



The impact of human development on fertility rates

An empirical analysis of fertility rates among the regions of Europe

Marcell Granát

April 12, 2021

Contents

Introduction	3
Theoretical consideration	3
Data	5
Human Development	6
A decent standard of living	7
Long and healthy life	7
Knowledge	7
Total fertility rates, family benefit and unemployment statistics	9
Explore the data	9
Pairwise comparison	9
Regression tree	9
Model building	13
Framework I: with unemployment	13
Methodology of model estimating	17
Interpreting the results	19
Framework II: without unemployment	19
Interpreting the results	19
Conclusion	21
Appendix: R codes	23

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 weak 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

1	Indicators of similarity between the Human Development Indices provided by UNDP and GDL	6
2	Indicators of similarity between the income component of the Human Development Indices provided by GDL and the estimation based on regional GDP	7
3	Indicators of similarity between the knowledge component of Human Development Indices provided by UNDP and the calculated principal components using educational attainment level	8
4	Models	13
5	Comparison of average values of the variables for incomplete and complete observations (Framework I)	16
6	Comparison of average values of the variables for incomplete and complete observations (Framework II)	19

List of Figures

1	Main steps of the study.	4
2	Negative correlation between gross domestic product and fertility rates based on nation level observations (2017)	5
3	HDI and its components based on the dataset from Global Data Lab (2017)	6
4	PCAs and the explained variance	8
5	Pairwise correlation among TFR and calculated human development indices	10
6	Pairwise correlation among TFR and youth unemployment rates	11
7	Correlation between TFR and family benefits	12
8	Regression tree explaining the TFR excluding youth unemployment rates ($cp = 0.02$)	12
9	Regression tree explaining the TFR using all the mentioned explanatory variables ($cp = 0.01$)	14
10	Panel models on the total fertility rates	15
11	Number of complete observations by countries when the model includes or excludes unemployment statistics	16
12	Performance of lasso regression models with different parameters	17
13	Estimated coefficient of the fixed panel model controlling for youth unemployment indicators	18
14	Estimated coefficient of the fixed panel model omitting youth unemployment indicators	20

Introduction

Total fertility rates have decreased significantly over the past few decades, and in most of the developed countries, they are far below the minimum level, which would ensure the reproducibility of the population (2.1 children/woman). The motivation why this demographic pattern frequently becomes the focus of scientific and political discussions is that it has a radical impact on the economy of the future.

The one which must be mentioned is that low fertility is one of the key determinants that lead a country to a high old-age dependency ratio. In extreme cases, today's poor childbearing tendency may put huge economic weight on the future working generation. An interesting fact is that this risk is typically faced by developed economies. Moreover, historical data show a clear correlation between the increasing trend of human development and falling childbearing willingness. The first research question I aim to answer in this study is the following: Does human development truly contribute to the fall in fertility rates, and if it does, are the three components (healthy life; access to knowledge; decent standard of living) included with the same weight?

The related literature contains many uncertainties about this question. Not even the existence of the relation, but the correct direction of the effect is also unclear. An important milestone is a piece of empirical evidence that the association between development and fertility might turn from negative to positive above a given threshold of development [Myrskylä et al., 2009].

Detecting the trend of falling fertility rates is not new. It is well documented that this issue was a key question in politics in many countries after the second world war, and so is nowadays. Governments around the world have already implemented several policies to increase fertility: Abortion ban, campaigns and various financial incentives. But the effectiveness (and their social benefit in many other important aspects) of these are often questioned. Since only the last one is frequently applied (and can be observed accurately) in the developed countries, my second research question is whether spending on family support increases the childbearing willingness significantly.

The contribution of the current paper to the literature is that it applies longitudinal econometrics on NUTS-2 level regional annual observations. Eurostat is the main source of the data. I use fixed effect panel regression to estimate the effect of the mentioned predictors on fertility. 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. Figure 1 visualize the main steps of the current paper.

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 weak evidence for the significant effect of expenditure on family on total fertility rates on the long run*.

To ensure fully reproducibility I enclose the used codes in the Appendix, but these are also available on the following GitHub repository: <https://github.com/MarcellGranat/fertilityEU>.

Theoretical consideration

To give a detailed answer to our research question I introduce the discussed indicators and the possible causal mechanisms among them.

Political debates and unprofessional public discourses usually focus on the total number of livebirth per year. If the target is to describe the effect of demographic changes on the economy, this perspective is usually right. However, the situation is different when the focus is on the causal mechanisms that affect the level of childbearing willingness. The reason is that the total number of births relies explicitly on the demographic trends in the past. This leads to the problem that the number of births decreases if the number of women at childbearing age reduces. Measuring the efficiency of family support expenditures and demographic policies requires eliminating these effects. The total fertility rate (TFR) is a suitable indicator to manage this.

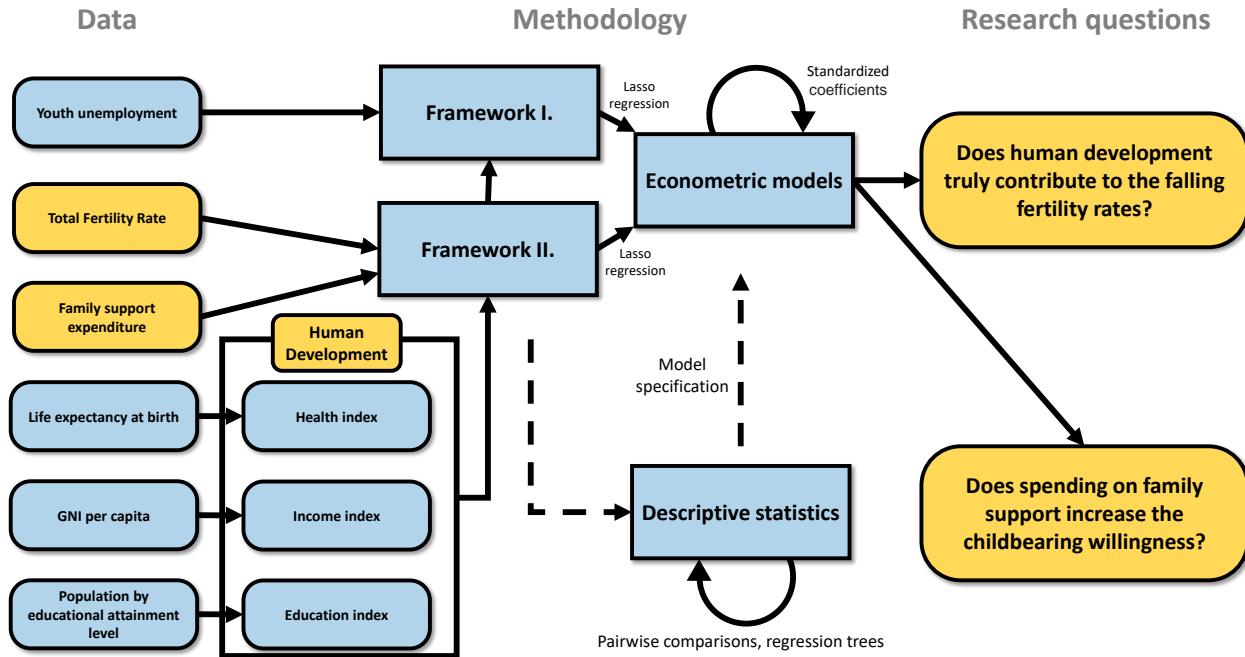


Figure 1: Main steps of the study.

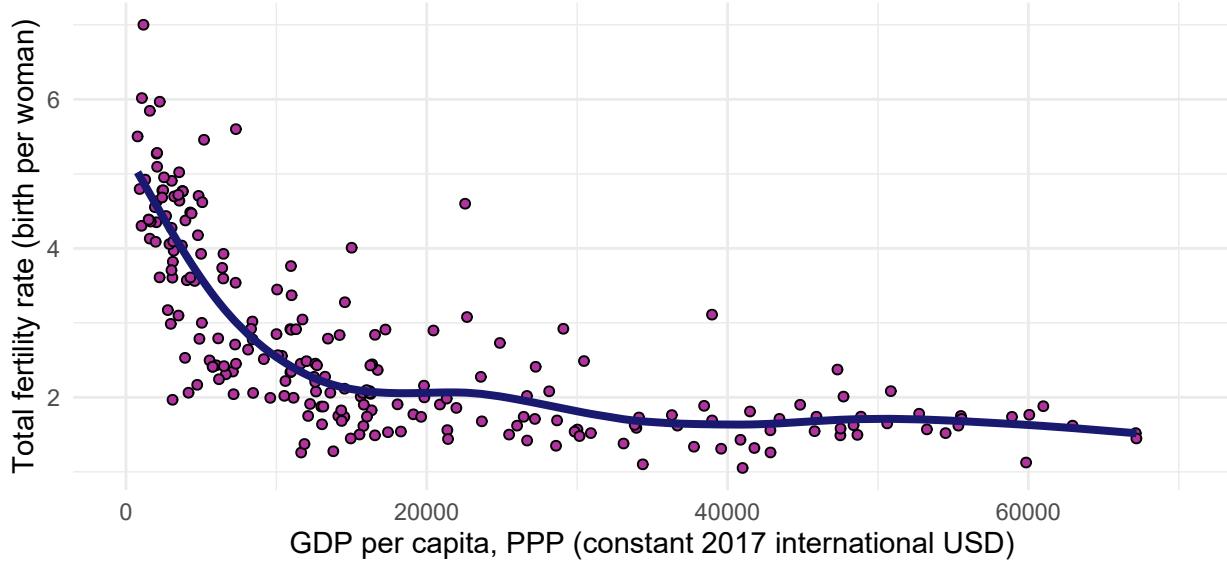
Total fertility rate (TFR) is “the average number of children born per woman over a lifetime given current age-specific fertility rates and assuming no female mortality during reproductive years. TFRs are computed as the sum of age-specific fertility rates defined over five-year intervals” [oec, 2018]. This calculation method ensures that TFR is insensitive to the existing demographic characteristics (robust to the change in the number of childbearing-aged women).

Total fertility rates have dropped significantly in the last decades. More precisely, in the OECD countries, the average has decreased from 2.8 children to 1.7 in the last 50 years [oec, 2019]. Another well-known fact is that fertility is lower in the more developed countries. This correlation is visualized in figure 2.

Figure 1 shows a clear pattern between income and fertility, but national incomes correlate with hundreds of indicators which can be the true reason for failing fertility. The presented relationship was an unquestionable rule of demography a few decades ago, but the empirical evidence of the “J-shaped” fertility trend changed this [eco, 2009]. Myrskyl et al. (2009) reported in their study that the relationship between development and fertility is reversing. The authors explain this with the innovation in family behavior and government policies that improve the compatibility between economic success and family life. This argument confirms the relevance of investigating the combined effect of human development and family support expenditures on fertility.

It is important to note that the mentioned study refers to the Human Development Index, which is a composite index of three key dimensions of human development. “Human development encompasses more than just economic development. The concomitant construction of the HDI offered a simple, yet multidimensional approach to comparatively evaluate the human development of various countries.” [Sagar and Najam, 1998]

To describe the motivation behind the decomposition of human development I refer to the findings of the closely related literature. A sizeable number of studies focus on the effect of education, health, income and fertility. The range of methodologies is wide: several studies use cross-sectional or longitudinal national-level indicators, while others refer to data from questionnaires (Generations- and Gender Surveys is a frequently used source). This paper belongs to the first group of studies but extends the research with regional level observations.



Own editing based on the Figure 5-2. from Kreiszné Hudák (2019).

The trend is drawn via splines.

Source of the data: World Bank.

Figure 2: Negative correlation between gross domestic product and fertility rates based on nation level observations (2017)

Harttgen and Vollmer (2014) write that pairwise investigation among the components of human development confirms that development leads to higher fertility in countries where the HDI is higher than 0.86. The authors used country-level data and they highlighted that proving the robustness of their findings with a different sample is needed. The current study performs this.

As mentioned previously, Myrskl et al. (2009) argue that innovations related to feasible work-life balance may be the reason for the increasing fertility in highly developed countries. Taking that into account, I extend the set of explanatory variables with employment statistics. The effect of unemployment on fertility is already investigated and proved using the observations of OECD nations [Adser, 2004].

The perception of rapidly declining fertility in Europe during deteriorated labor market conditions provides further evidence for the relevance of unemployment in the model framework [Matysiak et al., 2020]. Since the former findings state that the working and financial condition affects young people more sensitively [Frejka et al., 2016], I focus on the effect of youth unemployment rates on fertility.

Data

Eurostat is an abundant and reliable source of NUTS-2 level statistics, so it is a reasonable choice to use¹. However, human development is a hardly available indicator. “The Human Development Index (HDI) is a summary measure of achievements in three key dimensions of human development: a long and healthy life, access to knowledge and a decent standard of living.” [UNITED NATIONS DEVELOPMENT PROGRAMME, 2020] The disadvantage of the HDI and its components defined by the United Nations Development Programme is that they are reported officially only on the national level. To ensure comparability with studies that used national-level HDI, I aim to obtain or calculate equivalent indicators.

¹Technical note: To ensure easy reproducibility data was download with the dedicated Eurostat R package [Lahti et al., 2017].

Human Development

A possible source for this target is the database provided by the Global Data Lab. [Global Data Lab - Innovative Instruments for Turning Data into Knowledge] This website and the data are created by the Institute of Management Research at Radboud University. Their reported values for the national level are almost fully equivalent to the ones published by UNDP. Table 1 figures the similarity between the two datasets.

