# DEGGED - Dynamic Exploration of Greenhouse Gas Emissions and its Determinants using R and Shiny

CHOONG Shi Lian Selene
Singapore Management University
slchoong.2019@mitb.smu.edu.sg

JIANG Weiling Angeline
Singapore Management University
wljiang.2019@mitb.smu.edu.sg

WONG Wei Sheng Dylan
Singapore Management University
dylan.wong.2019@mitb.smu.edu.sg

## ABSTRACT

In December 2020, the European Union (EU) leaders committed to an ambitious goal of reducing greenhouse gas emissions levels by 55% by 2030 to tackle climate change. As the world's third largest emitter of greenhouse gases, it is important to determine the factors which are contributing significantly to greenhouse gas emissions. Most existing literatures focused largely on the relationship between drivers and greenhouse gas emission levels, without considering mitigation factors which also plays a role in reducing greenhouse gas emissions. Earlier research also often presented findings in static forms limiting the amount of data exploitation that can be performed. Hence, our research aims to study the relationship of both drivers and mitigation measures on greenhouse gas emission levels using Ordinary Least Squares regression and Panel Data regression. We designed and developed DEGGED, an interactive web-based dashboard, to allow policymakers and environmentalists to explore and analyse determinants of greenhouse gas emissions. DEGGED will be intuitive for non-technical users to perform fundamental data analysis and regression modeling without any coding needed from users.

## 1. INTRODUCTION

Global warming is expected to result in a rise of the average global temperature between 1.1 to 6.4 degree celsius over the century, if no interventions are taken to reduce the emissions of greenhouse gases [1]. Multiple drivers which include economic activities, such as electricity production, transportation and waste generation, contributes to greenhouse gases. To mitigate greenhouse gas emissions, measures like emission taxes and use of renewable energy sources were introduced. Understanding the relationship between greenhouse gas emissions and the determinants will provide greater insights on the influence of determinants on emission levels.

With large amount of environmental, social and economic statistics available on Eurostat, our aim is to provide policy makers and environmentalists with an interactive analytical tool for them to understand the influence of drivers and mitigation measures on greenhouse gas emission levels via Data Exploratory Analysis (EDA) and advanced statistical techniques namely Ordinary Least Square regression (OLS) and Panel data regression.

This paper documents our approach to design and implement a web-enabled client-based analytics tool and consists of 5 sections. Section 1 provides the general introduction of the paper. This is followed by the motivation and objectives behind our research, along with review of the existing works. Section 3 discusses on user interface design and the next section illustrates the data preparation and practical use of the interactive tool built using climate change data from Eurostat. Lastly, the paper concludes by highlighting value-add of this research and potential future work moving on from this research.

## 2. MOTIVATION AND OBJECTIVES

Marcotullio et al.(2013) [10] conducted exploratory data analysis on global urban greenhouse gas emissions to describe emissions by region and sector and to examine the distribution of emissions through the urban-to-rural gradient, with the descriptive findings presented in tabular forms. Kijewska & Bluszcz (2016) [9] explored including descriptive charts on the greenhouse gas emission levels over time and by countries prior to its clustering analysis. The static nature of how the data were presented in these researches restricted the amount of exploratory data analysis that users can perform.

To estimate the relationship between one or more independent variables and a dependent variable, statistical methods namely OLS regression is used to analyse at the cross-sectional level while panel data regression is employed to analyse from the longitudinal perspective. Studies from Budiono et al. (2019) [3], Grunewald & Martínez-Zarzoso (2009) [6] and Guo & Jiang (2011) [7] explored using OLS to estimate the relationship between various factors and greenhouse gas emissions. While studies from González-Sánchez & Martín-Ortega (2020) [5], Dogan & Seker (2016) [4] and Azevedo, Horta & Leal (2017) [2] investigated the determinants of greenhouse gas emissions in the European Union through panel data models.

Despite multiple past researches conducted to investigate the relationship between greenhouse gas emissions and the determinants, most research papers usually only published the final model used in their analysis and omitted the intermediate steps used to derive the final set of independent variables. Tests conducted to ensure that model assumptions were not violated were also seldom presented in the papers. In addition, tabular forms were often adopted when presenting the regression findings as well as the results of the various validation tests performed (if any).

