

# Exploring Factors Impacting Innovation Rates: an analysis on OECD Countries (2001-2015)

By Elia Landini  
Tra My Le  
Veronica Manusardi

Panthéon-Sorbonne Master in Economics  
Course: Econometrics

# TABLE OF CONTENT

1. INTRODUCTION	1
2. LITERATURE REVIEW	3
3. DATA ANALYSIS	7
4. EMPIRICAL APPROACH	12
5. RESULTS	16
6. CONCLUDING REMARKS	21
7. REFERENCES	22

---

## 1. INTRODUCTION

America is the land of possibilities and personal growth, under many aspects. The American dream, the idea that through hard work, courage and determination it is possible to achieve an improved quality of life and economic prosperity, is what pushed people to migrate in the US since the late 1880s. The US has soon become a symbol of entrepreneurial opportunities for people coming from different parts of the globe and it started to be seen as the ideal place to cultivate ideas and transform them into revolutionary inventions - from Edison's light bulb in 1878 till the invention of the internet and the *www*.

However, behind this captivating and fascinating vision, access to opportunities to become an inventor is not uniformly distributed, as there are challenges and social characteristics that may hinder the path for some.

There are indeed various factors, but those who have the strongest influence on the possibility of becoming an inventor are especially the socio-economic ones, as highlighted by Raj Chetty in his research paper "*Who becomes an inventor in America? The importance of exposure to innovation*", published in 2018 on The Magazine for Economic Performance.

The paper explores the main socio-economic barriers that some individual face when pursuing an innovation career, concluding that what is making individual miss the opportunity to become an inventor is their characteristic at birth, so gender, ethnicity, the social condition of their families and the environment they were born in.

The author coined the expression of “*Lost Einsteins*”, referring to all those people who could have made highly influential innovations if they had been provided with the opportunities they rightly earned, but that could not due their characteristics at birth.

In particular, the definition refers to students who are part of *social minorities*, such as women, ethnic minorities and low-income families, who due to misallocation and the limited resources of their families and the lack of a strong mentorship in the early stages of education, missed the possibility to delve into STEM educational fields and develop skills that may push them in creating something and becoming an inventor.

However, Chetty’s investigation was focused exclusively on the United States, mainly because he argued that relatively little was known about the individuals who become inventors in the modern era in the United States and he decided to delve more into this geographical area.

However, his research *opened us a door* to a crucial question: does the American dream of limitless opportunities reflect reality or do these socio-economic barriers cast shadows on this ideal? Behind the façade of the United States as the land of opportunities, we must consider whether certain social boundaries exert such a significant influence that alternative setting may prove more conducive to inventive pursuits.

This question has inspired our research, aimed at exploring whether the United States is indeed the place where anything is possible for aspiring inventors, even considering the different social obstacles that could influence the path to innovative success.

Through this inquiry, we seek to understand if there are other countries or social environments where the dream of becoming an inventor is more easily achievable than in the United States, narrowing our research exclusively on a selected group, countries within the Organisation for Economic Co-operation and Development (OECD)<sup>1</sup>, in order to ensure fair comparisons and draw meaningful insights specific to advanced economies.

Additionally, our research considered data from 2000 to 2015 and thought this is not an extended period of time, it was the base chosen for our research to avoid exogenous effects quite hard to be captured by our model. An example of it, imagining analysing data from the ‘70s, regardless their actual availability, we could have to capture also various generational changes that plays a determinant role in encouraging individual to apply for a patent, and perhaps do it in a specific country. Among them, more binding and widespread intellectual property laws, the strong diffusion of English as international language, or even the epiphany of social media. Hence, we consider these macro-trends as *ceteris paribus* over here, notwithstanding each country has its own factors that either

---

<sup>1</sup> To this date, OECD countries are 38, for the list of members visit: <https://www.oecd.org/about/members-and-partners/>

help or hinder innovation, for example economic structures, education systems, regulations and cultural attitudes. Therefore, we considered some of these factors there appeared to be common in the countries we analysed, such as GDP, trade openness, employment rate, gender gap and government expenditure (in R&D, education or healthcare).

To unravel the complexities of innovation dynamics and in order to understand those factors that set countries apart in terms of innovation, we have built a regression model with the country's rate of innovation as the dependent variable (y) and considering various independent variables representing representing economic, social, and institutional factors.

Our research aims to provide a contribution to the debate on innovation opportunities in an increasingly interconnected and competitive world, not only to highlight where is actually easier to become an inventor, but also to provide proof of those factors that have a high influence on this event.

## 2. LITERATURE REVIEW

Firstly, to understand this research, we consider that it is essential to recognize who the authors and literature consider as an inventor. An inventor is an individual that holds a patent, that gives exclusive rights to sell or profit from an invention and protects their intellectual property. The patent is the only element granting the right to exclude anyone else from the re-production or use of a specific new invention for a stated number of years (an average of 20 years) and to be such, it has to present three main characteristics: novelty (not previously published elsewhere in the world), non-obviousness (sufficient inventive step) and industrial applicability (commercial utility).

Even though the main role of a patent is to protect the intellectual property of the inventor for him to draw economic advantage from it<sup>2</sup>, a patent promotes innovation and technical progress firstly by providing a temporary monopoly for the inventor and secondly, by generating new ideas and encouraging new inventions.

However, the advantages of a patent and their characteristics are not related solely on the inventor, on a micro level, but they must be considered on a macro level, involving countries, firms and a globe wide economy in order to fully comprehend their impact.

Literature has debated on the importance of patents for a country<sup>3</sup>, examining their role in shaping market dynamics and their impact on market competitiveness. Patents not only show us a nation's

---

<sup>2</sup> Griliches, Zvi. "Patent statistics as economic indicators: a survey." *R&D and productivity: the econometric evidence*. University of Chicago Press, 1998. 287-343.

<sup>3</sup> Gambardella, Alfonso, Dietmar Harhoff, and Bart Verspagen. "The value of patents." *Universita Bocconi, Ludwig-Maximilians-Universitaet, and Eindhoven University, Working Paper*: [http://www.creiweb.org/activities/sc\\_conferences/23/papers/gambardella.pdf](http://www.creiweb.org/activities/sc_conferences/23/papers/gambardella.pdf) (2005)

commitment to encourage creativity and protect intellectual property, but they also generate an economic value and represents for a country the ideal mechanism for incentivizing research, competitiveness, and technological progress in a country.

<b>R&amp;D</b>	<ul style="list-style-type: none"> <li>- patent provides legal protection for innovations</li> <li>- incentive for companies to engage in research and development activities</li> </ul>
<b>Tech</b>	<ul style="list-style-type: none"> <li>- attracts FDI</li> <li>- stimulate domestic investment</li> <li>- contributes to technology evolution</li> </ul>
<b>Competitiveness</b>	<ul style="list-style-type: none"> <li>- Patented technologies can serve as a barrier to entry for competitors</li> <li>- Influence on market dynamics</li> </ul>
<b>Market Dynamics</b>	<ul style="list-style-type: none"> <li>- patenting strategies influence firms and new entrants</li> </ul>

#### *Research and Development (R&D)*

Investment in R&D encourages innovation, which in turn, spurs economic growth. Indeed, an economy with a higher degree of innovation - meant as original patented products - will tend to be more independent on the market and gain a competitive advantage. As a result, it will strengthen and enhances its connectivity, becoming more appealing to foreign investments. +

Literature has studied the relation existing between the number of patent applications and R&D expenditure and some<sup>4</sup> were able to build a regression model which showed a strong positive correlation between R&D expenditure and patent application, meaning that an increase of R&D expenditure will enlarge the frontier R&D activities, improve the probability of R&D success and eventually increase the number of patent applications. More literature<sup>5</sup> showed that both business and non-business R&D have a positive influence on patenting activity.

#### *Technological progress*

As emphasized throughout this study, patents provide inventors with exclusive rights to their invention for a certain timeframe. This exclusivity serves as an appealing incentive for companies to invest in innovation, explore new ideas and their subsequent development, and this incentive extends

<sup>4</sup> Prodan, Igor. "Influence of research and development expenditures on number of patent applications: selected case studies in OECD countries and central Europe, 1981-2001." *Applied Econometrics and International Development* 5.4 (2005).

<sup>5</sup> Westmore, B. (2013), "R&D, Patenting and Growth: The Role of Public Policy", OECD Economics Department Working Papers, No. 1047, OECD Publishing, Paris,

especially to foreign investors, making the country more attractive for attracting Foreign Direct Investment (FDI). Nations boasting well-established patent systems tend to be magnets for FDI, as foreign companies are more willing to invest in countries where their intellectual property is protected, as this minimizes the risk of competitors copying or infringing upon their innovations. The elevated presence of patent plays a role in fostering technology transfer among different companies, industries and countries. This means that a strong patent framework for a country does not only stimulate foreign and domestic investment, but also contributes to the collaborative evolution of technology across various sectors.

### *Market competitiveness and dynamics*

Protecting individuals' intellectual property, patents promote fair market generate a incentive for innovation, research and development, and since patents may lead to significant profits, firms and countries are interested and willing to research and innovate more and more.

Additionally, firms and countries with strong patent protection will be better positioned on the market the operate in and are able to establish a competitive advantage by offering innovative products and services, generating entry barriers for potential competitors.

As stated above, the presence of a well-established patent framework may be appealing not only to domestic and foreign investors, but also to other companies that may be more likely to engage in synergies and joint ventures that can drive innovation and market competitiveness.

Certainly, patents wield a significant impact on the innovation rate of a country. They are undeniable instrumental in shaping a nation's innovation landscape. The correlation between patents and a country's innovation level is unmistakable. Without any doubt, the loss of potential inventors – those *Lost Einsteins* discussed above- due to a weak approach to patent protection can have extremely negative effect on a country's innovation environment. Especially if these individuals are left out because of extreme inequality.

Therefore, having a strong patent system is not just a legal necessity but a strategic move for any country looking to boost and maintain its innovation capabilities.

Raj Chetty's research on "Lost Einsteins" has delve into and highlight the main differences among students that hinder them to becoming an inventor, focusing on socio-economic factors such as gender, race (differences between children from black or white families) or growing up in a low-income family or neighborhood. His research has shown what the impact of these factors is and how significant it is, concluding that losing this part of society may have a deep negative impact on US economic growth.

In fact, the critical point faced in his research is that the demographic groups found represent a significant portion of America population, meaning that US society may be missing out on a substantial number of potential inventors. As stated by D. Leonhardt<sup>6</sup>, US society appears to be missing out on most potential inventors from minorities group, both ethnic and gender (in the first place, African American and women), and these groups together make up most of the American population.

However, he did not explain if these factors are especially correlated to the US environment or not and so it is not completely clear if they are the only ones that may hinder potential inventors. Our conclusion was that there may be a possibility that the US is not the most suitable place and country to become an inventor.

The rest of the literature that focuses on this topic often revolves around two main points. Firstly, there is a focus on the significance and influence of parental income, particularly father's income, on a child's possibility and likelihood of becoming an inventor and those other influential factors, such as ethnicity and gender. Secondly, more recent literature explores the relation between misallocation and economic growth, examining how the loss of potential inventors within these marginalized groups could impact the overall welfare of the country.

The impact of innovation on economic growth is undeniable and missing out on a part of it will generate negative impacts on the economy of a country.

The rate of innovation, measured by the introduction of novel ideas, products, or processes, serves as a key indicator of a nation's competitive edge and it is fundamental to revolutionize industries, improve quality of life, and enhance global competitiveness.

However, it still seems difficult to measure the innovative output of a country. The main part of literature has chosen to use R&D spending in a country as a measure of technical change, due to availability and reliability of data, usually using a Cobb-Douglas approach. Studying how overall productivity grows, authors and research either use the amount of R&D capital in a total factor productivity regression or the intensity of R&D in a regression for changes in total factor productivity. These two methods are shown in the equations below.

$$1. \quad \log TFP_t = \log A + \gamma \log RDK_t + \beta_t$$

Where RDK is the stock of R&D capital in time t, the parameter  $\gamma$  is knowledge. This equation yields a measure of the elasticity of output with respect to knowledge.

---

<sup>6</sup> Lost Einsteins: The Innovations We're Missing, by D. Leonhardt, New York Times, 2017  
<https://www.nytimes.com/2017/12/03/opinion/lost-einsteins-innovation-inequality.html>

$$2. \quad d\log TFP_t = \rho \frac{RD_t}{Q_t} + \beta$$

Where RD is the flow of R&D in time t, the parameter  $\rho$  is return to knowledge. This equation yields a measure of the social gross (excess) rate of return to knowledge.

However, both these two approaches present problems. Firstly, it is not explained that knowledge is separable in the production function and there may be measurement problems.

Therefore, other methods to measure innovation rate exist, even though only a few authors have not considered the key role of R&D in measuring innovation.

For example, authors like Geroski (1989) and Budd and Hobbis (1989) conducted studies involving firms in the UK and their results indicated that patenting by UK firms had a significant and positive impact on productivity.

In any case, innovation and technical change has always been considered in literature as the biggest responsible for economic growth of a country, starting from the work of Solow, and overlooking the central role of innovation would represent a substantial loss for a country.

However, as Chetty states<sup>7</sup>, it is not possible to deduce that American economy would be better if these individuals had the possibility to become inventors, but it is possible to notice the failure in society to consent these individuals to become inventors.

### 3. DATA ANALYSIS

In this section, we describe the data sources and define the key variables used in our analysis. In total, 11 databases of the 36 OECD countries from 2001 to 2015 are used. They are merged to form a panel dataset that contains 540 observations of per country innovation rate over the mentioned period across the corresponding GDP, GDP growth, Trade openness, FDI, GDP expenditure on R&D and education among other key determinants of innovation.

#### *Data sources*

The 2 data sources for the datasets are OECD.Stat, which includes data/ metadata for OECD countries among others and the World Bank Open data.

##### 1. OECD.Stat

10/11 of our variables come from this platform, which consolidates data across OECD's many databases. The dependent variable, the "Regional Innovation rate" (specifically, the PCT patent

---

<sup>7</sup> Alex Bell, Raj Chetty, Xavier Jaravel, Neviana Petkova, John Van Reenen "Who Becomes an Inventor in America? The Importance of Exposure to Innovation", The Quarterly Journal of Economics, Volume 134, Issue 2, May 2019, Pages 647–713



applications per million inhabitants) is taken from the “Regional Innovation” dataset, whose metadata collection is undertaken by the Centre for Entrepreneurship, SMEs, Regions and Cities via an annual questionnaire sent to the delegates of the Working Party on Territorial Indicators and via access to websites of National Statistical Offices and Eurostat.

## 2. The World Bank Open Data

Only one variable, the “GDP yearly growth rate” in OECD countries, is extracted from The World Bank Database, which uses World Bank national accounts data and OECD National Accounts data files. The GDP data represents the sum of value added by all its producers (before accounting for consumption of fixed capital in production) and Growth rates of GDP are calculated using the least squares method and constant price data in the local currency.