The reason why this dataset does not fully meet the requirements to use it in this research is that the territorial units do not completely fit with the ones reported by Eurostat, so merging the two data sources is only possible with a high rate of mismatching observations. Although I do not use these data for model building, the matching observations are useful as benchmark points for calculating the index based on the available data from the Eurostat database. This usability is confirmed by its similarity to the officially reported national-level dataset by UNDP.

Table 1: Indicators of similarity between the Human Development Indices provided by UNDP and GDL

Indicator	Value
R^2	99.76%
Spearman R^2	99.71%
Mean absolute deviation	0.007
Mean absolute percentage deviation	1.20%



Figure 3: HDI and its components based on the dataset from Global Data Lab (2017)

A decent standard of living

To determine the standard of living dimension of the development, the gross national income (GNI) per capita would be required in PPP terms (constant 2017 PPP\$) according to the technical note of UNDP [UNITED NATIONS DEVELOPMENT PROGRAMME, 2020]. To transform it into an index, the income is put into the following equation:

$$\text{Income index} = \frac{\ln(\text{GNI}) - \ln(100)}{\ln(75,000) - \ln(100)} \quad (1)$$

This normalization process keeps the value between 0 (at 100 dollars annual income) and 1 (at 75,000 dollars annual income). Eurostat reports GNI in PPP terms (constant 2020 PPP) in euros for regional observations. To eliminate this difference, I divide the given GNI value by the corresponding annualized value of the EUR/USD exchange rate (similarly downloaded from the Eurostat database). Table 2 figures the similarity to the values reported by Global Data Lab for the matching observations (where the territorial unit corresponds with the one used by Eurostat). In this aspect, I would like to highlight the very high Spearman R^2 between them. It shows that it hardly ever happens that a region has a different rank in the two datasets.

Table 2: Indicators of similarity between the income component of the Human Development Indices provided by GDL and the estimation based on regional GDP

Indicator	Value
R^2	92.20%
Spearman R^2	94.18%
Mean absolute deviation	0.0520
Mean absolute percentage deviation	6.62%

Long and healthy life

The health dimension of HDI is determined based on the life expectancy at birth. The logic of the variable transformation is the same as previously written at the income index:

$$\text{Health index} = \frac{\text{life expectancy at birth} - 20}{85 - 20} \quad (2)$$

The natural minimum of life expectancy is at 20 years, while the maximum is estimated at 85 years, and with the given formula, the range of health index is also set between 0 and 1. Since life expectancy is reported for NUTS 2 levels by Eurostat, this transformation can be performed without any additional calculation.

Knowledge

Determining the education index is more complex compared to the previous ones. The required expected years of schooling are not available at a regional level. However, Eurostat reports several pieces of information on the topic of education. In this study, I substitute the original formula with an index based on the available variables. For this purpose, I use data about the population by educational attainment level and NUTS 2 regions. The available data are also grouped by age. Considering that the discussed decline in fertility rates is mainly related to women under the age of 30 [oec, 2019], statistics about younger generations seem more relevant. To manage this, I only use the age class 20-24 and 25-34.

The educational dimension is based on the International Standard Classification of Education (ISCED 2011). This means the following categories: Less than primary, primary and lower secondary education belongs to

DATA

levels 0-2, upper secondary and post-secondary non-tertiary education is labeled with levels 3 and 4 and tertiary education with levels 5-8.

These two dimensions lead to eight possible variables, and the target is to define one index from their values. To manage this dimensionality reduction, I use principal component analysis (PCA). PCA is a useful tool to identify independent sources of variance, and in many cases, economic interpretation can be found in the loadings [Maddala and Lahiri, 1992]. For this step, I impute missing values based on auxiliary regressions using mice R package, and I normalize the variables. Figure 4 shows the result.

The next step is to find the relevant principle component or components. In addition to the interpretability, I also focus on how well this principle describes the education index reported by Global Data Lab. Table 3 contains the similarity indicators.

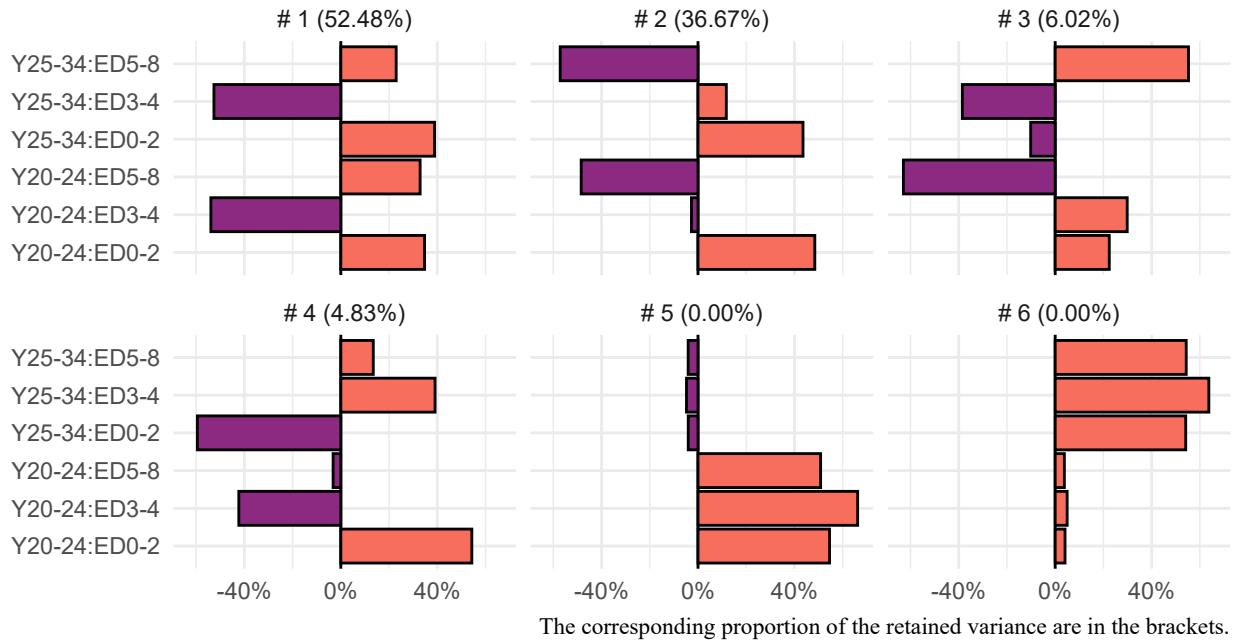


Figure 4: PCAs and the explained variance

The next step is to find the relevant principle component or components. In addition to the interpretability, I also focus on how well this principle describes the education index reported by Global Data Lab. Table 3 contains the similarity indicators.

Table 3: Indicators of similiarity between the knowledge component of Human Development Indices provided by UNDP and the calculated principal components using educational attainment level

Indicator	Comp 1	Comp 2	Comp 3	Comp 4	Comp 5	Comp 6
R^2	5.73%	31.18%	17.31%	12.33%	0.00%	0.05%
Spearman R^2	2.72%	26.35%	16.43%	11.23%	6.59%	0.86%

Since the R^2 and Spearman R^2 are high between the second principal component and the education index from Global Data Lab, and the component has good interpretability, I use the second principal component in the following. Based on the loading, its value increases if the proportion of lower educated people increases and decrease if the share of highly educated people increases. This tells the opposite of what the education index means. For this purpose, I determine the education index as the normalized value of the given PCA

score multiplied by minus one.

$$\text{Education index} = \frac{-\text{PCA score} + |\min(-\text{PCA score})|}{\max(-\text{PCA score})} \quad (3)$$

Total fertility rates, family benefit and unemployment statistics

Family benefit data can be found in the Eurostat database. It is reported on the national level, and in this paper, I choose its unit as a percentage of gross national income. The national-level family support corresponds to all regions of a country. The highly unbalanced distribution of family support within countries can lead to incorrect results. However, better statistics are not available at present. In this paper, I systematically refer to this indicator (expenditure on family/children benefits as a percentage of gross national income) as a family benefit or family support.

Youth unemployment rate and fertility statistics are reported for NUTS-2 levels by Eurostat. To perform this empirical analysis, the transformation of these variables is not required.

Explore the data

This section contains an exploratory analysis of the calculated indices and used variables to describe the easily identifiable patterns before the model building. For this purpose, I report distributions, pairwise comparisons and regression trees about the variables.

Pairwise comparison

Figure 5 shows the relationship between human development indices and total fertility rates. All the pairwise linear correlation coefficients are positive and significant, so it is reasonable to assume that spurious correlations exist in this framework. (An index can explain a significant proportion of the variance of TFR, but the indices also explain each other. As a consequence, the correct mechanism has to be identified with a multivariate model).

Similarly, correlations among youth unemployment statistics and fertility are shown in figure 6. Unemployment rates show a negative correlation with TFR, but the pattern here is noisy as well. In contrast, the unemployment statistics seem to move strongly together, foreshadowing the existence of high multicollinearity in our regression.

Unlike the previous relations, the share of family support in gross national income shows a clear pattern with the fertility rate. This is visualized in figure 7. The value of the linear correlation coefficient between the two indicators is 0.51, and one percentage point increase in the support goes with a 0.138 children/woman increase in the TFR. T-statistic is 16.0 in this bivariate regression, which indicates high statistical significance. Since values of family benefit have a wide range in the sample (from 0 to 4.5), the mentioned slope of the regression line shows significance from the economic aspect as well.

Regression tree

A regression tree is a useful statistical tool, where the process behind the regression is the sorting of the observations into as homogenous groups as possible concerning the response variable (total fertility rates in our case) [James et al., 2013]. This model does not have any longitudinal characteristics, and is by far not sufficient to answer the research question², but it can visualize and describe much descriptive statistical information with simple interpretations. To find the optimal cut-points, I use the CART algorithm.

The first step is to determine the set of predictors. At this point, it is important to note again that unemployment rates are available only for a critical short period in most cases. That being the case, I decide

²Hence I do not focus on tuning the hyperparameter of the models. I set the complexity parameter to return a well-visualisable number of nodes.

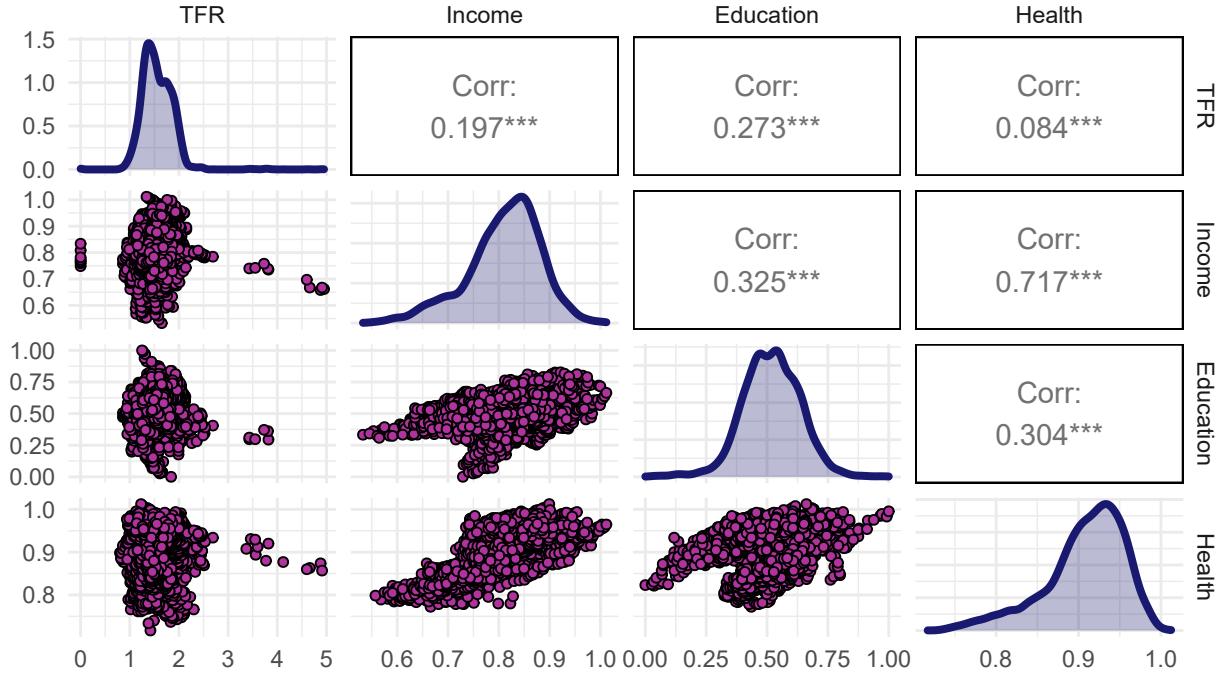


Figure 5: Pairwise correlation among TFR and calculated human development indices

to apply two different frameworks: one without unemployment, and one including it. The model contains each of the observations regardless of the time dimension. The result of the first frame is shown in figure 8.

