

The driving factors behind food prices in developing countries

Benjamin Jaidi¹, Stefan Ramakrishnan¹, Raiber Alkurdi¹

¹Ciência da Computação – Faculdade Integrada da Grande Fortaleza (FGF)
Av. Porto Velho, 410 – 60.510-040 – Fortaleza – Ce – Brasil

{elidiane,Edvan,hitalo}@fgf.edu.br

Contents

1	Introduction	2
1.1	Related Work	2
2	Methodology	3
2.1	Workflow	3
2.2	Dataset Descriptions	5
2.3	Function Descriptions	6
3	Implementation	8
3.1	Data Preparation	8
3.2	Data overview	10
3.3	Variable Selection	12
3.3.1	VIF based method	12
3.3.2	Lasso based method	14
3.3.3	Random Forest based method	16
3.4	Result Exploration	17
4	Content related analysis	18
4.1	Datasets	18
4.2	Results	20
4.2.1	Correlation / VIF based method	20
4.2.2	Lasso based method	29
4.2.3	Random Forest based method	31
4.2.4	Summary	34
5	Outlook	35
6	Conclusion	36

1. Introduction

The issue of food security is important in our times of more extreme weather occurrences, political instability and economic uncertainty. According to the Food and Agriculture Organization of the United Nations (FAO), the demand for food is growing while production to cover this demand is being surpassed. For example the global production of grain will have reached 2.1 billion tons by 2030 whereas the demand for grain will have increased up to 2.7 Billion tons. (FAO, 2016). Although population growth will most likely be slower in the future, even now 2 billion people in developing countries spent up to 70% of their disposable income on food (Erokhin, 2017). This problem is further exacerbated for low incoming families. Ören recognizes income level as the most decisive variable for food security (Ören, 2013).

As pointed out by Erokhin an important factor contributing to food security are food prices, as well as economic, climatic conditions and political stability. We leaned on the (Erokhin, 2017) paper in the development of the contributing factor framework. The OECD has further identified a set of factors influencing crop prices. This paper copes with the influence on crop prices on domestic markets in selected developing countries. The goal is to find out to which extent the OECDs findings are also applicable to India, Rwanda and the Philippines. This country selection aims to cover variety of economic factors as both Rwanda and the Philippines are considered developing countries whereas India is considered a lower-middle income country. Geographical and population differences between these countries allow us to form a more comprehensive conclusion if the OECD approach is reproducible for a variety of base conditions. Apart from the raw data selection process the Factor selection is approach in a variety of ways, with importance of factors varying widely in the different frameworks. We leaned on the related literature for the creation of contributing factor framework and adjusted it on granular level to fit our needs. Erokhin groups the factors through two parameters. One the one hand he focuses on physical availability of food (domestic production and import), the second being economic access (purchasing power, food inflation, distribution etc.) This was take into consideration but we choose the overarching parameters division of supply and demand factors as our structure format. Similar to (Smith, 1990) he groups factors into supply (weather, production, policy incentives, stocks and imports) and demand factors (population growth, income growth and distribution, and export revenue). Our finished framework has the two overarching categories of supply factors - demand factors. The subcategories for the supply side are climatic factors consisting of rain and temperature data, production factors consisting of the oil price and production amount and macro economic factors consisting of the Agricultural GDP, inflation and imports. On the demand side the subcategories are demographic factors consisting of GNI, per capita caloric intake and population growth and macro economic factors consisting of export of goods. This selection of factors will be discussed in more detail in the Methodology part of our paper.

1.1. Related Work

This works mainly leaned on the “Establishing Food Security and Alternatives for International Trade in Emerging Economies” by Erokhin for the creation of our methodology as highlighted priorly in the introduction section. The OECD work regarding the rise of food prices and the common consequences was used as a benchmark. In their work they take global look at food prices development and how specific factors impacted these prices movements. The argument being that tight market conditions for essential agricultural commodities need to be understood so that national governments cant create meaningful policy answers (OECD, 2008). The OECD paper aims to aid governments in the policy making process which necessitates an understanding on which factors impact food prices. While the OECD goal is to aid the policy formation of national governments our work builds upon the general framework of the

OECD's understanding of food price impacting factors. The OECD created a global framework, whereas our work takes these global notions to test if they are applicable on a country specific level. Our work builds upon this more global framework by testing the OECD global approach on a country specific level in order to access if these general notions hold true when used within a smaller case scenarios.

