

University of Duisburg-Essen  
Faculty of Business Administration and  
Economics  
Chair of Econometrics



# P-Approximation

Seminar in Econometrics

Term Paper

Submitted to the Faculty of  
Business Administration and Economics  
at the  
University of Duisburg-Essen

from:

Jens Klenke and Janine Langerbein

---

Reviewer: Christoph Hanck

Deadline: Jan. 17th 2020

---

Name:	Jens Klenke	Janine Langerbein
Matriculation Number:	3071594	3061371
E-Mail:	jens.klenke@stud.uni-due.de	janine.langerbein@stud.uni-due.de
Study Path:	M.Sc. Economics	M.Sc. Economics
Semester:	5 <sup>th</sup>	5 <sup>th</sup>
Graduation (est.):	Summer Term 2021	Summer Term 2021

# Contents

List of Figures	II
List of Tables	II
List of Abbreviations	II
1 Introduction	1
2 Bayer Hanck Test	1
3 Simulation	3
4 Models	4
4.1 Data Pre-Processing . . . . .	4
4.2 Polynomial Regression . . . . .	5
4.3 Lasso . . . . .	6
4.4 Other regression models . . . . .	7
5 Model evaluation	7
6 Package	7
References	III
Software-References	VI
A Appendices	VII

**List of Figures**

**List of Tables**

**List of Abbreviations**

# 1 Introduction

Meta tests have been shown to be a powerful tool when testing for the null of non-cointegration. The distribution of their test statistic, however, is mostly not available in closed form. This might pose difficulties when implementing the meta tests in econometric software packages, as one has to include the full null distribution for each combination of the underlying tests. Software package size limitations are therefore quickly exceeded.

In this paper we propose supervised Machine Learning Algorithms to approximate the p-values of the meta test by Bayer and Hanck (2012) which tests for the null of non-cointegration. This approach might reduce the size of associated software packages considerably. The algorithms are trained on simulated data for various specifications of the aforementioned test.

Ergebnis der Models (1-2 Sätze)

Inhalt Paper

## 2 Bayer Hanck Test

The choice as to which of the available cointegration tests to use is a recurrent issue in econometric time series analysis. Bayer and Hanck (2012) propose powerful meta tests which provide unambiguous test decisions. They combine several residual- and system-based tests in the manner of Fisher's (1932) Chi-squared test.

Bayer and Hanck build their paper on previous work from Pesavento (2004), who defines the underlying model as  $z'_t = [x'_t, y_t]$ , with  $x_t$  being an  $n_1 \times 1$  vector and  $y_t$  a scalar, which displays the cointegration relation. They can be written as

$$\Delta x_t = \tau_1 + v_{1t} \quad (2.1)$$

$$y_t = (\mu_2 - \gamma' \mu_1) + (\tau_2 - \gamma' \tau_1)t + \gamma' x_t + u_t, \quad (2.2)$$

$$u_t = \rho u_{t-1} + v_{2t}. \quad (2.3)$$

$\Delta x_t$  presents the regressor dynamics.  $\mu_1, \mu_2, \tau_1$  and  $\tau_2$  are the deterministic parts of the model. They are subject to the following restrictions: (i)  $\mu_2 - \gamma' \mu_1$  and  $\tau = 0$  which translates to no deterministics, (ii)  $\tau = 0$  which corresponds to a constant in the cointegrating vector, (iii)  $\tau_2 - \gamma' \tau_1 = 0$ , a constant plus trend.

$v_t = [v'_{1t} v_{2t}]'$  with  $\Omega$  the long-run covariance matrix of  $v_t$ . For derivation of  $v_t$  see Pesavento (2004). Pesavento shows that  $\{v_t\}$  satisfies an FCLT, i.e.  $T^{-1/2} \sum_{t=1}^{[T\cdot]} v_t \Rightarrow \Omega^{1/2} W(\cdot)$ . It is further assumed that the  $x_t$  are not cointegrated.

rho und p?

It clearly follows from (2.3) that  $z_t$  is cointegrated if  $\rho < 1$ . Hence the null hypothesis of no cointegration is  $H_0 : \rho = 1$ . Furthermore, Pesavento introduces two other parameters. First,  $R^2$  measures the squared correlation of  $v_{1t}$  and  $v_{2t}$ . It can be interpreted as the influence of the right-hand side variables in (2.2). It ranks between zero and one. When there is no long-run correlation between those variables and the errors from the cointegration regression,  $R^2$  equals zero. Secondly, the number of lags is approximated by a finite number  $k$ .

k geht hier nicht

### Assumptions (BH S. 84)?

Bayer and Hanck's (2012) meta test considers the test statistics of up to four stand-alone tests. Namely, these are the tests of Engle and Granger (1987), Johansen (1988), Boswijk (1994) and Banerjee, Dolado, and Mestre (1998). For the sake of brevity the detailed derivation of the underlying tests has been deliberately omitted here.