Including every datapoint leads to 3773 observations (see the top of the tree). Based on the regression tree, the first logical statement made to generate homogenous groups is that if the share of family benefit is below 1.85. This sorts the observations into two subgroups. The TFR in the group where the statement holds (the family benefit is higher) is higher by 0.2 children/woman on average. The second statement is whether the education index is higher than 0.51 (this generates again two similar-sized subgroups), and fertility is higher in the higher educated subgroup. The third pivotal question is also about family benefit. In the subgroup where the family support as a percentage of GDP is higher than 3.4 percentage, the average TFR is 1.8, and 1.5 in the subgroup where this statement does not hold. At this node, the model separates a group of homogeneous datapoints: where the health index is above 0.9, the TFR is extremely high. These few observations come from France.

These lead to the following interpretation: (1) Even if TFR is lower in low educated regions, there are observations with high family benefit rate and high fertility among these regions. (2) Childbearing willingness tends to increase related to a higher share of family support as a percentage of GDP, but under the mentioned circumstances, one can observe extremely high TFR data with modest family benefits. Integrating youth unemployment statics into the modeling framework decreases the number of complete observations significantly. The extended regression tree contains only 1039 data points. On the other hand, variable importance indicators assign youth unemployment as an important explanation to the variance of fertility among the regions. 49 percentage of prediction error reduction is related to the family benefit, but 16 percentage to the unemployment among young people aged 15 to 19.

The effect of youth unemployment is interesting in this regression tree. TFR is higher in the category where the unemployment rate among young people is lower at node 2, but at node 3 the conclusion is the opposite. The cutting-points has two main differences: observations, where the family benefit is below 2.1, belongs to node 2, and that cutting-point is based on the unemployment rate among young people aged 15-19, while node 3 belongs to the data points with family support above 2.1, and the model finds youth unemployment

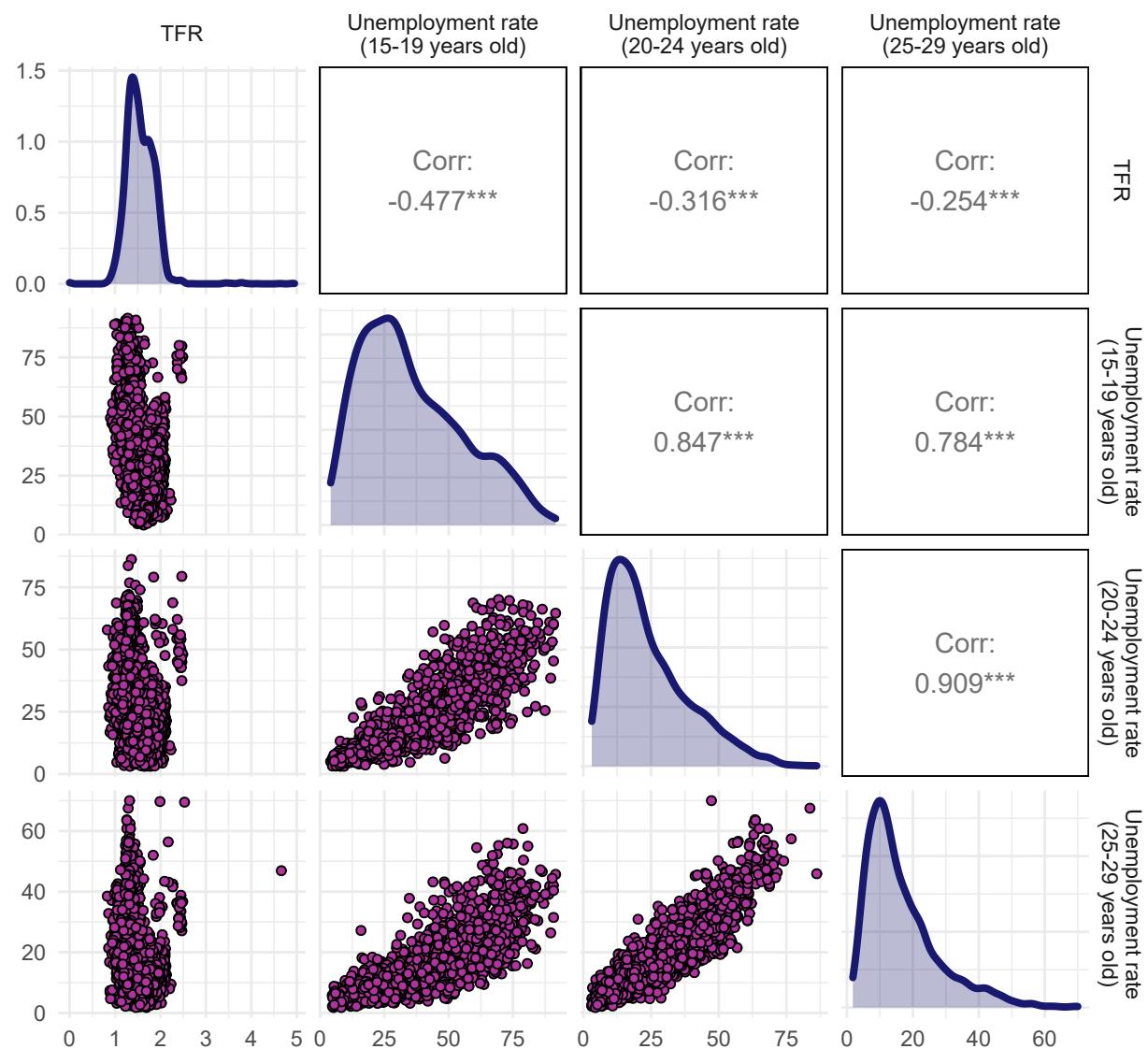


Figure 6: Pairwise correlation among TFR and youth unemployment rates

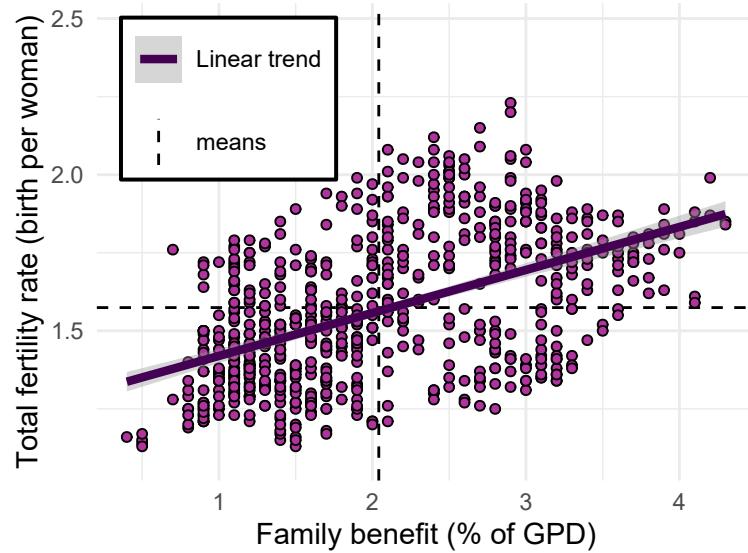


Figure 7: Correlation between TFR and family benefits

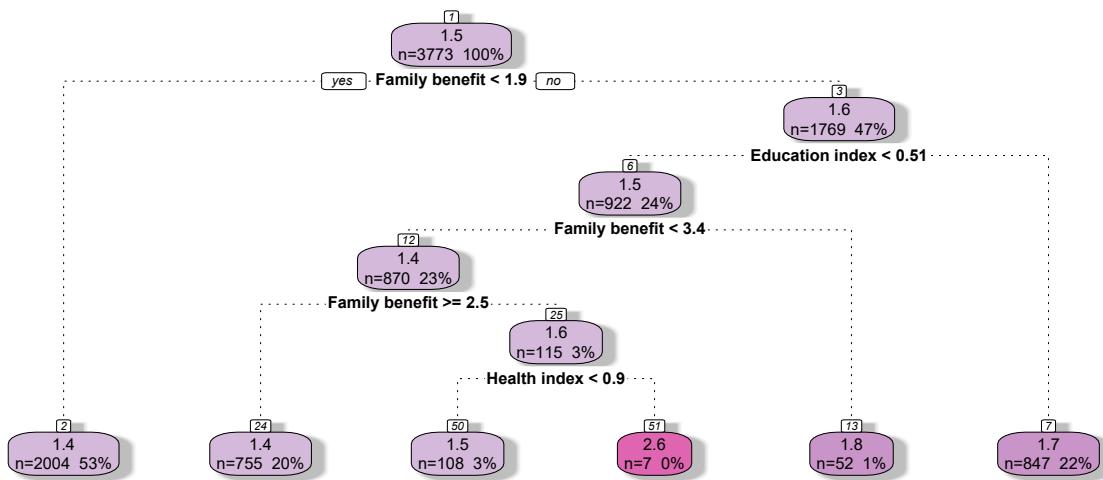


Figure 8: Regression tree explaining the TFR excluding youth unemployment rates ($cp = 0.02$)

rate among 20-24 years olds more important. Therefore, the difference in the direction of the effect can be explained by the amount of family support or the fact, that the effect-mechanism of unemployment differs in the two age categories.

Based on the literature and theoretical considerations, the first alternative is more plausible. In countries where the family support is higher, young people are probably less sensitive to their working situation, when they decide about childbearing. A high value of family support-to-GDP ratio can eliminate the aforementioned mechanism, that some young parents postpone their childbearing because they cannot afford it. Based on these findings, I extend the design-matrix with the interactions of youth unemployment rates and family benefits.

Model building

Following the findings in the previous chapters, I highlight the main ideas for the model building: (1) Youth unemployment is a significant factor in explaining the variance of fertility, but including it in the model reduces the number of complete observations, so I report one model with the youth unemployment included and one without that. (2) Extending the design matrix with the interactions of youth unemployment rates and fertility is justified. (This extension only concerns the framework which contains unemployment).

Additionally, concerning the nature of fertility (duration of pregnancy) and childbearing decisions (parents may interpolate their expected socio-economic situation from their past) extending the model with the lagged value of the predictors is also reasonable. This raises the issue that the model contains too many predictors and variable selection becomes difficult. I manage variable selection based on lasso selection.

The first step to estimate panel regression models is to identify the appropriate model type. Choosing among the pooled, within and random-effects model requires performing the Chow test and Hausman test. The null hypothesis in the former one is that whether a significant difference among the individual intercepts exists, and the test is performed based on an F-test. If the given H_0 cannot be rejected at any standard significance level, estimating a pooled model is suggested. In other cases, the outcome depends on the result of the Hausman test. The Hausman test (or Durbin-Wu-Hausman test) is more complex from mathematical aspects, but the interpretation is simple: if the given H_0 cannot be rejected, the within (fixed-effects) model is not as efficient as the random-effects, and estimating random-effects model is suggested. If the null hypothesis is rejected at any standard significance level, the fixed-effects model will be suitable [Wooldridge, 2016].

The optimal model may differ if the set of predictors change, but as a starting point, I estimate the one-one simplest model for the two mentioned model frameworks (including unemployment statistics or not). These models do not contain interaction but lagged effect and quadratic terms are included. The performed tests show that the fixed-effects model is optimal (see the results in Table 5), and the estimated coefficients from the fixed-effects models are presented in figure 10.

Table 4: Models

Indicator	Model I.	Model II.
Pooltest	0.00%	0.00%
Phtest	0.00%	0.00%
Adjusted R^2	20.19%	29.24%
Observations	699	3257

In the following, I perform lasso variable selection on these model frameworks to find the most relevant effects.

Framework I: with unemployment

As repeatedly mentioned in the previous section, including unemployment statistics in the regression analysis drastically reduces the number of complete observations. Moreover, there is a good reason to believe that

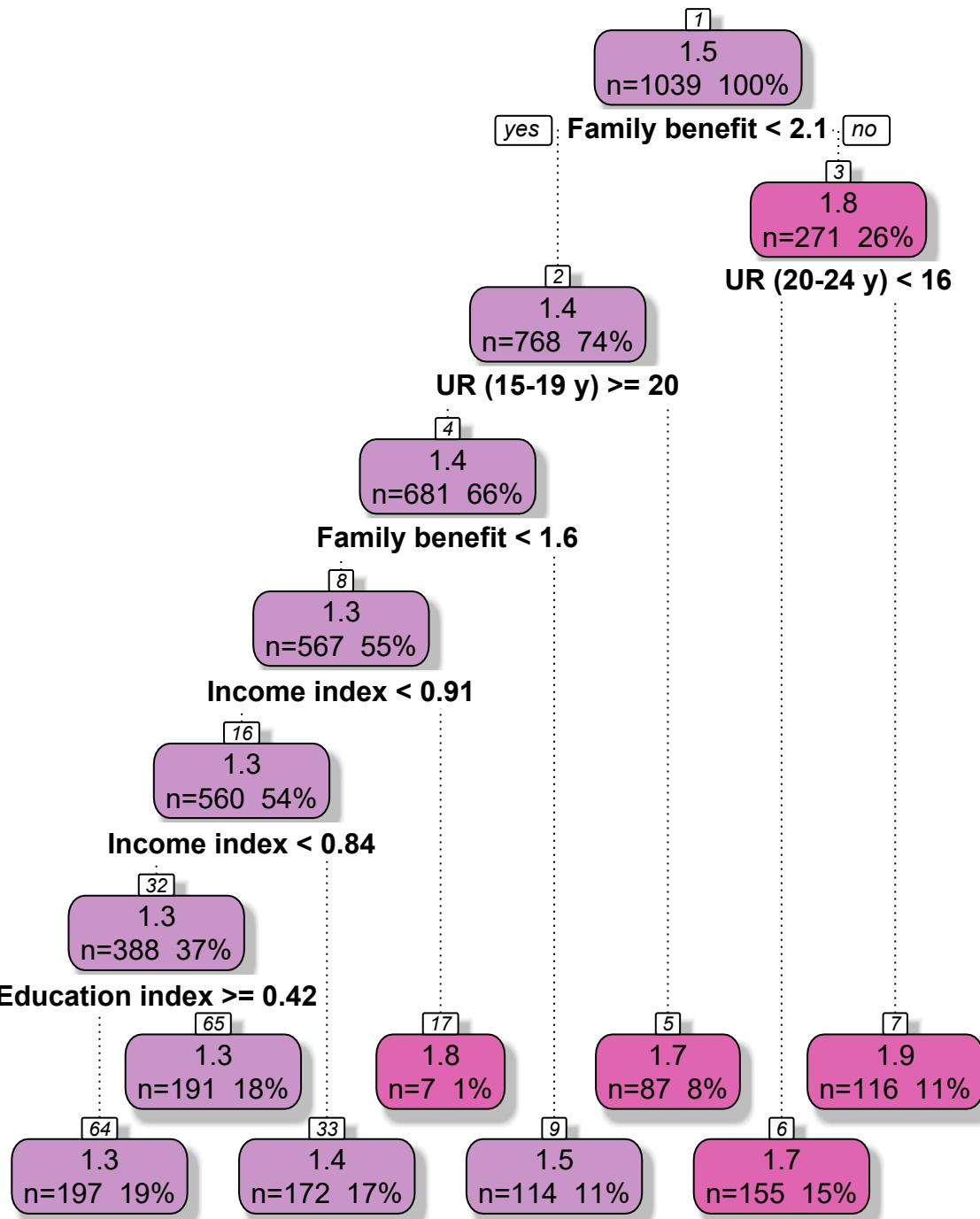


Figure 9: Regression tree explaining the TFR using all the mentioned explanatory variables (cp = 0.01)