### *Data description*

The definitions are highlighted below for the variables (per OECD country) of our research question. Specifically, we discuss those of ‘Regional Innovation rate’ as the dependent variable ( $Y_0 = i$ ) and 10 remaining independent variables detailed below ( $\{X_i\}_1^{10}$ ).

- *Regional Innovation rate* (PCT patent applications per million inhabitants,  $Y_0 = i$ ):

A patent is an exclusive right granted by a national or regional patent office for a limited period. Patent Co-operation Treaty (PCT), an international treaty, administered by the World Intellectual Property Organization (WIPO), enables applicants to seek patent protection simultaneously in each of many countries by filing one international patent application instead of several national or regional applications. The dependent variable, Innovation rate, is defined as PCT patent application data averaged per million inhabitants in a country.

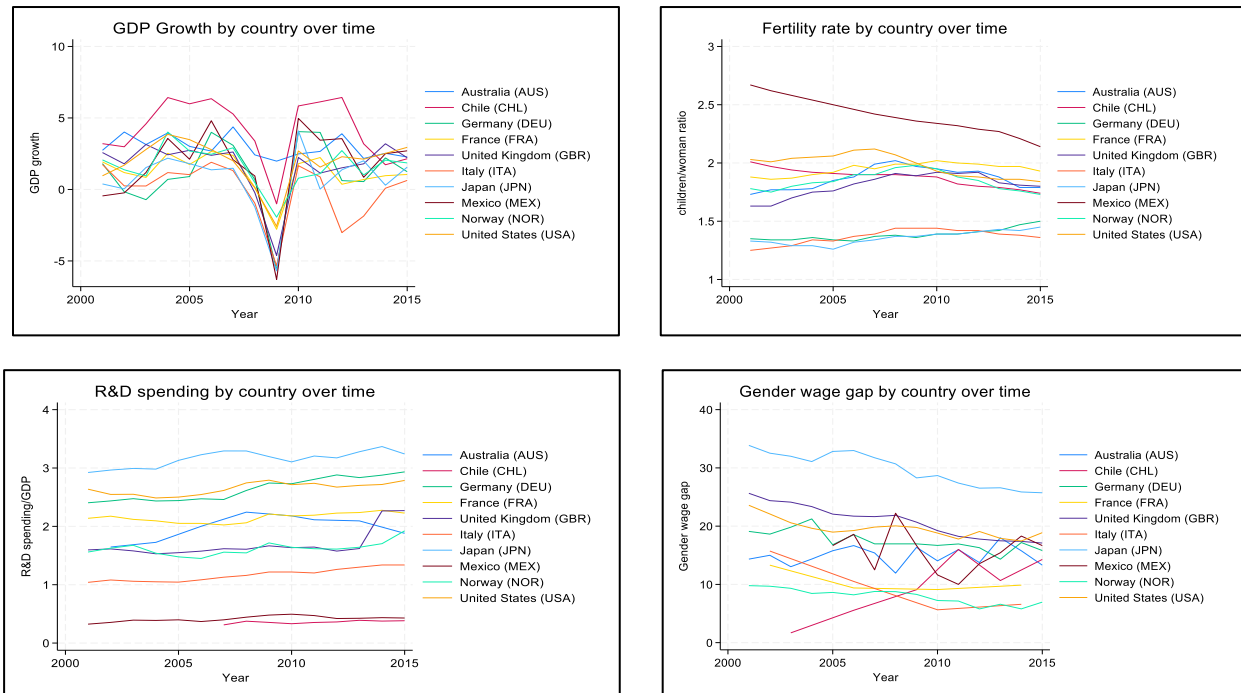
- *Gross Domestic Product* (GDP in million USDs,  $X_1 = y$ ):

GDP is the standard measure of the value added created via production of goods and services in a country during a period and the income earned from that production or the total spending on final goods and services (less imports). It is based on nominal GDP (GDP at current prices).

- *GDP Growth* (annual %,  $X_2 = g$ ):

GDP growth is the annual percentage growth rate of GDP at market prices based on constant local currency. Aggregates are based on constant 2015 prices, expressed in USD.

**Graph 1: raw data, variables' trend by country over time**



Source: our data elaboration on Stata (not all countries included)

- *Foreign direct Investment (FDI) (million USDs, X3=s):*

FDI flows record cross-border transaction value of direct investment in a period, often a quarter or a year. Financial flows include equity transactions, earnings reinvestment and intercompany debt transactions. Outward flows represent transactions that increase the investment that investors in the reporting economy have in enterprises in a foreign economy, less transactions that reduce such investment.

- *Trade Openness - Trade in Goods and Services (million USDs, X4=t):*

Trade in Goods and Services is the transactions in goods and services between residents and non-residents. It is measured in million USDs at 2015 constant prices and purchasing power for net trade (exports minus imports)

- *Employment Rate (thousands of persons or % of working age population, X5=e):.*

Employment rate is a measure of the extent to which available labor resources are being used, calculated as the ratio of the employed to the working age population. The employed are those aged 15 or over who report having worked in gainful employment for at least one hour in the previous week or who had a job but were absent from work during the reference week. The working age population refers to people aged 15-64. It is seasonally adjusted and measured in 'thousand persons' aged 15 and over.

- *Gender Wage Gap indicator (X6=w):*

Gender wage gap is the difference between median earnings of men and women relative to median earnings of men.

- *Fertility rate* (Total children per woman,  $X7=f$ ):

Fertility rate in a year is the total number of children that would be born to each woman if she were to live to the end of her child-bearing years and give birth to children in alignment with the prevailing age-specific fertility rates. It is calculated by totaling the age-specific fertility rates as defined over 5-year intervals.

- *Public Spending on R&D (%GDP) ( $X8=r$ ):*

GDP spending on R&D is the total expenditure on R&D carried out by all resident companies, research institutes, university and government laboratories, etc., in a country, including R&D funded from abroad, excluding domestic funds for R&D outside the domestic economy.

- *Public Spending on Education (% GDP,  $X9=u$ ):*

Public spending on education includes direct expenditure on educational institutions and educational-related public subsidies given to households and administered by educational institutions. It is shown as a percentage of GDP, divided by primary, primary to post-secondary non-tertiary and tertiary levels. OECD Database on Education Statistics

- *Public spending on healthcare (% GDP,  $X10=h$ ):*

Health spending measures the final consumption of health care goods and services including personal health care and collective services but excluding spending on investments. It is measured as a share of GDP, as a share of total health spending

#### Data description

Variable names (per country)	Units	Variables in regression (specified if log was taken)	Frequency	Years of available data	Source
<b>Regional Innovation rate (<math>Y0=i</math>)</b>	PCT patent applications per million inhabitants	Regional innovation rate (Natural logarithm of $i$ )	annually	1990-2015	OECD
<b>Gross domestic product (GDP) (<math>X1=y</math>)</b>	Million USDs	National GDP (Natural logarithm of $y$ )	annually	1970-2022	OECD
<b>GDP Growth (<math>X2=g</math>)</b>	Annual %	GDP yearly growth rate	annually	1961-2022	World Bank
<b>Foreign direct Investment (FDI) (<math>X3=s</math>)</b>	Million USDs	Foreign direct Investment (FDI) (Natural logarithm of $s$ )	quarterly	2005-2022	OECD
<b>Trade Balance (Trade in Goods and Services) (<math>X4=t</math>)</b>	Million USDs	Trade balance (Natural logarithm of $t$ )	quarterly	1995-2022	OECD

<b>Employment Rate</b> (X5=e)	thousands of persons	Employment rate (LFS) (Natural logarithm of e)	annually	1955-2021	OECD
<b>Gender Wage Gap indicator</b> (X6=w)	employees	Gender wage gap indicator <i>*Excluding the self-employed</i>	annually	1995-2022	OECD
<b>Fertility rate</b> (X7=f)	Total children per woman ratio	Fertility rate (children/woman)	annually	1970-2022	OECD
<b>Public Spending on R&amp;D</b> (X8=r)	% GDP	R&D spending	annually	2000-2022	OECD
<b>Public Spending on Education</b> (X9=u)	% GDP	Educational spending	annually	2000-2020	OECD
<b>Public spending on healthcare</b> (X10=h)	% GDP	Healthcare system spending	annually	1970-2022	OECD

### *Data Manipulation*

Firstly, disaggregated quarterly data was transformed into aggregated annually data for some variables.

Secondly, with regards to missing values, where existed for some variables, we employ a panel random effect estimator for each to generate a prediction on both possible past and future observations. As per Hausman test, this provides better estimation for our unadjusted data. Despite large number of variables, we could only employ those with a complete range of observations over time (otherwise the regression would continue predicting missing results). Notice that for some countries, data is complete for variables not listed in our model but to avoid over-noise-capturing bias which might affect only those countries, we opted for a uniformed prediction model, including those complete variables common to all countries.

Thirdly, in the case of log-transformed variables (specified in the table below), we proceed with retaking the log function for the new predicted values. We hold that doing a 2-stage log transformation makes the regressive forecast more robust to possible heteroskedasticity or high-skewness possible biases.

Next, in the case of log-transformed variables, we proceed with retaking the log function for the new predicted values. We hold that doing a 2-stage log transformation makes the regressive forecast more robust to possible heteroskedasticity or high-skewness possible biases.

At this point, stationarity tests are crucial to validate that statistical properties of mean and variance remain constant over time in our data frame. Ensuring stationarity is vital for accurate modelling, forecasting and reliable statistical inference as non-stationarity data may mislead models and compromise forecasting accuracy. Hence, stationarity verification is essential for maintaining the

stability and interpretability of the time series models. To investigate the heterogeneous range of non-stationarity biases, we performed 4 tests for each independent variable: Im-Pesaran-Shin (IPS) test, Breitung Panel Unit Root Test, Hadri M test and Levin, Lin and Chu (LLC) test. Where there existed non-consensus result, the test with the most reliable significant statistics would be taken. It is also worth noticing that in our case, traditional approaches such as Augmented Dickey Fuller test are not available given the multitude of panel belonging to the data frame. The results suggested that the variables 'y', 'g', 's', 't' and 'r' are likely non-stationary across the panel.

Given the results in terms of stationarity extrapolated from the employed stationarity test, we moved on to implement corrections to the non-stationary variables, by recurring to first differences, following the rationale of eliminating trends and making the series more amenable to statistical analyses. As a result, we found that the differenced variables 'y', 'log\_y', 'g', 's', 't', 'log\_t' and 'r', appear stationary across panels, supporting the effectiveness of the differencing transformation achieving stationarity.

## 4. EMPIRICAL APPROACH

### *Introduction*

It is common practice in current literature to present past models in order to give strength to possible assumptions or avoiding specifications already widely described by other sources. However, the current literature on this side was quite meagre, and largely focused on social factors, giving macroeconomic factors a backstage role or sometimes even excluding them.

Hence, in addressing our research question, i.e., how and to what extent macroeconomic, social and government-biased factors affect the innovation rate in a country, we were not really backed from any previous research which could propose similar or feasible models to be apply in our case. Moreover, when it comes to deal with regional innovation there exist factors that are largely affecting individuals on their own micro-sphere, related to both intrinsic characteristics at birth and characteristics developed over time, but that are out of the scope of our research demand<sup>8</sup>. This idea of a macro-level analysis is probably the key added value of our research to the current literature, which can be translated into a radical switch from a literature that investigates individuals' factors leading to holding a patent, to research, ours, that investigates country's factors that influence the chance of having new patent holders on its domestic soil. It is worth to notice that we are mentioning patent holders and not number of patents released, since our dataset refers to a density value in terms of millions of inhabitants of people granted of this right and this leads to the potential risk of having people holding more than one patent.

---

<sup>8</sup> Alex Bell, Raj Chetty, Xavier Jaravel, Neviana Petkova, John Van Reenen "Who Becomes an Inventor in America? The Importance of Exposure to Innovation", The Quarterly Journal of Economics, Volume 134, Issue 2, May 2019, Pages 647–713

It is paramount to underline that there exists always potential bias in the estimation of regional innovation rate, since we are aware of omitted variables that are not available, hard to measure and can be related to the economic development and not really to the innovation itself.

### *Model Specifications*

In order to find the most adherent model for our research problem without incurring in both endogeneity and high-multicollinearity issues, we have developed multiple linear panel regression models to estimate the marginal impact of each factor on the regional innovation rate (i). Aside from our dependent variable, independent variables have been clustered in three groups according to their nature and then progressively added to the model to verify their relevance and models' fit (R-squared). But before giving insight into these different model structure, it is mandatory to explicate the choice of adopting a random effect estimator rather than more common methods of estimations for panel data, among them fixed effects or pooled OLS.

This econometric approach, grounded in the random effects model, is chosen to efficiently capture unobserved heterogeneity across countries, recognizing the diverse influences shaping innovation trajectories within the dynamic context of OECD member states. Pooled OLS is appropriate when you treat the panel data as if it were a single cross-sectional dataset, ignoring any individual-specific or time-specific effects. This assumes that there is no correlation between the individual effects and the independent variables, which is not our case. On the other hand, to verify whether random effects assumption was more suitable than the fixed effects one, we conducted the Hausman test. The results of the test indicate that the p-value is considerably high (0.9992), suggesting that you fail to reject the null hypothesis. The null hypothesis in the Hausman test is that the difference in coefficients between the fixed-effects (FE) and random-effects (RE) models is not systematic. Since the p-value is critical, there is no strong evidence to suggest that the differences in coefficients between the two models are systematic. Therefore, we chose the random-effects model as it is consistent with the assumptions of the test.

**Table 4: Hausmann test**

Variable	FE (b)	RE (B)	Difference (b-B)	Std. Err. sqrt(diag(V_b-V_B))
log_y	142.44640	146.30180	-3.85538	-
g	-093209	-0.91844	-0.01365	-
log_s	67132.79	67984.52	-851.72360	-
log_t	-66392.85	-67234.6	84.74470	-
log_e	39.33977	40.10877	-0.76901	-
w	-0.87264	-0.71456	-0.15809	-
f	27.24597	29.78793	-2.54196	2.49745
r	1.99552	1.92523	0.07029	-
u	-626554	-6.06510	-0.20043	-
h	13.66471	14.19723	-0.53252	0.07353

Source: our data elaboration on Stata

The model specifications follow:

- *Macro-trends specification*

$$iit = \beta_0it + \beta_1i \ln(yit) + \beta_2igit + \beta_3i \ln(sit) + \beta_4i \ln(tit) + \beta_5i \ln(eit) + \varepsilon it$$

where  $i$  represents regional innovation rate by country over time,  $\ln(y)$  the logarithmic function of country's GDP,  $g$  the yearly GDP's growth rate,  $\ln(s)$  the logarithmic function of foreign direct investment (FDI),  $\ln(t)$  the logarithmic function of trade balance (trade openness),  $\ln(e)$  the logarithmic function of employment rate,  $\varepsilon it$  is the idiosyncratic error term with the assumption  $E(\varepsilon it) = 0$ .

