

Statistical Learning Final Project: Predicting Country HDI

Daniel Alonso

December 19th, 2020

Contents

Dataset of choice	2
Variables	2
The target variable	3
Data preprocessing	3
Methodology and steps	3
Exploratory Data Analysis	4
Looking at variables by themselves	4
Significantly right-skewed variables	4
Education variables	5
Other economic variables	5
Other interesting variables	6
Correlation matrix	7
Analysis	8
Modelling	9
Method 1: Linear Discriminant Analysis	9
Confusion matrix and variable importance	9
Method 2: Logistic Regression	9
Confusion matrix and variable importance	9
Method 3: Quadratic Discriminant Analysis	10
Confusion matrix and variable importance	10
Method 4: K-Nearest Neighbors	10
Confusion matrix and variable importance	10
Conclusion and final thoughts on the first part of the project	11

Dataset of choice

For this project I decided to pick 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 has lots of economic, education, energy use, and population specific metrics.

I find demographic data fascinating, and I think this dataset will be quite good for predicting country development measures along with providing quite interesting and relevant information.

Variables

NOTE:

- **blue** = used for training/predicting
- **red** = target variable
- **green** = ID variables
- **purple** = variables excluded as either they were *components* of HDI or they were 100% correlated to another variable (like GNI/GDP, which are both 100% correlated and GNI is a component of HDI)

Variables in the original dataset as constructed using the World Bank Databank tool (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
- **foreign_inv_inflows**: Foreign direct investment, net inflows (BoP, current US\$)
- **exports_perc_gdp**: Exports of goods and services (as a % of GDP)
- **inflation_perc**: Inflation, consumer prices (annual %)
- **education_years**: Compulsory education, duration (years)
- **education_perc_gdp**: Government expenditure on education, total (as a % of GDP)
- **gds_perc_gdp**: Gross domestic savings (as a % of GDP)
- **gross_savings_perc_gdp**: Gross savings (as a % of GDP)
- **int_tourism_arrivals**: International tourism, number of arrivals
- **int_tourism_receipts**: International tourism, receipts (in current US\$)
- **perc_internet_users**: Individuals using the Internet (as a % of population)
- **access_to_electricity**: Access to electricity (% of population)
- **agricultural_land**: Agricultural land (% of land area)
- **birth_rate**: Birth rate, crude (per 1,000 people)
- **gne**: Gross national expenditure (% of GDP)
- **mobile_subscriptions**: Mobile cellular subscriptions (per 100 people)
- **infant_mort_rate**: Mortality rate, infant (per 1,000 live births)
- **sex_ratio**: Sex ratio at birth (male births per female births)
- **greenhouse_gas_em**: Total greenhouse gas emissions (kt of CO2 equivalent)
- **urban_pop_perc**: Urban population (% of total population)
- **hdi**: human development index
- **hdi_cat**: Human development index as a category
- **life_exp**: Life expectancy at birth, total (years)
- **gdp**: GDP (current US\$)
- **gni**: GNI (current US\$)
- **fertility_rate**: Fertility rate, total (births per woman)

The target variable

As all these variables could perhaps tell us how developed a country is, we used a constructed categorized Human development index variable in order to classify the countries using the above variables (unless stated otherwise by their colour) as training variables.

The criteria for constructing the categorical variable *hdi_cat* was the following:

- Very high: HDI above 0.8
- High: HDI between 0.7 and 0.799
- Medium: HDI between 0.55 and 0.699
- Low: HDI under 0.55

This categorization is emulated from Wikipedia's construction and uses the same ranges as used in every Wikipedia article referencing HDI.

Data preprocessing

The data cleanup and preliminary feature engineering was done in Python and the imputation was then done in R within the same Jupyter Notebook (called *preprocessing.ipynb*).