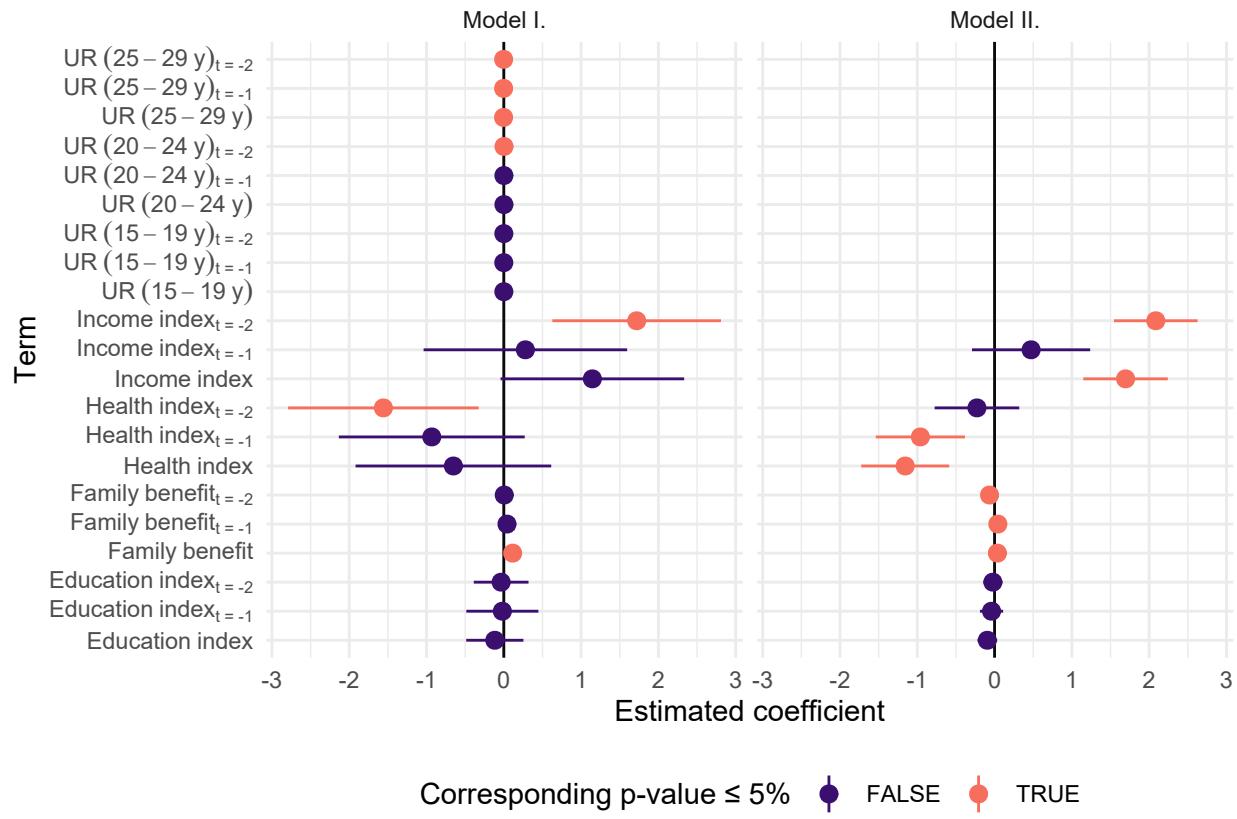


Figure 10: Panel models on the total fertility rates

this selection method is not random, and omitting the incomplete data points may cause bias in the results. Figure 10 shows the number of complete observations by territorial units related to the two frameworks.

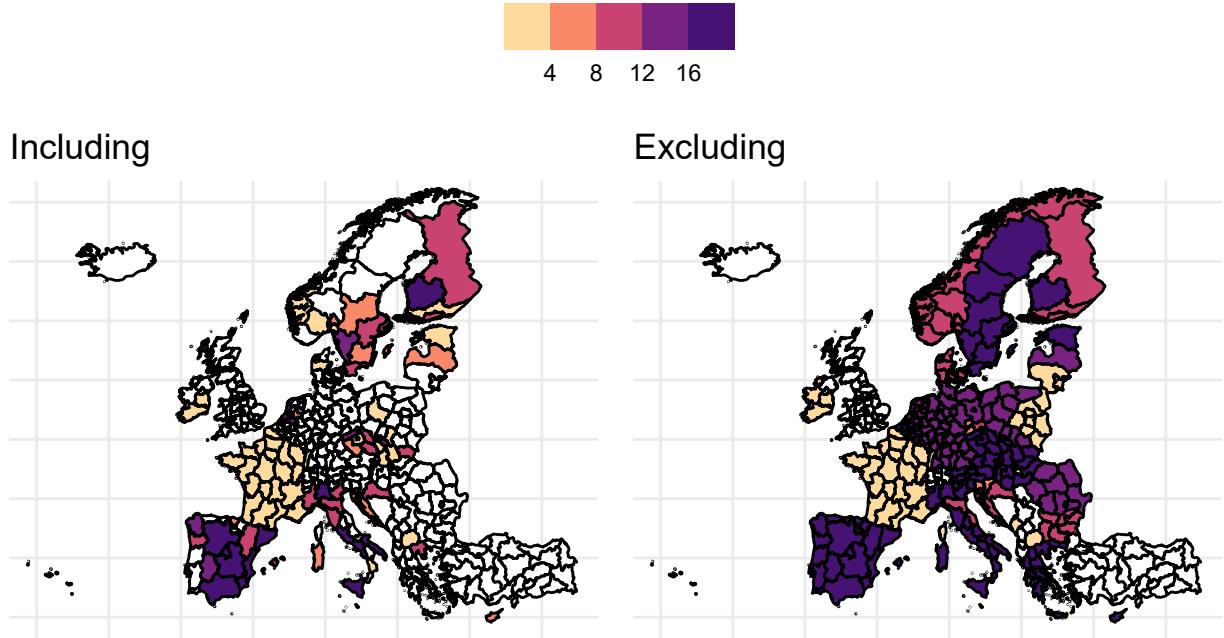


Figure 11: Number of complete observations by countries when the model includes or excludes unemployment statistics

Figure 10 reveals that there are unbalances among the regions concerning their representation in the dataset. A prominent amount of incomplete and thus unusable observations concerns the Central European region and the UK. In this paper, I do not impute these missing values, and therefore, the limited generalizability of the results is taken into account.

To measure the bias, I calculate the mean of the given variables in the total sample (including incomplete observations), and the sample that is used in this framework (excluding incomplete observations). The result is reported in table 5. The table describes that the health index and income index are higher, while the family benefit is lower in the used sample compared to the dataset containing the incomplete observations as well.

Table 5: Comparison of average values of the variables for incomplete and complete observations (Framework I)

Variable	Mean in total sample	Mean in used sample	Number of observations in the total sample
TFR	1.5575	1.4801	7249
Education index	0.5152	0.5088	5373
Health index	0.9089	0.9375	6983
Income index	0.8120	0.8298	4291
Family benefit	2.0759	1.6302	6155
UR (15-19 y)	36.6883	44.4313	1713
UR (20-24 y)	24.3994	25.5202	2935
UR (25-29 y)	17.0107	17.0332	2664

Methodology of model estimating

As described above, an optimal statistical tool for a high-dimensional dataset is the lasso regression. From mathematical aspect lasso regression means to add a $\lambda \sum_{j=1}^p |\beta_j|$ term to the target function of regression [James et al., 2013]. Intuitively, with this additional term, the value of the target function is lower (which has to be minimized), if more parameters are equal to zero, but the prediction error does not decrease significantly.

Lasso regression has a hyperparameter λ . Setting its value to zero leads to the unmodified OLS model. In contrast, a lasso regression with $\lambda = 1$ would lead to an empty model. Finding the optimal λ hyperparameter requires estimating the model with different parameters. In each case, the model contains a different number of variables. To determine the optimal value of λ , leave-one-out cross-validation is performed with 10 folds, then the model having the lowest mean squared error on the validation set is chosen. This process is visualized in figure 11.

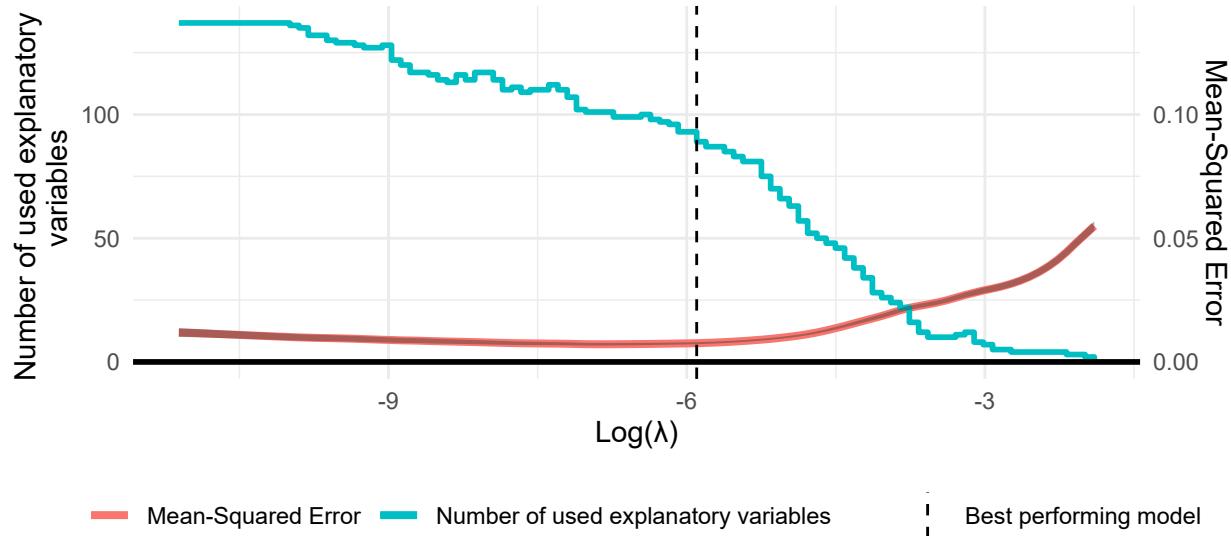


Figure 12: Performance of lasso regression models with different parameters

It is important to note that the above-described algorithm eliminated the insignificant individual intercept terms from the model. To adjust this property, I reestimate the within model including all the individual intercept terms and predictor variables from the best-performing lasso regression model.

In addition, the difference in the measurement of the variables causes complexity in interpretation. Interpreting the direct effect of a predictor is simple, but the coefficients are not sufficient to describe the importance of the variables, because they are measured on different scales (one percentage point change in the family benefit-to-GDP ratio would be extreme, but not as outstanding as one percentage point change in the youth unemployment rate). To manage this, I reestimate the model with standardized variables. The benefit of this model is that the explained variance of the regression model can be decomposed with it:

$$R^2 = \sum_{j=1}^p r_j \times \beta_{standardized,j} \quad (4)$$

Based on equation 4, the coefficient multiplied by the linear correlation coefficient (r) can be interpreted as the contribution to the explained variance. The result of this and the above-mentioned computations are reported in figure 12.