- *Female empowerment specification*

$$iit = \beta_0i + \beta_1i \ln(yit) + \beta_2igit + \beta_3i \ln(sit) + \beta_4i \ln(tit) + \beta_5i \ln(eit) + \beta_6iwit + \beta_7ifit + \varepsilon it$$

where  $i$  represents regional innovation rate by country over time,  $\ln(y)$  the logarithmic function of country's GDP,  $g$  the yearly GDP's growth rate,  $\ln(s)$  the logarithmic function of foreign direct investment (FDI),  $\ln(t)$  the logarithmic function of trade balance (trade openness),  $\ln(e)$  the logarithmic function of employment rate,  $w$  the gender wage gap,  $f$  the fertility rate,  $\varepsilon it$  is the idiosyncratic error term with the assumption  $E(\varepsilon it) = 0$ .

- *Government expenditure targets specification*

$$iit = \beta_0i + \beta_1i \ln(yit) + \beta_2igit + \beta_3i \ln(sit) + \beta_4i \ln(tit) + \beta_5i \ln(eit) + \beta_6iwit + \beta_7ifit + \beta_8irit + \beta_9iuit + \beta_{10}ihit + \varepsilon it$$

where  $i$  represents regional innovation rate by country over time,  $\ln(y)$  the logarithmic function of country's GDP,  $g$  the yearly GDP's growth rate,  $\ln(s)$  the logarithmic function of foreign direct investment (FDI),  $\ln(t)$  the logarithmic function of trade balance (trade openness),  $\ln(e)$  the logarithmic function of employment rate,  $w$  the gender wage gap,  $f$  the fertility rate,  $r$  the public expenditure on R&D,  $u$  the public expenditure for education,  $h$  the public expenditure for healthcare,  $\varepsilon it$  is the idiosyncratic error term assuming  $E(\varepsilon it) = 0$ .

Once independent variables have been included in the model, we must depurate our dependent variable from its own recursive bias, given by its past values: the autocorrelation test. The analysis undertaken on the variable  $i$  and its first lag, denoted as  $L.i$ , discloses a robust positive autocorrelation with a correlation coefficient of approximately 0.9938. This latter underscores a pronounced and consistent linear relationship between the innovation rate ( $i$ ) and its past values averagely for each panel. The strength of this positive autocorrelation signifies that the current level of innovation is strongly influenced by its historical performance, and the nearly perfect positive correlation suggests a persistent pattern in the innovation rate over time, indicating a substantial degree of inertia or momentum in the innovation process. This empirical finding leads us to have evidence of potential influence of past innovation dynamics on present outcomes, emphasizing the relevance of incorporating temporal dependencies in the modelling and analysis of innovation trends. However, including *lagged*  $i$  in the model to tackle autocorrelation and temporal dependencies would have brought to over capture the data noise and mislead eventual conclusions. Therefore, we have kept working with our original specifications, without including any lag variables, aware of the autocorrelation of innovation rate.

As regards the specification choice, the model fit statistics for the three distinct specifications, namely Specification 1 (Macro-trends), Specification 2 (Female Empowerment), and Specification 3 (Government Expenditure Targets), offer a comprehensive view of their respective explanatory prowess in elucidating variations in the innovation rate ( $i$ ) across the observed countries. While Specification 1 captures an overall R-squared of 0.1311, denoting a modest but notable explanatory power based on macroeconomic indicators, Specification 2, integrating variables that deal with women empowerment, demonstrates a slightly higher R-squared of 0.0718. The inclusion of the wage gap ( $w$ ) and fertility rate ( $f$ ) as significant contributors suggests their relevance in shaping innovation dynamics, however risk of correlation but not causality must be considered. In contrast, Specification 3 (Government Expenditure Targets) stands out with the highest overall R-squared of 0.2877, underscoring the substantial impact of government-related variables such as R&D spending ( $r$ ), education spending ( $u$ ), and healthcare spending ( $h$ ) on innovation rates. Reason why we then moved forward adopting this latter as the better explicative specification.

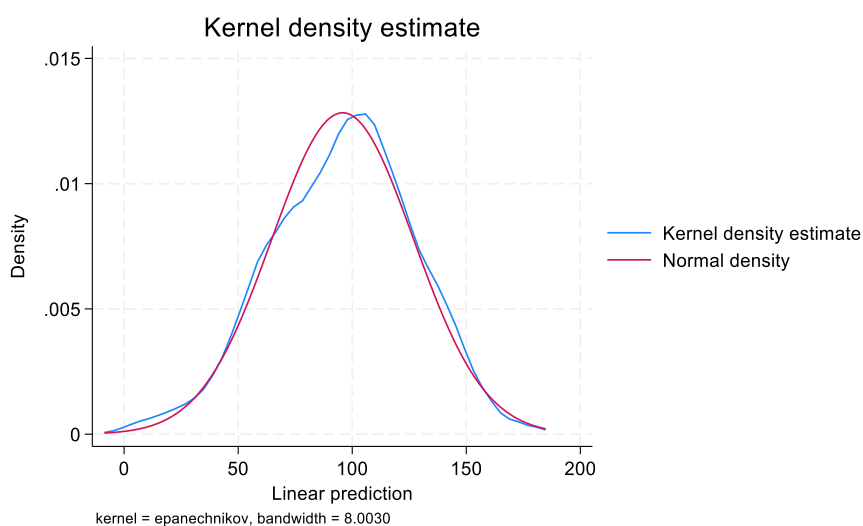
Moreover, and additional proof of the suitability of our choice was ultimately given by our checks on error terms' normality. For Specification 1, the Jacques-Bera test, Skewness and Kurtosis tests, and the Shapiro-Wilk test consistently reject the null hypothesis, suggesting a departure from normality for the error terms associated with the innovation rate ( $i$ ). Similarly, Specification 2 follows suit, with all normality tests indicating a significant departure from normal distribution. However, in Specification 3, the p-values from the Jacques-Bera test, Skewness and Kurtosis tests, and the



Shapiro-Wilk test are relatively higher. While the Jacques-Bera test's p-value is not notably elevated, the other two tests provide some support for the hypothesis of normality. Therefore, based on these normality tests, the error terms in Specification 3 appear to be relatively close to a normal distribution compared to the other specifications. Together with the better model fit, the relevance of normality in ensuring the validity of statistical inferences, Specification 3 is again chosen as the preferred model, where the error terms demonstrate a comparatively better fit to a normal distribution.

Eventually, remaining on Specification 3, the functional form of the model had to be evaluated to address whether our dependent variable was following a different trend structure with respect to the one we have predicted ex ante. Running the test, the p-value resulted quite small ( $p = 0.0000$ ), indicating that we reject the null hypothesis, suggesting that there might be a specification error in the model related to quadratic and cubic functional forms, which leads us to rule them out of the model specification. The test confirmed us that chosen functional form fits our data.

**Graph 2: Kernel density estimate of se for Specification 3**



*Source: our data elaboration on Stata*

## 5. RESULTS

Beginning with each specification's marginal effects estimates, notably, the natural logarithm of GDP ( $\log\_y$ ) exhibits distinctive impacts across specifications, with Specification 2, centred around female empowerment, displaying the highest coefficient (5.152.036), suggesting that when accounting for gender-related dynamics, the positive influence of GDP on innovation is more pronounced. In

contrast, Specification 3, emphasizing government expenditure targets, reports the smallest coefficient for  $\log\_y$  (1.463.018), which could lead to hold that the varying significance of  $\log\_y$  underscores the importance of considering gender-related factors in innovation processes, potentially indicating that fostering female empowerment contributes significantly to innovation, even though we have causality issue must be addressed.

Similarly, the GDP yearly growth rate exhibits nuanced relationships across specifications. Specification 3 stands out with the most substantial negative impact (-0.91844), indicating that higher GDP growth, within the context of government expenditure targets, is associated with a lower innovation rate. This may be attributed to policy choices and resource allocation strategies influenced by government expenditure targets, revealing a complex relationship between economic growth and innovation in this specific context, which actually goes against our expectations on this side since we were expecting that a favourable economic environment would have been largely more beneficial to the vitality of the entire research and development field, fostering patent grants. It is even worth to mention that the result could be biased by the fact that individuals are encouraged to apply for a patent during those periods characterized by solid economic growth and optimistic expectations, but that does not guarantee that the patent would be released in a year that preserve the same characteristics as at the application time (application and verification procedures can take years according to countries' bureaucracy).

Furthermore, the natural logarithm of FDI ( $\log\_s$ ) shows intriguing patterns. In Specification 1, emphasizing macro-trends,  $\log\_s$  has the highest coefficient (82044.77), indicating a robust positive relationship between Foreign Direct Investment and innovation. This aligns with the understanding that a favourable economic environment, as captured by macroeconomic trends, attracts higher levels of FDI, thereby fostering innovation.

Public expenditure variables in Specification 3 offer nuanced insights into their impact on innovation rates within OECD countries. The positive coefficient for government spending on research and development (R&D) (1.925.225) aligns with expectations, indicating that higher R&D investments relative to GDP positively contribute to innovation. Unexpectedly, higher public spending on education as a percentage of GDP exhibits a negative association with the innovation rate (-6.065101). This counterintuitive result warrants further investigation, and suggest the risk of possible unobserved variables, since from one side the global trend of this “patent rush” has been pushing regional innovation to its highest historical data, on the other, given external and procyclin macro-dynamics, such as crisis, trade sanctions and ultimately war tensions, the government public spending in R&D has been progressively reduced to often leave the space to economic recovery. Conversely, increased public spending on healthcare relative to GDP (14.19723) is positively linked to higher

innovation rates, suggesting a potential role of facilitating the access to the healthcare system and a healthy population in fostering innovation. These findings underscore the intricate relationship between public expenditure components and innovation dynamics within the OECD countries frame. The table below does not mention statistical significance data (p-value, z) since all the parameters result to be significant within our data frame and model structure.

**Table 2: Specification Analysis**

	SPECIFICATION 1		SPECIFICATION 2		SPECIFICATION 3	
	b	se	b	se	b	se
Natural logarithm of y (GDP)	2514.705	2677.466	5152.036	2583.135	1463.018	2542.915
GDP yearly growth rate	-0.3883	0.309913	-0.58297	0.301106	-0.91844	0.275109
Natural logarithm of s (FDI)	82044.77	37259.08	78259.32	35546.32	67984.52	32154.55
Natural logarithm of t (TB)	-81121.9	36772.13	-77377.9	35081.75	-67234.6	31734.32
Natural logarithm of e	4098.525	5268.398	4148.856	5022.998	4010.877	4546.576
Gender wage gap indicator	-	-	-149.357	0.316669	-0.71456	0.299045
Fertility rate (children/woman)	-	-	4601.328	1045.343	2978.793	9642.898
R&D spending	-	-	-	-	1925.225	2600.789
Education spending	-	-	-	-	-60.651	2837.728
Healthcare system spending	-	-	-	-	1419.723	1362.253
Constant	-394.487	6432.373	-453.711	6425.924	-522.374	5968.244
Observations	521	-	521	-	521	-

*Source: our data elaboration on Stata*

Chosen the model specification, we now display results for what were the potential estimators to employ to analyses our panels, i.e., Fixed Effects (FE), Fixed Effects with Clustered Standard Errors (FECL), Random Effects (RE), Random Effects with Robust Standard Errors (RERB), and lastly Pooled Ordinary Least Squares (pooled OLS). Notably,  $\log_y$  (GDP) exhibits minimal variance across RE and RERB, suggesting robustness in this variable's estimation method. GDP yearly growth rate shows consistency across all models, indicating its stable impact on innovation rates. However, for FDI ( $\log_s$ ) and Trade Balance ( $\log_t$ ), the choice between RE and RERB is crucial, as OLS provides notably different estimates, emphasizing the significance of accounting for unobserved heterogeneity. Employment Rate demonstrates a consistent impact across all models, with marginal differences.

**Table 3: Estimated Values**

Variable	FE	FECL	RE	RERB	OLS
<b>log_y</b>	142.44644	142.44644	146.30182	146.30182	198.48275
<b>g</b>	-0.93209	-0.93209	-0.91844	-0.91844	-0.31454
<b>log_s</b>	67132.795	67132.795	67984.518	67984.518	65787.939
<b>log_t</b>	-66392.851	-66392.851	-67234.596	-67234.596	-65060.566
<b>log_e</b>	39.33977	39.33977	40.10877	40.10877	81.41465
<b>w</b>	-0.87264	-0.87264	-0.71456	-0.71456	3.83101
<b>f</b>	27.24597	27.24597	29.78793	29.78793	79.19797
<b>r</b>	1.99552	1.99552	1.92522	1.92522	4.25235
<b>u</b>	-6.26554	-6.26554	-6.06510	-6.06510	-9.32652
<b>h</b>	13.66471	13.66471	14.19724	14.19724	29.19655
<b>_cons</b>	-500.92249	-500.92249	-522.37432	-522.37432	-1286.28070

*Source: our data elaboration on Stata*

In our analysis, the risk of heteroskedasticity introduces a concern regarding the assumption of constant variance in the residuals. Heteroskedasticity, or uneven variability of the error terms, can compromise the reliability of standard errors and, consequently, the validity of statistical inferences. In testing these intrinsic risks in our analysis, the Breusch-Pagan test outcomes, with missing statistics and p-values, left the assessment inconclusive, possibly due to peculiarities in the residuals or data. However, even though not such valuable, the test statistic (F-statistic) is 6.14 with a p-value of 0.0000. Since the p-value is less than the conventional significance level of 0.05, we would have rejected the null hypothesis. This provides potential preliminary evidence against the assumption of constant variance, suggesting the presence of heteroskedasticity in the model.

In addition, the subsequent White test offered deeper insights, revealing substantial evidence of heteroskedasticity in the residuals. The F-statistic of 6.14, coupled with a p-value of 0.0000, unequivocally indicated the rejection of the assumption of constant variance. Although specific test statistics (W0, W50, W10) and associated p-values were not available, the mean and standard deviation of squared residuals for each country ID underscored significant variability, bolstering the conclusion of heteroskedasticity. Acknowledging the implications of this issue on parameter estimates, we opted for robust standard errors to enhance the robustness of our panel data analysis, accounting for potential variations in error term variance across nations and reinforcing the validity of our statistical inferences.

