

Advanced Regression and Prediction Final Project: Predicting access to and reliance on technology

Daniel Alonso

May 9th, 2021

Contents

Dataset of choice	2
Variables	2
The target variable	3
Data preprocessing	4
Methodology and steps	4
Exploratory Data Analysis	5
Correlation	5
Very high correlation	5
Negative	5
Positive	5
Moderately correlated	6
Positive	6
Not correlated almost <i>at all</i>	6
Modelling: statistical tools	7
Simple models	7
Simple linear regression (univariate):	7
Prevalence of anemia in pregnancy to predict ART index	7
Infant mortality rate to predict ART index	7
Human development index to predict ART index	7
Simple linear regression (multivariate):	8
Using all variables to predict ART index	8
Using only variables with 20%+ importance	8
Robust linear model	9
Using OLS model selection to improve on the simple linear model with all vars	10
OLS forward based on p-val	10
OLS forward based on AIC	10
OLS step backward based on AIC	10
OLS stepwise AIC	11
Advanced Regression models	11
Backward regression	11
Lasso Regression	11
Stepwise Regression	12
Partial Least Squares	12
Ridge Regression	12
Principle Component Regression	12
Elastic Net	13

Variable importance for Elastic Net	14
Trying an elastic net model with only the selected variables in VarImp	14

Dataset of choice

For this project I picked a custom-built dataset obtained from the [World bank Databank](#), specifically the [World Development Indicators database](#). This is the “primary World Bank collection of development indicators” as stated on the database description. It is a database that specifically lists development variables, many of which also are related to health, economics, the environment and human development in a broad sense.

A second helper dataset that I also used and merged with the one obtained through the World Bank is a COVID-19 dataset constructed by [Our World in Data](#) and obtained from [their GitHub repository](#), or directly from [the direct link in the README](#).

The amount of historical data on both databases is massive, and the amount of possible predictors is immense. I have decided to stick to an amount of variables that is also quite large, but will be very interesting to reduce with the techniques learned in class.

Variables

I have made a really large selection of variables (over 45) initially, however, as many included NAs and a significant amount of countries and dependencies lacked a large amount of metrics, the amount of columns and countries has been reduced to 38 and 173 respectively. Nonetheless, this is still a reasonable amount of data, and also enough to make a reasonably acceptable prediction.

NOTE:

- **blue** = used for training/predicting
- **red** = target variable
- **green** = ID variables

Variables in the original dataset as constructed using the World Bank Databank tool and the COVID-19 dataset (all variables were renamed):

- **year**: year the data was obtained in
- **year_code**: code for the year as the world bank databank sets it
- **country_name**: name of the country
- **country_code**: alpha-3 ISO 3166 code for the country
- **access_to_electricity**: Access to electricity (% of population)
- **birth_rate**: Birth rate, crude (per 1,000 people)
- **life_expectancy**: Life expectancy at birth, total (years)
- **exports_perc_gdp**: Exports of goods and services (% of GDP)
- **education_years**: Compulsory education, duration (years)
- **gdp_per_capita_ppp**: GDP per capita, PPP (current international USD)
- **perc_internet_users**: Individuals using the Internet (% of population)
- **infant_mort_rate**: Mortality rate, infant (per 1,000 live births)
- **inflation_perc**: Inflation, consumer prices (annual %)
- **consumer_price_index**: Consumer price index (2010 = 100)
- **crop_production_index**: Crop production index (2014-2016 = 100)
- **goods_imports**: Goods imports (BoP, current USD)
- **prevalence_of_anemia_pregnancy**: Prevalence of anemia among pregnant women (%)
- **diabetes_prevalence**: Diabetes prevalence (% of population ages 20 to 79)
- **human_capital_index**: Human capital index (HCI) (scale 0-1)
- **net_enrollment_rate**: Adjusted net enrollment rate, primary (% of primary school age children)
- **measles_immunization**: Immunization, measles (% of children ages 12-23 months)
- **co2_emissions_gaseous_fuel**: CO2 emissions from gaseous fuel consumption (% of total)

