

Exploring Countries' Sustainability Using ESG Data

Kathy Wu, Nathan Lai, Ruoxin Wang, Yuchen Fang.

Author contributions

Kathy Wu contributed Data Cleaning, Principal Component Analysis Code and Result.

Nathan Lai contributed Datasets Description, Principal Components Analysis plots and Discussion.

Ruoxin Wang contributed Introduction, ESG score evaluation coding and Result (ESG score evaluation part).

Yuchen Fang contributed Abstract and Methods.

Abstract

This report mainly investigates the sustainability of countries through Environment, Social, and Governance aspects and analyzes the variation among the principal components of each country. The motive is to find out which country is most sustainable during the Covid-19 pandemic outbreak. According to Principle Component Analysis, there are no significant changes in the chosen ESG variables during the pandemic. Although the variation of GDP growth is thought to be an influential component, it doesn't affect the result a lot. Hence, a new sustainability score evaluation method is used and Vietnam is considered the most sustainable country in 2019 and 2020.

0. Introduction

-> Background

ESG is an abbreviation for Environmental, Social, and Governance ---- a combination of three categories of non-financial factors that are increasingly applied by investors as part of their analysis process to evaluate material risks and growth opportunities nowadays. However, to better align with the global goals, the World Bank Group rearranges it into a new data frame that further classifies 17 key sustainability themes based on the original environmental, social, and governance categories, which some of the key themes are shown in the image below (image credit: bondevalue.com).

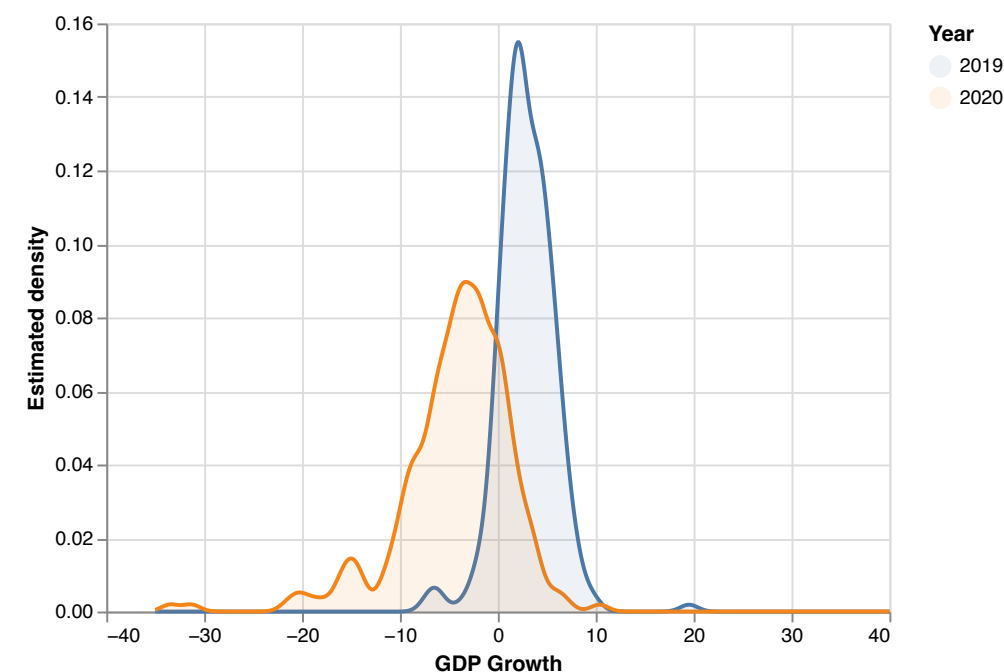


The World Bank Group believes that these themes are crucial for financial sector representatives to consider when assessing the contribution of investments or policies to sustainable development. Based on the [ESG Score Evaluation](#), an ESG score ranging from 0 to 100 will be assigned. Usually, a score of less than 50 will be regarded as poor performance and a score of more than 70 will be considered as excellent performance.

-> Aims

At the beginning of 2020, a worldwide pandemic hit hard across the world. Economic fallout, the unemployment rate remained high, people suffered from living hardships, and the decline in GDP growth is also inevitable. This decrease can be shown in the estimated density plot of GDP Growth below for the year 2019 and 2020.

Figure 1: The distribution of GDP growth of each country from 2019 to 2020 in the data.