Methodology and steps

1. Prior to importing there was a search within the .csv file with the following regex:

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

as a find and replace and removing 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, the dataset was filtered excluding these region codes (the codes were excluded as they were country aggregates, and we're only interested in the countries themselves).
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 (2020, 2019) was backfilled with data from previous years, as it still fits our modelling purposes. The data that goes the furthest back is from 18 years prior, so 2002.
5. Columns with more than 45 NA values were removed as this represents roughly 25
6. The life expectancy, GDP, GNI and fertility rate columns were removed as these are either components of HDI or (in the case of fertility rate) are 100
7. We scrape Wikipedia in order to obtain an updated metric on HDI estimates for each country. We also scrape it to obtain the alpha-3 3166 codes for each country. The data was obtained with the purpose of making it easier to join with the dataset obtained from the World Bank, as joining by country names is not reliable enough (misses 20 countries which the alpha-3 codes do not).
8. We add HDI to the main dataframe
9. Numerical columns were converted into floats and column names were simplified for easier future manipulation.
10. The categorical target variable is constructed as explained previously in this step.
11. The data is then exported as a csv, along with a json file that shows how the columns were renamed
12. 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.

Exploratory Data Analysis

Looking at variables by themselves

The plots of choice for visually analyzing the variables prior to any processing has been Boxplots and Histograms/density plots all looked at individually and categorized by *hdi_cat* which is our categorical target variable (HDI).

Given the space constraints of this report, if a plot for a specific variable is not shown here but talked about or mentioned, it will be shown in the *eda.ipynb* and *eda.html* notebook and notebook export used to perform the visual analysis. Only variables which are used in the modelling section will be visually analyzed and described.

Significantly right-skewed variables

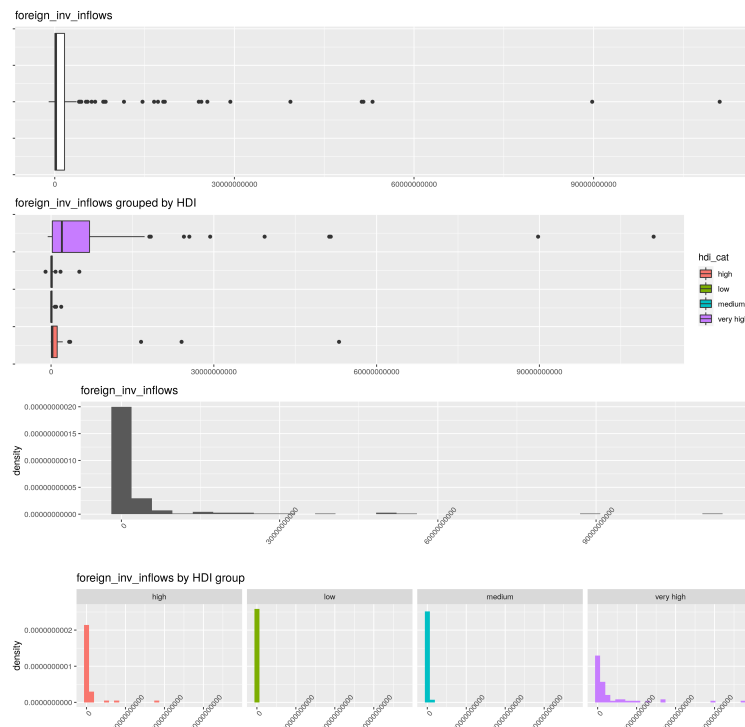


Figure 1: Foreign Investment Inflows

The following variables showed a similarly right-skewed shape with significant outliers:

- **foreign_inv_inflows:** Heavily right skewed, significant outliers (USA, UK, Germany, China...) etc... Given the histogram of *foreign_inv_inflows* when grouped by very high HDI vs low HDI, it would seem like the more developed a country is, the more foreign investment it should have. The grouped boxplot also shows the same trend.
- **int_tourism_receipts** and **int_tourism_arrivals:** Basically the same as with *foreign_inv_inflows*, we see that there's a significant tendency for higher development nations to receive significantly more tourism receipts and arrivals.
- **greenhouse_gas_em:** Just like with the previous two variables, fossil fuel use, meat production and economic activities like such produce huge amounts of greenhouse gases, and the more developed a country is, the more emissions it produces. However, this one is a bit less strongly inclined like the previous two. We can clearly see that there's many less developed countries with a high level of emissions. A notorious example of this is China, which is the 2nd country with the most emissions and

falls under the high category of HDI, by far exceeding the emission levels of most countries with very high HDI countries.