- **fossil_fuel_energy_consumption**: Fossil fuel energy consumption (% of total)
- **fuel_exports**: Fuel exports (% of merchandise exports)
- **fuel_imports**: Fuel imports (% of merchandise imports)
- **investment_inflows**: Foreign direct investment, net inflows (% of GDP)
- **investment_outflows**: Foreign direct investment, net outflows (% of GDP)
- **mobile_subscriptions**: Mobile cellular subscriptions (per 100 people)
- **agricultural_land**: Agricultural land (% of land area)
- **greenhouse_gas_em**: Total greenhouse gas emissions (% change from 1990)
- **age_dependency_ratio**: Age dependency ratio (% of working-age population)
- **imports_annual_growth**: Imports of goods and services (annual % growth)
- **int_tourism_arrivals**: International tourism, number of arrivals
- **total_covid_cases_per_million**: Total COVID-19 cases per 1 million inhabitants
- **total_covid_deaths_per_million**: Total COVID-19 deaths per 1 million inhabitants
- **population_density**: Population density
- **median_age**: Median Age
- **human_development_index**: Human Development Index (HDI, scale 0-1)
- **ar_tech**: Access to and reliance on technology (scale 0-1)

The target variable

From a very roughly inferential standpoint and looking at the data, we can somewhat see that clearly, both developed and developing countries were hit very hard. However, there's a pattern, where a large chunk of developing countries have barely felt the effects of the pandemic (either due to isolation or the fact that there's just other more serious issues, like large-scale wars or poverty).

Mixing development metrics and COVID data might help us see whether COVID data does indeed contribute to this, and also, how well it can perform when compared to those development metrics when it comes to predicting a particularly good metric to assess development of a nation.

In my statistical learning project, I predicted the HDI Group of a country based on many of these development metrics (and a few other differing ones). One of those metrics, and one that was particularly good at assessing a country's development were **access_to_electricity**: percentage of population with access to electricity, **perc_internet_users**: percentage of population which use the internet and **mobile_subscriptions**: mobile cellular subscriptions per 100 people.

These three were incredibly good at predicting HDI and could be a fantastic substitute for it. And reasonably so, we know that the more developed a country is, the more modern its infrastructure for electricity distribution is and the more it eventually relies on the internet to perform certain activities.

Having lived in both a developing and developed country myself, I can say that the consistency of services like electricity supply and how common the internet was used to provide services to the population (to internet users) is quite significant.

Therefore, either of these could be used as target variable, but, as this would be somewhat repetitive, I've decided to create a weighted metric to contain all three of them. This metric will be called: **ar_tech**.

This measure will represent the following: "*Access to and reliance on technology*", or *ART* for short.

The measure will be calculated as a weighted measure as follows:

$$\text{ar_tech} = \frac{\text{mobile_subscriptions} * 10 + \text{perc_internet_users} * 35 + \text{access_to_electricity} * 55}{100}$$

This highlights the importance of electricity, as none of the other two would be possible without it, and that clearly, electricity is significantly more important than use of the internet or mobile subscriptions. However, more people using the internet is a huge sign of development, this is also true for mobile subscriptions, but perhaps to a lesser extent.

Data preprocessing

The data cleanup and very basic feature selection (only removing those with too many NAs) was performed in Python. Then a small imputation was then performed in R (only for 5 remaining missing values) within the same Jupyter Notebook (called *preprocessing.ipynb*).

Methodology and steps

1. Prior to importing the data, I performed a find and replace within the .csv file with the following regex:

```
\s\[["\w\S"]{1,}\]
```


to remove every instance of it. This regex matches a metadata tag that the world bank uses in their dataset. Following this step, the data was imported.
2. Metadata at the end of the dataset was removed. This data corresponded with indexes and regions, whereas the rest of the dataset corresponded with data specific to the countries conforming these indexes/regions. We only keep the countries.
3. The year column was converted to integer and the '.' values (which represent NAs in the World bank databank) were replaced for numpy NaNs and then the values were sorted by year and country name.
4. Data missing in later years (2021, 2020) was backfilled with data from previous years (2000-2019), as it still fits our modelling purposes, many of the backfilled data points still correspond to metrics that fit our criteria. The oldest possible data point would go as far back as the year 2000, however, it's unlikely we should have data this old, most should have been backfilled.
5. The dataset which contains the COVID-19 cases, deaths, population density, population, median age, etc.. is imported, we make a short variable selection within it, only keeping values which correspond to May 1st, 2021. Then, we merge this dataset with the previous dataset utilizing ISO code and using an inner join, only keeping countries in common in both datasets.
6. Columns with more than 45 NA values were removed as this represents around 25
7. Numerical columns were converted into floats and column names were simplified for easier future manipulation and to call the column names in a more practical way.
8. We take the **mobile_subscriptions** variable, and convert it to a percentage of the population, in this case, we leave it in ratio form and do not multiply by 100, as we prefer a 0-1 scale.
9. The categorical target variable is constructed as explained previously to this step.
10. The data is then exported as a csv (the dataset itself, called *data.csv*), along with a json (*columns.json*) file which contains the renaming used for the columns in the originally preprocessed dataset up to that step.
11. The data is then imported in an R cell and a *mice* imputation is performed using the 'cart' method with $m = 5$ in order to impute the very few missing values remaining (about 5 missing values).