**Table 6: IV regressions**

Variable	Legend	IV1	IV2	IV3
w	b	0.16271	30.37978	46.68347
	se	7.31529	12.95725	32.83681
	p	0.9823	0.019	0.1551
log_y	b	38.73575	684.79440	987.49227
	se	123.46527	323.72974	652.96308
	p	0.7537	0.0344	0.1305
log_s	b	11.98115	-260256.28	-511968.85
	se	40.20147	367416.75	724887.6
	p	0.7657	0.4787	0.48
log_t	b	-2.28751	256646.49	504819.36
	se	8.84477	362503.95	715371.43
	p	0.7959	0.479	0.4804
f	b	39.49464	176.35456	242.39084
	se	25.94485	58.37765	134.56462
	p	0.1279	0.0025	0.0717
r	b	7.56697	-20.65011	-36.12017
	se	13.32579	29.5427	63.73141
	p	0.5701	0.4772	0.5709
u	b	5.71598	-2.46841	-2.05175
	se	10.15501	19.45304	32.62235
	p	0.5735	0.899	0.9499
h	b	20.59631	66.35703	86.96926
	se	4.67605	17.2541	41.66987
	p	0.0	0.0001	0.0369
_cons	b	-162.46829	-126.29357	-1825.6752
	se	223.65626	452.59054	1145.7025
	p	0.4676	0.0053	0.111

*Source: our data elaboration on Stata*

Our original random effect regressive model, “xtreg, re” in Stata notation, assumed exogeneity of the explanatory variables. However, concerns about endogeneity prompted the need for rigorous testing. A detailed exploration revealed intricate relationships within the dataset, raising doubts about the independence of variables, such as public expenditure on education ( $u$ ) and public expenditure on healthcare ( $h$ ). The potential correlation of these variables with the error term necessitated endogeneity tests to validate the assumptions of our model.

Upon scrutiny, both the Durbin-Wu Hausman and Hansen Sargan tests yielded low p-values, signaling strong evidence against exogeneity. Consequently, we opted for an instrumental variable (IV) approach to enhance the robustness of our regression analyses. The selection of public expenditure on education ( $u$ ) and public expenditure on healthcare ( $h$ ) as instrumental variables stems from their theoretical relevance and orthogonality with the error term.

Our rationale for choosing  $u$  and  $h$  as instruments lies in their economic significance: public expenditure on education ( $u$ ) reflects investments in human capital, and public expenditure on healthcare ( $h$ ) is a key determinant of societal well-being.

The effectiveness of  $u$  and  $h$  as instruments is affirmed by IV1, IV2, and IV3 methodologies, where their coefficients exhibit consistency across specifications. Furthermore, correlation analyses confirm the weak correlation between the instrumental variables ( $w$ ,  $g$ ,  $\log_e$ ) and the chosen instruments ( $u$ ,  $h$ ), indicating their orthogonality. This comprehensive validation process reinforces our confidence in  $u$  and  $h$  as high-quality instruments, supporting the robustness of our regression models in the presence of endogeneity challenges.

## 6. CONCLUDING REMARKS

In this paper we developed a random effect panel regression model with robust standard errors for a panel of 36 OECD countries over the period from 2001 to 2015 to investigate what are the pivotal factors influencing a country's capability of attracting innovation. We started by considering those factors, as suggested by literature and our knowledge, that could directly influence the innovation rate. However, it was not always trivial to measure and to observe each variable and the external macro-dynamics that could influence them (crises, trade sanctions, war tensions), so we are still aware of the evidence of bias in the estimation, even though also confident with the reliability of our results and their respective robustness checks. Additionally, the presence of bias could be linked to the likelihood of individuals to apply for patents during periods of solid economic growth, which may not guarantee the same characteristics at the time of patent release.

By considering three different model specifications, Macro-trends, Female empowerment and Government expenditure targets, we observed a drastic increase in fitting when social variables, as female empowerment were added, but the ultimate contribution has been given by what were government expenditure targets.

Firstly, we noticed that a higher GDP growth is associated with a lower innovation rate, which was for us an unexpected result. We conclude that this effect may be attributed to policy choices and resource allocation strategies influenced by government expenditure targets, revealing a complex relationship between economic growth and innovation. Also, the result describing the relation with public spending on education and innovation rate was unexpected, as we discovered a negative correlation that would require further investigation. The counterintuitivity of our result suggests the risk of possible unobserved variables, since from one side the global trend of this "patent rush" has been pushing regional innovation to its highest historical data, on the other, given external and

procyclic macro-dynamics, such as crisis, trade sanctions and ultimately war tensions, the government public spending in R&D has been progressively reduced to often leave the space to economic recovery.

However, we were delighted to notice that there exists a robust positive relationship between foreign Direct Investment (FDI) and innovation. In fact, both government spending on research and development (R&D) and public spending on healthcare positively contributes to a country innovation rate, aligning with expectations. In summary, the findings underscore the complex and multifaceted nature of the relationship between economic variables, government policies, and innovation within OECD countries, paving the way to future discussion, especially on those variables that could lift the veil on endogeneity issues. Ultimately, the analysis suggests the importance of considering gender-related factors and the need for further investigation into unexpected results, potential biases, and the intricate dynamics between different components of public expenditure and innovation.

## 7. REFERENCES

### *Bibliography*

Griliches, Zvi. *"Patent statistics as economic indicators: a survey."* R&D and productivity: econometric evidence. University of Chicago Press, 1998. 287-343.

Gambardella, Alfonso, Dietmar Harhoff, and Bart Verspagen. *"The value of patents."* Universita Bocconi, Ludwig-Maximilians Universitaet, and Eindhoven University, (2005).

Prodan, Igor. *"Influence of research and development expenditures on number of patents applications: selected case studies in OECD countries and central Europe, 1981-2001."* Applied Econometrics and International Development 5.4 (2005).

Westmore, B. (2013), *"R&D, Patenting and Growth: The Role of Public Policy"*, OECD Economics Department Working Papers, No. 1047, OECD Publishing, Paris.

Alex Bell, Raj Chetty, Xavier Jaravel, Neviana Petkova, John Van Reenen *"Who Becomes an Inventor in America? The Importance of Exposure to Innovation"*, The Quarterly Journal of Economics, Volume 134, Issue 2, May 2019, Pages 647–713.

Cameron, Gavin. *"Innovation and economic growth"*. No. 277. Centre for Economic Performance, London School of Economics and Political Science, 1996.

D. Leonhardt, *"Lost Einsteins: The Innovations We're Missing"*, New York Times, 2017  
<https://www.nytimes.com/2017/12/03/opinion/lost-einsteins-innovation-inequality.html>

Hana, Urbancová. *"Competitive advantage achievement through innovation and knowledge."* Journal of competitiveness 5.1 (2013): 82-96.

Raj Chetty, Nathaniel Hendren, Patrick Kline, Emmanuel Saez, *Where is the land of Opportunity? The Geography of Intergenerational Mobility in the United States*, The Quarterly Journal of Economics, Volume 129, Issue 4, November 2014, Pages 1553–1623.

Alex Bell, Raj Chetty, Xavier Jaravel, Neviana Petkova, John Van Reenen “*Who Becomes an Inventor in America? The Importance of Exposure to Innovation*”, The Quarterly Journal of Economics, Volume 134, Issue 2, May 2019, Pages 647–713.

Organization for Economic Co-operation and Development (OECD) (2005). Oslo manual: Guidelines for collecting and interpreting innovation data, Third Edition, Paris.

Akduğan, U.; Doğan, N. Factors Affecting Innovation in OECD Countries. EKOIST Journal of Econometrics and Statistics, [Publisher Location], v. 0, n. 36, p. 111-136, 2022.

### *Datasets*

OECD. "National Accounts at a Glance", OECD National Accounts Statistics (database), 2024.  
<https://doi.org/10.1787/data-00369-en>

OECD. "Main Science and Technology Indicators", OECD Science, Technology and R&D Statistics (database), 2024.  
<https://doi.org/10.1787/data-00182-en>

OECD. "Education at a glance: Educational finance indicators", OECD Education Statistics (database), 2024.  
<https://doi.org/10.1787/c4e1b551-en>

OECD. "Health expenditure and financing: Health expenditure indicators", OECD Health Statistics (database), 2024  
<https://doi.org/10.1787/data-00349-en>

OECD (REGPAT). Regional Innovation rate Data, 2010-2019.  
[https://stats.oecd.org/Index.aspx?DataSetCode=REGION\\_INNOVATION#](https://stats.oecd.org/Index.aspx?DataSetCode=REGION_INNOVATION#)

OECD. National GDP Data, 1970-2022.  
<https://data.oecd.org/gdp/gross-domestic-product-gdp.htm>

OECD. GDP yearly growth rate Data, 1961-2022.  
<https://data.worldbank.org/indicator/NY.GDP.MKTP.KD.ZG?locations=OE>

OECD. FDI Data, 2005-2022  
<https://data.oecd.org/fdi/fdi-flows.htm>

OECD. Trade Balance Data, 1995-2022.  
<https://data.oecd.org/trade/trade-in-goods-and-services.htm>

OECD. Employment rate Data, 1955-2021.  
[https://stats.oecd.org/Index.aspx?DataSetCode=LFS\\_D](https://stats.oecd.org/Index.aspx?DataSetCode=LFS_D)

OECD. Gender Wage Gap indicator Data, 1955-2022.  
<https://data.oecd.org/earnwage/gender-wage-gap.htm>

OECD (Family Indicators database). Fertility rate Data, 1970-2022.



<https://data.oecd.org/pop/fertility-rates.htm>

OECD (Database on Main Science and Technology Indicators). R&D Spending Data, 2000-2022.  
<https://data.oecd.org/rd/gross-domestic-spending-on-r-d.htm>

OECD. Education Spending Data, 2000-2020.  
<https://data.oecd.org/eduresource/public-spending-on-education.htm>

OECD. Healthcare System Spending Data, 1970-2022.  
<https://data.oecd.org/healthres/health-spending.htm>

```

1  *Who becomes an inventor today?-Experimental analysis-Do-file**
2
3  //////////////////////////////////////
4  //////////////////////////////////Introductory setup////////////////////////////////
5  //////////////////////////////////////
6
7  clear mata
8  capture log close
9  cd "D:\Sorbonne\Econometrics\Project PSME1\Datasets"
10 clear
11
12 //////////////////////////////////////
13 //////////////////////////////////Data import and manipulation////////////////////////////////
14 //////////////////////////////////////
15
16 *0) Y0=i--->Regional innovation rate (PCT patent applications per million inhabitants)
17 *-----
18 *Data Source: OECD (https://stats.oecd.org/Index.aspx?DataSetCode=REGION_INNOVATION)
19 import delimited RIRi
20 describe
21 *save in Stata format
22 save dfi, replace
23 *now, we do not need all the included variables, but only the only usefull ones for hour
  regression or functional to our merging process
24 keep value territoryleve~y reg_id region year
25 *equally, another specification must be introduced, i.e., we are only interested in country's
  level of regional innovation (patent application density per million inhabitants), given our data
  availability on macro variables
26 keep if territorylevelandtypology == "Country"
27 keep value reg_id year
28 *keep only target countries (OECD members)
29 keep if reg_id == "AUS" | reg_id == "AUT" | reg_id == "BEL" | reg_id == "CAN" | reg_id == "CHL" |
  reg_id == "CZE" | reg_id == "DNK" | reg_id == "EST" | reg_id == "FIN" | reg_id == "FRA" | reg_id
  == "DEU" | reg_id == "GRC" | reg_id == "HUN" | reg_id == "ISL" | reg_id == "IRL" | reg_id == "ISR"
  | reg_id == "ITA" | reg_id == "JPN" | reg_id == "KOR" | reg_id == "LVA" | reg_id == "LTU" |
  reg_id == "LUX" | reg_id == "MEX" | reg_id == "NLD" | reg_id == "NZL" | reg_id == "NOR" | reg_id
  == "POL" | reg_id == "PRT" | reg_id == "SVK" | reg_id == "SVN" | reg_id == "ESP" | reg_id == "SWE"
  | reg_id == "CHE" | reg_id == "TUR" | reg_id == "GBR" | reg_id == "USA"
30 *rename (and relabel) original variables' names to better fit the model
31 rename reg_id country_code
32 label var country_code "OECD member country code"
33 rename value i
34 label var i "Regional innovation rate"
35 *sort by country_code and year
36 sort country_code year -i
37 *we know check for duplicates or disaggregate observations not usefull for our analysis
38 duplicates drop country_code year, force
39 *normal works, take the log
40 generate logi = log(i)
41 label var logi "Natural logarithm of i"
42 *verify the better fit of the log transformation through Kernel density plots (example=USA)
43 kdensity i if country_code == "USA"
44 kdensity logi if country_code == "USA"
45 *declare the data as a panel dataset
46 egen country_id = group(country_code)
47 label var country_id "Time-serie country ID number"
48 xtset country_id year
49 *ultimately manipulate the variables order and sort them
50 order country_id country_code year i
51 sort country_code year
52 save datai, replace
53 browse
54 *display a bar chart to a better visualization of aggregate average value of i per each year
  within the OECD frame
55 egen mean_i = mean(i), by(year)
56 summ mean_i, detail
57 twoway (bar mean_i year, bargap(20)), ///
58       title("OECD average i per year") ///

```

```

59         xtitle("Year") ytitle("Mean i") ///
60         yscale(range(`r(min)' `r(max)'))
61
62 *1) X1=y--->Gross domestic product (GDP)
63 *-----
64 *Data Source: OECD (https://data.oecd.org/gdp/gross-domestic-product-gdp.htm)
65 clear
66 import delimited GDPy
67 describe
68 *save in Stata format
69 save dfy, replace
70 *now, we do not need all the included variables, but only the only usefull ones for hour
  regression or functional to our merging process
71 *first delete those observations going beyond our target time frame
72 keep if time >= 2001 & time <= 2015
73 *keep only OECD members countries
74 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
  == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
  location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
  | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
  "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
  == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
  location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
  | location == "TUR" | location == "GBR" | location == "USA"
75 *some national statistics are displayed in USD_CAP, and some others in MLN_USD depending on data
  availability. Since we hold that MLN_USD is a more wide-spread and immediate way of measurement
  and since we are committed to maintain an uniform unit of scale for each variable, we decided to
  pick this latter and henceforth we will give this assumption as granted when it comes to deal
  with GDP.
76 keep if measure == "MLN_USD"
77 *drop non-involved variables
78 drop indicator measure subject frequency flagcodes
79 *rename and relabel targetted variables
80 rename location country_code
81 label var country_code "OECD member country code"
82 rename time year
83 label var year "Year"
84 rename value y
85 label var y "National GDP"
86 *declare the data as a panel dataset
87 egen country_id = group(country_code)
88 label var country_id "Time-serie country ID number"
89 xtset country_id year
90 *normal works, take the log
91 generate logy = log(y)
92 label var logy "Natural logarithm of y"
93 *verify the better fit of the log transformation through Kernel density plots (example=USA)
94 kdensity y if country_code == "USA"
95 kdensity logy if country_code == "USA"
96 *sort by country_code and year
97 order country_id country_code year y
98 sort country_code year
99 save datay, replace
100 browse
101 *plot results through a line graph for each country
102 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
  "FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
  "MEX" | country_code == "NOR" | country_code == "USA"
103 xtline y, overlay i(country_code) t(year) ///
104         title("GDP by country over time") ///
105         xtitle("Year") ytitle("GDP") ///
106         legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
107                 label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
108                 label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
109                 label(10 "United States (USA)"))
110
111 *2) X2=g--->GDP growth
112 *-----