2. Methodology

This section gives an overview of every step in the workflow from raw data to interpretable results as well as involved utility functions and the datasets.

2.1. Workflow

The process to achieve the desired results was structured into four main stages as depicted in **Figure 1**. Each step includes working R scripts that were used to generate the desired output files for the following steps, as well as Quantlets. These are supposed to work as runnable examples of every step conducted in the process. They are functionally independent of the working scripts and of their resources since every resource required by each quantlet is stored in the quantlet's own folder. Working scripts are executable as well except for scripts located in folder `Processing_scripts`, due to missing resource files that could not be uploaded because of githubs's file size constraints.

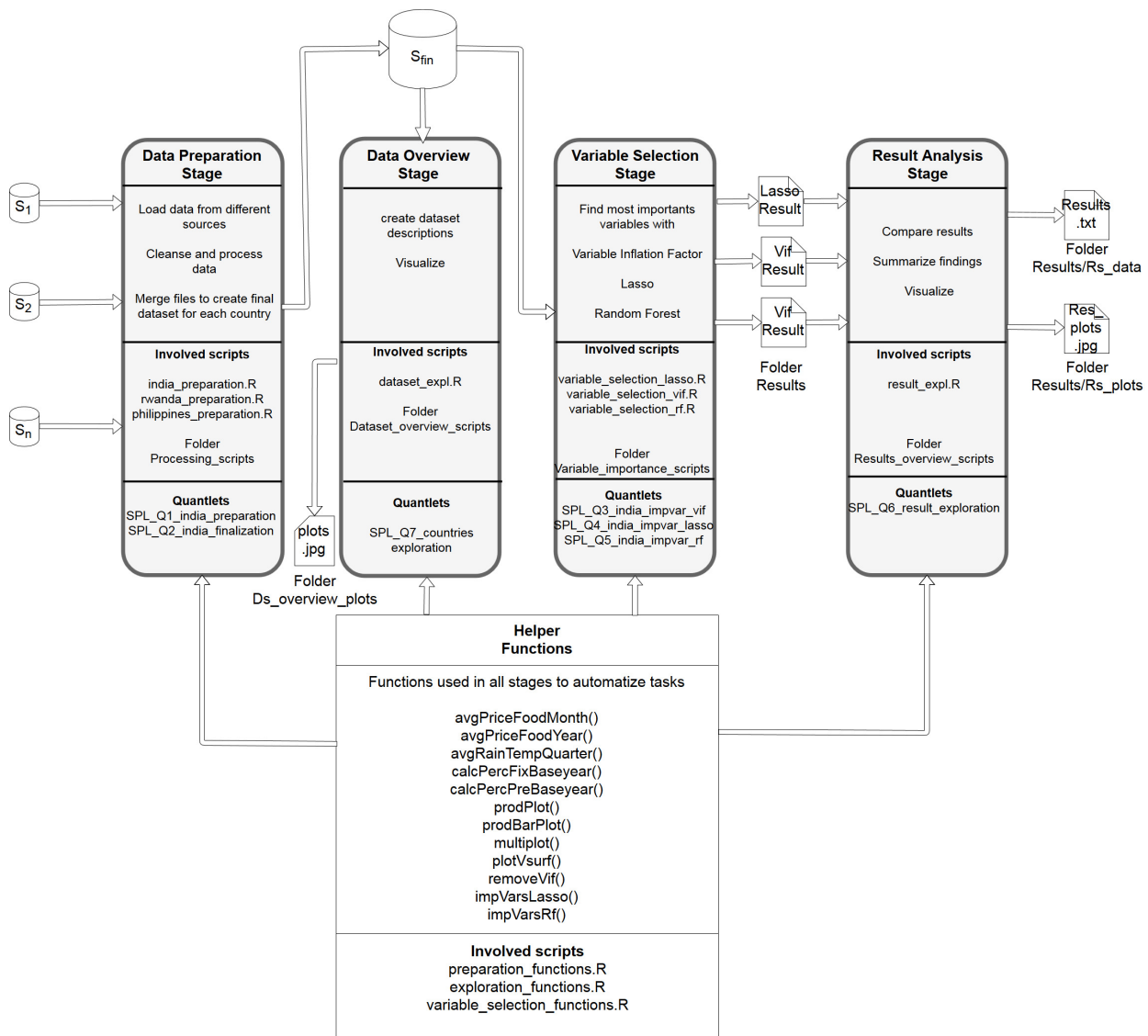


Figure 1. Workflow Structure.

In order to build a clean and valid dataset suitable for further analysis, first of all various different pieces of information needed to be taken from a range of different sources and blended together. This is done in the data preparation stage. There is a preparation script for each country we analyzed in the folder `Processing_scripts` in our github repository (https://github.com/jaidikam/sps_ws1718). The scripts make use of utility functions `avgPriceFoodMonth()`, `avgPriceFoodYear()` and `avgRainTempQuarter()` defined in the script `preparation_functions` in folder `Helper_functions`. Quantlets, `SPL_Q1_india_preparation` and `SPL_Q2_india_finalization` serve as a runnable more compact example of all the steps executed in this stage for the country india, independent from our working scripts. The subsequent data overview stage includes all the steps taken in the explorative analyza-tion of the datasets created in the previous stage. The script `dataset_expl.R` in the folder `Dataset_overview_scripts` holds all the code required to generate graphs such as the