Engle and Granger (1987) propose a two-step procedure to test the null hypothesis of no cointegration against the alternative of at least one cointegrating vector. First, the long-run relationship between  $y_t$  and  $\mathbf{x}_t$  is estimated by least squares regression. The obtained residuals  $\hat{u}_t$  are then tested for a unit root. For this, Engle and Granger suggest the use of the  $t$ -statistic  $t_{\gamma}^{\text{ADF}}$  in the Augmented Dickey-Fuller (ADF) regression:

$$\Delta \hat{u}_t = \gamma \hat{u}_{t-1} + \sum_{i=1}^k \pi_i \Delta \hat{u}_{t-i} + \varepsilon_t. \quad (2.4)$$

The rejection of a unit root points to a cointegration relationship.

Johansen's (1988) maximum eigenvalue test is a system-based test that allows for several cointegration relationships. Take the vector error correction model (VECM)<sup>1</sup>

$$\Delta \mathbf{z}_t = \Pi \mathbf{z}_{t-1} + \sum_{i=1}^k \Gamma_i \Delta \mathbf{z}_{t-i} + \mathbf{d}_t + \varepsilon_t. \quad (2.5)$$

blabla Johansen test statistic

Banerjee and Boswijk

---

<sup>1</sup>Due to practical reasons we omit the derivation of the VECM which is presumed to be known.

To combine the results from the underlying tests Bayer and Hanck draw upon Fisher’s combined probability test (Fisher, 1932). It merges the tests using the formula

$$\tilde{\chi}_{\mathcal{I}}^2 := -2 \sum_{i \in \mathcal{I}} \ln(p_i). \quad (2.6)$$

Let  $t_i$  be the  $i^{th}$  test statistic. If test  $i$  rejects for large values, take  $\xi_i := t_i$ . If test  $i$  rejects for small values, take  $-\xi_i := t_i$ . With  $\Xi_i(x) := \Pr_{\mathcal{H}_i}(\xi_i \geq x)$  the p-value of the  $i^{th}$  test is  $p_i := \Xi_i(\xi_i)$ .

Fisher (1932) shows that under the assumption of independence the null distribution of  $\tilde{\chi}_{\mathcal{I}}^2$  follows a chi-squared distribution with  $2\mathcal{I}$  degrees of freedom. If this assumption is violated the null distribution is less evident. Here, the latter case occurs, as the  $\xi_i$  are not independent. The  $\tilde{\chi}_{\mathcal{I}}^2$ , however, have well-defined asymptotic null distributions  $F_{\mathcal{F}_{\mathcal{I}}}$ , as  $\tilde{\chi}_{\mathcal{I}}^2 \rightarrow_d \mathcal{F}_{\mathcal{I}}$  under  $\mathcal{H}_0$  if  $T \rightarrow \infty$ , with  $\mathcal{F}_{\mathcal{I}}$  some random variable. It is therefore feasible to simulate the joint null distribution of the  $\xi_i$  to obtain the distribution  $F_{\mathcal{F}_{\mathcal{I}}}$  of (2.6). The  $F_{\mathcal{F}_{\mathcal{I}}}$  depend on which and how many tests are combined. The distributions of the  $\xi_i$  depend on  $K - 1$  and the deterministic case.

### 3 Simulation

In this section, we describe the simulation of the null distribution of the Bayer Hanck meta test. The objective is to obtain data for training machine learning algorithms on approximating the p-values of the aforementioned test. In consideration of the different forms of the meta test we generate six data sets. These vary according to the specific combinations of the underlying tests and also account for the above-mentioned restrictions on the deterministic parts of the model.

The following approach relies largely on previous work by Pesavento (2004). For calculating the Bayer Hanck test statistic we require the p-values of the underlying tests. For this, we simulate their null distributions. It can be shown that asymptotically these are functions of standard Brownian motions. Here, the latter are constructed by step functions using Gaussian random walk of size  $N = 1000$ . The number of repetitions is set to 1,000,000. Furthermore, we consider  $R^2 \in \{0, 0.05, 0.1, \dots, 0.95\}$ , the maximum number of lags  $K = 11$  and  $c = 0^2$  (c mal definieren).