Education variables

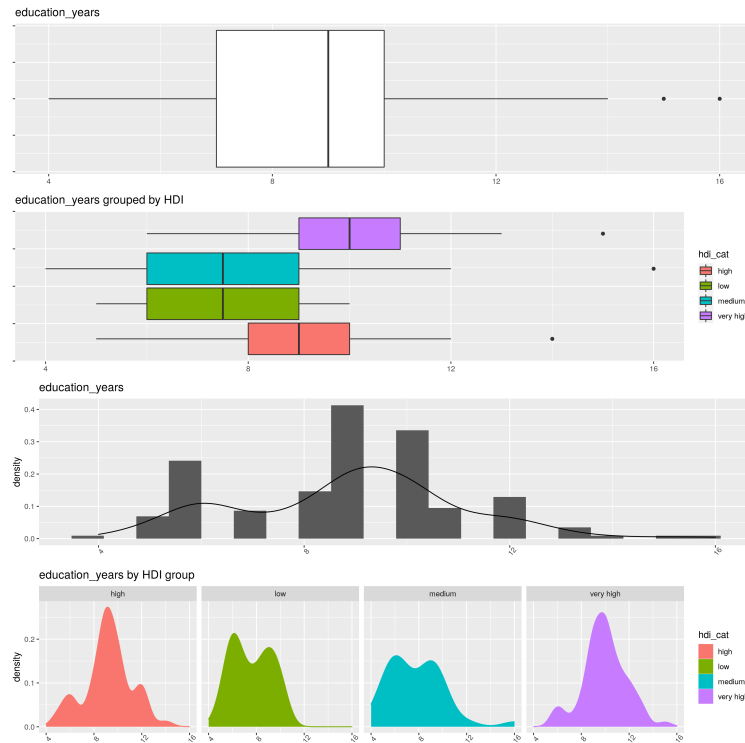


Figure 2: Years of compulsory education

- **education_years:** There's a clear tendency, both shown by our grouped boxplots and grouped density plots, for very high HDI countries to have significantly longer periods compulsory education, with few exceptions in each group, however a clear tendency that peaks between 8 and 12 years of education for both high and very high HDI countries. Low and medium HDI countries could be clumped together with similar compulsory education times, but with more variability for medium HDI countries.
- **education_perc_gdp:** How much are countries spending on education as a percentage of their GDP is a very interesting variable, as we can see that the consistency for higher and lower spendings respectively for low and very high HDI countries is quite solid, while medium HDI countries show an incredible variability versus other groups. High HDI countries show significant variability versus that of very high and low HDI countries, however, still aligning more iwth very high HDI countries than the rest.

Other economic variables

These variables might have some skewness towards more tails or longer tails than a normally distributed variable should have. Therefore they're hard to classify as particularly right or left skewed.

- **exports_perc_gdp:** This variable has a long left tail and it is right-skewed. We can see the box for medium developed countries is larger than others while very high and high HDI countries tend to have higher export amounts as percentage of GDP. Low HDI countries lag behind, as expected. It is interesting to see how medium HDI countries completely bridge the gaps between high, very high and low HDI countries in terms of exports, in the sense that there's plenty of countries with medium HDI with exports just as high as others with significantly higher HDIs.