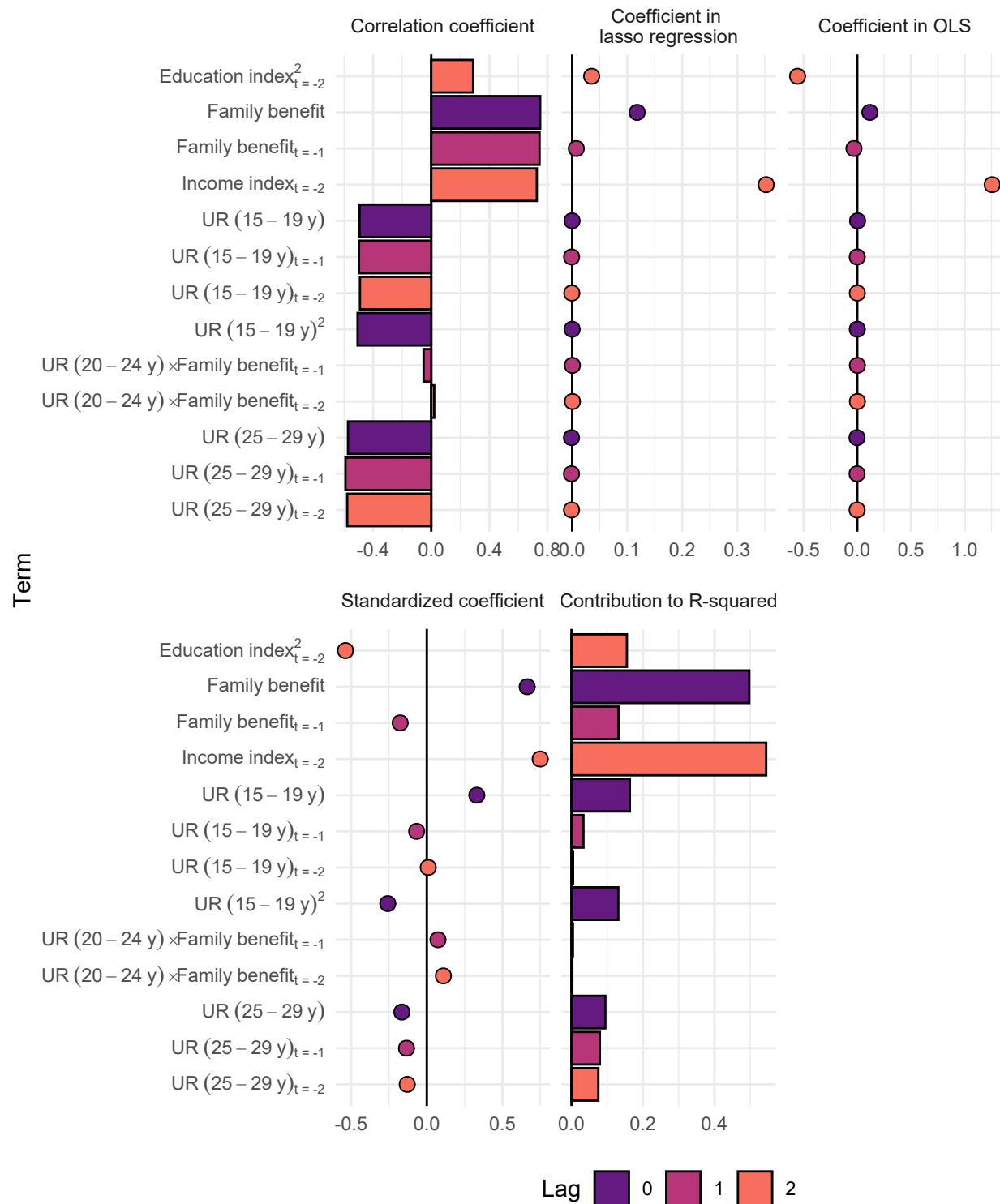


Figure 13: Estimated coefficient of the fixed panel model controlling for youth unemployment indicators

Interpreting the results

The high contribution of income index and family benefit to the explained variance revealed by figure 12. Both of them seem to have a positive effect on fertility. Increasing income per capita leads to higher total fertility rate based on the model parameters (positive coefficient correspond to each lagged variable). In contrast, coefficients related to the different lagged values of family support indicate a complex mechanism: family support has a high instantaneous effect on fertility, but the lagged negative effect implies that this birth-surplus disappears in the following years.

Based on the results the answer to our research question is that in the developed world (1) *income is far the most important component of human development influencing fertility (highest contribution to R²)* and (2) *family support also has a significant instantaneous effect on childbearing willingness, but it seems weaker on the long run.*

Youth unemployment rates among 25-29 year-olds have a negative effect on fertility, but its total effect is lagged. Among 15-19 year-olds this effect is different. In their case, the increase in unemployment causes an instantaneous increase in fertility. But in the case of a permanent increase in unemployment, this increase disappears (*ceteris paribus*). Extending the model with the interactions of youth unemployment and family benefit was truly beneficial comparing the standardized effect of the unemployment rates and the interactions.

Education index is also detected as a significant explanation of the variance of fertility, but its interpretation is more complex. The lagged quadratic terms are represented in the model with negative coefficients. This reflects that the higher education index leads to lower fertility, and this effect is stronger in the case of those regions, where the education index is higher. This confirms the former findings in the literature that highly educated women tend to have less children [Martin, 1995], but recent research found empirical evidence, that “higher educated German women, who already decided to have a child despite their high opportunity costs are more family oriented” [König, 2011].

Framework II: without unemployment

I continue my study reporting the results from the model excludes unemployment statistics. This framework omits significantly fewer data points, so the probability of contra selection-caused bias is reduced. Table 6 describes the average difference comparing the used sample and the values from the incomplete observations.

The methodology of the model estimation is equivalent to the one described in the first model framework. The results are presented in figure 15. The H₀ of the Chow test ($p = 0.00$) and the Hausman test ($p = 0.00$) are rejected, so fixed effect model is adequate. The R² of this model is 14.85%³.

Table 6: Comparison of average values of the variables for incomplete and complete observations (Framework II)

Variable	Mean in total sample	Mean in used sample	Number of observations in total sample
TFR	1.5575	1.4935	7249
Education index	0.5152	0.4971	5373
Health index	0.9089	0.9202	6983
Income index	0.8120	0.8228	4291
Family benefit	2.0759	1.9699	6155

Interpreting the results

Figure 13 shows that the outstandingly high contribution to the R² of the income index did not change, so the answer for my first research question is robust to the framework: *income index explains significantly more*

³The heterogeneity of the countries containing complete observation is higher in this setup, that is comparing the R² of the two frameworks is not suggested.

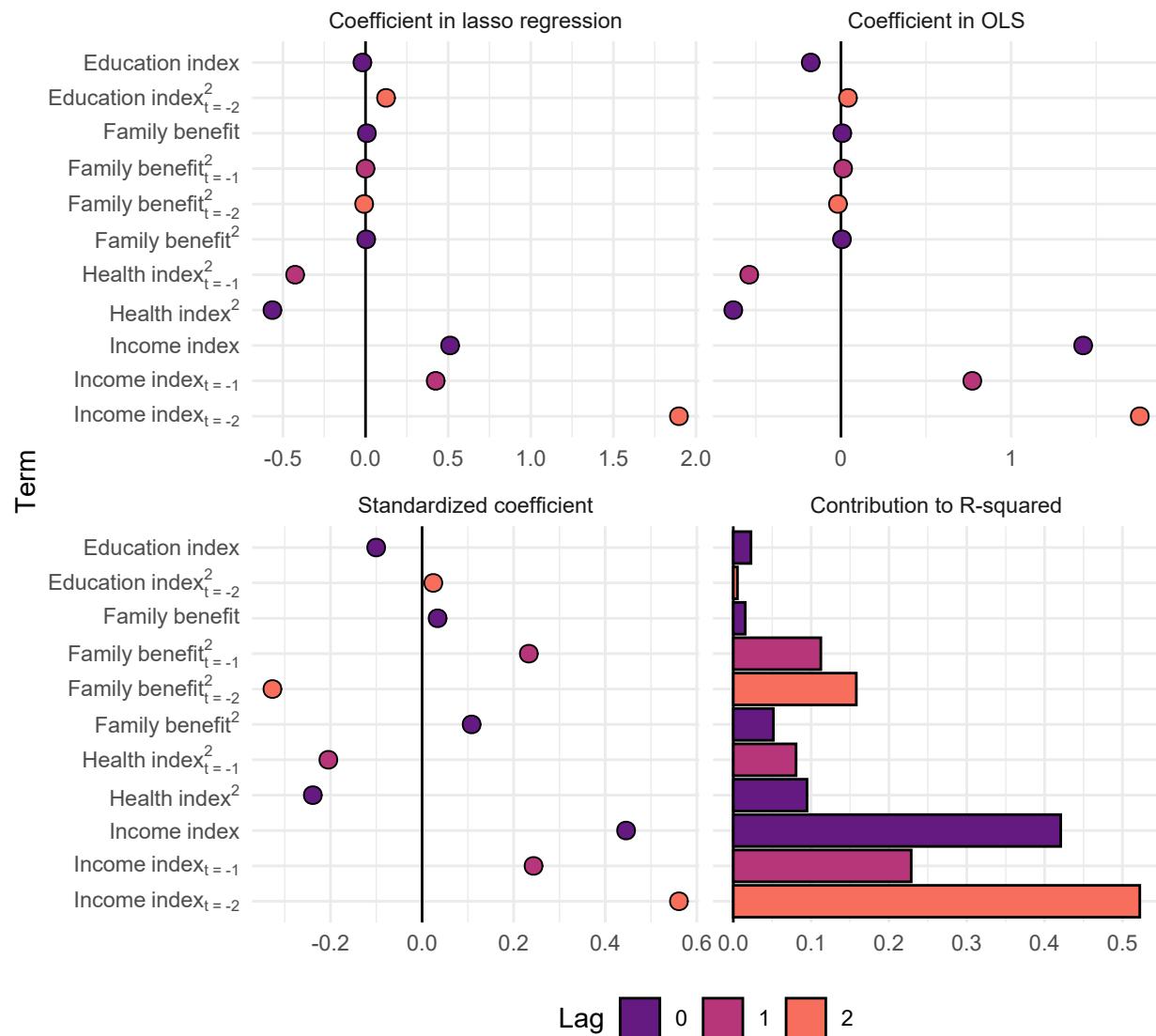


Figure 14: Estimated coefficient of the fixed panel model omitting youth unemployment indicators

CONCLUSION

of the variance of the fertility than the other components of human development, and it has a positive effect on childbearing willingness.

The estimated effect of family benefit differs in this framework compared to the one including unemployment statistics (and excluding observations where youth unemployment is not available). The new observations come from Central-European regions and the estimated structure of the effect of family benefit became significantly different. The reversal effect in this framework is close to the instantaneous effect of the family spending. This leads to the interpretation that family support has only an instantaneous impact on fertility, but on the long run it can not significantly influence the fertility.

The main difference between the results reported by the two frameworks is that the health index is signed as a significant variable in the second one. Since all of its transformed terms has a negative coefficient, a higher health index leads to lower fertility. The frequently mentioned reason for this is the changing lifestyle and women in the EU are having their first child later. One possible explanation why health index was not relevant in the previous framework is that the share of observations from Central-European countries are much higher in this model (due to the lack of unemployment statistics from the early 2000s). Extension of average childbearing age led to failing fertility rates in this area [Berde and Németh, 2014]. But many articles suggest that an adjusted TFR should be considered [Bongaarts and Feeney, 1998], because the drastically low fertility during the time of this mechanism. However, these indices are currently not available for regional datasets, but this could be a possible further research direction.

Conclusion

This paper focused on the effect of two key factors on fertility rates in Europe: human development and expenditures on family support. To quantify the relationships, econometric modeling on regional dataset were performed. This requires to generate or estimate the human development indices for regional level. From statistical aspect the key issue is the trade off between limited complete observations or omitting possibly important explanatory variables. This is due to the lack of data about youth unemployment. To manage this a I set up two model frameworks: one including unemployment statistics and one excluding it. Based on theoretical considerations and the previously performed analysis the design-matrix was extended with lagged variables (duration of pregnancy), quadratic terms (Theory of “J-shaped” effect) and interactions between youth unemployment and family benefit (based on the previously performed regression tree). This lead to a high-dimensional dataset, so lasso-based feature selection was the bases of the reported longitudinal econometric models.

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 framework. In contrast, the research come up *only with weak evidence for the significant effect of expenditure on family on total fertility rates*.