Since the estimated density of GDP growth is observed to be a significant difference for two years, this project will focus on comparing various aspects of Environment, Social, and Governance (ESG) categories for most countries before and after the existence of a pandemic and how each variable drives the variation of the data. Using Principal Components Analysis to discover which ESG key factors would have a larger influence, as well as, whether they cause any effects on the accuracy of the ESG Score Evaluation Method. However, the result reveals that there are no significant changes in the distribution of the selected ESG variables, which suggests that the accuracy of the ESG Scoring Method is relatively objective and will not be easily swayed by such unexpected incidents as COVID-19. Furthermore, a new scoring method adapted based on the original system is created to study the sustainability of each country during 2019 and 2020, which finds out the most sustainable one with the highest sum of variable scores.

1. Dataset Description

In order to shift financial flows so that they are better aligned with global goals, the World Bank Group (WBG) is working to provide financial markets with improved data and analytics that shed light on countries' sustainability performance. This dataset is classified as Public under the Access to Information Classification Policy. It provides information on sustainability themes spanning environmental, social, and governance categories. Along with new information and tools, the World Bank can develop research on the correlation between countries' sustainability performance and the risk and return profiles of relevant investments.

These data are publicly available:

[Environment, Social and Governance Data, The World Bank](#) and is licensed under [Creative Commons Attribution 4.0.](#))

After cleaning up the raw data from the source, the final dataset contains 2-year ESG information from 239 countries all over the world with a total of 478 observations and 11 ESG variables.

-> Sample and measurement information

For the collection method, since this data is census data, the values in the topic Governance and Social are obtained from surveys, and most of the data on the topic Environment is collected by using scientific equipment.

Furthermore, for the sampling design and scope of inference this, all countries reporting environment, social and governance data is the sampling frame, the census is the sampling mechanism, and the scope of inference is none.

Table 1: variable descriptions, topic, data type and units for each variable in the dataset.

Name	Variable Description	Topic	Type	Units of measurement
fore_area	Forest area	Environment	Numeric	% of land area
pop_denst	Population density	Environment	Numeric	people per sq. km of land area
rate_labor	Ratio of female to male <i>labor force participation rate</i>	Governance	Numeric	% (modeled ILO estimate)
gdp_grow	GDP growth	Governance	Numeric	annual %
parliment_women_seat	Proportion of seats held by women	Governance	Numeric	% in national parliaments
unemp_rate	Unemployment, total	Social	Numeric	% of total labor force (modeled ILO estimate)
life_exp	Life expectancy at birth, total	Social	Numeric	years
acce_electr	Access to <i>electricity</i>	Social	Numeric	% of population
mortal_rate	Mortality rate, under-5	Social	Numeric	per 1,000 live births
acce_fuel_tech	Access to <i>clean fuels</i> and <i>technologies</i> for cooking	Social	Numeric	% of population
pop_65	Population ages 65 and above	Social	Numeric	% of total population

In tiding up the data, the *Mortality Rate* was converted to the percentage base, and *Population Density* was converted using logarithm to obtain a smaller variance. The first few rows of the ESG data are shown in Table 2 below.

Table 2: example rows of data.

Row	Country Name	Country Code	Year	acce_fuel_tech	acce_electr	fore_area	gdp_grow	life_exp	mortal_rate	pop_65	parliment_women_seat	rate_labor	unemp_rate	log_pop_denst
0	Afghanistan	AFG	2019	31.9	97.7	1.850994087	3.9116034	64.833	0.0601	2.615794213	27.86885245	30.00988032	11.21700000	1.76544050640
1	Afghanistan	AFG	2020	33.2	97.7	1.850994087	-2.351100673	65.173	0.058	2.64906965	27.01612903	24.68587789	11.7100000	1.7754458352
2	Albania	ALB	2019	80.7	100	28.79197080	2.113419981	78.573	0.0097	14.20263062	29.50819672	77.8291194	11.47000026	2.0177324695
3	Albania	ALB	2020	81.3	100	28.79197080	-3.955397926	78.686	0.0098	14.70458131	29.50819672	75.85744815	13.32900047	2.01523872040
4	Algeria	DZA	2019	99.7	99.5	0.814110350	0.9999999999	76.88	0.0233	6.552777881	25.75757575	24.87786005	10.51299953	1.25710943107