```

```

113 *Data Source: World Bank (https://data.worldbank.org/indicator/NY.GDP.MKTP.KD.ZG?locations=OE)
114 clear
115 import delimited GDPg
116 describe
117 *save in Stata format
118 save dfg, replace
119 *keep only OECD members countries
120 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
== "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
| location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
"KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
== "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
| location == "TUR" | location == "GBR" | location == "USA"
121 *drop non-involved variables
122 drop indicator measure subject frequency flagcodes
123 *rename and relabel targetted variables
124 rename location country_code
125 label var country_code "OECD member country code"
126 rename time year
127 rename value g
128 label var g "GDP yearly growth rate"
129 *unlike our previous datasets, this dataframe contains not only yearly data, but also
disaggregated quarterly data. With respect to our research purpose, we hold that quarterly data
would not show statistically significant variations if compared to yearly ones (Chetty, Raj, and
Nathaniel Hendren, 2018), hence, we will keep on working with the latter.
130 *avoiding complicated drops function conditional to having a string variable rather than a
numeric one, as in the case of quarterly data expressed as YYYY-Qn, we can just easily drop those
observations containing a Q letter in their respective year's value.
131 drop if strpos(year, "Q") > 0
132 *destring year to make it a numeric variable
133 destring year, generate(year_n)
134 drop year
135 rename year_n year
136 label var year "Year"
137 *first delete those observations going beyond our target time frame
138 keep if year >= 2001 & year <= 2015
139 *declare the data as a panel dataset
140 egen country_id = group(country_code)
141 label var country_id "Time-serie country ID number"
142 xtset country_id year
143 *sort by country_code and year
144 order country_id country_code year g
145 sort country_code year
146 save datag, replace
147 browse
148 *plot results through a line graph for each country
149 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
"FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
"MEX" | country_code == "NOR" | country_code == "USA"
150 xtline g, overlay i(country_code) t(year) ///
151     title("GDP Growth by country over time") ///
152     xtitle("Year") ytitle("GDP growth") ///
153     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
154         label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
155         label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
156         label(10 "United States (USA)"))
157
158 *3) X3=s--->Foreign direct investement (FDI)
159 *-----
160 *Data Source: OECD (https://data.oecd.org/fdi/fdi-flows.htm)
161 clear
162 import delimited FDI
163 describe
164 *save in Stata format
165 save dfs, replace
166 *As for g (GDP growth), this dataframe contains not only yearly data, but also disaggregated

```

```

quarterly data, for the same reasons mentioned above, we will refuse quarterly data, and keep on
working on yearly ones.
167 *avoiding complicated drops function conditional to having a string variable rather than a
numeric one, as in the case of quarterly data expressed as YYYY-Qn, we can just easily drop those
observations containing a Q letter in their respective year's value.
168 drop if strpos(time, "Q") > 0
169 *destring year to make it a numeric variable
170 destring time, generate(year_n)
171 drop time
172 rename year_n year
173 label var year "Year"
174 *now, we do not need all the included variables, but only the only usefull ones for hour
regression or functional to our merging process
175 *first delete those observations going beyond our target time frame
176 keep if year >= 2001 & year <= 2015
177 *keep only OECD members countries
178 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
== "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
| location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
"KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
== "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
| location == "TUR" | location == "GBR" | location == "USA"
179 *drop non-involved variables
180 drop indicator measure subject frequency flagcodes
181 *rename and relabel targetted variables
182 rename location country_code
183 label var country_code "OECD member country code"
184 rename value s
185 *the dataset still appears not perfectly structured since, even though the time variable (year)
should present yearly data for s (FDI), the same year is repeated 4 time, leading us to deduce
that this repetition may refer to data with a quarterly frequency. To validate this hypothesis we
have cross-checked this data with the ones presented by the World Bank (BX.KLT.DINV.CD.WD), and
we actually confirmed our deductions, making a correction to aggregate results a necessary step
to be undertaken. This action must be retained mandatory also considering the nature of our
dataframe (panel). Stata would not be able to declare this latter as a panel df if it is not able
to uniquelyvocably associate each observation to both a country_code and a year.
186 collapse (sum) s, by(country_code year)
187 label var s "Foreign direct investement (FDI)"
188 *declare the data as a panel dataset
189 egen country_id = group(country_code)
190 label var country_id "Time-serie country ID number"
191 xtset country_id year
192 *normal works, take the log
193 *unlike previous variables, s (IDF) can potentially assume both positive and negative values,
making critical to deal with log since the logarithmic function is undefined for non-positive
numbers. To approach this issue we added a constant to "s" before taking the logarithm, small
enough to ensure us to be able to take the log without bringing any significant bias to the
analysis. When, we say "small enough", we mean a portion of the minimum value taken by s in our
data range.
194 egen min_s = min(s)
195 generate log_s = s + abs(min_s) + 1
196 generate logs = log(log_s)
197 label var logs "Natural logarithm of s"
198 drop min_s
199 drop log_s
200 *verify the better fit of the log transformation through Kernel density plots (example=USA)
201 kdensity s if country_code == "USA"
202 kdensity logs if country_code == "USA"
203 *plot results through a line graph for each country
204 tsline s, by(country_code) ///
205     title("") ///
206     xtitle("Year") ytitle("FDI")
207 *sort by country_code and year
208 order country_id country_code year s
209 sort country_code year
210 save datas, replace

```

```

211  browse
212
213  *4) X4=t--->Trade Openess (Trade in Good and Services)
214  *-----
215  *Data Source: OECD (https://data.oecd.org/trade/trade-in-goods-and-services.htm)
216  clear
217  import delimited TGSt
218  describe
219  *save in Stata format
220  save dft, replace
221  *now, we do not need all the included variables, but only the only usefull ones for hour
  regression or functional to our merging process
222  *first delete those observations going beyond our target time frame
223  keep if time >= 2001 & time <= 2015
224  *keep only OECD members countries
225  keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
  == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
  location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
  | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
  "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
  == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
  location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
  | location == "TUR" | location == "GBR" | location == "USA"
226  *since we hold that MLN_USD is a more wide-spread and immediate way of measurement and since we
  are committed to maintain an uniform unit of scale for each variable, we decided to pick this
  latter and henceforth we will give this assumption as granted when it comes to deal with GDP. In
  this sense, we have refused other measurement units such as %GDP.
227  keep if measure == "MLN_USD"
228  *data are also classified by their subject, i.e., if the value represents either exports or
  imports, but it also already displays the net difference between these latter according to the
  2007 notation (exp-imp). As measure of trade openess, we hold that the trade balance is a
  representative and robust index to adress this measure.
229  *keep only trade balances' results
230  keep if subject == "NTRADE"
231  *drop non-involved variables
232  drop indicator measure subject frequency flagcodes
233  *rename and relabel targetted variables
234  rename location country_code
235  label var country_code "OECD member country code"
236  rename time year
237  label var year "Year"
238  rename value t
239  label var t "Trade Balance (TB)"
240  *declare the data as a panel dataset
241  egen country_id = group(country_code)
242  label var country_id "Time-serie country ID number"
243  xtset country_id year
244  *as for s (IDF), t (TB) can potentially assume both positive and negative values, making critical
  to deal with log since the logarithmic function is undefined for non-positive numbers. To
  approach this issue we added a constant to "t" before taking the logarithm, small enough to
  ensure us to be able to take the log without bringing any significant bias to the analysis. When,
  we say "small enough", we mean a portion of the minimum value taken by t in our data range.
245  egen min_t = min(t)
246  generate log_t = t + abs(min_t) + 1
247  generate logt = log(log_t)
248  label var logt "Natural logarithm of t"
249  drop min_t
250  drop log_t
251  *verify the better fit of the log transformation through Kernel density plots (example=USA)
252  kdensity t if country_code == "USA"
253  kdensity logt if country_code == "USA"
254  *sort by coutry_code and year
255  order country_id country_code year t
256  sort country_code year
257  save datat, replace
258  browse
259  *plot results through a line graph for each country
260  keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==

```



```

    "FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
    "MEX" | country_code == "NOR" | country_code == "USA"
261 xtline t, overlay i(country_code) t(year) ///
262     title("Trade openness by country over time") ///
263     xtitle("Year") ytitle("Trade Balance") ///
264     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
265             label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
266             label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
267             label(10 "United States (USA)"))
268
269 *5) X5=e--->Employment rate (LFS)
270 *-----
271 *Data Source: OECD (https://stats.oecd.org/Index.aspx?DataSetCode=LFS\_D)
272 clear
273 import delimited LFS_e
274 describe
275 *save in Stata format
276 save dfe, replace
277 *now, we do not need all the included variables, but only the only usefull ones for hour
278 *regression or functional to our merging process
279 *first delete those observations going beyond our target time frame
280 keep if time >= 2001 & time <= 2015
281 *keep only OECD members countries
282 rename country location
283 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
284 == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
285 location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
286 | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
287 "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
288 == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
289 location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
290 | location == "TUR" | location == "GBR" | location == "USA"
291 *data are also disaggregated by gender, but at this tage we are only interested in aggregated
292 results (MW).
293 *keep only aggregated results
294 keep if sex == "MW"
295 *equally we are interested in the whole labourforce (LF), without furthering micro-level
296 consideration in different age-range clusters
297 keep if v6 == "Total"
298 *same for our target series (E=LFS employment)
299 keep if series == "E"
300 *drop non-involved variables
301 drop v2 sex v4 age v6 series v8 frequency v10 v12 unit unitcode powercodecode powercode
302 referenceperiodcode referenceperiod flagcodes flags
303 *rename and relabel targetted variables
304 rename location country_code
305 label var country_code "OECD member country code"
306 rename time year
307 label var year "Year"
308 rename value e
309 label var e "Employment rate (LFS)"
310 *declare the data as a panel dataset
311 egen country_id = group(country_code)
312 label var country_id "Time-serie country ID number"
313 xtset country_id year
314 *normal works, take the log
315 generate loge = log(e)
316 label var loge "Natural logarithm of e"
317 *verify the better fit of the log transformation through Kernel density plots (example=USA)
318 kdensity e if country_code == "USA"
319 kdensity loge if country_code == "USA"
320 *sort by country_code and year
321 order country_id country_code year e
322 sort country_code year
323 save datae, replace
324 browse
325 *plot results through a line graph for each country
326 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==

```

```

"Fra" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
"MEX" | country_code == "NOR" | country_code == "USA"
316 xtline e, overlay i(country_code) t(year) ///
317 title("Employment rate by country over time") ///
318 xtitle("Year") ytitle("Employment rate (LFS)") ///
319 legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
320 label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
321 label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
322 label(10 "United States (USA)"))
323
324 *6) X6=w--->Gender wage gap indicator (difference between median earnings of men and women
relative to median earnings of men)
325 *-----
-----
326 *Data Source: OECD (https://data.oecd.org/earnwage/gender-wage-gap.htm)
327 clear
328 import delimited GWGw
329 describe
330 *save in Stata format
331 save dfw, replace
332 *now, we do not need all the included variables, but only the only usefull ones for hour
regression or functional to our merging process
333 *first delete those observations going beyond our target time frame
334 keep if time >= 2001 & time <= 2015
335 *keep only OECD members countries
336 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
== "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
| location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
"KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
== "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
| location == "TUR" | location == "GBR" | location == "USA"
337 *moreover, we hold that self-employed women represent outliers in the overall frame of female
labour force, bringing on the table a considerable risk of bias if included in the model with
suitable specifications. Given the limitedness of our data availability, for most pf the
countries listed within our dataframe, data regarding the share of self-employed women over the
whole female labour force were not available, making impossible to run a weighted average to also
include these values in our estimation. Hence, we retain more representative and robust to
consider the employee component, eventhough we signal that these ruled out observations presents
the highest wage gap between male and female individuals.
338 drop if subject == "SELFEMPLOYED"
339 *drop non-involved variables
340 drop indicator subject measure frequency flagcode
341 *rename and relabel targetted variables
342 rename location country_code
343 label var country_code "OECD member country code"
344 rename time year
345 label var year "Year"
346 rename value w
347 label var w "Gender wage gap indicator (difference between median earnings of men and women
relative to median earnings of men)"
348 *declare the data as a panel dataset
349 egen country_id = group(country_code)
350 label var country_id "Time-serie country ID number"
351 xtset country_id year
352 *sort by country_code and year
353 order country_id country_code year w
354 sort country_code year
355 save dataw, replace
356 browse
357 *plot results through a line graph for each country
358 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
"FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
"MEX" | country_code == "NOR" | country_code == "USA"
359 xtline w, overlay i(country_code) t(year) ///
360 title("Gender wage gap by country over time") ///
361 xtitle("Year") ytitle("Gender wage gap") ///

```



```

362     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
363            label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
364            label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
365            label(10 "United States (USA)"))
366
367 *7) X7=f--->Fertility rate (children/woman ratio)
368 *-----
369 *Data Source: OECD (https://data.oecd.org/pop/fertility-rates.htm)
370 clear
371 import delimited FERf
372 describe
373 *save in Stata format
374 save dff, replace
375 *now, we do not need all the included variables, but only the only usefull ones for hour
376 regression or functional to our merging process
377 *first delete those observations going beyond our target time frame
378 keep if time >= 2001 & time <= 2015
379 *keep only OECD members countries
380 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
381 == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
382 location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
383 | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
384 "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
385 == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
386 location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
387 | location == "TUR" | location == "GBR" | location == "USA"
388 *drop non-involved variables
389 drop indicator subject measure frequency flagcode
390 *rename and relabel targetted variables
391 rename location country_code
392 label var country_code "OECD member country code"
393 rename time year
394 label var year "Year"
395 rename value f
396 label var f "Fertility rate (children/woman)"
397 *declare the data as a panel dataset
398 egen country_id = group(country_code)
399 label var country_id "Time-serie country ID number"
400 xtset country_id year
401 *sort by coutry_code and year
402 order country_id country_code year f
403 sort country_code year
404 save dataf, replace
405 browse
406 *plot results through a line graph for each country
407 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
408 "FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
409 "MEX" | country_code == "NOR" | country_code == "USA"
410 xtline f, overlay i(country_code) t(year) ///
411     title("Fertility rate by country over time") ///
412     xtitle("Year") ytitle("children/woman ratio") ///
413     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
414            label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
415            label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
416            label(10 "United States (USA)"))
417
418 *8) X8=r--->R&D spending (public spending on research and development as percentage of GDP)
419 *-----
420 *Data Source: OECD (https://data.oecd.org/rd/gross-domestic-spending-on-r-d.htm)
421 clear
422 import delimited RDSr
423 describe
424 *save in Stata format
425 save dfr, replace
426 *now, we do not need all the included variables, but only the only usefull ones for hour
427 regression or functional to our merging process
428 *first delete those observations going beyond our target time frame
429 keep if time >= 2001 & time <= 2015