References

- [eco, 2009] (2009). The best of all possible worlds? page 1.
- [oec, 2018] (2018). *SF2.1 Fertility rates*. OECD.
- [oec, 2019] (2019). *Society at a Glance 2019*. OECD.
- [Adser, 2004] Adser, A. (2004). Changing fertility rates in developed countries. the impact of labor market institutions. *Journal of Population Economics*, 17(1):17–43.
- [Adsera, 2004] Adsera, A. (2004). Changing fertility rates in developed countries. the impact of labor market institutions. *Journal of population economics*, 17(1):17–43.
- [Berde and Németh, 2014] Berde, É. and Németh, P. (2014). Az alacsony magyarországi termékenység új megközelítésben. *Statisztikai Szemle*, 92(3):253–274.
- [Bongaarts and Feeney, 1998] Bongaarts, J. and Feeney, G. (1998). On the quantum and tempo of fertility. *Population and development review*, pages 271–291.
- [Frejka et al., 2016] Frejka, T., Gietel-Basten, S., Abolina, L., Abuladze, L., Aksyonova, S., Akrap, A., Antipova, E., Bobic, M., Čipin, I., Fakheyeva, L., et al. (2016). Fertility and family policies in central and eastern europe after 1990. *Comparative Population Studies*.
- [James et al., 2013] James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An introduction to statistical learning*, volume 112. Springer.
- [König, 2011] König, S. (2011). Higher order births in germany and hungary: Comparing fertility intentions in a national context. *Sozialforschung MZES Mannheim*, pages 1–18.
- [Lahti et al., 2017] Lahti, L., Huovari, J., Kainu, M., and Biecek, P. (2017). eurostat r package. Version 3.6.84.
- [Maddala and Lahiri, 1992] Maddala, G. S. and Lahiri, K. (1992). *Introduction to econometrics*, volume 2. Macmillan New York.
- [Martin, 1995] Martin, T. C. (1995). Women’s education and fertility: results from 26 demographic and health surveys. *Studies in family planning*, pages 187–202.
- [Matysiak et al., 2020] Matysiak, A., Sobotka, T., and Vignoli, D. (2020). The great recession and fertility in europe: A sub-national analysis. *European Journal of Population*, 37(1):29–64.
- [Myrskylä et al., 2009] Myrskylä, M., Kohler, H.-P., and Billari, F. C. (2009). Advances in development reverse fertility declines. *Nature*, 460(7256):741–743.
- [Sagar and Najam, 1998] Sagar, A. D. and Najam, A. (1998). The human development index: a critical review. *Ecological economics*, 25(3):249–264.
- [UNITED NATIONS DEVELOPMENT PROGRAMME, 2020] UNITED NATIONS DEVELOPMENT PROGRAMME (2020). Technical notes: Calculating the human development indices.
- [Wooldridge, 2016] Wooldridge, J. M. (2016). *Introductory econometrics: A modern approach*. Nelson Education.

Appendix: R codes

```
1 # Set up -----
2
3 ## Packages =====
4
5 library(tidyverse)
6 library(patchwork)
7 library(knitr)
8 library(broom)
9 library(eurostat)
10
11 ## Gg theme =====
12
13 update_geom_defaults("point", list(fill = "#B1339E",
14                               shape = 21,
15                               color = "black",
16                               size = 1.4))
17 update_geom_defaults("line",
18                      list(color = "midnightblue", size = 1.4))
19
20 update_geom_defaults("smooth", list(color = "red4", size = 1.4))
21
22 update_geom_defaults("density",
23                      list(color = "midnightblue", fill = "midnightblue", alpha = .3,
24                           size = 1.4))
25
26 extrafont::loadfonts(device="win")
27
28 theme_set(theme_minimal() + theme(
29   legend.direction = "vertical",
30   # text = element_text(family = "Impact"),
31   plot.caption = element_text(family = "serif"),
32   legend.key=element_blank()
33 ))
34
35 # https://data.worldbank.org/indicator/SP.DYN.TFRT.IN
36
37 WB_fertility <- read_csv("WB_fertility.csv", skip = 4)
38
39 # https://data.worldbank.org/indicator/NY.GDP.PCAP.PP.KD
40
41 WB_GDP <- read_csv("WB_GDP.csv", skip = 4)
42
43 merge(WB_fertility %>%
44       select('Country Name', '2017') %>%
45       rename(tfr = '2017'),
46       WB_GDP %>%
47       select('Country Name', '2017') %>%
48       rename(GDP = '2017')) %>%
49 ggplot(aes(GDP, tfr)) + geom_point() +
50   ggfformula::geom_spline() +
51   scale_x_continuous(limits = c(0, 7e+4)) +
52   labs(y = "Total fertility rate (birth per woman)",
```

```
53     x = "GDP per capita, PPP (constant 2017 international USD)",  
54     caption = "Own editing based on the Figure 5-2. from Kreiszné Hudák (2019).  
55     The trend is drawn via splines.  
56     Source of the data: World Bank."  
57 )  
58  
59 # Global Data Lab =====  
60  
61 #### Nation-level data #####  
62  
63 GDL_nat <- read_csv("GDL-Sub-national-HDI-data.csv") %>%  
64   filter(Level == "National") %>%  
65   select(Country, 6:34) %>% pivot_longer(-1, names_to = "time", values_to = "values") %>%  
66   mutate(time = as.numeric(time)) %>%  
67   na.omit()  
68  
69 #### Data from UNDP #####  
70  
71 # source: http://hdr.undp.org/  
72  
73 HDI_UNDP <- read_csv("Human Development Index (HDI).csv",  
74   skip = 5) %>%  
75   select(!starts_with("X"), - 'HDI Rank') %>%  
76   mutate_at(-1,  
77     function(x) {as.numeric(ifelse(x == "...", NA, x))}) %>%  
78   pivot_longer(-1, names_to = "time", values_to = "values") %>%  
79   mutate(time = as.numeric(time)) %>%  
80   na.omit()  
81  
82 #### Comparison ---> Table 1 #####  
83  
84 merge(GDL_nat %>% rename(GDL = values),  
85       HDI_UNDP %>% rename(UNDP = values)) %>%  
86 {  
87   c(  
88     scales::percent(cor(x = .\$GDL, y = .\$UNDP)^2, accuracy = .01),  
89     scales::percent(cor(x = .\$GDL, y = .\$UNDP, method = "spearman")^2, accuracy = .01),  
90     as.character(format(mean(abs(.\$GDL - .\$UNDP)), digits = 1)),  
91     scales::percent(mean(abs(.\$GDL - .\$UNDP)/.\$UNDP), accuracy = .01)  
92   )  
93 } %>%  
94 {tibble(  
95   Indicator = c("$R^2$", "Spearman $R^2$",
96                 "Mean absolute deviation", "Mean absolute percentage deviation"),
97   Value = .
98 )} %>%  
99 kable(  
100   caption = "Indicators of similarity between the Human Development Indices  
101   provided by UNDP and GDL"
102 )  
103  
104 #### Sub-national data #####
```

```
106 # source of csv files: https://globaldatalab.org/
107
108 GDL_import <- function(x) {
109   get_eurostat_geospatial(nuts_level = 2) %>%
110   data.frame() %>%
111   tibble %>%
112   mutate(
113     ISO_Code = countrycode::countrycode(CNTR_CODE, origin = "iso2c", "iso3c"),
114     ISO_Code = ifelse(CNTR_CODE == "UK", "GBR", ISO_Code),
115     ISO_Code = ifelse(CNTR_CODE == "EL", "GRC", ISO_Code),
116   ) %>%
117   select(ISO_Code, NUTS_NAME, geo) %>%
118   merge(read_csv(x), by = "ISO_Code") %>%
119   mutate(
120     z = stringdist::stringsim(NUTS_NAME, Region)
121   ) %>%
122   arrange(desc(z)) %>%
123   filter(!duplicated(Region)) %>%
124   filter(!duplicated(NUTS_NAME)) %>%
125   filter((z > .5 | Country %in% c("Greece", "Turkey", 'Romania',
126                                     'Malta', 'Italy')) & NUTS_NAME != "Dresden") %>%
127   select(geo, '1990':'2018') %>%
128   pivot_longer(-1, names_to = "time", values_to = "values") %>%
129   mutate(time = as.numeric(time))
130 }
131
132 GDL_subnat <- GDL_import("GDL-Sub-national-HDI-data.csv") %>% rename(HDI = values) %>%
133   merge(
134     GDL_import("GDL-Educational-index--data.csv") %>% rename(education = values)
135   ) %>%
136   merge(
137     GDL_import("GDL-Health-index-data.csv") %>% rename(health = values)
138   ) %>%
139   merge(
140     GDL_import("GDL-Income-index-data.csv") %>% rename(income = values)
141   )
142
143 plot_NUTS2 <- function(df, viridis_c = T, ..., all.x = F) {
144   p <- df %>%
145     {merge(eurostat::get_eurostat_geospatial(nuts_level = 2), ., all.x = all.x)} %>%
146     ggplot(aes(fill = values)) +
147     geom_sf(color = "black") +
148     theme(
149       axis.text = element_blank()
150     ) +
151     xlim(c(-30, 44)) +
152     ylim(c(35, 70)) +
153     labs(fill = NULL)
154
155   if (viridis_c) {
156     p <- p + scale_fill_viridis_c(option = "magma", ...,
157                                   guide = guide_colorbar(frame.colour = "black",
158                                             ticks.colour = "black")),
159   }
160 }
```

```
159                               na.value = "white")
160 }
161 p
162 }
163
164 GDL_subnat %>%
165   filter(time == 2017) %>%
166   select(-time) %>%
167   pivot_longer(-1, names_to = "var", values_to = "values") %>%
168   filter(!is.na(var)) %>%
169   mutate(
170     var = str_to_title(var),
171     var = str_replace(var, "Hdi", "HDI"),
172     var = factor(var,
173       levels = c("HDI", "Income", "Education", "Health"),
174       ordered = T)
175   ) %>%
176   plot_NUTS2() + facet_wrap(~ var, ncol = 2) +
177   labs(caption = "Source: https://globaldatalab.org") +
178   scale_fill_viridis_c(guide = guide_colorsteps(), option = 'magma')
179
180 GDP_index <- get_eurostat("nama_10r_2gdp", time_format = "num") %>%
181   filter(unit == "PPS_EU27_2020_HAB") %>% # Purchasing power standard (PPS) per inhabitant
182   select(-unit) %>%
183   rename(GDP = values) %>% merge(
184     get_eurostat("ert_bil_eur_a", time_format = "num") %>% # EUR/USD annual avg exc r
185       filter(currency == "USD" & statinfo == "AVG") %>%
186       select(time, e = values)
187   ) %>%
188 {
189   GDP <- .\$GDP/.\$e # mutate to USD
190   mutate(.,
191     GDPindex = (log(GDP) - log(100))/
192       (log(75000) - log(100)))
193   )
194 }
195
196 merge(GDL_subnat, GDP_index) %>%
197 {
198   c(
199     scales::percent(cor(x = .\$income, y = .\$GDPindex)^2, accuracy = .01),
200     scales::percent(cor(x = .\$income, y = .\$GDPindex, method = "spearman")^2,
201       accuracy = .01),
202     as.character(format(mean(abs(.\$income - .\$GDPindex)), digits = 1, nsmall = 4)),
203     scales::percent(mean(abs(.\$income - .\$GDPindex)/.\$GDPindex), accuracy = .01)
204   )
205 } %>%
206 {tibble(
207   Indicator = c("$R^2$", "Spearman $R^2$",
208                 "Mean absolute deviation", "Mean absolute percentage deviation"),
209   Value = .
210 )} %>%
211 kable()
```

```
212     caption = "Indicators of similarity between the income component of the  
213     Human Development Indices provided by GDL and the estimation based on regional GDP"  
214 )  
215  
216 health_index <- get_eurostat("demo_r_mlifexp", time_format = "num") %>%  
217   filter(age == "Y_LT1" & sex == "T") %>%  
218   select(geo, time, le = values) %>%  
219   mutate(  
220     health_index = (le - 20) / (85 - 20)  
221   )  
222  
223 edu_wide <- get_eurostat("edat_lfse_04", time_format = "num") %>%  
224   filter(sex == "T" & !str_detect(isced11, "GEN") &  
225     !str_detect(isced11, "VOC") & isced11 != "ED3-8" & age != 'Y25-64' &  
226     age != 'Y30-34' & !str_detect(geo, "TR")  
227 ) %>%  
228   mutate(  
229     var = str_c(age, ":", isced11)  
230   ) %>%  
231   select(geo, time, var, values) %>%  
232   pivot_wider(names_from = var, values_from = values) %>%  
233 {  
234   x <- .  
235   names(x) <- letters[1:length(x)]  
236   x <- cbind(  
237     .[, 1:2],  
238     mice::complete(mice::mice(select(x, -a,-b), printFlag = F))  
239   )  
240   names(x) <- names(.)  
241   x  
242 }  
243  
244 edu_comps <- edu_wide%>%  
245   select(-time, -geo) %>%  
246   na.omit() %>%  
247   {princomp(scale(.))}  
248  
249 edu_comp_vars <- edu_comps %>%  
250   summary() %>%  
251   {.$sdev^2/sum(.\$sdev^2)} %>%  
252   scales::percent(accuracy = .01) %>%  
253   {str_c("# ", 1:length(.), " (", ., ")")}  
254  
255 edu_comps %>%  
256   .$loadings %>%  
257   unclass() %>%  
258   data.frame() %>%  
259   rownames_to_column() %>%  
260   pivot_longer(-1) %>%  
261   mutate(  
262     name = as.numeric(str_remove(name, 'Comp.')),  
263   ) %%  
264   arrange(name) %>%
```

```

265   mutate(
266     name = edu_comps %>%
267       summary() %>%
268       {$.sdev^2/sum($.sdev^2)} %>%
269       scales::percent(accuracy = .01) %>%
270       {str_c("# ", 1:length(.), " (", ., ")")} %>%
271       .[name],
272     rowname = str_replace_all(rowname, '_ ', '-')
273   ) %>%
274   ggplot +
275   aes(rowname, value, fill = value < 0) +
276   geom_hline(yintercept = 0) +
277   geom_col(color = 'black') +
278   coord_flip() +
279   scale_fill_viridis_d(guide = F, option = "magma", begin = .4,
280                         end = .7, direction = -1) +
281   scale_y_continuous(labels = scales::percent) +
282   facet_wrap(~name, ncol = 3) +
283   labs(x = NULL, y = NULL, caption =
284         "The corresponding proportion of the retained variance are in the brackets.")
285
286 edu_comps %>% .$scores %>%
287   cbind(edu_wide) %>% merge(GDL_subnat) %>%
288   select(starts_with('Comp'), education) %>%
289   {
290     x <- .$education
291     apply(select(., -education), 2, function(y) {
292       c(
293         scales::percent(cor(x = x, y = y)^2, accuracy = .01),
294         scales::percent(cor(x = x, y = y, method = "spearman")^2, accuracy = .01)
295       )
296     })
297   } %>%
298   data.frame() %>%
299   mutate(Indicator = c("$R^2$", "Spearman $R^2$")) %>%
300   rename_all(function(x) str_replace(x, "p.", "p ")) %>%
301   select(Indicator, everything()) %>%
302   kable(
303     caption = "Indicators of similarity between the knowledge component of Human
304     Development Indices provided by UNDP and the calculated principal components using
305     educational attainment level"
306   )
307
308 edu_index <- edu_comps %>%
309   .$scores %>%
310   data.frame() %>%
311   select(2) %>%
312   cbind(edu_wide) %>%
313   select(geo, time, edu_index = Comp.2) %>%
314   mutate(
315     edu_index = -edu_index,
316     edu_index = edu_index + abs(min(edu_index)),
317     edu_index = edu_index/max(edu_index)