Exploratory Data Analysis

Correlation

Here we obtain the correlation matrix maximized per correlation coefficient (using all three: Kendall, Pearson and Spearman) for our dataset.

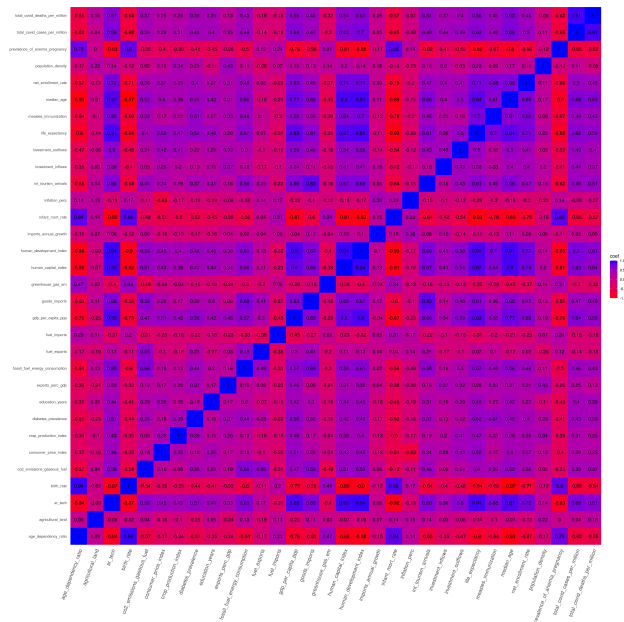


Figure 1: Correlation matrix maximizing correlation between variables

There's several things we can notice here that could preliminarily tell us how some of these variables could be excellent predictors.

The most relevant values to compare with our target variable are:

Very high correlation

Negative

- *age_dependency_ratio*: -0.84
- *birth_rate*: -0.87
- *infant_mort_rate*: -0.92
- *prevalence_of_anemia_pregnancy*: -0.93

Positive

- *gdp_per_capita_ppp*: 0.88
- *human_capital_index*: 0.92
- *human_development_index*: 0.94
- *life_expectancy*: 0.91
- *median_age*: 0.87

All these metrics are clearly all significantly descriptive of how developed a country is, those positive ones tell us that the higher their value, the more developed a country is, and the opposite for the negative ones.

A reasonable hypothesis can be proposed, that our constructed target variable is a reasonably good descriptor of development as a whole. These high correlation variables could be instrumental to our predictions

Moderately correlated

Positive

- *total_covid_cases_per_million* and *total_covid_deaths_per_million* at 0.69 and 0.61 respectively
- *net_enrollment_rate*: 0.72
- *int_tourism_arrivals* and *fossil_fuel_energy_consumption* at 0.66 and 0.65 respectively
- *goods_imports*: 0.68
- *investment_outflows*: 0.6

These moderately correlated variables could also be interesting at predicting, perhaps for better or for worse depending on how we treat them. But it's interesting to see how correlated COVID-19 cases and deaths are to our constructed target index. We could go back to a previous inference we made where we speculated that the most developed and some developing countries were hit the hardest when it comes to the pandemic, while other developing countries, while being hit hard economically, also handled the pandemic much better than most developed countries.

The rest of the metrics in this category we can infer that are related to each country's development measure as well, but maybe not as much as those previous ones we saw.

Not correlated almost *at all*

- *agricultural_land*
- *imports_annual_growth*
- *fuel_exports*
- *population_density*
- *inflation_perc*

These variables maybe relate with development, maybe don't, but for the most part, we can't comfortably say they're very significantly related to it.

They will be used, but I suspect these will be the first ones to be ditched by variable selection algorithms or models themselves.