This research was motivated by the general lack of interactive web-enabled client-based analytics tool for uncovering greenhouse gas emission patterns and influence of drivers and mitigation measures on emissions. It aims to deliver a R-Shiny app that provides interactive user-interface design to: 1) Apply fundamental data analysis to understand the factors affecting greenhouse gas emissions, including by country and time period 2) Identify main factors contributing to greenhouse gas emissions from cross-sectional view using OLS regression. Step-by-step guidance for users to detect any violation of OLS assumptions. 3) Identify main factors contributing to greenhouse gas emissions from panel data view using Panel data regression and list of validation tests available for panel analysis.

## 3. DESIGN FRAMEWORK
Development of the DEGGED was done on Shiny, an R package for building interactive web apps. Our design considerations when developing the application includes:

- Detailing all data preparation and modelling within R for reproducibility
- Using R packages supported on the Comprehensive R Archive Network (CRAN) for supportability
- Offering selection options and responsive visualisations for interactivity

DEGGED supports three main categories of statistical analysis, namely, (a) Exploratory data analysis, (b) OLS regression model and (b) Panel data regression model. Every interface has similar design where the left panel consists of the input variables and the output is on the right panel. The design of DEGGED follows the taxonomy of interactive dynamics from Heer & Shneiderman (2012) [8] (Data and View Specification, View Manipulation, and Process and Provenance). Detailed functionalities of each interface, including data input, parameter configuration, output views will be described in the following sub-sections. On top of these three main interfaces, DEGGED also offers users an Data Overview module to explore the data used for the application.

## 3.1 Exploratory Data Analysis
The EDA module consists of 4 sub-modules (Descriptive Analysis, Time Trend, Correlation matrix and Scatterplot) and aims to overcome the limitation of static visualisation in data exploration by providing users with an interactive and user-friendly display of the data, which also serve as a ground work in the variables selection for OLS and Panel data regression analysis. The plotly package in R was used for chart plotting as it has the pan and zoom features which enhances user experience.

### 3.1.1 Descriptive Analysis
To understand the distribution of the variables in the dataset, users can click on the Descriptive analysis module. The interface includes a slider bar for the selection of years of analysis, a drop-down list (including the option to select all or deselect all) for multiple countries selection and a drop-down list for users to select the variables of interest. (Figure. 1 - Left) Upon clicking on the Apply changes button, the summary statistics tab provides users with the measures of central tendency (Mean and Median) and measures of spread (min-max, 25th and 75th percentile) in tabular form, while the histogram tab shows the distribution of the values of the variables. By hovering over the bars in the histograms, users could read off the count of each value. (Figure. 1- Right)



Figure. 1: User-interface for Descriptive analysis and Summary output

### 3.1.2 Time Trend
Users who are interested in finding patterns on how greenhouse gas emission levels or levels of drivers and mitigation factors have changed over time can click on the Time trend sub-module. Unlike the Descriptive analysis sub-module, the variable selection has been split into two drop-down lists. The Individual option allows users to compare the trend of a single variable across selected countries (Figure. 2), while the Grouped option groups the countries and displays the trend of all the variables across time (Figure. 3). Users can hover over to the data point and the actual or mean value will be displayed.

### 3.1.3 Correlation Matrix and Scatterplot
The Correlation matrix sub-module aims to allow users to analyse the relationship between two variables. Users can select the year of analysis from the single-select drop-down list and countries and variables selection are the same as other sub-modules where multiple selections are allowed. In addition, users have the choice on the correlation method and significance level of their preferred choice by selecting from the Type (Pearson, Spearman and Robust) and Significance