2. Methods

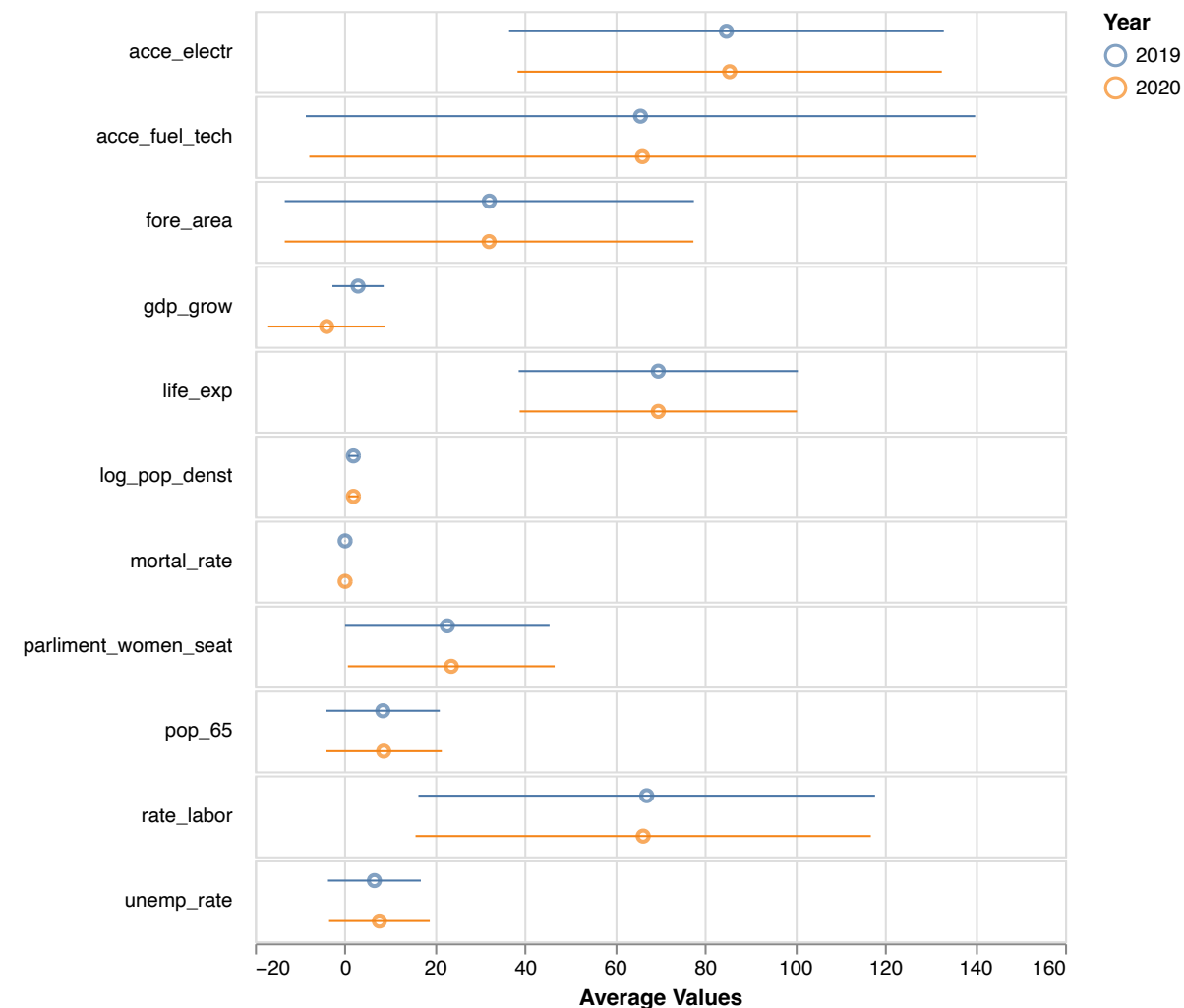
Exploratory analysis on multiple variables is applied while dealing with the dataset, such as examining correlation structure and computing and selecting principal components. A heatmap is made among all variables from the original dataset to explore the correlation and discard some unnecessary components. Then, 11 variables are chosen among 3 big categories from the dataset and ran PCA analysis. Principal Components Analysis(PCA) is to find variable combinations that capture large portions of the variation and covariation in our dataset. PCA can also capture the changes in the variation and covariation of the components in specific years such as 2019 to 2020 when there was a global Covid-19 outbreak, and see whether there are any changes in the weighted components affected by the pandemic. To further compare the sustainability of each country, all evaluated variables are ranked and assigned a score from 1 to 5 based on their positions of ranking quantiles in the overall sorted numeric value list. The sum of scores for all variables will be calculated for each country ranging from 5 to 55 and then, based on the original ESG Score Evaluation, the countries with the highest sum of scores will be the ones that have the most sustainable performance generally.