devel-opment of prices over the years and production rates. The results are stored in .jpeg file format in folder `Dataset_overview_scripts/Plots`. Utility functions used are `calcPercFixBaseyear()`, `calcPercPreBaseyear()`, `prodPlot()` and `multiplot()` in the script `exploration_functions` in folder `Helper_functions`. The Quantlet `SPL_Q7_countries_exploration` shows an example of how the explorative graphs were created. The next step is to find out the set of important variables for

each dataset. This is achieved in the variable selection stage. For each technique, there is a script in the folder Variable_importance_scripts. Scripts make use of the functions removeVif(), impVarsLasso(), impVarsRf() defined in the script variable_selection_functions.R in the folder Helper_functions. The resulting Files including the important variables and additional information are stored in folder Results. The Quantlets SPL_Q3_india_impvar_vif, SPL_Q4_india_impvar_lasso, SPL_Q5_india_impvar_rf show a functionally independent example of how each technique was used in our process. The final result analysis stage contains all the steps related to display and analyzation of the result files produced in the previous step. Code is included in the working script result_expl.R in folder Results_overview_scripts. For demonstrating purposed there is an independent Quantlet SPL_Q6_result_exploration.

2.2. Dataset Descriptions

Dataset description: The following tables will showcase the variables used in our datasets. We have three country specific datasets (India, Rwanda and the Philippines) containing information on supply and demand related factors influencing food prices. For the sake of clarity we have split the variables in tables below showing which data is shared by the individual sets and which data is country specific. Information on the variable names, type, range, factor group as well as a description is included.

The initial table show cases information of the data shared by all three countries: India / Philippines / Rwanda.