```

```

419 *keep only OECD members countries
420 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
    == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
    location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
    | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
    "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
    == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
    location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
    | location == "TUR" | location == "GBR" | location == "USA"
421 *henceforth to evaluate public spending in each sector we will consider percentage of GDP as unit
    of scale, since our demand scope is to evaluate whether a greater either government or private
    (aggreagate) focus on a specific public sector may impact an individual's chances of being
    granted of a patent, and not whether this impact is due to higher absolute values where obviously
    more populated and richer countries would be advantaged.
422 keep if measure == "PC_GDP"
423 *drop non-involved variables
424 drop indicator subject measure frequency flagcode
425 *rename and relabel targetted variables
426 rename location country_code
427 label var country_code "OECD member country code"
428 rename time year
429 label var year "Year"
430 rename value r
431 label var r "R&D spending (public spending on research and development as percentage of GDP)"
432 *declare the data as a panel dataset
433 egen country_id = group(country_code)
434 label var country_id "Time-serie country ID number"
435 xtset country_id year
436 *sort by coutry_code and year
437 order country_id country_code year r
438 sort country_code year
439 save datar, replace
440 browse
441 *plot results through a line graph for each country
442 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
    "FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
    "MEX" | country_code == "NOR" | country_code == "USA"
443 xtline r, overlay i(country_code) t(year) ///
    title("R&D spending by country over time") ///
    xttitle("Year") ytitle("R&D spending/GDP") ///
    legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
    label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
    label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
    label(10 "United States (USA)"))
450
451 *9) X9=u--->Education spending (public spending on education as percentage of GDP)
452 *-----
453 *Data Source: OECD (https://data.oecd.org/eduresource/public-spending-on-education.htm)
454 clear
455 import delimited EDUu
456 describe
457 *save in Stata format
458 save dfu, replace
459 *now, we do not need all the included variables, but only the only usefull ones for hour
    regression or functional to our merging process
460 *first delete those observations going beyond our target time frame
461 keep if time >= 2001 & time <= 2015
462 *keep only OECD members countries
463 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
    == "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
    location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
    | location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
    "KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
    == "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
    location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
    | location == "TUR" | location == "GBR" | location == "USA"
464 *we take in account the PRY_NTRY index, since according to the current literature, it is the more
    reliable and widespread pre-processed index to measure public spending in education.

```

```

465 keep if subject == "PRY_NTRY"
466 *drop non-involved variables
467 drop indicator subject measure frequency flagcode
468 *rename and relabel targetted variables
469 rename location country_code
470 label var country_code "OECD member country code"
471 rename time year
472 label var year "Year"
473 rename value u
474 label var u "Education spending (public spending on education as percentage of GDP)"
475 *declare the data as a panel dataset
476 egen country_id = group(country_code)
477 label var country_id "Time-serie country ID number"
478 xtset country_id year
479 *sort by coutry_code and year
480 order country_id country_code year u
481 sort country_code year
482 save datau, replace
483 browse
484 *plot results through a line graph for each country
485 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
"Fra" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
"MEX" | country_code == "NOR" | country_code == "USA"
486 xtline u, overlay i(country_code) t(year) ///
487     title("Education spending by country over time") ///
488     xtitle("Year") ytitle("Education spending/GDP") ///
489     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
490         label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
491         label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
492         label(10 "United States (USA)"))
493
494 *10) X10=h--->Healthcare system spending (public spending on healthcare as percentage of GDP)
495 *-----
496 *Data Source: OECD (https://data.oecd.org/healthres/health-spending.htm)
497 clear
498 import delimited HEAh
499 describe
500 *save in Stata format
501 save dfh, replace
502 *now, we do not need all the included variables, but only the only usefull ones for hour
regression or functional to our merging process
503 *first delete those observations going beyond our target time frame
504 keep if time >= 2001 & time <= 2015
505 *keep only OECD members countries
506 keep if location == "AUS" | location == "AUT" | location == "BEL" | location == "CAN" | location
== "CHL" | location == "CZE" | location == "DNK" | location == "EST" | location == "FIN" |
location == "FRA" | location == "DEU" | location == "GRC" | location == "HUN" | location == "ISL"
| location == "IRL" | location == "ISR" | location == "ITA" | location == "JPN" | location ==
"KOR" | location == "LVA" | location == "LTU" | location == "LUX" | location == "MEX" | location
== "NLD" | location == "NZL" | location == "NOR" | location == "POL" | location == "PRT" |
location == "SVK" | location == "SVN" | location == "ESP" | location == "SWE" | location == "CHE"
| location == "TUR" | location == "GBR" | location == "USA"
507 *we consider the total aggregate value over national GDP
508 keep if subject == "TOT"
509 *as for r and u (R&D spending and education spending over GDP) we take %GDP as unit of scale for
the same reasons above mentioned
510 keep if measure == "PC_GDP"
511 *drop non-involved variables
512 drop indicator subject measure frequency flagcode
513 *rename and relabel targetted variables
514 rename location country_code
515 label var country_code "OECD member country code"
516 rename time year
517 label var year "Year"
518 rename value h
519 label var h "Healthcare system spending (public spending on healthcare as percentage of GDP)"
520 *declare the data as a panel dataset
521 egen country_id = group(country_code)

```

```

522 label var country_id "Time-serie country ID number"
523 xtset country_id year
524 *sort by country_code and year
525 order country_id country_code year h
526 sort country_code year
527 save datah, replace
528 browse
529 *plot results through a line graph for each country
530 keep if country_code == "AUS" | country_code == "CHL" | country_code == "DEU" | country_code ==
    "FRA" | country_code == "GBR" | country_code == "ITA" | country_code == "JPN" | country_code ==
    "MEX" | country_code == "NOR" | country_code == "USA"
531 xtline h, overlay i(country_code) t(year) ///
532     title("Healthcare spending by country over time") ///
533     xtitle("Year") ytitle("Healthcare spending/GDP") ///
534     legend(label(1 "Australia (AUS)") label(2 "Chile (CHL)") label(3 "Germany (DEU)") ///
535         label(4 "France (FRA)") label(5 "United Kingdom (GBR)") label(6 "Italy (ITA)") ///
536         label(7 "Japan (JPN)") label(8 "Mexico (MEX)") label(9 "Norway (NOR)") ///
537         label(10 "United States (USA)"))
538
539 ///////////////////////////////////////////////////
540 ///////////////////////////////////////////////////Final merging////////////////////////////////////
541 ///////////////////////////////////////////////////
542
543 *in this section we progressively merge each dataset to work with a unique dataframe enclosing
    all our variables and observations in a panel format
544 use datai, clear
545 merge 1:1 country_id year using "datay"
546 dropmerge
547 save dataa, replace
548 *we repeat the same procedure for all the disaggregate dataframes referred to each variable
549 use dataa, clear
550 merge 1:1 country_id year using "datag"
551 dropmerge
552 save data2, replace
553
554 use data2, clear
555 merge 1:1 country_id year using "datas"
556 dropmerge
557 save dataa3, replace
558
559 use dataa3, clear
560 merge 1:1 country_id year using "datat"
561 dropmerge
562 save dataa4, replace
563
564 use dataa4, clear
565 merge 1:1 country_id year using "datae"
566 dropmerge
567 save dataa5, replace
568
569 use dataa5, clear
570 merge 1:1 country_id year using "dataw"
571 dropmerge
572 save dataa6, replace
573
574 use dataa6, clear
575 merge 1:1 country_id year using "dataf"
576 dropmerge
577 save dataa7, replace
578
579 use dataa7, clear
580 merge 1:1 country_id year using "datar"
581 dropmerge
582 save dataa8, replace
583
584 use dataa8, clear
585 merge 1:1 country_id year using "datau"
586 dropmerge

```

```

587 save data9, replace
588
589 use data9, clear
590 merge 1:1 country_id year using "datah"
591 drop _merge
592 save merged_data, replace
593 describe
594 browse
595 *declare the data as a panel dataset
596 xtset country_id year
597
598 ///////////////////////////////////////////////////
599 //Missing values regressive forecasts ///////////////////////////////////////////////////
600 ///////////////////////////////////////////////////
601
602 *since we have missing values for some variables, we employ a panel random effect estimator for
each of these latter to generate a prediction on both possible past and future observations,
which, according to Hausman test, provides the better estimation for our unadjusted data.
Notwithstanding our large number of involved variables, we are only able to employ those
presenting a complete range of observations over time, otherwise the regression would keep
predicting missing results. It is worth to notice that for some countries data are complete also
for variables not listed in our estimation model, but in order to avoid over-noise-capturing
bias, that would then possibly affect only these exceptional countries, we have opted for a
uniformed prediction model, exclusively including those complete variables common to all countries.
603
604 use merged_data, clear
605 xtset country_id year
606
607 *1) predicted s
608 *-----
609 xtreg s f logt g logy, re
610 predict pred_s
611 *Replace missing values in 's' with predicted values
612 replace s = pred_s if missing(s)
613 drop pred_s
614 describe
615 *in case of log transformed variables, we also need to retake the log function for the new
predicted values. We hold that doing a 2-stage log transformation makes the regressive forecast
more robust to possible heteroschedacity or high-skewness possible biases.
616 *Redo log for new predicted values
617 egen max_s = max(s)
618 generate log_s1 = t + abs(max_s) + 1
619 generate log_s = log(log_s1)
620 label var log_s "Natural logarithm of e"
621 drop max_s
622 drop log_s1
623 drop logs
624
625 *2) predicted e
626 *-----
627 xtreg e f logt g logy, re
628 predict pred_e
629 *Replace missing values in 'e' with predicted values
630 replace e = pred_e if missing(e)
631 drop pred_e
632 describe
633 *Redo log for new predicted values
634 egen max_e = max(e)
635 generate log_e1 = t + abs(max_e) + 1
636 generate log_e = log(log_e1)
637 label var log_e "Natural logarithm of e"
638 drop max_e
639 drop log_e1
640 drop loge
641
642 *3) predicted w
643 *-----
644 xtreg w f logt g logy, re

```

```

645 predict pred_w
646 *Replace missing values in 'w' with predicted values
647 replace w = pred_w if missing(w)
648 drop pred_w
649
650 *4) predicted r
651 *-----
652 xtreg r f logt g logy, re
653 predict pred_r
654 *Replace missing values in 'w' with predicted values
655 replace r = pred_r if missing(r)
656 drop pred_r
657
658 *5) predicted u
659 *-----
660 xtreg u f logt g logy, re
661 predict pred_u
662 *Replace missing values in 'w' with predicted values
663 replace u = pred_u if missing(u)
664 drop pred_u
665
666 *uniformly rename variables
667 rename logi log_i
668 rename logy log_y
669 rename logt log_t
670 *adjust the order of variables to better reflect our format
671 *sort by country_code and year
672 order country_id country_code year i log_i y log_y g s log_s t log_t e log_e w f r u h
673 sort country_code year
674 save pred_data, replace
675 *declare the data as a panel dataset
676 xtset country_id year
677 browse
678
679 ///////////////////////////////////////////////////////////////////
680 ///////////////////////////////////////////////////////////////////Stationarity Tests/////////////////////////////////////////////////////////////////
681 ///////////////////////////////////////////////////////////////////
682
683 *At this point, stationarity tests are crucial to validate that statistical properties like mean
and variance remain constant over time throughout our dataframe. Ensuring stationarity is vital for
accurate modeling, forecasting, and reliable statistical inference. Non-stationary data may
mislead models and compromise forecasting accuracy. Overall, stationarity verification is
essential for maintaining the stability and interpretability of time series models. To
investigate the heterogeneous range of non-stationarity biases we have involved four different
tests: Im-Pesaran-Shin (IPS) test, Breitung Panel Unit Root Test, Hadri LM Test and Levin, Lin,
and Chu (LLC) Test. It is worth to notice that in our case, traditional approaches, such as
Augmented Dickey-Fuller test, are not available given the multitude of panel belonging to the
dataframe.
684
685 *1) stationarity analysis for y
686 *-----
687 *The Im-Pesaran-Shin (IPS) test is assessing whether the target variable is stationary across
individual time series. It provides insights into the common stochastic trend shared by the panel.
688 xtunitroot ips y
689 *The Im-Pesaran-Shin (IPS) unit-root test results for the variable "y" indicate a positive test
statistic of 8.8557 and a p-value of 1.0000. These values fail to reject the null hypothesis of
unit roots, suggesting that "y" is likely non-stationary across the panel.
690
691 *The Breitung Panel Unit Root Test examines stationarity in the target variable, considering
cross-section dependence. It's robust in the presence of correlated data among entities.
692 xtunitroot breitung y
693 *The test yielded a test statistic of 13.5466 and a p-value of 1.0000. These results indicate a
failure to reject the null hypothesis, suggesting that y is likely non-stationary across the panel.
694
695 *The Hadri LM Test assesses stationarity in the target variable for panel data, particularly
considering the impact of cross-sectional dependence on the results.
696 xtunitroot hadri y
697

```



```

698 *The Levin, Lin, and Chu (LLC) Test evaluates stationarity in the target variable in panel data
699 while accounting for both individual effects and cross-section dependence.
700 xtunitroot llc y
701 *The Hadri LM test yielded a highly significant test statistic of 47.9938 and a p-value of
702 0.0000. This strong evidence led to the rejection of the null hypothesis that all panels are
703 stationary, indicating that variable "y" is likely non-stationary across the panel
704
705 *2) stationarity analysis for g
706 *-----
707 *IPS test
708 xtunitroot ips g
709 *The test resulted in a test statistic of -6.6116 and a p-value of 0.0000. These findings lead to
710 the rejection of the null hypothesis, providing evidence that the variable "g" is likely
711 non-stationary across the panel.
712
713 *Breitung Panel Unit Root Test
714 xtunitroot breitung g
715 *The test yielded a test statistic of -9.6198 and a p-value of 0.0000. These results indicate a
716 rejection of the null hypothesis, providing evidence that the variable "g" is likely
717 non-stationary across the panel.
718
719 *Hadri LM Test
720 xtunitroot hadri g
721 *The test yielded a test statistic of 4.6573 and a p-value of 0.0000. These results indicate a
722 rejection of the null hypothesis, providing evidence that the variable "g" is likely
723 non-stationary across the panel.
724
725 *LLC test
726 xtunitroot llc g
727 *These results indicate a rejection of the null hypothesis, providing strong evidence that the
728 variable "g" is likely non-stationary across the panel
729
730 *3) stationarity analysis for s
731 *-----
732 *IPS test
733 xtunitroot ips s
734 *The test resulted in a test statistic of -3.5428 and a p-value of 0.0002. These findings lead to
735 the rejection of the null hypothesis, providing evidence that the variable "s" is likely
736 non-stationary across the panel.
737
738 *Breitung Panel Unit Root Test
739 xtunitroot breitung s
740 *The test yielded a test statistic of -2.5754 and a p-value of 0.0050. These results indicate a
741 rejection of the null hypothesis, providing evidence that the variable "s" is likely
742 non-stationary across the panel.
743
744 *Hadri LM Test
745 xtunitroot hadri s
746 *The test yielded a test statistic of 19.2455 and a p-value of 0.0000. These results strongly
747 support the rejection of the null hypothesis, indicating that the variable "s" is likely
748 non-stationary across the panel.
749
750 *LLC test
751 xtunitroot llc s
752 *These results strongly indicate a rejection of the null hypothesis, providing evidence that the
753 variable "s" is likely non-stationary across the panel.
754
755 *4) stationarity analysis for t
756 *-----
757 *IPS test
758 xtunitroot ips t
759 *The test resulted in a test statistic of -1.1581 and a p-value of 0.9861. The critical values
760 for rejection are -1.830 (1%), -1.740 (5%), and -1.690 (10%). Given that the test statistic is
761 less extreme than the critical values, we fail to reject the null hypothesis. This indicates that
762 variable "t" is likely non-stationary across the panel.
763
764 *Breitung Panel Unit Root Test
765 xtunitroot breitung t