```

```

318 )
319
320 # Variables without transformation =====
321
322 f_data <- get_eurostat("demo_r_find2", time_format = "num") %>% # TFR
323   select(geo, time, var = indic_de, values) %>%
324   filter(!str_detect(geo, "TR"))
325
326 FAM_df <- get_eurostat("spr_exp_sum", time_format = "num") %>% # Family benefit
327   filter(spdeps == "FAM" & unit == "PC_GDP") %>%
328   rename(FAM = values, country = geo) %>%
329   select(-(spdeps:unit))
330
331 yth_empl_byage <- get_eurostat("yth_empl_110", time_format = "num") %>%
332   # youth unemployment
333   filter(unit == "PC" & sex == "F") %>%
334   filter(age %in% c("Y15-19", "Y20-24", "Y25-29")) %>%
335   select(-unit, -sex) %>%
336   pivot_wider(names_from = age, values_from = values) %>%
337   rename(
338     "uY15" = "Y15-19",
339     "uY20" = "Y20-24",
340     "uY25" = "Y25-29"
341   )
342
343 # Merging the data.frames =====
344
345 dat <- f_data %>%
346   pivot_wider(names_from = var, values_from = values) %>%
347   merge(edu_index, all = T) %>%
348   merge(health_index, all = T) %>%
349   merge(GDP_index, all = T) %>%
350   filter(!str_detect(geo, "TR")) %>%
351   mutate(country = str_sub(geo, end = 2)) %>%
352   merge(FAM_df, all.x = T, all.y = F) %>%
353   merge(yth_empl_byage, all.x = T, all.y = F) %>%
354   filter(!str_detect(geo, "TR"))
355
356 f.clean_names <- function(v, Tosparse = F) {
357   v <- str_replace_all(v, "GDPindex", "Income index") %>%
358     str_replace_all("health_index", "Health index") %>%
359     str_replace_all("edu_index", "Education index") %>%
360     str_replace_all("TOTFERRT", "TFR") %>%
361     str_replace_all("GDPindex", "Income index") %>%
362     str_replace_all("FAM", "Family benefit") %>%
363     str_replace_all("uY15", "UR (15-19 y)") %>%
364     str_replace_all("uY20", "UR (20-24 y)") %>%
365     str_replace_all("uY25", "UR (25-29 y)") %>%
366     str_replace_all("Fu15", "UR (15-19 y)*Family benefit") %>%
367     str_replace_all("Fu20", "UR (20-24 y)*Family benefit") %>%
368     str_replace_all("Fu25", "UR (25-29 y)*Family benefit")
369
370   if(Tosparse) v <- str_replace_all(v, " ", "~") %>%

```

```
371     str_replace_all('\\\\*', '%\\\\*%')
372   v
373 }
374
375 # Explore the data -----
376
377 # Pairwise correlations =====
378
379 dat %>%
380   filter(str_length(geo) == 4) %>%
381   select(TOTFERRRT, GDPindex, edu_index, health_index) %>%
382   {set_names(., f.clean_names(names(.)))} %>%
383   rename_all(.funs = function(x) str_remove_all(x, ' index')) %>%
384   GGally::ggpairs()
385
386
387 dat %>%
388   filter(str_length(geo) == 4) %>%
389   select(TOTFERRRT, uY15, uY20, uY25) %>%
390   set_names("TFR", "Unemployment rate\n(15-19 years old)",
391             "Unemployment rate\n(20-24 years old)",
392             "Unemployment rate\n(25-29 years old)") %>%
393   GGally::ggpairs()
394
395
396 dat %>%
397   filter(str_length(geo) == 2) %>%
398   {ggplot(., aes(FAM, TOTFERRRT)) +
399    geom_vline(aes(xintercept = mean(.\$FAM, na.rm = T), linetype = "means")) +
400    geom_hline(yintercept = mean(.\$TOTFERRRT, na.rm = T), linetype = 2) +
401    geom_point() +
402    geom_smooth(method = "lm", aes(color = "Linear trend"), size = 1.5) +
403    scale_color_viridis_d() +
404    scale_linetype_manual(values = c("means" = 2)) +
405    labs(y = "Total fertility rate (birth per woman)",
406          x = "Family benefit (% of GPD)",
407          linetype = NULL, color = NULL
408        ) +
409    theme(
410      legend.position = c(.2, .8),
411      legend.box.background = element_rect(
412        colour = "black",
413        size = .7,
414        fill = "white"
415      ),
416      legend.spacing.y = unit(0.05, 'cm')
417    )
418 }
419
420 dat %>%
421   filter(str_length(geo) == 2) %>%
422   select(FAM, TOTFERRRT) %>%
423   na.omit() %>%
424   {print(cor.test(.\$FAM, .\$TOTFERRRT)); broom::tidy(lm(TOTFERRRT ~ FAM, data = .))}
```

```
424  
425 # Regression trees -----  
426  
427 m_part <- dat %>%  
428   filter(str_length(geo) == 4) %>%  
429   select(TOTFERRT, GDPindex, edu_index, health_index, FAM) %>%  
430   {set_names(., f.clean_names(names(.)))} %>%  
431   na.omit() %>%  
432   rpart::rpart(formula = TFR ~ ., cp = .02)  
433  
434 rattle::fancyRpartPlot(m_part, palettes = 'PuRd', sub = NULL)  
435  
436 summary(m_part)  
437  
438  
439 m_part2 <- dat %>%  
440   filter(str_length(geo) == 4) %>%  
441   select(TOTFERRT, GDPindex, edu_index, health_index, uY15, uY20, uY25, FAM) %>%  
442   {set_names(., f.clean_names(names(.)))} %>%  
443   na.omit() %>%  
444   rpart::rpart(formula = TFR ~ ., cp = .01)  
445  
446 rattle::fancyRpartPlot(m_part2, palettes = 'PuRd', sub = NULL)  
447  
448 summary(m_part2)  
449  
450 # Model building -----  
451  
452 # Transform the data for panel modeling -----  
453  
454 dat_plm <- dat %>%  
455   select(  
456     geo, time, TOTFERRT, edu_index, health_index, GDPindex, FAM, uY15, uY20, uY25  
457   ) %>%  
458   mutate( # TODO new element  
459     Fu15 = FAM*uY15,  
460     Fu20 = FAM*uY20,  
461     Fu25 = FAM*uY25  
462   ) %>% # TODO end of new  
463   filter(str_length(geo) == 4 & !is.na(TOTFERRT))  
464  
465 dat_plm <- dat_plm %>%  
466   select(-TOTFERRT) %>%  
467   mutate(time = time + 1) %>%  
468   {  
469     set_names(., ifelse(names(.) == 'geo' | names(.) == 'time', names(.),  
470                 paste0(names(.), '_1')))  
471   } %>%  
472   merge(dat_plm, all.x = F, all.y = T)  
473  
474 dat_plm <- dat_plm %>%  
475   select(!ends_with("_1")) %>%  
476   select(-TOTFERRT) %>%  
477   mutate(time = time + 2) %>%
```

```

478   {
479     set_names(., ifelse(names(.) == 'geo' | names(.) == 'time', names(.),
480                 paste0(names(.), '_1_1')))
481   } %>%
482   merge(dat_plm, all.x = F, all.y = T)
483
484 dat_plm <- dat_plm %>%
485   select(-TOTFERRRT) %>%
486   mutate_at(-(1:2), function(x) x^2) %>%
487   {
488     set_names(., ifelse(names(.) == 'geo' | names(.) == 'time', names(.),
489                 paste0(names(.), '_2')))
490   } %>%
491   merge(dat_plm, all.x = F, all.y = T) %>%
492   select(-c("Fu15_1_1_2", "Fu20_1_1_2", "Fu25_1_1_2", "Fu15_1_2", "Fu20_1_2",
493           "Fu25_1_2", "Fu15_2", "Fu20_2", "Fu25_2"))
494
495 ## Initial models =====
496
497 library(plm)
498
499 m_panels <- c(
500   'TOTFERRRT ~ edu_index_1_1 + health_index_1_1 + GDPindex_1_1 + FAM_1_1 +
501   uY15_1_1 + uY20_1_1 + uY25_1_1 + edu_index_1 + health_index_1 + GDPindex_1 + FAM_1 +
502   uY15_1 + uY20_1 + uY25_1 + edu_index + health_index + GDPindex + FAM + uY15 + uY20 +
503   uY25',
504   'TOTFERRRT ~ edu_index_1_1 + health_index_1_1 + GDPindex_1_1 + FAM_1_1 + edu_index_1 +
505   health_index_1 + GDPindex_1 + FAM_1 + edu_index + health_index + GDPindex + FAM'
506 ) %>%
507   lapply(function(formula) {
508     pooling <- plm(eval(formula), data = dat_plm, model = "pooling")
509     within <- plm(eval(formula), data = dat_plm, model = "within")
510     random <- plm(eval(formula), data = dat_plm, model = "random")
511
512     list(
513       tests = c(
514         pooltest(pooling, within)$p.value,
515         phtest(within, random)$p.value,
516         plm::r.squared(within, dfcor = T)),
517       model = within,
518       OLS = formula %>%
519         str_replace_all('\\~', ' ', ' ') %>%
520         str_replace_all('\\+', ' ', ' ') %>%
521         str_split(' ', ' ') %>%
522         .[[1]] %>%
523         trimws() %>%
524         {c(., 'geo')} %>%
525         {select(dat_plm, .)} %>%
526         na.omit() %>%
527         group_by(geo) %>%
528         summarise_all(.funns = function(x) mean(x)) %>%
529         merge(
530           plm::fixef(within) %>%

```

```

531     {tibble(geo = names(.), a = .)}
532   ) %>%
533   mutate(
534     TOTFERRT = TOTFERRT - a
535   ) %>%
536   lm(formula = formula)
537   )
538 })
539
540 ### Plot coefficients #####
541
542 m_panels %>%
543   lapply(function(output) {
544     output$model %>% broom::tidy(conf.int = T) %>%
545       rownames_to_column()
546   }) %>%
547   reduce(rbind) %>%
548   mutate(
549     rowname = paste0("Model ", as.roman(cumsum(rowname == 1)), "."),
550     term = f.clean_names(term, Tosparse = T),
551     term = gsub("_2", "^2", term),
552     term = gsub("_1_1", '["t = -2"]', term),
553     term = gsub("_1", '["t = -1"]', term),
554   ) %>%
555   ggplot() +
556   aes(estimate, term, color = p.value <= .05) +
557   geom_vline(xintercept = 0, color = "gray4") +
558   geom_point() +
559   geom_pointrange(aes(xmin = conf.low, xmax = conf.high)) +
560   facet_wrap(~rowname, nrow = 1) +
561   labs(x = "Estimated coefficient", y = "Term",
562         color = "Corresponding p-value \u2264 5%") +
563   scale_color_viridis_d(option = "magma", begin = .2, end = .7) +
564   scale_y_discrete(labels = scales::parse_format()) +
565   theme(
566     legend.position = "bottom",
567     legend.direction = "horizontal"
568   )
569
570 ### Model descriptions #####
571
572 m_panels %>%
573   lapply(function(output) {
574     c(output$tests, nrow(augment(output$model)))
575   }) %>%
576   reduce(rbind) %>%
577   t() %>%
578   data.frame() %>%
579   mutate_all(function(x) c(scales::percent(x[1:3], accuracy = .01),
580                           as.character(x[4]))) %>%
581   {set_names(., paste0("Model ", as.roman(1:ncol(.)), ".")) } %>%
582   mutate(
583     Indicator = c("Pooltest", "Phtest", "Adjusted $R^2$", "Observations")