Modelling: statistical tools

Onto the modelling part. We will test several models, simple and robust ones, however, all will be done utilizing 5 repeats of 10-fold cross validation. Defined as follows using the *caret* package:

```
ctrl <- trainControl(method = "repeatedcv",  
                     number = 10, repeats = 5)
```

Simple models

Simple linear regression (univariate):

We will test simple and multiple regression model for each set of variables, of those that have the highest correlations with our target variable.

Prevalence of anemia in pregnancy to predict ART index

Utilizing prevalence of anemia in pregnancy to predict ART index, we obtain an Rsquared of ~0.69 and a relatively low MAE of ~0.099. It's a very simple model, but we can say we're off to a good start.

Results table

Table 1: ART index ~ Prevalence of anemia in pregnancy

RMSE	MAE	Rsquared
0.1279857	0.0988103	0.6908135

Infant mortality rate to predict ART index

Testing with a different strongly negatively correlated variable, in this case infant mortality rate, a metric that's usually quite high in countries with low ART index, we can see that the model is reasonably better than that of prevalence of anemia in pregnancy ~ ART index.

We obtain an Rsquared of ~0.79 and a MAE of ~0.079, which is in itself an acceptable model. Perhaps this would be one of those variables that contributes to more robust and multivariate models.

Results table

Table 2: ART index ~ Infant mortality rate

RMSE	MAE	Rsquared
0.1072522	0.0785578	0.7967469

Human development index to predict ART index

As an analogous index to HDI, it is reasonable to check whether it is acceptable at predicting it in a simple linear regression model.

And we find that yes!, it's actually quite good, with an Rsquared of ~0.88 and MAE of ~0.06. This variable is probably going to contribute somewhat well in more robust and multiple models.

Results table

Table 3: ART index ~ HDI

RMSE	MAE	Rsquared
0.080499	0.0605895	0.8797171

Simple linear regression (multivariate):

For multivariate linear regression we first test out with all the variables, utilizing the same trainControl defined as before:

Using all variables to predict ART index

Utilizing all the variables yields a reasonably good model.

The Rsquared is of ~0.89 and the MAE is of ~0.061. This is only slightly better than when using only HDI.

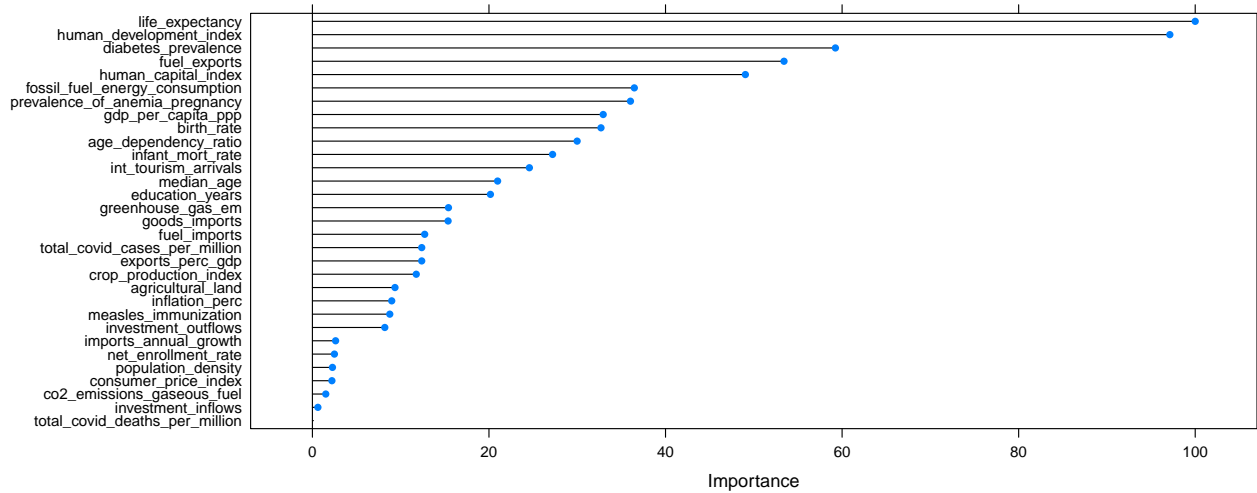
Results table

Table 4: ART index ~ .