VariableName	DataType	Range(India)	Range(Philippines)	Range(Rwanda)	FactorGroup	Description
Year	int	2001-2015	1998-2015	1991-2015	x	Starting and end year of data collection
prod_name	char	x	x	x	x	name of the selected product
prod_price	num	8.392 - 37.02	24.70-452.10	37.9-1095.9	supply/production	?
tas_q1	num	19.86-21.64	22.826-52.1206	18.54-21.13	supply/climatic	average temperature in celsius for quartile 1
tas_q2	num	28.59-30.25	11.5366-21.1364	19.03-20.90	supply/climatic	average temperature in celsius for quartile 2
tas_q3	num	26.67-27.36	15.3575-24.5690	19.00-21.26	supply/climatic	average temperature in celsius for quartile 3
tas_q4	num	21.07-22.34	15.9073-48.1226	19.17-21.42	supply/climatic	average temperature in celsius for quartile 4
pr_q1	num	7.039-23.189	13.4319-61.9415	72.86-198.95	supply/climatic	average rain fall in mm for quartile 1
pr_q2	num	54.93-109.00	16.2071-62.5248	45.45-126.96	supply/climatic	average rain fall in mm for quartile 2
pr_q3	num	161.3-246.1	19.1587-37.5307	19.39-135.49	supply/climatic	average rain fall in mm for quartile 3
pr_q4	num	21.64-49.27	77.865-49.4250	85.2-186.9	supply/climatic	average rain fall in mm for quartile 4
avg_p_barrel	num	23.12-109.45	12.28-109.45	12.28-109.45	supply/production	annual oil price (U.S. dollars per barrel)
population	num	1.071e+09-1.309e+09	74694-101716	5928-11630	demand/demographic	annual population level
prod_amount	int	41555-362333	4106698-28376518	2300-3547200	supply/production	produced amount in (1k tons)
gni_pc	num	778.4-1737.8	1784-3163	203.8-696.8	demand/demographic	per capita income(GNI per capita in constant 2010 US\$)
cp_inflation	num	3.685-11.992	1.434-9.235	-2.406-56.000	supply/macroeconomic	inflation consumer prices (annual %)
agri_gdp	num	2.092e+11-3.281e+11	1.636e+10-2.670e+10	4.256e+08-2.098e+09	supply/macroeconomic	GDP in agriculture value added (constant 2010 US\$)
daily_caloric_supply	num	2256-2459	2292-2595	1723-2270	demand/demographic	per capita calorie intake (kcal)
imp_cer	int	25911-992977	27159-117853	7625-142335	supply/macroeconomic	annual imports in thousands of dollars

Table 1. My caption

The following table will showcase information of the data shared by India and Rwanda

VariableName	DataType	Range(India)	Range(Rwanda)	FactorGroup	Description
imp_veg	num	1140825-6317271	1954-23112	supply/macroeconomic	annual vegetable imports in thousands of dollars
exp_cer	num	3.210e+07-2.248e+09	0.51-54746.39	demand/macroeconomic	annual cereal exports in thousands of dollars
exp_veg	num	851222-3559119	43.56-8891.15	demand/macroeconomic	exported vegetables in US\$

Table 2. My caption

The next table will showcase information of the data shared by the Philippines and Rwanda.

VariableName	DataType	Range(Rwanda)	Range(Philippines)	FactorGroup	Description
GDP	num	7.536e+08-8.261e+09	7.221e+10-2.928e+11	supply/macroeconomic	
exchange_rate	num	125.2-721.0	39.09-56.04	supply/macroeconomic	exchange rate in relation to 2010 US\$
population_unit	num	1000	1000	demand/demographic	Unit for population size (1000 Persons)
flag	chr	x	x	x	Country name
flag.description	chr	x	x	x	Description of source for country name

Table 3. My caption

The following table will showcase information of the data specific to India

VariableName	DataType	Range(India)	FactorGroup	Description
country	Factor	x	x	Country Name
imp_sug	num	26034-1419642	supply/macroeconomic	annual imports in thousands of dollars
exp_sug	num	91273-2247911	demand/macroeconomic	annual exports in thousands of dollars

Table 4. My caption

The final table will showcase information specific to the Philippines

VariableName	DataType	Range(Philippines)	FactorGroup	Description
exp_agri	num	36.30-142.32	demand/macroeconomic	

