = str_wrap(str_c(
280         "Points are jittered around their true y value for better ",
281          "visualisation (the number of cointegrated vectors is interger). ",
282           "Date of recession is from the National Bureau of Economic Research ",
283           "(https://www.nber.org/cycles.html)."), 50)
284     ) +
285     theme(
286       panel.grid.minor.y = element_blank()
287     ) +
288     scale_fill_manual(values = c("recession" = "#FF5B6B"))
289
290  johansen_tests_rw %>%
291    select(-1) %>%
292    gather() %>%
293    mutate(
294      key = case_when(
295        key == "pct1" ~ "1%",
296        key == "pct5" ~ "5%",
297        key == "pct10" ~ "10%"
298      ),
299      key = factor(key, levels = c("10%", "5%", "1%"))
300    ) %>%
301    group_by(key, value) %>%
302    tally() %>%
303    ggplot() +
304    geom_bar(aes(x = key, y = n, fill = factor(value, levels = 2:0)), position = "fill",
305             stat = "identity", color = "black") +
306    scale_y_continuous(labels = scales::percent_format(accuracy = 1), expand = c(0, 0),
307                       breaks = seq(from = 0, to = 1, by = .1)) +
308    scale_fill_grey() +
309    labs(
310      title = "Distribution of the Johansen-test results with rolling window",
311      x = "Alpha",
312      y = "Proportion",
313      fill = "Number cointegrated vectors (r)",
314      subtitle = "Size of window = 250"
315    ) +
316    theme(
317      legend.title = element_text(),
```

```
318    legend.position = "bottom"
319    )
320
321 johansen_tests_rw %>%
322    pivot_longer(-1) %>%
323    mutate(
324      name = factor(name, levels = c("pct1", "pct5", "pct10")),
325      t = as.Date(Bankdata$Date)[t + 125],
326      t = ifelse(t > as.Date("2007-12-01") & t < as.Date("2009-12-01"), "recession",
327                 "expansion")
328    ) %>% filter(t == "expansion") %>% group_by(name) %>% count(value) %>% pivot_wider(
329      id_cols = value, values_from = n, names_from = name
330    )  %>% mutate(
331      pct1 = scales::percent(pct1/sum(pct1, na.rm = T), accuracy = .01),
332      pct5 = scales::percent(pct5/sum(pct5, na.rm = T), accuracy = .01),
333      pct10 = scales::percent(pct10/sum(pct10, na.rm = T), accuracy = .01)
334    ) %>% rename_all(funs(c("# cointegrated vectors", "1%", "5%", "10%")))
335
336 johansen_tests_rw %>%
337    pivot_longer(-1) %>%
338    mutate(
339      name = factor(name, levels = c("pct1", "pct5", "pct10")),
340      t = as.Date(Bankdata$Date)[t + 125],
341      t = ifelse(t > as.Date("2007-12-01") & t < as.Date("2009-12-01"), "recession",
342                 "expansion")
343    ) %>% filter(t == "recession") %>% group_by(name) %>% count(value) %>%
344    pivot_wider(
345      id_cols = value, values_from = n, names_from = name
346    )  %>% mutate(
347      pct1 = scales::percent(pct1/sum(pct1, na.rm = T), accuracy = .01),
348      pct5 = scales::percent(pct5/sum(pct5, na.rm = T), accuracy = .01),
349      pct10 = scales::percent(pct10/sum(pct10, na.rm = T), accuracy = .01)
350    ) %>% rename_all(funs(c("# cointegrated vectors", "1%", "5%", "10%")))
```