- **inflation_perc**: For inflation we can see that independently of the HDI of a country, inflation could, perhaps, be an inevitable event of economic/political management or mismanagement and uncertainty. Either way, we can clearly see that the higher the HDI, the less uncertain such inflation rate will be. The boxes for medium and low HDI countries are significantly larger than very high and high therefore telling us that there's less consistency and high inflation events, while universal, are unpredictable, but less unpredictable and probably less common the higher the HDI of a country is.
- **gds_perc_gdp**: For gross domestic savings we can see that there's a clear difference between groups, with much more overlapping on the high end of HDI than the medium to low end of HDI. However, clearly, once again, the higher the HDI the higher the GDS. Much more variability again on those middle HDI groups (high and medium).
- **gross_savings_perc_gdp**: Gross savings as a percentage of GDP is an interesting variable, where the only defining feature of high and very high HDI countries being, again, relative consistency versus other groups, there might be seemingly no predictive capability in it due to the extremely low variability of the values among groups however, it definitely differentiates very high and low HDI countries quite well.

Other interesting variables

These variables seemed like quite appropriate to me for predicting HDI or any other development measure, as they could potentially be huge differentiators between the different HDI groups.

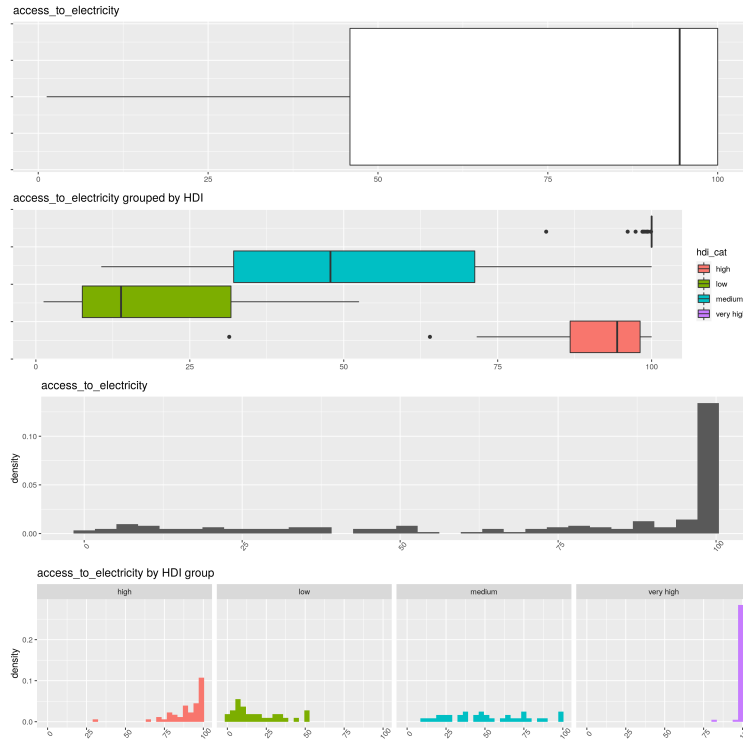


Figure 3: Percentage of population with access to electricity

- **perc_internet_users**: The percentage of internet users in each group is a surprisingly good differentiator as mentioned previously. We can see that, sure there's countries in each group with percentages that resemble other groups, but the difference in variability among groups is massive. Where the consistency of low HDI countries vs very high and high HDI countries is significantly different. We can see the standard deviations (and means) per group follow an interesting ladder: $\sigma_{low} \approx 0.714$, $\sigma_{medium} \approx 2.532$, $\sigma_{high} \approx 4.908$, $\sigma_{very_high} \approx 20.382$.
- **access_to_electricity**: another Hugely differentiating variable with similar properties to

perc_internet_users. We can see a clear division where if the country is in the very high HDI category, it will almost certainly have near full coverage in its electricity supply to its population. Countries in the high category are also clearly differentiated with the rest while also maintaining quite high rates. Medium HDI countries show a huge variability in their coverage, lower than high HDI countries but higher than low HDI countries. Low HDI countries show a significantly lower but clearly divided coverage with medium HDI countries.