```

```
584 ) %>%
585   select(Indicator, everything()) %>%
586   knitr::kable(caption = "Models", align = c("l", "c", "c", "c"))
587
588 theme_update(legend.direction = 'horizontal', legend.position = 'bottom')
589
590 ggpubr::ggarrange(
591   dat_plm %>%
592     na.omit() %>%
593     group_by(geo) %>%
594     summarise(
595       values = n()
596     ) %>%
597     plot_NUTS2(all.x = T) +
598     scale_fill_viridis_b(option = "magma", direction = -1, begin = .2,
599                           na.value = "white") +
600     ggtitle('Including'),
601   dat_plm %>%
602     select(!starts_with('u') & !starts_with('Fu')) %>%
603     na.omit() %>%
604     group_by(geo) %>%
605     summarise(
606       values = n()
607     ) %>%
608     plot_NUTS2(all.x = T) +
609     scale_fill_viridis_b(option = "magma", direction = -1, begin = .2,
610                           na.value = "white") +
611     ggtitle('Excluding'),
612   common.legend = T
613 )
614 # Framework I. =====
615
616 ## Bias of framework #####
617
618 names(dat_plm) %>%
619   { ifelse(
620     . %in% c("geo", "time") | str_detect(., "_1") | str_detect(., "_2") |
621     str_detect(., 'Fu')
622     , NA, .
623   )} %>%
624   na.omit() %>%
625   lapply(function(variable){
626     x <- pull(dat_plm, variable) %>% na.omit()
627     y <- pull(na.omit(dat_plm), variable)
628     tibble(
629       Variable = f.clean_names(variable),
630       'Mean in total sample' = mean(x),
631       'Mean in used sample' = mean(y),
632       'Number of observations in the total sample' = length(x)
633     )
634   }
635 ) %>% reduce(rbind) %>%
636   knitr::kable(caption =
```

```

637     'Comparison of average values of the variables for incomplete
638     and complete observations (Framework I)', digits = 4,
639     align = c('l', 'c', 'c', 'c', 'c'))
640
641 ### Run lasso regression #####
642
643 library(glmnet)
644
645 y <- na.omit(dat_plm)$TOTFERRT
646 X <- model.matrix(TOTFERRT ~ ., data = select(na.omit(dat_plm), -time))
647 LASSO <- cv.glmnet(X, y)
648
649 tidy(LASSO) %>%
650   ggplot() +
651   aes(log(lambda), ymin = conf.low*1000, ymax = conf.high*1000) +
652   geom_line(aes(log(lambda), estimate*1000, color = "Mean-Squared Error")) +
653   geom_step(aes(y = nonzero, color = "Number of used explanatory variables"),
654             size = 1) +
655   geom_ribbon(alpha = .4) +
656   geom_hline(yintercept = 0, size = 1) +
657   geom_vline(aes(xintercept = log(LASSO$lambda.1se), linetype = "Best performing model"),
658             color = 'black') +
659   scale_y_continuous(
660     name = "Number of used explanatory \n variables",
661     sec.axis = sec_axis( trans=~./1000, name = "Mean-Squared Error")
662   ) +
663   labs(y = "Mean-Squared Error", x = "Log(\u03bb)", color = NULL, linetype = NULL) +
664   scale_linetype_manual(values = c(2)) +
665   theme(legend.position = "bottom", legend.direction = "horizontal")
666
667 lasso_coefs <- capture.output(
668   coef(LASSO, LASSO$lambda.1se)
669 ) %>%
670   .[-(1:2)] %>%
671   {tibble(x = .)} %>%
672   mutate(term = gsub(" .*", "", x), coef = gsub(".* ", "", x)) %>%
673   select(-x) %>%
674   filter(!str_detect(term, "geo") & coef != "" & term != "(Intercept)")
675
676 m_panels2 <- paste("TOTFERRT ~", paste(lasso_coefs$term, collapse = " + ")) %>%
677   lapply(function(formula) {
678     pooling <- plm(eval(formula), data = dat_plm, model = "pooling")
679     within <- plm(eval(formula), data = dat_plm, model = "within")
680     random <- plm(eval(formula), data = dat_plm, model = "random")
681     list(
682       tests = c(
683         pooltest(pooling, within)$p.value,
684         phtest(within, random)$p.value,
685         plm::r.squared(within, dfcor = T)),
686       model = within,
687       OLS = formula %>%
688         str_replace_all('\\~', ' ', ') %>%
689         str_replace_all('\\+', ' ', ') %>%
```

```
690 str_split('!',!) %>%
691 .[[1]] %>%
692 trimws() %>%
693 {c(., 'geo')} %>%
694 {select(dat_plm, .)} %>%
695 na.omit() %>%
696 group_by(geo) %>%
697 summarise_all(.funns = function(x) mean(x)) %>%
698 merge(
699   plm::fixef(within) %>%
700   {tibble(geo = names(.), a = .)}
701 ) %>%
702 mutate(
703   TOTFERRT = TOTFERRT - a
704 ) %>%
705 lm(formula = formula)
706 )
707 }
708 )

709
710 m_panels2 %>%
711 lapply(function(output) {
712   output$tests
713 })
714
715 standard_beta <- m_panels2[[1]]$OLS %>%
716 QuantPsyc::lm.beta() %>%
717 {tibble(term = names(.), beta = .)} %>%
718 filter(!str_detect(term, "geo"))

719
720 standard_beta <- augment(m_panels2[[1]]$OLS) %>%
721 select(TOTFERRT:.fitted) %>%
722 select(-.fitted) %>%
723 cor() %>%
724 data.frame() %>%
725 select(1) %>%
726 rownames_to_column() %>%
727 rename(term = rowname) %>%
728 merge(standard_beta) %>%
729 mutate(explain = abs(TOTFERRT*beta)) %>%
730 select(term, cor = TOTFERRT, standard_beta = beta, explain)

731
732 lasso_coefs %>%
733 rename(lasso = coef) %>%
734 merge(tidy(m_panels2[[1]]$OLS, conf.int = T)) %>%
735 merge(standard_beta) %>%
736 mutate_at(-1, function(x) as.numeric(x)) %>%
737 pivot_longer(c(lasso:estimate, cor:explain)) %>%
738 mutate(
739   conf.low = ifelse(name != "estimate", NA, conf.low),
740   conf.high = ifelse(name != "estimate", NA, conf.high),
741   lag = ifelse(str_detect(term, "_1"),
742               ifelse(str_detect(term, "_1_1"), 2, 1), 0),
```

```

743 bar = ifelse(name %in% c("cor", "explain"), value, NA),
744 value = ifelse(!(name %in% c("cor", "explain")), value, NA),
745 term = f.clean_names(term, Tosparse = T),
746 term = gsub("_2", "^2", term),
747 term = gsub("_1_1", '[t = -2]', term),
748 term = gsub("_1", '[t = -1]', term),
749 name = case_when(
750   name == "cor" ~ 'Correlation coefficient',
751   name == "estimate" ~ 'Coefficient in OLS',
752   name == "lasso" ~ 'Coefficient in \nlasso regression',
753   name == "standard_beta" ~ 'Standardized coefficient',
754   name == "explain" ~ 'Contribution to R-squared'
755 ),
756 name = factor(
757   name, levels = c(
758     'Correlation coefficient',
759     'Coefficient in \nlasso regression',
760     'Coefficient in OLS',
761     'Standardized coefficient',
762     'Contribution to R-squared'
763   )
764 ),
765 ) %>%
766 {
767   ggplot(.) +
768     geom_vline(xintercept = 0) +
769     geom_point(aes(value, term, fill = factor(lag)), size = 3) +
770     geom_col(aes(bar, term, fill = factor(lag)), color = "black") +
771     scale_fill_viridis_d(option = "magma", begin = .3, end = .7) +
772     scale_y_discrete(labels=scales::parse_format(), limits = rev) +
773     facet_wrap(~ name, scales = "free_x") +
774     labs(x = NULL, y = "Term", fill = "Lag") +
775     theme(legend.position = "bottom", legend.direction = "horizontal")
776 }
777
778 dat_plm <- dat_plm %>%
779   select(!starts_with("uY") &! starts_with("Fu"))
780
781 names(dat_plm) %>%
782   {ifelse(
783     . %in% c("geo", "time") | str_detect(., "_1") | str_detect(., "_2"), NA, .
784   )} %>%
785   na.omit() %>%
786   lapply(function(variable){
787     x <- pull(dat_plm, variable) %>% na.omit()
788     y <- pull(na.omit(dat_plm), variable)
789     tibble(
790       Variable = f.clean_names(variable),
791       'Mean in total sample' = mean(x),
792       'Mean in used sample' = mean(y),
793       'Number of observations in total sample' = length(x)
794     )
795   })

```

```

796 ) %>% reduce(rbind) %>%
797 knitr::kable(caption = 'Comparison of average values of the variables for incomplete and complete obs',
798 align = c('l', 'c', 'c', 'c'))
799
800 LASSO <- cv.glmnet(model.matrix(TOTFERRT ~ ., data = select(na.omit(dat_plm), -time)),
801 na.omit(dat_plm)$TOTFERRT)
802
803 lasso_coefs <- capture.output(
804   coef(LASSO, LASSO$lambda.1se)
805 ) %>%
806   .[-(1:2)] %>%
807   {tibble(x = .)} %>%
808   mutate(term = gsub(" .*", "", x), coef = gsub(".* ", "", x)) %>%
809   select(-x) %>%
810   filter(!str_detect(term, "geo") & coef != "" & term != "(Intercept)")
811
812 m_panels2 <- paste("TOTFERRT ~", paste(lasso_coefs$term, collapse = " + ")) %>%
813   lapply(function(formula) {
814     pooling <- plm(eval(formula), data = dat_plm, model = "pooling")
815     within <- plm(eval(formula), data = dat_plm, model = "within")
816     random <- plm(eval(formula), data = dat_plm, model = "random")
817     list(
818       tests = c(
819         pooltest(pooling, within)$p.value,
820         phtest(within, random)$p.value,
821         plm::r.squared(within, dfcor = T)),
822       model = within,
823       OLS = formula %>%
824       str_replace_all('\\~', ' ', ' ') %>%
825       str_replace_all('\\+', ' ', ' ') %>%
826       str_split(' ', ' ') %>%
827       .[[1]] %>%
828       trimws() %>%
829       {c(., 'geo')} %>%
830       {select(dat_plm, .)} %>%
831       na.omit() %>%
832       group_by(geo) %>%
833       summarise_all(.funs = function(x) mean(x)) %>%
834       merge(
835         plm::fixef(within) %>%
836         {tibble(geo = names(.), a = .)}
837       ) %>%
838       mutate(
839         TOTFERRT = TOTFERRT - a
840       ) %>%
841       lm(formula = formula)
842     )
843   }
844 )
845
846 m_panels2 %>%
847   lapply(function(output) {
848     output$tests

```

```
849 })
850
851 standard_beta <- m_panels2[[1]]$OLS %>%
852   QuantPsyc::lm.beta() %>%
853   {tibble(term = names(.), beta = .)} %>%
854   filter(!str_detect(term, "geo"))
855
856 standard_beta <- augment(m_panels2[[1]]$OLS) %>%
857   select(TOTFERRT:.fitted) %>%
858   select(-.fitted) %>%
859   cor() %>%
860   data.frame() %>%
861   select(1) %>%
862   rownames_to_column() %>%
863   rename(term = rowname) %>%
864   merge(standard_beta) %>%
865   mutate(explain = abs(TOTFERRT*beta)) %>%
866   select(term, cor = TOTFERRT, standard_beta = beta, explain)
867
868 lasso_coefs %>%
869   rename(lasso = coef) %>%
870   merge(tidy(m_panels2[[1]]$OLS, conf.int = T)) %>%
871   merge(standard_beta) %>%
872   mutate_at(-1, function(x) as.numeric(x)) %>%
873   pivot_longer(c(lasso:estimate, cor:explain)) %>%
874   mutate(
875     conf.low = ifelse(name != "estimate", NA, conf.low),
876     conf.high = ifelse(name != "estimate", NA, conf.high),
877     lag = ifelse(str_detect(term, "_1"),
878                 ifelse(str_detect(term, "_1_1"), 2, 1), 0),
879     bar = ifelse(name %in% c("cor", "explain"), value, NA),
880     value = ifelse(!(name %in% c("cor", "explain")), value, NA),
881     term = f.clean_names(term, Tosparses = T),
882     term = gsub("_2", "^2", term),
883     term = gsub("_1_1", '[t = -2]', term),
884     term = gsub("_1", '[t = -1]', term),
885     name = case_when(
886       name == "cor" ~ 'Correlation coefficient',
887       name == "estimate" ~ 'Coefficient in OLS',
888       name == "lasso" ~ 'Coefficient in lasso regression',
889       name == "standard_beta" ~ 'Standardized coefficient',
890       name == "explain" ~ 'Contribution to R-squared'
891     ),
892     name = factor(
893       name, levels = c(
894         'Correlation coefficient',
895         'Coefficient in lasso regression',
896         'Coefficient in OLS',
897         'Standardized coefficient',
898         'Contribution to R-squared'
899       )
900     )
901   ) %>%
```

```
902 filter(name != 'Correlation coefficient') %>%
903 {
904   ggplot(.) +
905     geom_vline(xintercept = 0) +
906     geom_point(aes(value, term, fill = factor(lag)), size = 3) +
907     geom_col(aes(bar, term, fill = factor(lag)), color = "black") +
908     scale_fill_viridis_d(option = "magma", begin = .3, end = .7) +
909     scale_y_discrete(labels=scales::parse_format(), limits = rev) +
910     facet_wrap(~ name, scales = "free_x") +
911     labs(x = NULL, y = "Term", fill = "Lag") +
912     theme(legend.position = "bottom", legend.direction = "horizontal")
913 }
```