RMSE	MAE	Rsquared
0.0810849	0.0632482	0.8863016

Variable importance

As we suspected, we can see that the most important variables are headed by *life expectancy*, *HDI*, and a few others like *diabetes prevalence*, *fuel exports* and *human capital index*. Interestingly, some of these we did not expect to see up here, but seem to yield decent results along with the rest.



From our variable importance plot, we can then test a model where the variable importance only exceeds 25% importance.

Using only variables with 20%+ importance

We filter the previously shown data used to plot variable importance, and we only keep those with 20% or more importance filtering the table.

And looking at our results, we can see there's been a reasonable improvement in model quality, these apparently small incremental improvements when we have a very high Rsquared are quite welcome, and as

we have done utilizing the following set of variables:

The model yields an Rsquared of ~ 0.915 and a MAE of ~ 0.054 , both improving on previous model attempts.

Unfortunately, the COVID-19 metrics do not really contribute as much as we had expected before, yielding a quite low variable importance in the previous all variable model.

Results table

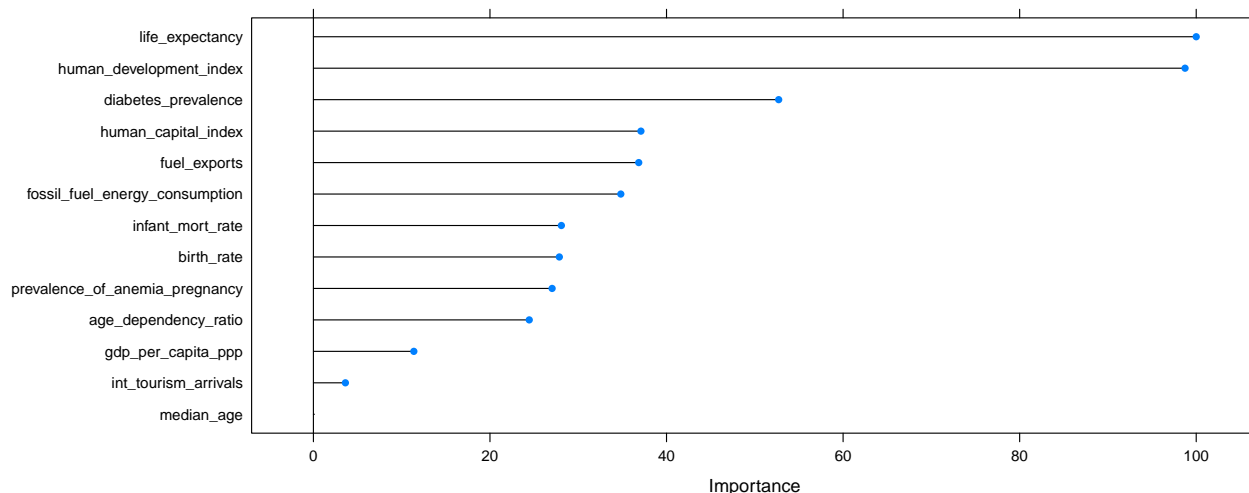
Table 5: ART index $\sim 20\%$ + importance variables from all variable model

RMSE	MAE	Rsquared
0.0692502	0.0543052	0.9132416

Variable importance

We also want to see variable importance after clearing up variables that are less important, and we see that the top variables remain the same with a few others going down in importance, much less so than in the previous models.

Before removing these, we want to test a model with all possible interactions, just to determine whether these serve a good purpose in the model.



Robust linear model

For the robust models, we test out once again with all the variables.

We first perform a train-test split with 75/25 proportions. and then run *rlm* using the previously defined formula with that preliminary variable selection.

Running the model we obtain the following RMSE:

```
#> [1] 1.166179
```

Which in an of itself is not a bad metric, but certainly inferior to the previous models, although those had the advantage of cross validation, which might have significantly improved results.

And Rsquared:

```
#> [1] 0.8818609
```

As for the Rsquared we get a quite decent result at ~ 0.89 , somewhat on par with the simple model using only HDI to predict ART index.

Using OLS model selection to improve on the simple linear model with all vars

OLS forward based on p-val

This method yields a model which performs a similar variable selection to that made through the variable importance plot, with the addition of a few other variables, however, it achieves a better model.