- **agricultural_land**: Agricultural land is an interesting variable that might not contribute a great deal to classifying, however, it is interesting to visualize how lower HDI countries tend to have a significantly higher percentage of agricultural land. However the rest of the variables all hover around the same values.
- **birth_rate**: Birth rate is a clearly defining variable, where we can see the lower HDI countries rise above the rest, along with the higher HDI countries jumping way down. Most very high to high HDI countries have already experienced a population growth plateau where the birth rate goes down significantly and usually either stays between 0 and 15 or very slightly above it (as it is necessary for demographic sustainability).
- **gne**: Gross national expenditure shows slightly lower values for very high HDI countries than the rest of groups, with higher variability for high HDI countries and low to medium HDI countries slightly higher. The general histogram for this value shows a long left tail with a normal-like distribution.
- **mobile_subscriptions**: An amazingly interesting variable with very similar properties to that of access to electricity and percentage of internet users. We see a very clear spread where the HDI groups are very clearly separated with very high HDI having an incredibly high variability while still being significantly higher than the rest of groups and the high HDI group lagging surprisingly behind but still a lot higher than low and medium HDI countries. The variable is quite right-skewed with most of the countries clumping around the low-medium area as there isn't an incredibly large difference between low and medium HDI countries for this one.
- **-infant_mort_rate**: Quite left skewed as both high and very high HDI countries have most of their values closer to zero while low and medium HDI countries have a much larger spread with low HDI countries taking the top position given the poverty conditions. Medium HDI countries lag behind low HDI countries while still maintaining a really high infant mortality rate in general. The variable also clearly defines groups and might be quite good at classifying.
- **-sex_ratio**: seemingly somewhat normally distributed with a long left tail. However quite surprisingly well defined for very high and high HDI countries. The differences between groups are interesting to see but maybe not quite as differentiating as other variables under this section.
- **-urban_pop_perc**: Population of people living in urban areas is a quite spread out variable as we know the more developed a country is, the more its population tends to live in urban areas. Urban areas are usually higher income areas and tend to offer more opportunities for workers, therefore, highly developed nations, as it is clear, tend to have quite good job offerings in urban areas. We clearly see the groupings here, very high HDI countries have the largest urban populations, while low to medium HDI countries have the lowest. Interestingly so, low and medium HDI groups are not as well separated.

Correlation matrix

This correlation matrix was constructed by calculating the correlation coefficient for the variables using *Kendall*, *Spearman* and *Pearson* correlation coefficients and then inputting the values into a dataframe where only the highest correlation of the three calculated is added. This way we are able to really illustrate how correlated the variables truly are.

The idea came from observing significant differences between correlation calculated for a few pairs of variables, where the correlation coefficients nearly doubled when moving from *Pearson* to *Spearman* correlation coefficients.