level (0.01, 0.05 and 0.10) drop-down list, respectively. The correlation matrix plot or correlogram displays the pairwise correlation coefficients. The color of each variables pair tells users whether the pairs are positively (green) correlated, negatively (orange) correlated, or not (white) correlated at all. The plot also crosses out the ones that are statistically insignificant at the chosen significance level. (Figure. 4)
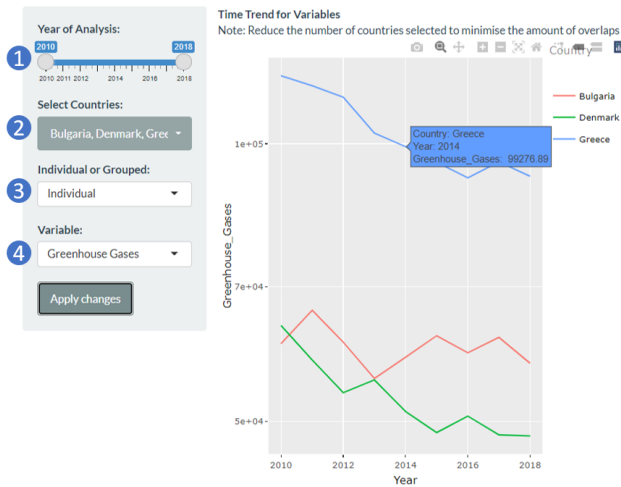
For deeper analysis, users can also click on the Scatterplot sub-module as it provides users the option to plot multiple years and countries of data within a plot. In addition, the scatterplot includes a loess (locally weighted smoothing) curve to help users see the general relationship between the two selected variables. (Figure. 5)
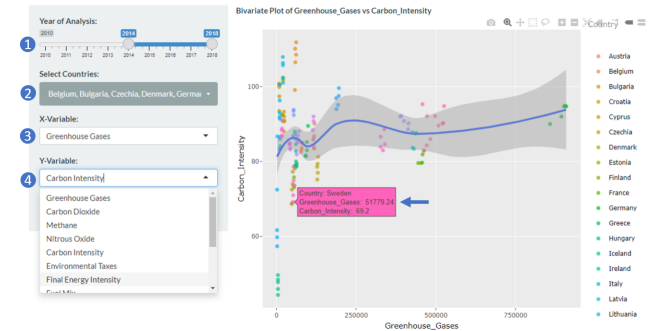


Figure. 2: User-interface for Time series (Individual analysis)



Figure. 3: User-interface for Time series (Group analysis)



Figure. 4: User-interface for Correlation Matrix



Figure. 5: User-interface for Scatterplot

## 3.2 Ordinary Least Square Regression

The interface of the OLS module consists of 2 sub-modules (Variable Selection and Model Selected), and provides a step-by-step guide in building OLS regression models.

### 3.2.1 Variable Selection

The first step is for users to explore and identify the best subset of independent variables to include in the model using the variable selection sub-module. Users will select the dependent and independent variables to be included in the model, and the year of analysis. Thereafter, users will click on a drop-down list consisting of different stepwise OLS methods (Stepwise AIC forward/backward/both and Stepwise BIC forward/backward/both) which are available in the MASS package in R and reactive in Shiny.(Figure. 6 - Left) Upon submitting the query, the ANOVA and regression summary will show a quick summary of the steps taken at each search and the parameter estimates of the final model. A user guide to interpreting the output was also included after each output. (Figure. 6 - Right) To have a clearer visualisation in the change of the criterion value after each search step, users can click on the Plot tab. (Figure. 7) For more detailed steps in the model ran, users can click on the Detailed output tab.