---

<sup>2</sup>Since we solely aim at simulating the distribution of the null of no cointegration we

From the mass of test statistics we build the cumulative distribution function of each underlying test and calculate the respective p-values. These are inserted into (2.6) to eventually obtain the Bayer Hanck test statistics. Analogous to the previous approach, we deduce the associated null distribution and the p-values.

## 4 Models

We now use the generated data for training machine learning algorithms on predicting the approximated empirical CDF of the Bayer Hanck test. We work with the values of the test statistic and the number of lags  $k$  as predictors. As it is our objective to describe the null distribution with a less memory-intensive model we will only consider linear methods. For the same objective we compare the models according to their in-sample RMSE. The threat of overfitting is thus of no particular relevance here. For this reason, and to reduce computation time, we use no cross-validation.

As the empirical CDF is typically known to be curved in an S-shape we skip the classic linear regression in favor of a more flexible model. We stay with least squares regression, but try various combinations of polynomial functions and interaction terms of the aforementioned regressors. The search for the best model is carried out via brute-force.

### 4.1 Data Pre-Processing

One approach for improving a model's predictive ability is the pre-processing of the training data. Some models, like linear regression, react sensitively to certain characteristics of the predictor or response data. Those characteristics include, inter alia, distributional skewness and outliers and there exist several methods to lower their potentially bad impact on the model's performance.

Since we simulated our training data under the null of non-cointegration we expect the distribution of the test statistic to be rather right skewed. Plot also reveals it to have a long right tail. If we train our regression model on this raw data it can possibly have difficulties predicting from high values of the test statistic.

One of the aforementioned methods to deal with such issues are power transforms. One might decide freely which transformation to apply. Alternatively, will not consider any further values of  $c$  here.

Hier Bayer und Hanck zitieren, da Sie selbst den Test ja anders nennen

klingt aus meiner Sicht so, als wäre memory concerns entscheidend ob in-sample RMSE

klingt widersprüchlich, verstehe was du meinst, aber da müssen wir mal mit dem Wording gucken

$X^2$  distribution ist von 0 bis unendlich. vllt mehr auf die Theorie abstellen

there exist statistical methods to determine an appropriate transformation. A well-known family of transformations to un-skew data is the Box-Cox transformation (Box & Cox, 1964). They aim at transforming the data so that it closely resembles the normal distribution. The exact transformation depends on the parameter  $\lambda$ , whose optimal value can be empirically estimated:

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \lambda \neq 0 \\ \log(y), & \lambda = 0 \end{cases} \quad (4.1)$$

würde eher vorschlagen, dass Box/Cox das ursprünglich nur für die response vorgestellt hat und Kuhn dann gezeigt hat, dass es auch für Regressors funktioniert

It is visible from (4.1) that Box and Cox (1964) developed these transformations for the dependent variable. Kuhn and Johnson (2013), however, report that it proves as effective for transforming individual regressors.

Plots hier einfügen, Test Stat all und ej mit und ohne bc.

We estimate lambda for the values of the test statistics of the Bayer-Hanck test and transform them according to (4.1). This forces their distribution into a more symmetric form.

Response transformieren

k als dummy/factor

## 4.2 Polynomial Regression

Aus meiner Sicht relativ lang und doppelt zum oberkapitel oder?

Due to the reasons given above we restrict ourselves to linear models. The empirical CDF, which we aim to predict, is known to have a curved shape. For this reason, a simple linear regression model is very unlikely to provide a satisfactory fit to the data. We are in need of a more flexible model to predict the response as accurately as possible.

Polynomial Regression extends the classic linear regression model by fitting a polynomial equation of arbitrary order to the data. A polynomial regression with  $n$  degrees thus takes the form

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_n x_i^n + \varepsilon_i, \quad (4.2)$$

where  $\varepsilon_i$  is the error term. Quelle?

Here, we calculate orthogonal polynomials of the test statistic of the Bayer-Hanck Test, considering up to 15 degrees. We estimate the parameters with ordinary least squares. To potentially increase the predictive performance



Efficiency concerns  
ansprechen (Varianz der  
geschätzten Koef nehmen  
zu wenn viele Koef)

of our model we also add interaction terms and different transformations of the regressor  $k$ . **Appendix** lists all calculated models. Since there is no need to prevent overfitting we expect higher order polynomials to perform best, as they are highly flexible. These polynomials, however, tend to show a wiggly behaviour at the boundaries. This makes extrapolation beyond the limits of our simulated data a risky endeavour. We will address and fix this issue later on.