The model gives up some AIC points for a lower RMSE and higher Rsquared, yielding a resulting model with an Rsquared of ~ 0.925 and an RMSE of ~ 0.0642 .

```
#>
#>
#> Selection Summary
#> -----
#> Variable Entered R-Square Adj. R-Square C(p) AIC RMSE
#> -----
#> 1 human_development_index 0.8866 0.8857 48.6391 -307.0716 0.0745
#> 2 life_expectancy 0.9110 0.9096 12.6394 -337.0462 0.0662
#> 3 fossil_fuel_energy_consumption 0.9204 0.9185 0.0017 -349.7749 0.0629
#> 4 greenhouse_gas_em 0.9228 0.9204 -1.7517 -351.8321 0.0622
#> 5 infant_mort_rate 0.9250 0.9220 -3.1036 -353.5639 0.0615
#> 6 education_years 0.9266 0.9230 -3.6081 -354.4228 0.0611
#> 7 fuel_exports 0.9279 0.9238 -3.6895 -354.8467 0.0608
#> 8 fuel_imports 0.9292 0.9246 -3.7114 -355.2446 0.0605
#> 9 crop_production_index 0.9299 0.9247 -2.7368 -354.4777 0.0605
#> -----
```

OLS forward based on AIC

Here we get a similar result to the previous one but with a more simplified set of variables. The quality of the model is very similar if not basically identical with an Rsquared of ~ 0.9246 and a similar AIC.

```
#>
#>
#> Selection Summary
#> -----
#> Variable AIC Sum Sq RSS R-Sq Adj. R-Sq
#> -----
#> human_development_index -307.072 5.639 0.721 0.88660 0.88573
#> life_expectancy -337.046 5.794 0.566 0.91099 0.90961
#> fossil_fuel_energy_consumption -349.775 5.853 0.506 0.92039 0.91853
#> greenhouse_gas_em -351.832 5.869 0.491 0.92280 0.92037
#> infant_mort_rate -353.564 5.882 0.477 0.92495 0.92197
#> education_years -354.423 5.893 0.467 0.92656 0.92304
#> fuel_exports -354.847 5.901 0.459 0.92790 0.92383
#> fuel_imports -355.245 5.909 0.450 0.92919 0.92459
#> -----
```

OLS step backward based on AIC

On the backward based AIC approach, we end up with a similarly high quality model, only very slightly higher than the previous ones, resulting in an Rsquared of ~ 0.931 and AIC lower than the previous ones.

```
#>
#>
#> Backward Elimination Summary
#> -----
#> Variable AIC RSS Sum Sq R-Sq Adj. R-Sq
#> -----
#> Full Model -322.171 0.408 5.951 0.93580 0.91590
#> investment_inflows -324.160 0.408 5.951 0.93579 0.91672
#> co2_emissions_gaseous_fuel -326.111 0.408 5.951 0.93577 0.91751
#> infant_mort_rate -328.072 0.409 5.951 0.93575 0.91829
#> consumer_price_index -330.010 0.409 5.951 0.93572 0.91903
#> measles_immunization -331.968 0.409 5.951 0.93570 0.91978
#> median_age -333.883 0.409 5.951 0.93566 0.92048
#> agricultural_land -335.736 0.410 5.950 0.93559 0.92114
#> net_enrollment_rate -337.515 0.410 5.949 0.93548 0.92174
```

```
#> human_capital_index      -339.296    0.411    5.949    0.93537    0.92233
#> population_density      -341.007    0.412    5.948    0.93523    0.92287
#> crop_production_index   -342.720    0.413    5.947    0.93509    0.92339
#> total_covid_cases_per_million -344.296    0.414    5.946    0.93488    0.92383
#> total_covid_deaths_per_million -346.167    0.415    5.945    0.93482    0.92443
#> inflation_perc          -347.675    0.416    5.944    0.93457    0.92482
#> gdp_per_capita_ppp      -348.933    0.418    5.941    0.93421    0.92505
#> imports_annual_growth   -350.025    0.421    5.938    0.93375    0.92518
#> investment_outflows     -351.227    0.424    5.936    0.93335    0.92537
#> goods_imports           -351.907    0.428    5.932    0.93268    0.92526
#> int_tourism_arrivals    -353.187    0.430    5.929    0.93231    0.92549
#> exports_perc_gdp        -353.838    0.435    5.925    0.93162    0.92535
#> fuel_imports            -353.911    0.441    5.918    0.93061    0.92488
#> prevalence_of_anemia_pregnancy -354.189    0.447    5.913    0.92970    0.92451
#> -----
```