### 3.2.2 Model Selected

Based on the set of selected independent variables which the variable selection module has determined as the best fit model, users proceed to the Model selected module to obtain the parameter estimations and run the necessary OLS assumption tests. Similarly, users will select the dependent and independent variables to be included in the model, and the year of analysis. Thereafter, users can also choose from a drop-down list their preferred choice of test to run to check for any violation of OLS assumptions and check for model fit. List of tests are (i) Heteroskedasticity test (Bvensch Pagan Test; Score test; F-test), (ii) Residual diagnostics (QQ-plot; Residual vs Fitted value; Residual Histogram), (iii) Collinear diagnostics and (iv) Model fit assessment
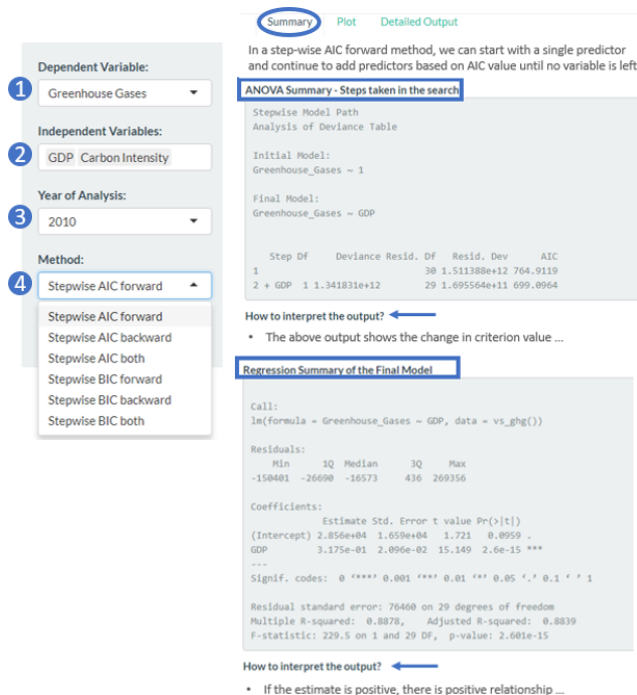
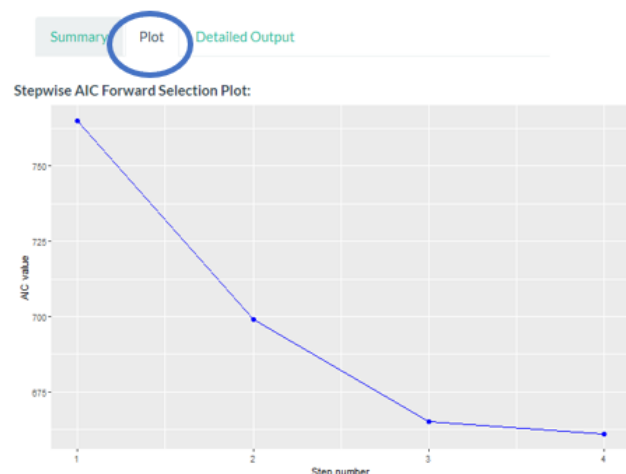Figure. 6: User-interface for Variable Selection module and Summary output



Figure. 7: Graphical summary output

(Residual Fit Spreadplot; Observed vs Predicted plot). (Figure. 8 - Left) In the output, DEGGED presents the abundant results in parameter estimation in tabular and dot plots (Figure. 8 - Right) for easy interpretation. The results from heteroskedasticity test and collinearity diagnostics are displayed in tabular form. While, results from residual diagnostics and model fit assessment will be shown using scatter plot as users will be able to easily detect any violation of the assumption by looking at the plotted points. (Figure. 9) For example, if the data is normally distributed, the points in the QQ-normal plot will lie on a straight diagonal line. The olsrr package in R was used to build the OLS assumption tests.



Figure. 8: User-interface for Model Selected module and Parameter estimations plot



Figure. 9: OLS Assumption tests

## 3.3 Panel Data Regression

Panel data regression extends beyond OLS to consider the longitudinal aspect of the data, where users can examine the variations along both the individual and time dimensions. The complexity involved in constructing panel models and ensuring its robustness will be addressed by DEGGED which simplifies the process into 2 sub-modules (Model Building and Validation Tests), and provides a step-by-step guide in building and evaluating panel regression models. As OLS is explored as a separate module in the earlier section, the panel regression module only covers fixed effects and random effects models.

The fixed effects model considers individual differences over time, and assumes correlation between the individuals' error term and predictor variables, with no correlation between time-invariant and individuals' other characteristics. The random effects model considers unique errors to be uncorrelated with the regressors, and assumes that the variations across individuals to be random and uncorrelated with the predictor variables. The plm package in R was used to build the panel data models.

### 3.3.1 Model Building

Apart from the dynamic time and country dimensions, users can select the independent and dependent variables of interest to be included in the model. Users can subsequently decide the type of panel data models (Fixed effect model and Random effect model) as well as the type of effect (Individual, Time and Two-ways). For random effect model, the application allows for four different estimators of the parameter, namely Swamy and Arora, Wallace and Hussain, Amemiya and Nerlove. Users can also select the type of Lagrange Multipler test (Breush-Pagan, Honda and King & Wu) to perform for the subsequent Model Fit Assessment. (Figure. 10 - Left) Running the selected model will produce a coefficient estimate table (Figure. 10 - Right) along with a dot-and-whisker plot of the regression coefficient estimates (Figure. 11) for clearer visualisation. For more statistical details, users can also refer to the Detailed output tab.