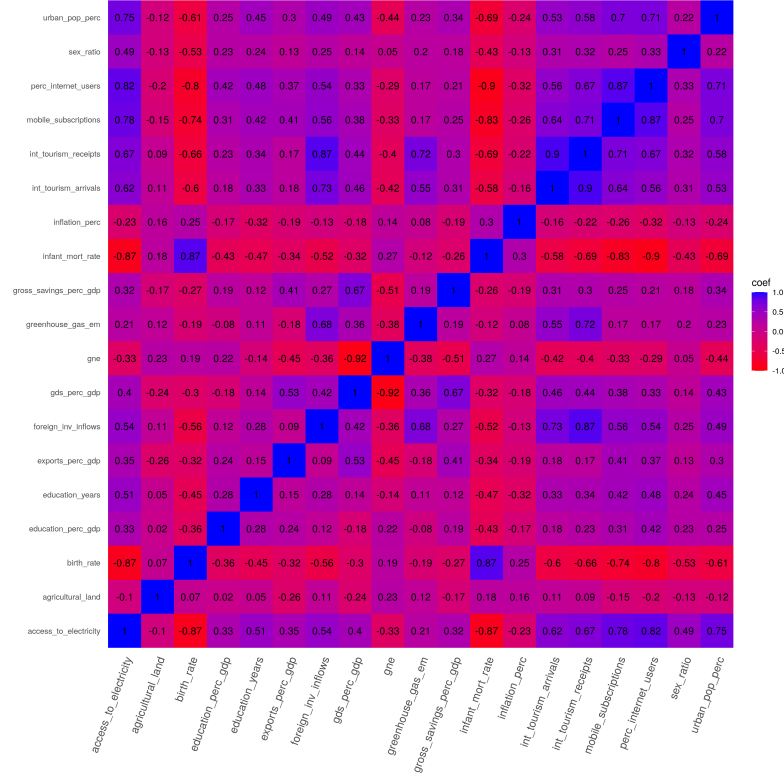


Figure 4: Correlation heatmap

Analysis

Looking at the correlations we highlight some important details:

- **int_tourism_receipts** and **int_tourism_arrivals** are extremely highly correlated at around 0.9 . These consideration will be looked at later on as they may or may not provide the model with interesting information when both are included in it.
- **access_to_electricity** and **perc_internet_users** are quite highly correlated at around 0.82 along with **perc_internet_users** and **mobile_subscriptions** at around 0.87 which also showed very similar trends in the individual variable analysis.
- **perc_internet_users** along with **access_to_electricity** and **infant_mort_rate**, both pairs with -0.9 and -0.87 respectively. The higher the infant mortality rate, the lower the percentage of internet users and people with access to electricity.
- **infant_mort_rate** and **birth_rate** at around 0.87 which is also quite high.

There's plenty of other interesting correlations to look at, however the general idea is that there's quite high correlation among some key variables. This makes a lot of sense, however, given how well some of these variables separate groups, this might not be a particular problem when predicting. Later on we will see whether this hurts the model in general or not.

When plotting a *pairs* plot (scatter plots) between the variables when grouped by HDI group, we see that there's a clear distinction similar to what we could appreciate in the boxplots, there's clear separation between groups in several of the plots.

Modelling

All models were prepared and ran in R using 3 repeats of 10-fold cross validation:

```
control <- trainControl(method = "repeatedcv",
                        repeats = 3,
                        number = 10 )
```

Method 1: Linear Discriminant Analysis

LDA is a generalization of fisher's linear discriminant. A method used to obtain a linear combination of deatures that classifies elements in two or more groups.

For our LDA model we used the following parametrization:

```
ldafit <- train(hdi_cat ~ ., method = "lda", data = data, preProcess = c("center", "scale"),
               tuneGrid = expand.grid(alpha = seq(0, 2, 0.1), lambda = seq(0, .1, 0.01)),
               metric = "Accuracy", trControl = control)
```

Confusion matrix and variable importance

```
#>           Reference
#> Prediction  high low medium very high
#>   high      48   0    7         9
#>   low        0  29    5         0
#>  medium      3   3   24         0
#>  very high   1   0    0        55
```

We can see in the confusion matrix for this model that the model classifies quite decently, achieving an accuracy of ~**84.78%**, with a *95% confidence interval* of (**78.76%**, **89.64%**).

As for variable importance we see that our LDA model uses with somewhat similarly high importance the variables *infant_mort_rate*, *access_to_electricity*, *birth_rate*, *perc_internet_users*, *mobile_subscriptions* and *urban_pop_perc* for most variables, therefore confirming what we had suspected previously, that these variables would be strong group separators.

Method 2: Logistic Regression

A logistic/logit model has the purpose of modelling hte propability of a certain class or event with 2 or more categories/groups.

For our LR model we used the following parametrization:

```
lrfit <- train(hdi_cat ~ ., method = "glmnet", data = data, preProcess = c("center", "scale"),
              tuneGrid = expand.grid(alpha = seq(0, 2, 0.1), lambda = seq(0, .1, 0.01)),
              metric = "Accuracy", trControl = control)
```

As shown above, we performed a grid search for both hyperparameters of the model, *lambda* and *alpha* with their respective intervals between 0 and 0.1 with steps of 0.01 for *lambda* and between 0 and 2 with steps of 0.1.