### 4.3 Lasso

As mentioned above our polynomial regression models are likely to perform best with higher order polynomials. With each added polynomial, however, we increase the complexity of our model and potentially add redundant regressors. Although, still, overfitting plays no major role here, we generally prefer sparser models in case of equal results. One way to deal with this is the use of variable selection methods. A well-known example of such methods is the least absolute shrinkage and selection operator (LASSO).

The lasso estimate is defined as

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \text{ s.t. } \sum_{j=1}^p |\beta_j| \leq t, \quad (4.3)$$

where the first term describes the residual sum of squares, subject to a term known as L1 penalty. In its Lagrangian form this can be rewritten as

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (4.4)$$

$\lambda$  is a tuning parameter which defines the degree of regularisation. The lasso penalty shrinks the coefficients and, for  $\lambda$  sufficiently large, can set them to zero. The value of  $\lambda$  is data dependent and is usually estimated with cross-validation. ausführlicher? Quelle?

We plan on fitting a LASSO model to polynomials of grade 15. We consider the same transformations and interaction terms as in earlier steps. We therefore fit a total of **Anzahl** models.

## 4.4 Other regression models

We also considered various other regression models. For different reasons they were not too suitable for our use case. Conventional non-linear methods, like Generalized Additive Models or Multivariate Adaptive Regression Splines, might have provided a decent prediction. However, the fitted models take up more memory space than the aforementioned linear methods. For the same reason refrain from using tree based methods. In addition, the latter tend to perform poorly with such a small amount of regressors. Given these limitations, we decided to stick solely with linear regression models. LASSO fehlt

## 5 Model evaluation

We estimate all models for two different combinations of the underlying tests. Namely, these are a combination of the Engle-Granger and Johansen test (EJ) and a combination of all four underlying tests (all). Furthermore, we estimate one model per specification of the model deterministics. Altogether, this results in a total of six different models.

## 6 Package

## References

- Banerjee, A., Dolado, J., & Mestre, R. (1998). Error-correction mechanism tests for cointegration in a single-equation framework. *Journal of Time Series Analysis*, 19(3), 267–283. Retrieved from <https://EconPapers.repec.org/RePEc:bla:jtsera:v:19:y:1998:i:3:p:267-283>
- Bayer, C., & Hanck, C. (2012). Combining non-cointegration tests. *Journal of Time Series Analysis*.
- Boswijk, H. P. (1994). Testing for an unstable root in conditional and structural error correction models. *Journal of Econometrics*, 63(1), 37–60. Retrieved from <https://EconPapers.repec.org/RePEc:eee:econom:v:63:y:1994:i:1:p:37-60>
- Box, G. E. P., & Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 26(2), 211–252. Retrieved from <http://www.jstor.org/stable/2984418>
- Engle, R., & Granger, C. W. (1987). Co-integration and error correction: Representation, estimation and testing. *Econometrica*, 55, 251–276.
- Fisher, R. A. (1932). *Statistical methods for research workers*. Oliver, Boyd, Edinburgh, and London.
- Johansen, S. (1988). Statistical analysis of cointegration vectors. *Journal of Economic Dynamics and Control*, 12(2), 231–254. doi:[https://doi.org/10.1016/0165-1889\(88\)90041-3](https://doi.org/10.1016/0165-1889(88)90041-3)
- Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling*. SpringerLink : Bücher. Springer New York. Retrieved from <https://books.google.de/books?id=xYRDAAAQBAJ>
- Pesavento, E. (2004). Analytical evaluation of the power of tests for the absence of cointegration. *Journal of Econometrics*, 122(2), 349–384.