OLS stepwise AIC

This approach performs similar to the p-value approach, and results in also a similarly quality model, with an Rsquared of ~0.923 and also a similar AIC, lower than that of the backward model, but similar to the others.

```
#>
#>
#> Stepwise Summary
#> -----
#> Variable          Method      AIC      RSS      Sum Sq      R-Sq      Adj. R-Sq
#> -----
#> human_development_index addition -307.072    0.721    5.639    0.88660    0.88573
#> life_expectancy      addition -337.046    0.566    5.794    0.91099    0.90961
#> fossil_fuel_energy_consumption addition -349.775    0.506    5.853    0.92039    0.91853
#> greenhouse_gas_em     addition -351.832    0.491    5.869    0.92280    0.92037
#> infant_mort_rate      addition -353.564    0.477    5.882    0.92495    0.92197
#> education_years       addition -354.423    0.467    5.893    0.92656    0.92304
#> fuel_exports          addition -354.847    0.459    5.901    0.92790    0.92383
#> fuel_imports          addition -355.245    0.450    5.909    0.92919    0.92459
#> -----
```

Advanced Regression models

For these models we stick with the variables we had previously selected as important with the addition of all 2-way interactions among them, as these advanced models will optimize on the variable selection as they go.

Backward regression

For the backward regression model, we obtain a very strong model with a low ~0.0546 MAE and a decent ~0.91 Rsquared. I reckon this could be improved significantly, but it's already better than most of the previous ones seen.

Table 6: ART index ~ 20%+ importance variables from all variable model and 2-way interactions

RMSE	MAE	Rsquared
0.0695646	0.0535982	0.9120786

Lasso Regression

As for Lasso Regression, we somewhat improve upon the previous result from backward regression. Here we obtain an Rsquared of ~0.921 and a MAE of ~0.052.

Table 7: ART index ~ 20%+ importance variables from all variable model and 2-way interactions

RMSE	MAE	Rsquared
0.0662578	0.0514261	0.9190318

Stepwise Regression

In stepwise regression we obtain a similar result to the previous one, with also quite good metrics and a significantly less computationally intensive model, although at our scale, this does not pose a noticeable difference in runtime.

The Rsquared is of about ~0.91-0.915 and the MAE lands between that of the Lasso and Backward regression models, at ~0.053

Table 8: ART index ~ 20%+ importance variables from all variable model and 2-way interactions

RMSE	MAE	Rsquared
0.0676918	0.05314	0.9140535

Partial Least Squares

As for partial least squares, we exclude interactions, as these land a significantly worse model. We stick with simply including the variable selection we had done prior in the simple multivariate regression. Here we obtain an Rsquared lower than the previous ones, at around ~0.887 and a MAE of ~0.0619. The model is good, but not good enough compared to the previous ones.

Table 9: ART index ~ 20%+ importance variables from all variable model

RMSE	MAE	Rsquared
0.0790439	0.061777	0.8880343

Ridge Regression

The Ridge regression model is also very computationally intensive, therefore we do without the interactions. But we also obtain good results, similar to that of the stepwise model in terms of quality. An Rsquared of ~0.916 and a MAE of ~0.0536

Table 10: ART index ~ 20%+ importance variables from all variable model

RMSE	MAE	Rsquared
0.0681362	0.0535956	0.9162932

Principle Component Regression

Similar to the partial least squares model we obtain a similar result in terms of Rsquared at ~0.89 and a MAE of ~0.062. Testing with all the variables, like with the previous models, yields worse results, therefore we stick to the previous variable selection and here we also do not include interactions, this yields a bad model with an extremely low Rsquared and 4-5x higher MAE compared to the rest of models in this section.

Table 11: ART index $\sim 20\%$ + importance variables from all variable model

RMSE	MAE	Rsquared
0.0791306	0.0628011	0.8846109

Elastic Net

Using 2-way interactions to improve upon the model, here in the Elastic Net model we get a lot more to choose from, and ends up resulting in quite good models, with a somewhat small variance among runs.

