

README for GitHub repository accompanying the paper:

# Augmenting the availability of historical GDP per capita estimates through machine learning

by Philipp Koch, Viktor Stojkoski, & César A. Hidalgo

## Version updates

v2: Now including source data for European regions in 1900 and 1950 provided by Rosés & Wolf.<sup>1,2</sup> We rescale their estimates such that they match country-level data by the Maddison project. This leads to slightly different out-of-sample estimates in 1900 and 1950 than in the original version. The repository, thus, also includes an Excel file with the original version 'historicalGDPs\_augmented\_dataset\_v1.xlsx' and information in the code '00\_master.R' on how to replicate v1.

The purpose of this document is twofold. First, we want to provide explanations for the final dataset and results that we publish. Second, we want to provide detailed information on the code to ensure replicability.

For all questions regarding the repository, please contact Philipp Koch ([philipp.koch@ecoaustria.ac.at](mailto:philipp.koch@ecoaustria.ac.at)).

## Dataset & results

The final dataset (*historicalGDPs\_augmented\_dataset.xlsx*) contain 5,700 observations with the following columns:

- *ID*: location\_year
- *location*: denotes the code for the respective country (iso-3) or region. The classification of regions follows the NUTS-2 (2021 edition) classification in Europe, metro- and micropolitan areas in the US, metropolitan areas in Canada, and oblasts in Ukraine and Russia.
- *country\_0*: This column is relevant, if the location is a region. Then country\_0 denotes the country the region is in.
- *location\_name*: denotes the name of the respective country or region.
- *year*: in 50-year intervals
- *GDPpc*: GDP per capita levels denoted in 2011 USD.

---

<sup>1</sup> Rosés, J. R., & Wolf, N. (2021). Regional growth and inequality in the long-run: Europe, 1900–2015. *Oxford Review of Economic Policy*, 37(1), 17–48. <https://doi.org/10.1093/oxrep/graa062>

<sup>2</sup> Rosés, J. R., & Wolf, N. (2019). *The economic development of Europe's regions: A quantitative history since 1900*. Routledge.

- *flag*: denotes whether an observation is based on source data, or an out-of-sample estimate based on our machine learning models. Also, it denotes the specific source.
- *GDPpc\_lower* & *GDPpc\_upper*: 90% confidence interval for out-of-sample estimates obtained by bootstrapping.

The dataset contains 5,700 observations, with 1,574 observations from source data and 4,126 out-of-sample estimates. All references for source data observations are described in detail in the manuscript and the SI Appendix.

The folder *genfiles* contains all figures included in the manuscript and the SI Appendix:

- Folder *figures*: contains raw versions of the subplots of Fig. 2, Fig. 3, and Fig. 4 H
- Folder *figures\_SI*: contains several figures of the robustness checks, barcharts for Shapley values from all years, and results from the elastic net model (in the folder *LASSO\_results*).
- Folder *maps*: contains several maps of the GDP per capita levels for countries and regions, e.g. those in Fig. 4 A-E in the manuscript.

## Replication

The folder *scripts* contains all R scripts necessary to replicate our findings and data. The file *00\_master.R* calls all other scripts and reproduces the results.

It is necessary to download another folder with the data on famous individuals from here: <https://zenodo.org/doi/10.5281/zenodo.11546192>.

All files included in the zenodo repository must be saved in a subfolder called *./raw\_data/famous\_individuals* for the scripts to find it.

Below you find all details on the used packages and R versions:

R version 4.2.3 Patched (2023-03-29 r84137)  
Platform: aarch64-apple-darwin20 (64-bit)  
Running under: macOS 14.5

Matrix products: default  
LAPACK: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/lib/libRlapack.dylib

locale:  
[1] en\_US.UTF-8/en\_US.UTF-8/en\_US.UTF-8/C/en\_US.UTF-8/en\_US.UTF-8

attached base packages:  
[1] parallel stats graphics grDevices utils datasets methods base

other attached packages:  
[1] doParallel\_1.0.17 iterators\_1.0.14 foreach\_1.5.2 RColorBrewer\_1.1-3 treemap\_2.4-4 rgdal\_1.6-7  
[7] boot\_1.3-28.1 sandwich\_3.0-2 modelsummary\_1.4.1 data.table\_1.14.8 BMA\_3.18.17 rrcov\_1.7-4  
[13] inline\_0.3.19 robustbase\_0.99-0 leaps\_3.1 survival\_3.5-3 BMS\_0.3.5 caret\_6.0-94  
[19] lattice\_0.20-45 ipred\_0.9-14 rpart.plot\_3.1.1 rpart\_4.1.19 dplyr\_1.1.4 rsample\_1.1.1  
[25] ds4psy\_0.9.0 forcats\_1.0.0 htmlwidgets\_1.6.2 pals\_1.7 plyr\_1.8.8 plyr\_1.8.8 car\_3.1-2  
[31] carData\_3.0-5 ggpmisc\_0.5.3 ggpp\_0.5.3 countrycode\_1.5.0 fixest\_0.11.1 plm\_2.6-3  
[37] xlsx\_0.6.5 plotly\_4.10.2 DataVisualizations\_1.3.1 raster\_3.6-23 sp\_2.0-0 sf\_1.0-14  
[43] eurostat\_3.8.2 modi\_0.1.2 EconGeo\_2.0 cowplot\_1.1.1 flextable\_0.9.5 officer\_0.6.5  
[49] kableExtra\_1.4.0 knitr\_1.46 stringr\_1.5.0 ICC\_2.4.0 future.apply\_1.11.0  
[55] future\_1.33.0 irr\_0.84.1 lpSolve\_5.6.20 haven\_2.5.3 reshape2\_1.4.4 glmnet\_4.1-7  
[61] Matrix\_1.6-1.1 ggplot2\_3.5.0