Table 5. My caption

2.3. Function Descriptions

- avgPriceFoodMonth
 - description:
 - * calculates the average price per month for each product for a given dataframe.
 - * Adds column avg_price_prod_month.
 - input parameters:
 - * ds: food price data as a dataframe
 - * cm_name: the name of the column holding the product name as a String
 - * mp_price : the name of the column holding the product price as a String
 - * mp_year: the name of the column holding the year as a String
 - * mp_month: the name of the column holding the month as a String
 - output parameters
 - * ds: the input dataframe including an additional column avg_price_prod_month
 - typical issues:
 - * rong column name specified: Error Column name not passed as String Error
- avgPriceFoodYear
 - description:
 - * calculates the average price per year for each product for a given dataframe
 - * Adds column avg_price_prod_year.
 - input parameters:
 - * Ds: food price data as a dataframe
 - * cm_name: the name of the column holding the product name as a String
 - * mp_year: the name of the column holding the year as a String
 - * avg_price_prod_month: the name of the column holding the average food price per month as a String
 - output parameters
 - * ds: the input dataframe including an additional column avg_price_prod_year
 - typical issues:
 - * Wrong column name specified: Error
 - * Column name not passed as String: Error
- avgRainTempQuarter
 - description:

- * Calculates the average temperature and average amount of rain per quarter year for a given dataframe
- input parameters:
 - * ds: rain and temperature data for every month in each year as a dataframe
 - * month: the name of the column holding the month as a String
 - * mp_year: the name of the column holding the year as a String
 - * pr: the name of the column holding the amount of rain per month as a string
 - * tas: the name of the column holding the average temperature per month as a string
- output parameters
 - * ds: the input dataframe including an additional
 - * columns: tas_q1, tas_q2, tas_q3, tas_q4, pr_q1, pr_q2, pr_q3, pr_q4
- typical issues:
 - * Wrong column name specified: Error
 - * Column name not passed as String: Error
- plotVsurf
 - description:
 - * Plots VSURF objects for thresholding and interpretation step
 - input parameters:
 - * iVsurfOb: A VSURF object
 - * iStep: the step for which results are to be plotted as a string
 - * iCountry: the country for which results are to be plotted as a string
 - typical issues:
 - * Other object than VSRUF object passed to function
 - * Other string than "thresh" or "interp" passed as value for iStep
- removeVif
 - description:
 - * Removes multicorrelated numeric variables from a given dataframe based on variance inflation factor
 - input parameters:
 - * explan_vars: the numeric variables as a dataframe
 - * cutoffval: the maximum allowed vif for remaining variables as a number
 - output parameters
 - * tempresults: the remaining variable names and their corresponding vif as a dataframe
 - typical issues:
 - * non numeric variables in input dataframe
- impVarsLasso
 - description:
 - * Identifies important variables for a given dataframe based on the lasso method
 - input parameters:
 - * ds: the variables as a dataframe
 - * targ: the name of the target variable column in ds as a String
 - output parameters
 - * resultset: the lasso model, fitted model and the label for display in a graph as a vectorlist

- typical issues:
 - * Column name not passed as String: Error
- impVarsRf
 - description:
 - * Identifies important variables for a given dataframe based on MSE error minimization in OOB samples of random forest
 - input parameters:
 - * ds: the variables as a dataframe
 - * targ: the name of the target variable column in ds as a String
 - output parameters
 - * resultset: names of important variables and mean OOB rate as a vectorlist
 - typical issues:
 - * Column name not passed as String: Error

3. Implementation

What is this section about? Only about the code we created! We shall explain every step in our workflow extensively. Code examples shall be taken from the Quantlets. One country as an example should be enough. The helper functions implementation shall be part of the Quantlets, as well as the way we use them.

This section gives detailed information about the actual implementation of every step featured in **Section 2** including r code and theoretical background of every variable technique that was applied.