Upon running the selected panel regression model, users will also be given the Model Fit Assessment to assess which of panel model (pooled OLS, fixed effect or random effect model) is most suitable for the selected variables and data used. To assess if fixed or random effect model should be used, DEGGED conducts the Hausman Test with the null hypothesis being the preferred model to be random effects. To decide if fixed effect model should be used over pooled OLS, Chow Test of Poolability will be performed with the null hypothesis being the pooled effect model. Lagrange Multiplier Test, with null hypothesis being no panel effect, is used to determine if pooled OLS or random effect model should be used. Relevant tests will be shown based on the Model Type selected. Outputs from each of the model selection tests are presented in tabular forms due to the limitations of clutteredness when approaching the assessments using visual representations. (Figure. 12)

### 3.3.2 Tests for Assumptions

To ensure the robustness of the model, additional diagnostic tests were performed to ensure that the dataset complies with the assumptions of panel regression model. Users are offered the option to choose from a drop-down list their preferred choice of test to check for any violation of panel regression assumptions (Figure. 13 - Left). The list of diagnostic tests included in DEGGED are (i) Test of Normality (Residual QQ Plot; Residual vs Fitted Value Plot; Residual Histogram), (ii) Test of Serial Correlation (Unobserved Effect Test; Locally Robust Tests; Breusch-Godfrey/Wooldridge Test) and (iii) Test for Heteroskedasticity (Breusch-Pagan Test; Residual vs Fitted Value Plot). DEGGED presents the results from the various tests in either tabular form or charts depending on outputs (Figure. 13 - Right). Intepretation reference is included in the application to help guide users interpret the results.

## 4. USE CASE: GHG EMISSIONS IN EU

In December 2020, the European Union (EU) leaders committed to an ambitious goal of reducing greenhouse gases by 55% by 2030 to tackle climate change. As the world's third biggest emitter of greenhouse gases, it is important to determine the factors which are contributing significantly to greenhouse gas emissions.

A study conducted by González-Sánchez & Martín-Ortega (2020) [5] on the determinants of greenhouse gas emissions
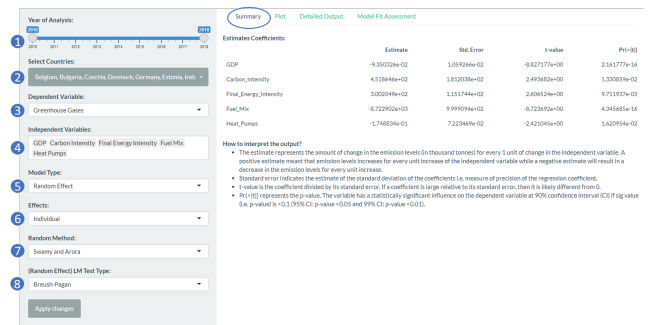


Figure. 10: User-interface for Model Building module and Summary output
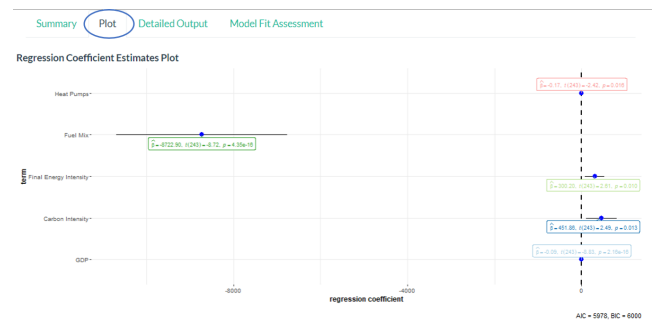