Confusion matrix and variable importance

```
#>           Reference
#> Prediction  high low medium very high
#>   high      48   0    9         8
#>   low        0  31    5         0
#>  medium      1   1   22         0
#>  very high   3   0    0        56
```

We can see in the confusion matrix for this model that the model classifies on par with the LDA model, achieving an accuracy of ~**83.15%**, with a *95% confidence interval* of (**76.95%**, **88.26%**).

The LR model seems significantly more selective with variables, essentially showing us that yes, the variables we previously suspected to be good classifying variables (*infant_mort_rate*, *access_to_electricity*, *birth_rate*, *perc_internet_users*, *mobile_subscriptions* and *urban_pop_perc*) had definitely a big role in classifying very high, high and medium HDI countries, however, this doesn't come without surprises, where a very relevant variable when classifying low HDI countries was *agricultural_land*, a variable that we didn't think would play a particularly big role in classification. However, while the model is good, it doesn't achieve the level of accuracy that we aim for (**greater than or equal to 90%**).

Method 3: Quadratic Discriminant Analysis

QDA is related to LDA where measurements per class are assumed to be normally distributed, however, it's not assumed that the covariance of each of the classes is identical.

QDA seems very appropriate for our purpose given our measurements and it shows its strengths in our model summary.

For our QDA model we used the following parametrization:

```
qdafit <- train(hdi_cat ~ ., method = "qda", data = data, preProcess = c("center", "scale"),
               metric = "Accuracy", trControl = control)
```

Confusion matrix and variable importance

```
#>           Reference
#> Prediction high low medium very high
#> high      45  0    1         5
#> low       0 32    1         0
#> medium    5  0   34         0
#> very high 2  0    0        59
```

We can see in the confusion matrix for this model that the model classifies very well, never placing medium or high HDI countries 2 categories above or below and reasonably classifying countries. The model achieves an accuracy of **~92.39%**, with a *95% confidence interval* of **(87.56%, 95.78%)** which is significantly higher than our LR and LDA models.

For QDA, our best performing model, we notice that the consideration of variables is near-equal to that of LDA, however, the model achieves better performance.

Method 4: K-Nearest Neighbors

For our KNN model we used the following parametrization:

```
knnfit <- train(hdi_cat ~ ., method = "knn", data = data, preProcess = c("center", "scale"),
               metric = "Accuracy", trControl = control)
```

Confusion matrix and variable importance

```
#>           Reference
#> Prediction high low medium very high
#> high      43  0   11         8
#> low       0 30   11         0
#> medium    2  2   14         0
#> very high 7  0    0        56
```

We can see in the confusion matrix for this model that the model classifies worse than previous models, making several mistakes in the middle groups (high and medium), achieving an accuracy of **~77.17%**, with a *95% confidence interval* of **(70.42%, 83.03%)**.

KNN variable importance strongly resembles that of LDA and QDA, and we could, perhaps, go as far as to say that for our conditions, it seems essentially identical. However, due to its performance versus the former 2, we do not use KNN for this analysis. At least not with the parametrization we have chosen. When testing

using grid search we didn't really achieve a hugely significant difference in results, and even when higher, we didn't quite rival the excellent performance of QDA for our dataset.

Conclusion and final thoughts on the first part of the project

This has been a highly interesting project so far, we looked at data that I personally find profoundly interesting and even culturally useful. We have found a fantastically well performing model for our purposes and learned a lot in the process.

We can conclusively say that the variables that best serve us a purpose when classifying HDI in groups are those that speak about the lifestyle and utilities of individuals in countries, things like percentage of people subscribed to a mobile data plan, or percentage of individuals with access to internet and electricity, birth rate, fertility rates, and infant mortality rate. While we obtained a sad conclusion, which is that countries with low HDI basically have the worst end of what quality of life services could be described as, we see and know that there is a lot to improve in the world.

The best performing model for our analysis was QDA, while I did test several other models using similar parametrizations and such, however, either these were out of the scope of this part of the project, or they were incredibly computationally intensive.

I look forward to part 2! On the next and last page the variable importance of the models is showcased.