3. Results

-> Averages and Variabtions on each ESG variables

The exploratory analysis here focuses on how each value of ESG variables shifted from 2019 to 2020. Figure 2 explains the averages and variabilities of value of each ESG variables for year 2019 and 2020.

Figure 2: Average relative abundance and variability by ESG variable for year 2019 and 2020; the error bars represent two standard deviations in either direction relative to the mean value across the same year.

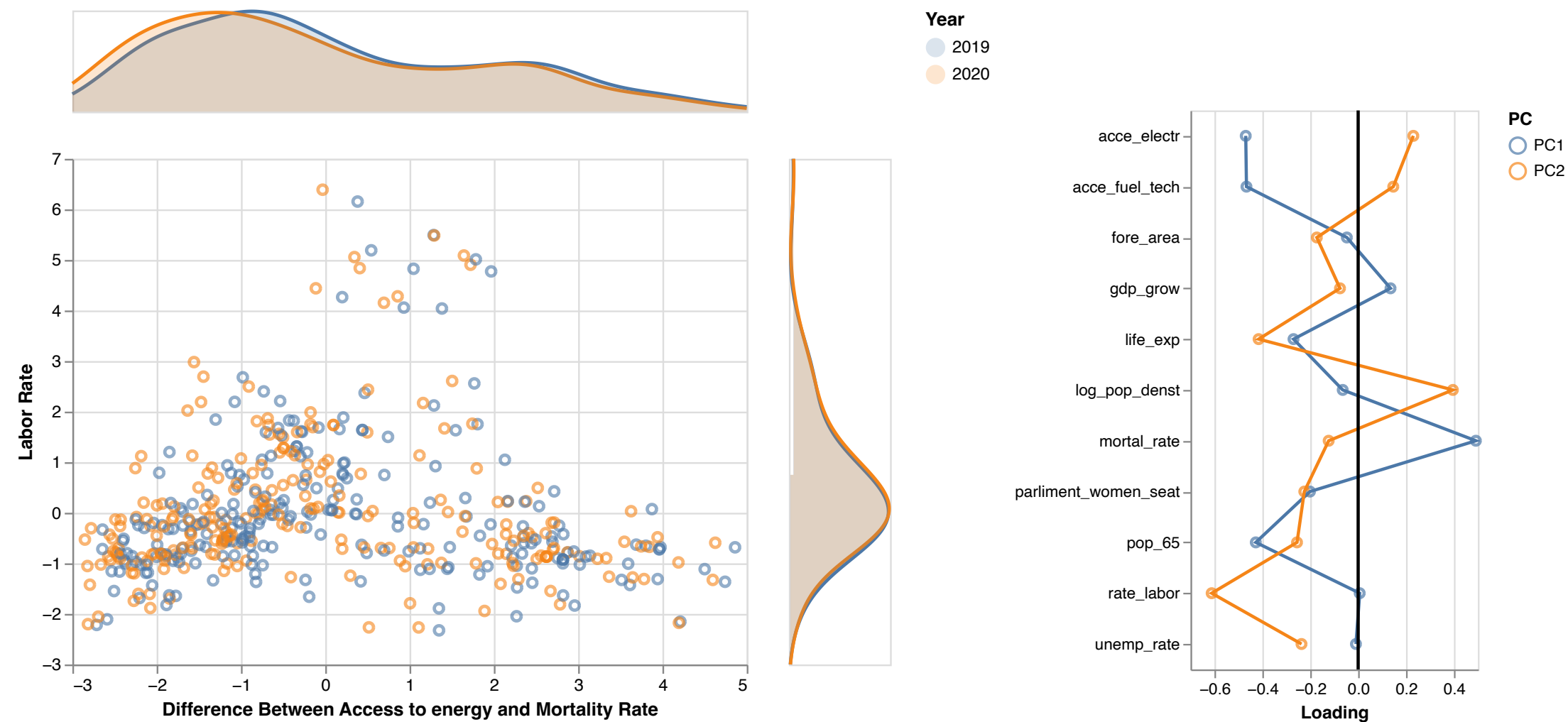


Only the *GDP Growth* shows noticable shift from 2019 to 2020. The other ESG variables reflect minimal changes between year 2019 and 2020. One possible reason of this shift might be the occurrence of COVID-19 at the beginning 2020.