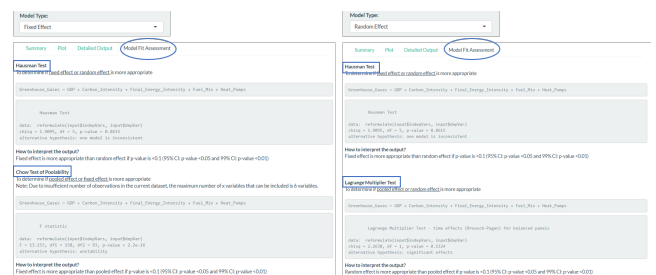Figure. 11: Graphical output of coefficient estimates



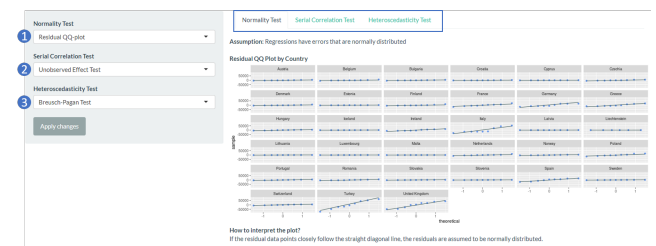Figure. 12: Model Fit Assessment for constructed panel model



Figure. 13: User interface for Tests of Assumption module

growth in Europe considered majority of the driver factors available on Eurostat database and found that GDP and final energy intensity are the main drivers for the reduction of greenhouse gas emissions in Europe. However, mitigation factors which also plays a part in reducing greenhouse gas emissions were not considered in their study. Hence, our research further explored the Eurostat database and extended the list with mitigation factors which could potentially play an important role in the reduction of greenhouse gas emissions in Europe.

The data used for our research contains the overall greenhouse gas emissions, excluding Land Use and Land Use Change and Forest (LULUCF) and the emission levels of common greenhouse gases (i.e. carbon dioxide, methane and nitrous oxide) from 33 European countries between the period 2010 to 2018. Data from four drivers (Gross Domestic Product, Final energy intensity, Fuel mix and Carbon intensity) and five mitigation factors (Renewable energy, Environmental taxes by economic activity, Liquid biofuels production capacities, Solar thermal collectors' surface and Heat pumps - technical characteristics by technologies) were also included.

As the downloaded datasets from Eurostat database had irrelevant rows and columns, data preparation steps were performed to tidy the data prior to combining into one dataset for our reserach use. Additional computation were also performed to obtain the final energy intensity as data for the indicator was not readily available. While Eurostat has distinguished the production capacities by different liquid biofuel types and the installed thermal capacity by different heat pump technologies, our research will only be analyzing these factors at the collective overall level. Hence, aggregation was performed to get the total value of the respective factors for each country at each year. The combined dataset was finally transposed to obtain a country-year record level data for analyzing.

Our assessment showed that DEGGED has several advantages over traditional exploratory and regression analysis. Firstly, upon detecting any unusual patterns from the overall view at the data exploratory step, users can pan and zoom to focus on specific data points with more information displayed by hovering over to the data point, without having the need to filter through the raw dataset. (Figure. 14) The pan and zoom features are from plotly package in R.
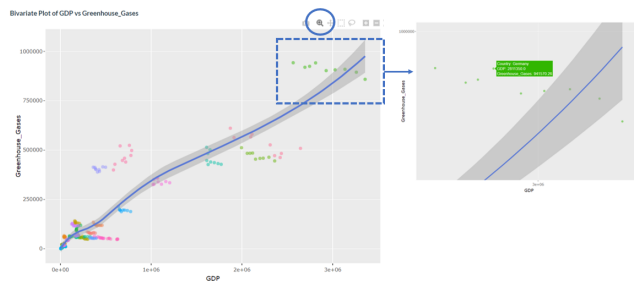