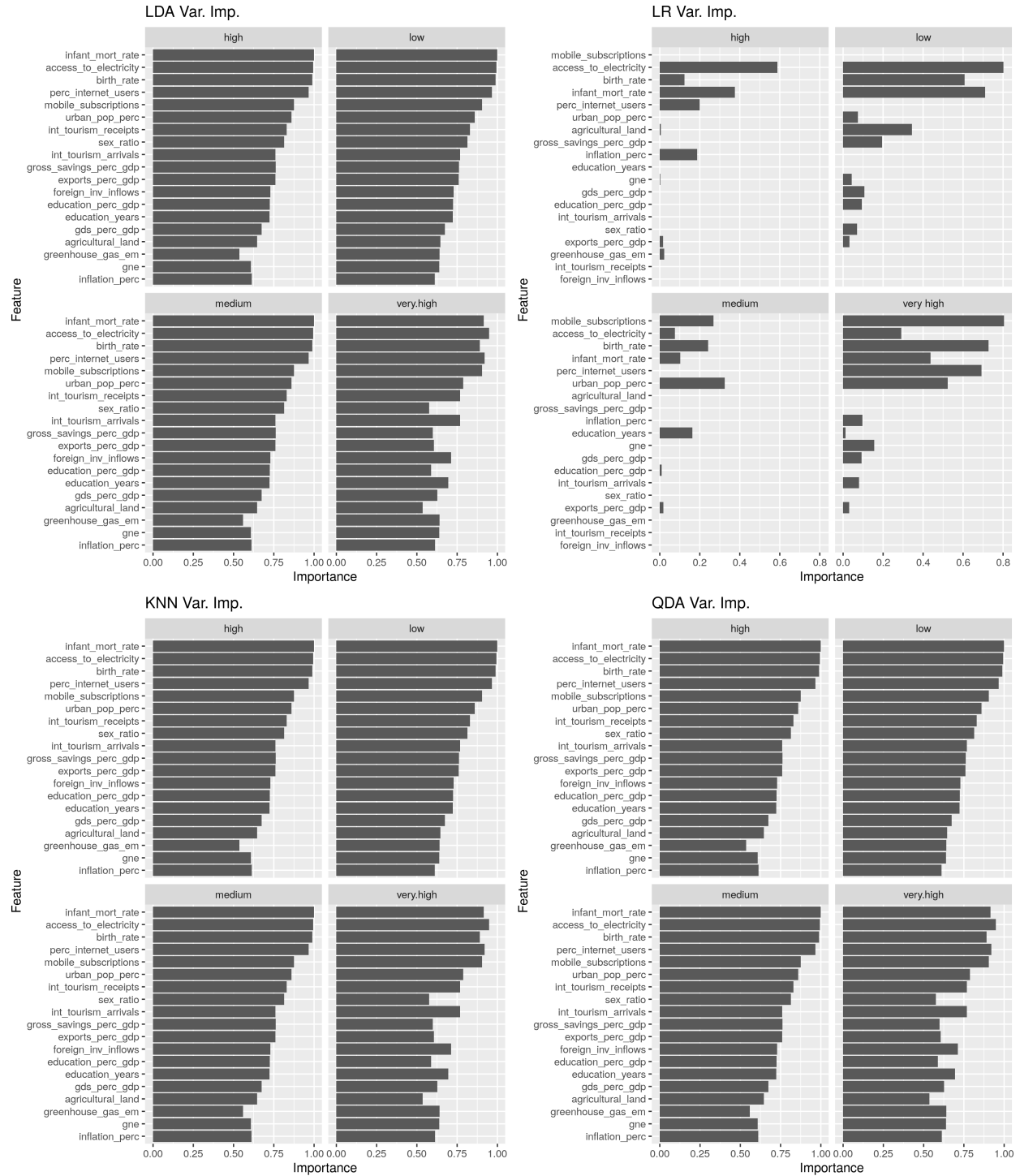


Figure 5: Variable importance per model