```

```

746  *The test statistic yielded a value of 1.3425 with a corresponding p-value of 0.9103. Given that
the test statistic is not extreme and the p-value is high, there is insufficient evidence to
reject the null hypothesis. This indicates that variable "t" is likely non-stationary across the
panel.
747
748  *Hadri LM test
749  xtunitroot hadri t
750  *The test statistic yielded a value of 24.9452 with a corresponding p-value of 0.0000. The high
test statistic and very low p-value provide strong evidence to reject the null hypothesis. This
suggests that variable "t" is likely non-stationary across the panel.
751
752  *LLC test
753  xtunitroot llc t
754  *The test statistic resulted in a value of -9.1361 with a corresponding p-value of 0.0000. The
negative test statistic and very low p-value provide strong evidence to reject the null
hypothesis, indicating that variable "t" is likely non-stationary across the panel.
755
756  *5) stationarity analysis for e
757  *-----
758  *IPS test
759  xtunitroot ips e
760  *The test statistics did not provide sufficient evidence to reject the null hypothesis, as the
p-value associated with Z-t-tilde-bar was 0.6817, exceeding the common significance levels.
Consequently, there is insufficient evidence to conclude that variable "e" is stationary based on
the IPS test, implying the potential presence of unit roots.
761
762  *Breitung Panel Unit Root Test
763  xtunitroot breitung e
764  *The test statistic lambda was -0.0312, and the p-value was 0.4876. The p-value exceeds common
significance levels, indicating that there is insufficient evidence to reject the null
hypothesis. Consequently, based on the Breitung test, variable "e" may contain unit roots,
implying non-stationarity across panels.
765
766  *Hadri LM test
767  xtunitroot hadri e
768  *The test statistic, lambda, yielded a value of -0.0312, with an associated p-value of 0.4876.
Given that the p-value exceeds common significance levels (e.g., 0.05), there is insufficient
evidence to reject the null hypothesis. Consequently, based on the Breitung test, variable "e"
may exhibit unit roots, indicating non-stationarity across panels.
769
770  *LLC test
771  xtunitroot llc e
772  *The test results yielded a significant test statistic, with an unadjusted t-value of -9.2830 and
an adjusted t-value of -4.2159, both implying strong evidence against the null hypothesis.
Therefore, based on the LLC test, variable "e" appears to be stationary across panels, suggesting
a lack of unit roots.
773
774  *In evaluating the reliability of the unit-root tests for variable "e," the Levin-Lin-Chu (LLC)
test stands out as more trustworthy. The LLC test incorporates a lagged regression and addresses
heteroskedasticity through the Bartlett kernel, providing a comprehensive assessment of
stationarity. The notable and negative adjusted t-value in the LLC test signals a rejection of
the null hypothesis, supporting the argument for stationarity.
775
776  *6) stationarity analysis for w
777  *-----
778  *IPS test
779  xtunitroot ips w
780  *The Im-Pesaran-Shin (IPS) unit-root test results for variable "w" indicate a significant test
statistic, with a Z-t-tilde-bar value of -6.9155 and a p-value of 0.0000. This suggests evidence
against the null hypothesis, implying non-stationarity for the variable.
781
782  *Breitung Panel Unit Root Test
783  xtunitroot breitung w
784  *The test statistic, lambda, was calculated as -0.4973 with a p-value of 0.3095. This result does
not provide sufficient evidence to reject the null hypothesis, indicating that variable "w" may
have unit roots, suggesting non-stationarity.
785
786  *Hadri LM test

```



```

787 xtunitroot hadri w
788 *The test yielded a significant statistic, with a z-value of 5.5229 and a p-value of 0.0000. This
strong evidence against the null hypothesis suggests that variable "w" is likely stationary
across panels, indicating a lack of unit roots.

789
790 *LLC test
791 xtunitroot llc w
792 *The test yielded a significant unadjusted t-value of -7.3412 and an adjusted t-value of -3.2765,
providing strong evidence against the null hypothesis. This implies that variable "w" appears to
be stationary across panels, indicating a lack of unit roots.

793
794 *Given the results of the Im-Pesaran-Shin (IPS), Breitung, Hadri LM, and Levin-Lin-Chu (LLC)
unit-root tests for variable "w," the significant p-value in the Hadri LM test (p-value = 0.0000)
and the adjusted t-value in the LLC test (t = -3.2765) suggest strong evidence against the null
hypothesis of unit roots, indicating that the variable "w" is likely stationary across panels.
Therefore, based on these findings, it is reasonable to conclude that variable "w" is stationary
in the panel dataset.

795
796 *7) stationarity analysis for f
797 *-----
798 *IPS test
799 xtunitroot ips f
800 *Based on the results of the Im-Pesaran-Shin (IPS) unit-root test for variable "f," the test
statistic (Z-t-tilde-bar) is 2.5333 with a p-value of 0.9944. These values suggest weak evidence
against the null hypothesis of unit roots, indicating that variable "f" is likely non-stationary
across panels.

801
802 *Breitung Panel Unit Root Test
803 xtunitroot breitung f
804 *The Breitung unit-root test for variable "f" yielded a test statistic (lambda) of 2.6238 with a
p-value of 0.9957. These results provide insufficient evidence to reject the null hypothesis,
suggesting that variable "f" likely contains unit roots and is non-stationary across panels.

805
806 *Hadri LM test
807 xtunitroot hadri f
808 *The Hadri LM test for variable "f" produced a test statistic of 29.8304 with a p-value of
0.0000. The highly significant p-value suggests strong evidence against the null hypothesis,
indicating that variable "f" is likely stationary across panels.

809
810 *LLC test
811 xtunitroot llc f
812 *The Levin-Lin-Chu unit-root test (LLC) for variable "f" resulted in a test statistic of -4.9840
and an adjusted t-value of -0.2989, with a p-value of 0.3825. The unadjusted and adjusted
t-values both fail to reject the null hypothesis of unit roots, suggesting a lack of stationarity.

813
814 *The Hadri LM test accounts for potential violations of homoskedasticity assumptions, making it
suitable for this dataset. Additionally, the highly significant test statistic and low p-value
reinforce its credibility in rejecting the null hypothesis of unit roots, leading to the
conclusion of stationarity for variable "f".

815
816 *8) stationarity analysis for r
817 *-----
818 *IPS test
819 xtunitroot ips r
820 *The test results reveal a test statistic of -1.3947 and a corresponding p-value of 0.9999. The
high p-value fails to provide evidence against the null hypothesis, suggesting a lack of
stationarity for variable "r."

821
822 *Breitung Panel Unit Root Test
823 xtunitroot breitung r
824 *The test yielded a statistic of 2.9776 and a corresponding p-value of 0.9985. The high p-value
fails to provide evidence against the null hypothesis, suggesting a lack of stationarity for
variable "r".

825
826 *Hadri LM test
827 xtunitroot hadri r
828 *The test resulted in a significant statistic of 29.1043 with a p-value of 0.0000, providing
strong evidence against the null hypothesis. This indicates that variable "r" is likely to be

```

non-stationary across panels, suggesting the presence of unit roots.

829

830 \*LLC test

831 `xtunitroot llc r`

832 \*The Levin-Lin-Chu (LLC) unit-root test for variable "r" in the panel dataset resulted in a non-significant test statistic, with an adjusted t-value of -0.9005 and a p-value of 0.1839. This suggests insufficient evidence to reject the null hypothesis that panels contain unit roots, implying a lack of clear stationarity for variable "r" across panels.

833

834 \*9) stationarity analysis for u

835 \*-----

836 \*IPS test

837 `xtunitroot ips u`

838 \*The test results yielded a significant test statistic, with a Z-t-tilde-bar value of -3.4928 and a p-value of 0.0002, providing strong evidence against the null hypothesis. Therefore, "u" appears to be stationary across panels, suggesting a lack of unit roots.

839

840 \*Breitung Panel Unit Root Test

841 `xtunitroot breitung u`

842 \*The test results produced a highly significant test statistic, with a lambda value of -5.7768 and a p-value of 0.0000, providing strong evidence against the null hypothesis. Therefore, "u" appears to be stationary across panels, indicating a lack of unit roots.

843

844 \*Hadri LM test

845 `xtunitroot hadri u`

846 \*The test results revealed a significant test statistic, with a z-value of 17.4116 and a p-value of 0.0000, providing robust evidence against the null hypothesis. Therefore, "u" is deemed to be stationary across panels, implying the absence of unit roots.

847

848 \*LLC test

849 `xtunitroot llc u`

850 \*The test results yielded a significant test statistic, with an unadjusted t-value of -11.6567 and an adjusted t-value of -4.8094, both implying strong evidence against the null hypothesis. Therefore, "u" appears to be stationary across panels, suggesting a lack of unit roots.

851

852

853 \*10) stationarity analysis for h

854 \*-----

855 \*IPS test

856 `xtunitroot ips h`

857 \*The test did not yield a statistically significant result, as the test statistic did not exceed the critical values at conventional significance levels (1%, 5%, 10%). Therefore, we do not reject the null hypothesis, suggesting that variable "h" may exhibit unit roots across the panels.

858

859 \*Breitung Panel Unit Root Test

860 `xtunitroot breitung h`

861 \*The test statistic, represented by lambda, was found to be 3.0187 with a p-value of 0.9987. As the p-value exceeds conventional significance levels, we fail to reject the null hypothesis. Therefore, "h" is likely to have unit roots across the panels, indicating non-stationarity.

862

863 \*Hadri LM test

864 `xtunitroot hadri h`

865 \*The test statistic, represented by "z," was found to be 37.4329 with a p-value of 0.0000, indicating strong evidence against the null hypothesis. Therefore, "h" is likely to be stationary across all panels, suggesting the absence of unit roots.

866

867 \*LLC test

868 `xtunitroot llc h`

869 \*The test results yielded a significant test statistic, with an unadjusted t-value of -7.2781 and an adjusted t-value of -3.5677, both implying strong evidence against the null hypothesis. Therefore, based on the LLC test, variable "h" appears to be stationary across panels, suggesting a lack of unit roots.

870

871 \*Im-Pesaran-Shin (IPS) and Breitung tests do not provide conclusive evidence in this case to clearly assert whether the variable can be retained stationary or not. Therefore, based on the more reliable Hadri LM and LLC tests, we can infer that variable "h" is likely stationary across panels, as the significance of their results shows.

872

```

873 ///////////////////////////////////////////////////////////////////
874 ///////////////////////////////////////////////////////////////////Stationarity Correction/////////////////////////////////////////////////////////////////
875 ///////////////////////////////////////////////////////////////////
876
877 *Given the results in terms of stationarity extrapolated from the employed stationarity tests, we
want now to implement corrections to those variables that were identified as non-stationary. To
do so, we will recur to first differences, following the rationale of eliminating trends and
making the series more amenable to statistical analyses.

878
879 *recall the unadjusted dataset
880 use pred_data, clear
881 xtset country_id year
882
883 *1) y adjusted
884 *-----
885 *generate first differences
886 gen y_diff = y - L.y
887 *drop the original variable if you want to keep only the differenced variable
888 drop y
889 *rename the differenced variable to the original variable name
890 rename y_diff y
891 *LLC test
892 xtunitroot llc y
893 *The differenced variable 'y' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
894 *repredict missing values to have complete data from 2001 to 2015
895 xtreg y g log_s log_t log_e w f r, re
896 predict pred_y
897 *Replace missing values in with predicted values
898 replace y = pred_y if missing(y)
899 drop pred_y
900 label var y "GDP"
901
902 *2) log_y adjusted
903 *-----
904 *generate first differences
905 gen log_y_diff = log_y - L.log_y
906 *drop the original variable if you want to keep only the differenced variable
907 drop log_y
908 *rename the differenced variable to the original variable name
909 rename log_y_diff log_y
910 *LLC test
911 xtunitroot llc log_y
912 *The differenced variable 'log_y' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
913 *repredict missing values to have complete data from 2001 to 2015
914 xtreg log_y g log_s log_t log_e w f r u h, re
915 predict pred_log_y
916 *Replace missing values in 's' with predicted values
917 replace log_y = pred_log_y if missing(log_y)
918 drop pred_log_y
919 label var log_y "Naural logarithm of y (GDP)"
920
921 *3) g adjusted
922 *-----
923 *generate first differences
924 gen g_diff = g - L.g
925 *drop the original variable if you want to keep only the differenced variable
926 drop g
927 *rename the differenced variable to the original variable name
928 rename g_diff g
929 *LLC test
930 xtunitroot llc g
931 *The differenced variable 'g' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
932 *repredict missing values to have complete data from 2001 to 2015
933 xtreg g log_y log_s log_t log_e w f r u h, re
934 predict pred_g