3.1. Data Preparation

In order to get comparable results, it was necessary to work with datasets containing comparable features. After specifying which variables were supposed to be included in the final datasets, those information needed to be put together from a range of heterogeneous sources. The source dataset containing food prices, for an instance, included the monthly price of each crop of interest for each marketplace in a certain country. That made it necessary to calculate the average price per crop per month for the whole country which is done by the following function:

```

1  #Define the function for food price per month across all markets
2  avgPriceFoodMonth = function(ds, cm_name, mp_price, mp_year, mp_month) {
3    dsfoods = unique(ds[[cm_name]])
4    ds$avg_price_prod_month = 0
5    for (k in min(ds[[mp_year]]):max(ds[[mp_year]])) {
6      print(paste('year is ', k))
7      for (j in 1:NROW(dsfoods)) {
8        a <- ds[ds[[mp_year]] == k & ds[[cm_name]] == dsfoods[j], ]
9        print(paste('food is ', dsfoods[j]))
10       for (i in 1:12) {
11         b <- a[a[[mp_month]] == i, ]
12         if (NROW(ds[ds[[mp_year]] == k & ds[[cm_name]] == dsfoods[j] & ds[[
13           mp_month]] == i, ]$avg_price_prod_month) > 0) {
14           ds[ds[[mp_year]] == k & ds[[cm_name]] == dsfoods[j] & ds[[mp_
15             month]] == i, ]$avg_price_prod_month = sum(b[[mp_price]]) /
16             NROW(b)
17           print(paste('month is ', i))
18         }
19       }
20     }
21   }
22 }
```



```

18 |     }
19 |     return(ds)
20 | }
21 | india = avgPriceFoodMonth(india,"cm_name","price","year","month")

```

Listing 1. SPL_Q1_india_preparation

Function avgPriceFoodMonth takes the dataframe containing price information and allows the user to explicitly pass the names of the columns holding relevant information as strings. First, possible duplicates in the input dataset are removed and the name of every crop is stored in a variable dsfoods. Then a column avg_price_prod_month to hold the result is added to the dataset with value 0. The first most outer loop iterates through the year. The second loop goes through every crop. The most inner loop goes through every month for every possible year/crop pair. If rows exist for a certain month/year/crop combination, the sum of prices for that combination are divided by the number of occurrences resulting in the desired average. The if – condition is necessary to avoid the possibility to divide by 0. Finally, the dataframe plus the newly created column is returned.

Later on in the process we decided to look at crop prices per year, so we created another similar function to calculate the average price per food per year:

```

1 | avgPriceFoodYear = function(ds,cm_name,mp_year,avg_price_prod_month) {
2 |   dsfoods = unique(ds[[cm_name]])
3 |   ds$avg_price_prod_year = 0
4 |   for (x in min(ds[[mp_year]]):max(ds[[mp_year]])) {
5 |     print(paste('year is ',x))
6 |     for(y in 1:NROW(dsfoods)){
7 |       print(paste('food is ',dsfoods[y]))
8 |       if(NROW(ds[ds[[mp_year]] == x & ds[[cm_name]] == dsfoods[y],]$avg_
9 |         price_prod_year) > 0){
10 |         ds[ds[[mp_year]] == x & ds[[cm_name]] == dsfoods[y],]$avg_price_
11 |           prod_year = sum(ds[ds[[mp_year]] == x & ds[[cm_name]] == dsfoods
12 |             [y],][[avg_price_prod_month]]) / NROW(ds[ds[[mp_year]] == x & ds
13 |               [[cm_name]] == dsfoods[y],][[avg_price_prod_month]])
14 |       }
15 |     }
16 |   }
17 |   return(ds)
18 | }
19 | india = avgPriceFoodYear(india,"cm_name","year","avg_price_prod_month")

```

Listing 2. SPL_Q2_india_preparation

The function takes a dataframe including crop prices on a monthly level and allows the user to specify the columns holding crop name, year and price per month. We pass the dataset that was created by the previous function and the column for the average price per year is added. The outer loop iterates through the years, the inner loop goes through the list of unique crop names. For every existing year/crop combination the average is price calculated.

Another dataset that required processing was the climate set, containing the amount of rain and the temperature for every month in a year, each in a column of its own. Instead of one row per month, we needed the average amount of rain / temperature for each quarter of the year to be a row identified by the year. A function was written to achieve that:

```

1 | avgRainTempQuarter = function(ds,month,mp_year,pr,tas) {
2 |   if (is.na(ds[[pr]]) || is.na(ds[[tas]])) {
3 |     message(paste("No missing values allowed!"))