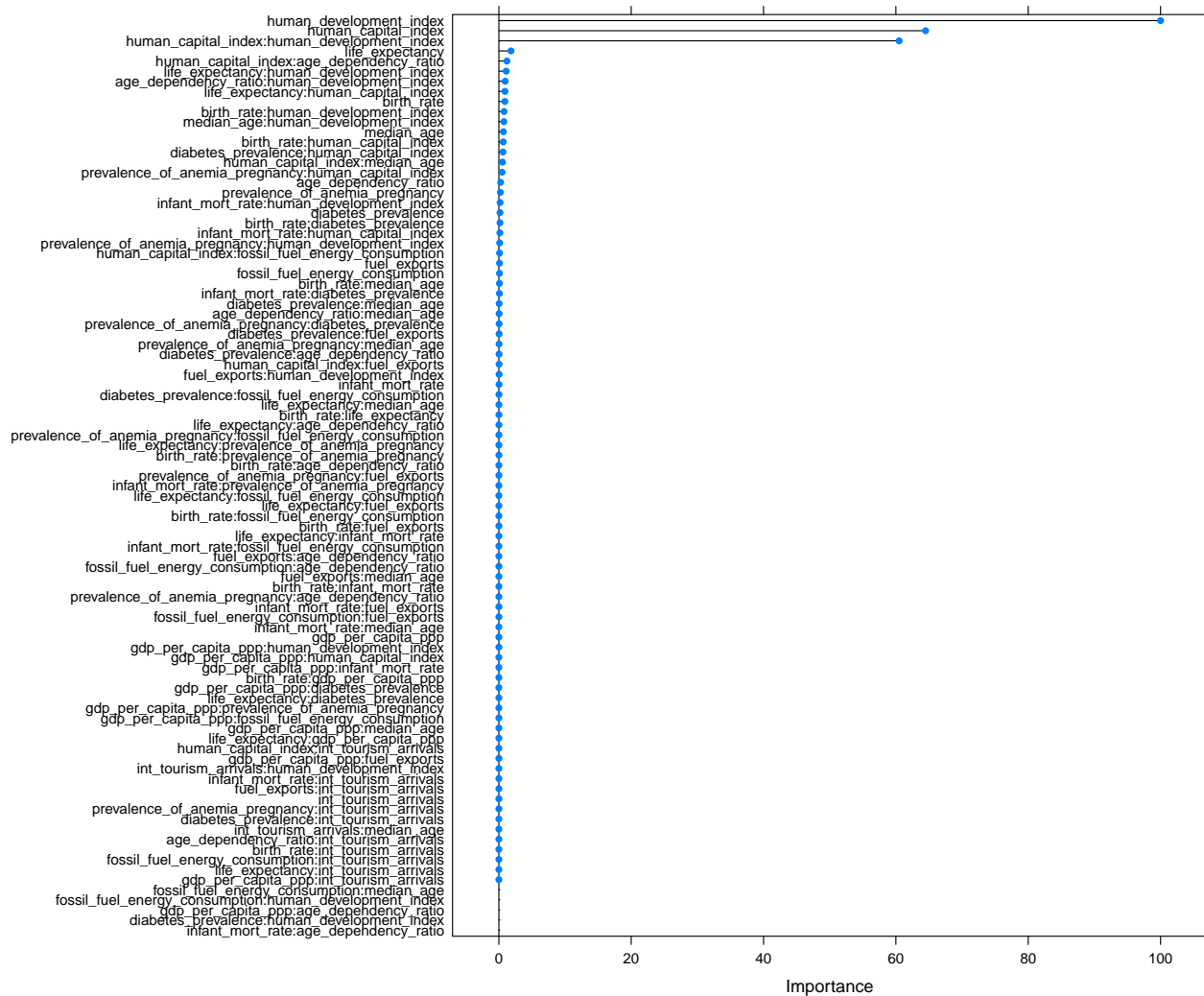
Here we obtain an Rsquared of ~ 0.924 and a MAE of ~ 0.051 . This model seems to provide the best results for our purposes so far.

Table 12: ART index $\sim 20\%$ + importance variables from all variable model and 2-way interactions

RMSE	MAE	Rsquared
0.0642156	0.0502014	0.9237852

Variable importance for Elastic Net

We can see the top variables for Elastic Net are HDI, HCI and its interaction followed by life expectancy in a much lower tier.



Trying an elastic net model with only the selected variables in VarImp

Here we get a very very good model, very simple, which only uses HDI, HCI and its interaction as prioritized by Elastic Net. This would be our final model with an Rsquared of ~0.9236 and a MAE of ~0.05

Table 13: ART index ~

RMSE	MAE	Rsquared
0.0839563	0.0649935	0.8676679