```

```

935 *Replace missing values in with predicted values
936 replace g = pred_g if missing(g)
937 drop pred_g
938 label var g "GDP yearly growth rate"
939
940 *4) s adjusted
941 *-----
942 *generate first differences
943 gen s_diff = s - L.s
944 *drop the original variable if you want to keep only the differenced variable
945 drop s
946 *rename the differenced variable to the original variable name
947 rename s_diff s
948 *LLC test
949 xtunitroot llc s
950 *The differenced variable 's' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
951 *repredict missing values to have complete data from 2001 to 2015
952 xtreg s log_y g log_t log_e w f r u h, re
953 predict pred_s
954 *Replace missing values in with predicted values
955 replace s = pred_s if missing(s)
956 drop pred_s
957 label var s "Foreign direct investement (FDI)"
958
959 *5) log_s adjusted
960 *-----
961 *generate first differences
962 gen log_s_diff = log_s - L.log_s
963 *drop the original variable if you want to keep only the differenced variable
964 drop log_s
965 *rename the differenced variable to the original variable name
966 rename log_s_diff log_s
967 *LLC test is not applicable in this case, since it requires strongly balanced data. We will then
try to deploy IPS test
968 *xtunitroot ips log_s
969 *None of the previous, and further, tests actually produces statistically significant results for
this variable analysis, leading to uncertain conclusions on the effectiveness of the differencing
transformation in achieving stationarity.
970 *repredict missing values to have complete data from 2001 to 2015
971 xtreg log_s g log_y log_t log_e w f r u h, re
972 predict pred_log_s
973 *Replace missing values in 's' with predicted values
974 replace log_s = pred_log_s if missing(log_s)
975 drop pred_log_s
976 label var log_s "Naural logarithm of s (FDI)"
977
978 *6) t adjusted
979 *-----
980 *generate first differences
981 gen t_diff = t - L.t
982 *drop the original variable if you want to keep only the differenced variable
983 drop t
984 *rename the differenced variable to the original variable name
985 rename t_diff t
986 *LLC test
987 xtunitroot llc t
988 *The differenced variable 't' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
989 *repredict missing values to have complete data from 2001 to 2015
990 xtreg t log_y g log_s log_e w f r u h, re
991 predict pred_t
992 *Replace missing values in with predicted values
993 replace t = pred_t if missing(t)
994 drop pred_t
995 label var t "Trade Balance (TB)"
996
997 *7) log_t adjusted

```

```

998 *-----
999 *generate first differences
1000 gen log_t_diff = log_t - L.log_t
1001 *drop the original variable if you want to keep only the differenced variable
1002 drop log_t
1003 *rename the differenced variable to the original variable name
1004 rename log_t_diff log_t
1005 *LLC test
1006 xtunitroot llc log_t
1007 *The differenced variable 'log_t' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
1008 *repredict missing values to have complete data from 2001 to 2015
1009 xtreg log_t g log_y log_s log_e w f r u h, re
1010 predict pred_log_t
1011 *Replace missing values in 's' with predicted values
1012 replace log_t = pred_log_t if missing(log_t)
1013 drop pred_log_t
1014 label var log_t "Naural logarithm of t (TB)"
1015
1016 *8) r adjusted
1017 *-----
1018 *generate first differences
1019 gen r_diff = r - L.r
1020 *drop the original variable if you want to keep only the differenced variable
1021 drop r
1022 *rename the differenced variable to the original variable name
1023 rename r_diff r
1024 *LLC test
1025 xtunitroot llc r
1026 *The differenced variable 'r' appears to be stationary across panels, supporting the
effectiveness of the differencing transformation in achieving stationarity.
1027 *repredict missing values to have complete data from 2001 to 2015
1028 xtreg r log_y g log_s log_t log_e w f u h, re
1029 predict pred_r
1030 *Replace missing values in with predicted values
1031 replace r = pred_r if missing(r)
1032 drop pred_r
1033 label var r "R&D spending (public spending on research and development as percentage of GDP)"
1034
1035 *reorder the final dataframe according to the new changes
1036 *sort by coutry_code and year
1037 order country_id country_code year i log_i y log_y g s log_s t log_t e log_e w f r u h
1038 sort country_code year
1039 save adj_data, replace
1040 *declare the data as a panel dataset
1041 xtset country_id year
1042 browse
1043
1044 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1045 ///////////////////////////////////////////////////////////////////Model Specification/////////////////////////////////////////////////////////////////
1046 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1047
1048 use adj_data, clear
1049 xtset country_id year
1050 ssc install outreg2
1051 ssc install estout
1052
1053 *-----
1054 *Specification_1----->Macro-trends
1055 xtreg i log_y g log_s log_t log_e, re
1056 predict pred_i1
1057 outreg2 using "Specification_1", replace ctitle(Baseline)
1058 est store spec1
1059
1060 *Specification_2----->Female empowerment
1061 xtreg i log_y g log_s log_t log_e w f, re
1062 predict pred_i2
1063 outreg2 using "Specification_2", replace ctitle(Baseline)

```

```

1064 est store spec2
1065
1066 *Specification_3----->Government expenditure targets
1067 xtreg i log_y g log_s log_t log_e w f r u h, re
1068 predict pred_i3
1069 outreg2 using "Specification_3", replace ctitle(Baseline)
1070 est store spec3
1071 *-----
1072 *export results
1073 esttab spec1 spec2 spec3 using "raw_res.csv", se label replace wide plain
1074
1075 //////////////////////////////////////
1076 //////////////////////////////////////Panel data specifications////////////////////////////////////
1077 //////////////////////////////////////
1078
1079 use adj_data, clear
1080 xtset country_id year
1081
1082 *1) Introductory setup
1083 *-----
1084 *we deploy this test to investigate which panel estimator works better with our data
1085 *FIXED EFFECTS----->without clustered SEs
1086 xtreg i log_y g log_s log_t log_e w f r u h , fe
1087 estimates store FE
1088 *FIXED EFFECTS----->with clustered SEs
1089 xtreg i log_y g log_s log_t log_e w f r u h , fe cluster(country_id)
1090 estimates store FECL
1091 *RANDOM EFFECTS----->without robust SEs
1092 xtreg i log_y g log_s log_t log_e w f r u h , re
1093 estimates store RE
1094 *RANDOM EFFECTS----->with robust SEs
1095 xtreg i log_y g log_s log_t log_e w f r u h , re robust
1096 estimates store REB
1097 *POOLED OLS
1098 reg log_y g log_s log_t log_e w f r u h
1099 estimates store OLS
1100 *Pooled OLS is appropriate when you treat the panel data as if it were a single cross-sectional
dataset, ignoring any individual-specific or time-specific effects. This assumes that there is no
correlation between the individual effects and the independent variables, which is not our case.
1101
1102 *comparision of different estimated results
1103 estimates table FECL REB REBS
1104
1105 *2) Hausman test
1106 *-----
1107 hausman FE RE
1108 *The results of the Hausman test indicate that the p-value is very high (0.9992), suggesting that
you fail to reject the null hypothesis. The null hypothesis in the Hausman test is that the
difference in coefficients between the fixed-effects (FE) and random-effects (RE) models is not
systematic. Since the p-value is high, there is no strong evidence to suggest that the
differences in coefficients between the two models are systematic. Therefore, we choose the
random-effects model as it is consistent with the assumptions of the test.
1109
1110 *hence, generalizing these findings to our previous estimations, we can assert that RE is the
best estimator we can deploy within our dataframe for this regression model.
1111
1112 *3) First-order autocorrelation check
1113 *-----
1114 sort country_id year
1115 gen l1gi
1116 by country_id year i l1gi
1117 correlate
1118 *The correlation analysis reveals a strong positive autocorrelation between the variable i and
its first lag L.i, with a correlation coefficient of approximately 0.9938. This indicates a
substantial linear relationship between the variable and its past values.
1119
1120 //////////////////////////////////////
1121 //////////////////////////////////////Functional form test////////////////////////////////////

```



```

1122 ///////////////////////////////////////////////////
1123
1124 use adj_data, clear
1125 xtset country_id year
1126
1127 *1) non-normality of error terms checks
1128 *-----
1129 ///////////////////////////////////////////////////SPEC1/////////
1130 xtreg i log_y g log_s log_t log_e, re
1131 predict pred_i1
1132 *Kernel density distribution vs normal
1133 kdensity pred_i1, normal
1134 *Jacques_Bera test
1135 sktest pred_i1
1136 *Shapiro-Wilk Test
1137 swilk pred_i1
1138 *The Jacques-Bera test for skewness and kurtosis, the Skewness and Kurtosis tests, as well as the
  Shapiro-Wilk test, all indicate a departure from normality. The p-values associated with these
  tests are all very low (close to or equal to zero), leading to the rejection of the null
  hypothesis that the data follows a normal distribution.
1139
1140 ///////////////////////////////////////////////////SPEC2/////////
1141 xtreg i log_y g log_s log_t log_e w f, re
1142 predict pred_i2
1143 *Kernel density distribution vs normal
1144 kdensity pred_i2, normal
1145 *Jacques_Bera test
1146 sktest pred_i2
1147 *Shapiro-Wilk Test
1148 swilk pred_i2
1149 *The Jacques-Bera test for skewness and kurtosis, the Skewness and Kurtosis tests, as well as the
  Shapiro-Wilk test, all provide evidence against the hypothesis that the data follows a normal
  distribution. The p-values associated with these tests are very low (close to or equal to zero),
  leading to the rejection of the null hypothesis. Therefore, based on these normality tests, data
  do not appear to be normally distributed.
1150
1151 ///////////////////////////////////////////////////SPEC3/////////
1152 xtreg i log_y g log_s log_t log_e w f r u h, re
1153 predict pred_i3
1154 *Kernel density distribution vs normal
1155 kdensity pred_i3, normal
1156 *Jacques_Bera test
1157 sktest pred_i3
1158 *Shapiro-Wilk Test
1159 swilk pred_i3
1160 *The p-values from the Jacques-Bera test for skewness and kurtosis, the Skewness and Kurtosis
  tests, and the Shapiro-Wilk test are relatively high. While the p-value from the Jacques-Bera
  test for skewness and kurtosis is not very elevate, the other two tests provide some support for
  the hypothesis of normality. Therefore, based on these normality tests, data appears to be
  relatively close to a normal distribution.
1161
1162 *Lower values of AIC and BIC indicate a better fit. Therefore, Specification 3 (Government
  expenditure targets) appears to have the best fit among the three models also according to
  R-squared. Moreover, R-squared values increase from Spec1 to Spec3, indicating a better ability
  to explain the variance in the dependent variable.
1163
1164 *2) Ramsey reset test
1165 *-----
1166 xtreg i log_y g log_s log_t log_e w f r u h, re
1167 predict pred_i
1168 gen form2 = pred_i^2
1169 gen form3 = pred_i^3
1170 *fit test including form2 and form3
1171 xtreg i log_y g log_s log_t log_e w f r u h form2 form3, re
1172 *RESET test over the new experimental variables
1173 test form2 form3
1174 *The p-value for the RESET test is very small (p = 0.0000), indicating that we reject the null
  hypothesis. This suggests that there might be a specification error in the model related to form2

```

```

and form3, leading us to rule them out of the model specification- The test confirm us that we
are working with the right functional form.
1175
1176 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1177 ///////////////////////////////////////////////////////////////////Heteroscedasticity robustness test/////////////////////////////////////////////////////////////////
1178 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1179
1180 use adj_data, clear
1181 xtset country_id year
1182 ssc install whitetst
1183
1184 *1) Breusch-Pagan test
1185 *-----
1186 xtreg i log_y g log_s log_t log_e w f r u h, re
1187 predict uhat, u
1188 *Manually run Breusch-Pagan LM test
1189 robvar uhat, by(country_id)
1190
1191 *The result of the Breusch-Pagan test for heteroscedasticity indicates that the test statistics
for all three weighting matrices (W0, W50, and W10) are missing (denoted by .), and the
corresponding p-values are also missing. This suggests that the test results are inconclusive or
invalid due to some issue with the residuals or the data.
1192
1193 *The test statistic (F-statistic) is 6.14 with a p-value of 0.0000. Since the p-value is less
than the conventional significance level of 0.05, you would reject the null hypothesis. This
provides evidence against the assumption of constant variance, suggesting the presence of
heteroskedasticity in the model.
1194
1195 *2) White test
1196 *-----
1197 *Manually run White test
1198 gen uhat_sq = uhat^2
1199 robvar uhat_sq, by(country_id)
1200 *The results indicate the presence of heteroscedasticity in the residuals. The test statistics
(W0, W50, W10) are not reported, and the associated p-values are also not available (denoted as
"."), suggesting potential issues with the estimation or the grouping variable. However, the mean
and standard deviation of squared residuals for each country ID reveal substantial variability,
supporting the conclusion of heteroscedasticity
1201
1202 *given these results, to avoid highly probable heteroskedasticity biases, we now want Stata to
estimate the standard errors using a method that is less sensitive to the presence of
heteroskedasticity. To do so, we employ robust standard errors
1203 xtreg i log_y g log_s log_t log_e w f r u h, re
1204 est store BASE
1205 xtreg i log_y g log_s log_t log_e w f r u h, robust
1206 est store ROBUST
1207 *export results
1208 esttab BASE ROBUST using "Spec_3_Base&Robust.csv", se label replace wide plain
1209
1210 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1211 ///////////////////////////////////////////////////////////////////Endogeneity robustness checks/////////////////////////////////////////////////////////////////
1212 ///////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
1213
1214 *To ensure the robustness of our analysis, we employ exogeneity tests, which scrutinize whether
our chosen explanatory variables are free from correlation with the error term. This test is
instrumental in validating the assumptions underlying our model, providing insights into the
potential bias introduced by endogeneity.
1215
1216 use adj_data, clear
1217 xtset country_id year
1218 ssc install estout
1219
1220 *1) IV regression
1221 *-----
1222 *pretend to use model variables as instrumental variables (IV validity)
1223 ivregress 2sls i log_y log_s log_t f r u h ( w=g), first
1224 est store IV1
1225 ivregress 2sls i log_y log_s log_t f r u h ( w= g log_e), first

```



```
1226 est store IV2
1227 ivregress 2sls i log_y log_s log_t f r u h ( w=log_e), first vce(robust)
1228 est store IV3
1229 *display results
1230 est table IV1 IV2 IV3, se p
1231 esttab IV1 IV2 IV3, se
1232
1233 *2) Durbin Wu Hausman test
1234 *-----
1235 estat endogenous
1236 *Both the robust score chi-squared test and the robust regression F-test have low p-values
(0.0000), indicating strong evidence against the null hypothesis. Therefore, we reject the null
hypothesis, suggesting that the model variables are endogenous. Hence, w and log_e must be
included in our model.
1237
1238 *3) Hansen Sargan test
1239 *-----
1240 estat overid
1241 *The robust score chi-squared statistic is 57.7889 with a p-value of 0.0000, and the robust
regression F-statistic is 37.0072 with a p-value of 0.0000. These findings provide robust support
for the rejection of the hypothesis that the included variables are exogenous, suggesting that
they are likely endogenous in your model. Hence, g and log_e must be included.
1242
1243 *4) verify the quality of our instrumental variables
1244 *-----
1245 corr w g log_e
1246 reg w g log_e log_y log_s log_t f r u h, robust
1247 test g log_e
1248 *the test indicates that there is no significant evidence to reject the null hypothesis that both
g and log_e coefficients are zero. The F-statistic is 1.59, and the p-value is 0.2044. This
suggests that these variables do not contribute significantly to explaining the variation in the
dependent variable.
1249
```