loaded via a namespace (and not attached):  
[1] SparseM\_1.81 ModelMetrics\_1.2.2.2 maxLik\_1.5-2 ragg\_1.2.5 tidyr\_1.3.0 hardhat\_1.3.0  
[7] generics\_0.1.3 terra\_1.7-39 proxy\_0.4-27 tzdb\_0.4.0 xml2\_1.3.5 lubridate\_1.9.2  
[13] httpuv\_1.6.11 assertthat\_0.2.1 gower\_1.0.1 xfun\_0.43 hms\_1.1.3 rJava\_1.0-6  
[19] evaluate\_0.21 promises\_1.2.0.1 DEoptimR\_1.1-1 fansi\_1.0.6 readxl\_1.4.3 igrapht\_1.5.1  
[25] DBI\_1.2.2 reshape\_0.8.9 stats4\_4.2.3 purrr\_1.0.1 ellipsis\_0.3.2 backports\_1.4.1  
[31] fontLiberation\_0.1.0 insight\_0.19.3 gridBase\_0.4-7 fontBitstreamVera\_0.1.1 vctrs\_0.6.5 quantreg\_5.96  
[37] here\_1.0.1 abind\_1.4-5 withr\_3.0.0 collapse\_1.9.6 bdsmatrix\_1.3-6 checkmate\_2.2.0  
[43] svglite\_2.1.1 pacman\_0.5.1 lazyeval\_0.2.2 crayon\_1.5.2 crul\_1.4.0 recipes\_1.0.6  
[49] pkgconfig\_2.0.3 labeling\_0.4.3 units\_0.8-2 nlme\_3.1-162 nnet\_7.3-18 rlang\_1.1.3  
[55] globals\_0.16.2 lifecycle\_1.0.4 MatrixModels\_0.5-2 fontquiver\_0.2.1 httpcode\_0.3.0 dichromat\_2.0-0.1  
[61] cellranger\_1.1.0 rprojroot\_2.0.3 lmtest\_0.9-40 datawizard\_0.8.0 ISOweek\_0.6-2 zoo\_1.8-12  
[67] viridisLite\_0.4.2 parameters\_0.21.1 KernSmooth\_2.23-20 pROC\_1.18.4 shape\_1.4.6 classInt\_0.4-9  
[73] parallelly\_1.36.0 readr\_2.1.5 scales\_1.3.0 magrittr\_2.0.3 bibtex\_0.5.1 compiler\_4.2.3  
[79] RefManageR\_1.4.0 miscTools\_0.6-28 cli\_3.6.2 unkn\_0.8.0 listenv\_0.9.0 Formula\_1.2-5  
[85] mgcv\_1.8-42 MASS\_7.3-58.2 tidyselect\_1.2.1 stringi\_1.7.12 textshaping\_0.3.6 askpass\_1.1  
[91] grid\_4.2.3 polynom\_1.4-1 tools\_4.2.3 timechange\_0.2.0 rstudioapi\_0.15.0 uuid\_1.1-0  
[97] tables\_0.9.17 prodlim\_2023.03.31 farver\_2.1.1 digest\_0.6.33 shiny\_1.7.4.1 lava\_1.7.2.1  
[103] gfonts\_0.2.0 Rcpp\_1.0.12 performance\_0.10.4 later\_1.3.1 httr\_1.4.7 gdtools\_0.3.7  
[109] Rdpack\_2.4 colorspace\_2.1-0 splines\_4.2.3 confintr\_1.0.2 xlsxjars\_0.6.1 mapproj\_1.2.11  
[115] systemfonts\_1.0.4 xtable\_1.8-4 dreamerr\_1.2.3 jsonlite\_1.8.8 timeDate\_4022.108 R6\_2.5.1  
[121] pillar\_1.9.0 htmltools\_0.5.5 mime\_0.12 glue\_1.7.0 fastmap\_1.1.1 class\_7.3-21  
[127] regions\_0.1.8 codetools\_0.2-19 maps\_3.4.1 pcaPP\_2.0-3 mvtnorm\_1.2-2 furr\_0.3.1  
[133] utf8\_1.2.4 tibble\_3.2.1 numDeriv\_2016.8-1.1 curl\_5.2.0 zip\_2.3.0 openssl\_2.1.0  
[139] rmarkdown\_2.23 munsell\_0.5.0 e1071\_1.7-13 gtable\_0.3.4 bayestestR\_0.13.1 rbibutils\_2.2.14