Figure. 14: Data exploration using pan and zoom features

Secondly, the interactive features available in the OLS (Figure. 15) and Panel Regression (Figure. 16) modules enable users to choose their list of independent and dependent variables and run their preferred model without any coding or having to code multiple code chunks to cater for different specifications. The regression output, both tabular and graphical format, are almost updated instantaneously. DEGGED also has a pre-loaded list of necessary assumption tests to detect any violation of the model assumptions each time a model is ran. (Figure. 17) (Figure. 18) To give users more options, users can also select their preferred test from the drop-down list of each assumption test. For example, users can choose between Breusch Pagan test, Score-test or F-test to detect heteroskedasticity problem. As a value-add,
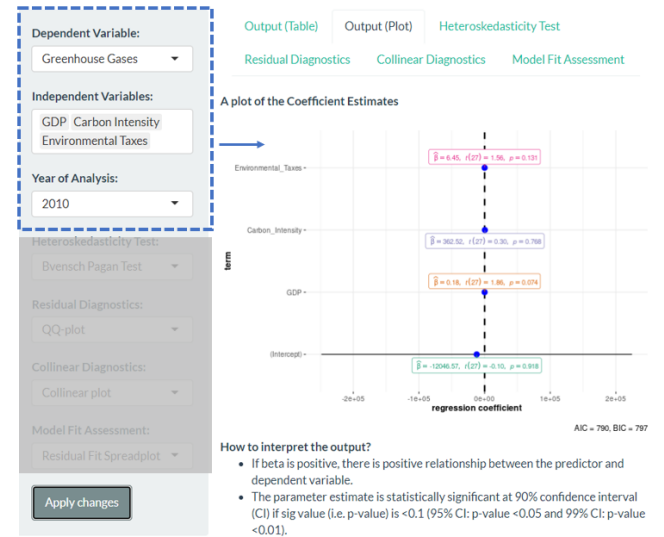


Figure. 15: Interactive parameter estimations for OLS
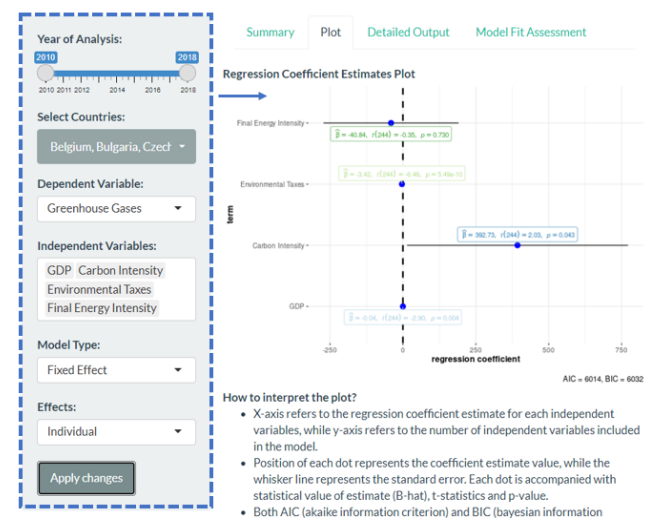


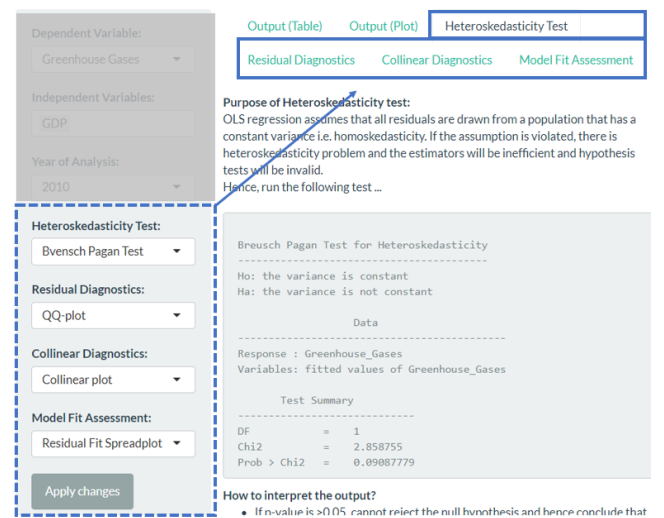Figure. 16: Interactive parameter estimations for Panel