-> Principal Components Analysis

By performing the analysis of the Principal components on a normalized basis, the resulting measures reflect the relative values for each ESG variable. These components together captured about fifty percent of the total variation in relative values of ESG variables from 2019 to 2020. The first Principal Component primarily explains the variables related to the Social Topic which are *Access to Electricity*, *Access to Clean Fuel*, and *Mortality Rate*. The second explains both Governance and Social Topics, however, the absolute loading value of *Labor rate* is much higher, hence it predominantly explains the Governance Topic.

Figure 3: Scatterplot of the principal components where Social variables on the x-axis and Governance variables on the y-axis while the points are colored according to the year (2019-2020), and univariate distributions of each measure is shown adjacent to the scatterplot accordingly; on the right, the high PC1 and PC2 loadings indicate the variables that drive the variation in the data.



According to the univariate distribution panels, the center and spread of each measure do not change noticeably between 2019 and 2020. It demonstrates that although the GDP growth has a predominantly changed from 2019 to 2020, it does not have much effect on the shape of the scatter plot from 2019 to 2020 and the composition of ESG variables.

-> Sustainability Score Evaluation

Lastly, to further explore which country performs the best in sustainable development in 2019 and 2020, a sustainable score evaluation method is used, which is adapted from the original ESG evaluation system. The final results with the most sustainable country for 2019 and 2020 are shown in the table below.

Table 3: The result table of country with the highest sustainability score for Year 2019 and 2020

	Country Name	Country Code	Year	score_sum
468	Vietnam	VNM	2019	47

	Country Name	Country Code	Year	score_sum
259	Luxembourg	LUX	2020	47
469	Vietnam	VNM	2020	47

From the table above, it is easy to see that Vietnam got the highest sustainability score of 47 out of 55 for both 2019 and 2020, and Luxembourg got the highest sustainability score of 47 out of 55 for 2020 as well.

4. Discussion

This project analyzes the sustainability from 2019 to 2020, to see how each country's sustainability changed affected by the COVID. Since the GDP is the most changeable variable through the years, the analysis first focused on the GDP growth before and after the pandemic, which is 2019 and 2020(Figure 1). Next, to identify both individual ESG variables that reflect corresponding shifts in average values(Figure 2) to locate which variable has the greatest changes. By applying and plotting the PCs(Figure 3), the graph shows that GDP growth doesn't have mainly effect on the sustainability. Conversely, the variables Access to Electricity, Access to Clean Fuel and Mortality Rate and Ratio of female to male labor force participation rate are the greatest affection variables to the sustainability. Moreover, the scatterplot shows that from 2019 to 2020, there are no significant changes to these weighted variables. Thus, based on the original evaluation method, we set up a new method(table 3) that is based on a currently existing well-established ESG Scores calculation. To analyze which is the most sustainable country before and after the existence of a pandemic.

The analysis suggests that from 2019 to 2020, the relationship between each ESG variable did not have any changes. Access to energy (electricity and clean fuels and technologies for cooking), which represent a country's technology and development level, is more relative to sustainability compared to other ESG variables(PC1 is typically slightly negative). In contrast, for such a well-developed country, its mortality rate will be lower than other developing countries. For example, based on the ESG Scores calculation, Luxembourg, which is a developed county, is more sustainable than other countries. Furthermore, the labor rate is also relative to sustainability compared to other ESG variables(PC2 is typically slightly negative). In comparison, the population density will be lower since if the supply of labor doesn't change, when population density goes up, the demand for labor will also go up, which will create a shortage in the labor market and lower the labor rate. For example, Vietnam. Many of the manufacturers are moving away from China and consider Vietnam as a lower-cost, reliable and quality source of parts, materials, and manufactured components. This will provide more job opportunities for Vietnam's people and increase the labor rate of the country, which also matches the result of the ESG Scores calculation.

Although it is not discussed here, this dataset can probably be used to predict the GDP Growth for each country using the other non-financial key factors in this ESG dataset. Accessing the whole ESG dataset would have an adequate amount of data to construct a model for prediction.