## Software-References

- Breiman, L., Cutler, A., Liaw, A., & Wiener, M. (2018). *Randomforest: Breiman and cutler's random forests for classification and regression*. R package version 4.6-14. Retrieved from <https://CRAN.R-project.org/package=randomForest>
- Croissant, Y., Millo, G., & Tappe, K. (2019). *Plm: Linear models for panel data*. R package version 2.1-0. Retrieved from <https://CRAN.R-project.org/package=plm>
- Friedman, J., Hastie, T., Tibshirani, R., Simon, N., Narasimhan, B., & Qian, J. (2019). *Glmnet: Lasso and elastic-net regularized generalized linear models*. R package version 2.0-18. Retrieved from <https://CRAN.R-project.org/package=glmnet>
- Greenwell, B., Boehmke, B., Cunningham, J., & Developers, G. (2019). *Gbm: Generalized boosted regression models*. R package version 2.1.5. Retrieved from <https://CRAN.R-project.org/package=gbm>
- Henry, L., & Wickham, H. (2019). *Purrr: Functional programming tools*. R package version 0.3.2. Retrieved from <https://CRAN.R-project.org/package=purrr>
- Hlavac, M. (2018). *Stargazer: Well-formatted regression and summary statistics tables*. R package version 5.2.2. Retrieved from <https://CRAN.R-project.org/package=stargazer>
- Izrailev, S. (2014). *Tictoc: Functions for timing r scripts, as well as implementations of stack and list structures*. R package version 1.0. Retrieved from <https://CRAN.R-project.org/package=tictoc>
- Kuhn, M., Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., ... Hunt, T. (2019). *Caret: Classification and regression training*. R package version 6.0-84. Retrieved from <https://CRAN.R-project.org/package=caret>
- Lumley, T., & Miller, A. (2017). *Leaps: Regression subset selection*. R package version 3.0. Retrieved from <https://CRAN.R-project.org/package=leaps>
- Mevik, B.-H., Wehrens, R., & Liland, K. H. (2019). *Pls: Partial least squares and principal component regression*. R package version 2.7-1. Retrieved from <https://CRAN.R-project.org/package=pls>

- Milborrow, S. (2019a). *Plotmo: Plot a model's residuals, response, and partial dependence plots*. R package version 3.5.5. Retrieved from <https://CRAN.R-project.org/package=plotmo>
- Milborrow, S. (2019b). *Rpart.plot: Plot 'rpart' models: An enhanced version of 'plot.rpart'*. R package version 3.0.7. Retrieved from <https://CRAN.R-project.org/package=rpart.plot>
- R Core Team. (2019). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Ripley, B. (2019a). *Class: Functions for classification*. R package version 7.3-15. Retrieved from <https://CRAN.R-project.org/package=class>
- Ripley, B. (2019b). *Mass: Support functions and datasets for venables and ripley's mass*. R package version 7.3-51.4. Retrieved from <https://CRAN.R-project.org/package=MASS>
- Ripley, B. (2019c). *Tree: Classification and regression trees*. R package version 1.0-40. Retrieved from <https://CRAN.R-project.org/package=tree>
- RStudio Team. (2019). *Rstudio: Integrated development environment for r*. Version 1.2.1541. RStudio, Inc. Boston, MA. Retrieved from <http://www.rstudio.com/>
- Rushworth, A. (2019). *Inspectdf: Inspection, comparison and visualisation of data frames*. R package version 0.0.4. Retrieved from <https://CRAN.R-project.org/package=inspectdf>
- Sievert, C., Parmer, C., Hocking, T., Chamberlain, S., Ram, K., Corvellec, M., & Despouy, P. (2019). *Plotly: Create interactive web graphics via 'plotly.js'*. R package version 4.9.0. Retrieved from <https://CRAN.R-project.org/package=plotly>
- Therneau, T., & Atkinson, B. (2019). *Rpart: Recursive partitioning and regression trees*. R package version 4.1-15. Retrieved from <https://CRAN.R-project.org/package=rpart>
- Ushey, K., Allaire, J., Wickham, H., & Ritchie, G. (2019). *Rstudioapi: Safely access the rstudio api*. R package version 0.10. Retrieved from <https://CRAN.R-project.org/package=rstudioapi>

- Wickham, H. (2019). *Stringr: Simple, consistent wrappers for common string operations*. R package version 1.4.0. Retrieved from <https://CRAN.R-project.org/package=stringr>
- Wickham, H., François, R., Henry, L., & Müller, K. (2019). *Dplyr: A grammar of data manipulation*. R package version 0.8.0.1. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- Wickham, H., & Henry, L. (2019). *Tidyr: Easily tidy data with 'spread()' and 'gather()' functions*. R package version 0.8.3. Retrieved from <https://CRAN.R-project.org/package=tidyr>
- Xie, Y. (2019). *Knitr: A general-purpose package for dynamic report generation in r*. R package version 1.23. Retrieved from <https://CRAN.R-project.org/package=knitr>

## A Appendices

### **Eidesstattliche Versicherung**

Ich versichere an Eides statt durch meine Unterschrift, dass ich die vorstehende Arbeit selbständig und ohne fremde Hilfe angefertigt und alle Stellen, die ich wörtlich oder annähernd wörtlich aus Veröffentlichungen entnommen habe, als solche kenntlich gemacht habe, mich auch keiner anderen als der angegebenen Literatur oder sonstiger Hilfsmittel bedient habe. Die Arbeit hat in dieser oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen.

Essen, den \_\_\_\_\_

\_\_\_\_\_  
Jens Klenke and Janine Langerbein