Figure. 17: Interactive assumption tests for OLS

DEGGED also includes interactive notes at the bottom of each output to guide and explain how each output should be interpreted. (Figure. 19)
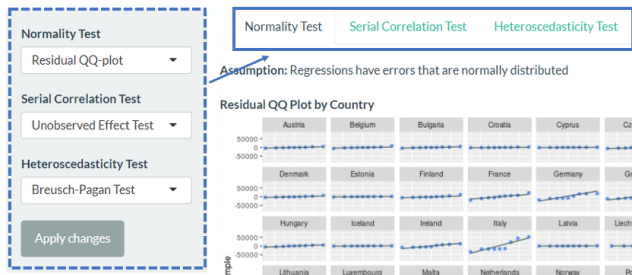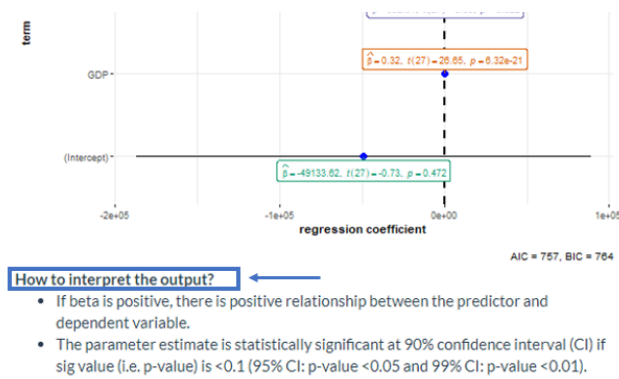


Figure. 18: Interactive asumption tests for Panel



Figure. 19: Interactive interpretation guide

# 5. CONCLUSION & FUTURE WORK

The demonstration of DEGGED clearly shows the benefits that policy makers and environmentalists, with or without technical background, can reap as they will be able to perform on the fly data exploratory analysis and advanced statistical techniques (OLS and Panel data regression) to gain quick and valuable insights for policy making. Nonetheless, instead of using MASS package in R which outputs only AIC criterion value at the stepwise modeling step, future work could be explored to modify the olsrr package in R for it to incorporate reactive dataset and hence able to output a panel of fit criterion (R-square, Adjusted R-square, Mallow's Cp, AIC, SBIC and SBC). In addition, the scope of DEGGED could be extended to cover greenhouse gas emissions in other countries such as U.S. as the U.S. aims to cut greenhouse gas emissions in half by 2030 as part of its new commitment to the Paris climate agreement.

# 6. ACKNOWLEDGEMENT

# References

[1]    Amanatidis, G. 2020. Combating climate change. European Parliament.

[2]    Azevedo, I. et al. 2017. Analysis of the relationship between local climate change mitigation actions and greenhouse gas emissions – Empirical insights. Energy policy. 111, (Dec. 2017), 204–213.

[3]    Budiono, R. et al. 2019. Modeling and analysis of CO2 emissions in million tons of sectoral greenhouse gases in Indonesia. IOP Conference Series. Materials Science and Engineering. 621, 1 (Oct. 2019).

[4]    Dogan, E. et al. 2016. An investigation on the determinants of carbon emissions for OECD countries - empirical evidence from panel models robust to heterogeneity and cross-sectional dependence. Environmental science and pollution research international. 23, 14 (Jul. 2016), 14646–14655.

[5]    González-Sánchez, M. and Martín-Ortega, J.L. 2020. Greenhouse Gas Emissions Growth in Europe - A Comparative Analysis of Determinants. Sustainability (Basel, Switzerland). 12, 3 (Jan. 2020), 1012.

[6]    Grunewald, N. and Martínez-Zarzoso, I. 2010. Driving Factors of Carbon Dioxide Emissions and the Impact from Kyoto Protocol.

[7]    Guo, H. and Jiang, Y. 2011. The Relationship between CO2 Emissions, Economic Scale, Technology, Income and Population in China. Procedia Environmental Sciences. 11 (2011), 1183–1188.

[8]    Heer, J. and Shneiderman, B. 2012. Interactive Dynamics for Visual Analysis. ACM Queue. 55, 4 (Feb. 2012), 45–54.

[9]    Kijewska, A. and Bluszcz, A. 2016. Analysis of greenhouse gas emissions in the European Union member states with the use of an agglomeration algorithm. Journal of Sustainable Mining. 15, 4 (Nov. 2016), 133–142.

[10]   Marcotullio, P.J. et al. 2013. The geography of global urban greenhouse gas emissions - an exploratory analysis. Climatic change. 121, 4 (Dec. 2013), 621–634.