```

```

4   } else {
5     ds$tas_q1 = 0
6     ds$tas_q2 = 0
7     ds$tas_q3 = 0
8     ds$tas_q4 = 0
9     ds$pr_q1  = 0
10    ds$pr_q2  = 0
11    ds$pr_q3  = 0
12    ds$pr_q4  = 0
13
14    for(z in min(ds[[mp_year]]):max(ds[[mp_year]])) {
15
16      ds[ds[[mp_year]] == z ,]$pr_q1 = sum(ds[ds[[mp_year]] == z & ds[[
17        month]] %in% c("1","2","3"),)[[pr]]/3
18      ds[ds[[mp_year]] == z ,]$pr_q2 = sum(ds[ds[[mp_year]] == z & ds[[
19        month]] %in% c("4","5","6"),)[[pr]]/3
20      ds[ds[[mp_year]] == z ,]$pr_q3 = sum(ds[ds[[mp_year]] == z & ds[[
21        month]] %in% c("7","8","9"),)[[pr]]/3
22      ds[ds[[mp_year]] == z ,]$pr_q4 = sum(ds[ds[[mp_year]] == z & ds[[
23        month]] %in% c("10","11","12"),)[[pr]]/3
24    }
25    return(ds)
26  }
27 }
28 raintemp = avgRainTempQuarter(raintemp,"month","year","pr","tas")
29 #load the india dataset and the remaining variables
30 india_wip = readRDS("../Qfolder1\\Q1_india_wip.rds")
31 rest      = readRDS("../Qfolder2\\Q2_india_rest.rds")
32 #merge india with wheater data
33 india_wip = merge(india_wip,unique(raintemp[c("tas_q1","tas_q2","tas_q3","
34   tas_q4","pr_q1","pr_q2","pr_q3","pr_q4","year")]),by=c("year"))
35 #join the datasets
36 india_fin = merge(india_wip, rest, by=c("prod_price")) #prod_price is
   unique, therefore it can be used as a key for the merge

```

Listing 3. SPL_Q2_india_finalization

The function takes a the climate dataset and allows the user to pass strings to identify the columns holding the month, year, rain and temperature respectively. Given that the dataset has no missing values in the rain and the temperature column, additional columns are added for every quarter of the year and the filled with the sum of each quarter divided by the number of month a quarter year has. The dataset with the added columns is then returned. Afterwards, the datasets are merged to create the final dataset used in the next step.

3.2. Data overview

In this section we will explain the process and the code which we used to produce four of our graphs, which provide us with some important statistical comparison. All the produced graphs are mmentioned in **Section 4.1** The resulted graphs were:

- Product Price Change

- Population Change
- Price Index
- Consumer Price Development

We have collected our databases from different resources, and for the values to be comparable between countries an world wide, we had to develop some functions to calculate the percentage of the change

```

1  #calculating the percentage of the change in a column's value on a fixed
   base year
2  calcPercFixBaseyear = function(ds, areacol, areaname, yearcol, baseyear,
   valuecol, perccol){
3    base = ds[ds[[yearcol]] == baseyear & ds[[areacol]] %in% areaname, ][[
   valuecol]]
4    if(is.null(ds[[perccol]])){
5      ds[[perccol]] = 0
6    }
7    #
8    for(i in baseyear: max(ds[[yearcol]])){
9      later = ds[ds[[yearcol]]==i & ds[[areacol]] %in% areaname, ][[valuecol]]
10     sub = later - base
11     ds[ds[[yearcol]] == i & ds[[areacol]] %in% areaname, ][[perccol]] = (sub
       / later) * 100
12   }
13   return(ds)
14 }

```

Listing 4. [SPL_Q7_countries_exploration.R](#)

Function calcPercFixBaseyear calculate the precentage change in a value basd on a specifed base year during a timeframe, where the prcentage will always represent the difference between the cvalue column and the value stored in the base year, beside the base year the function takes also targeted dataframe, area column, area name, year column, value column and the name of the produced new column. We firstly store the given year in the base variable then we check if the goal new column have been already identified or not. After that we go through a loop in value column depending on the targeted year column and store the resulted prcentage change in the new column and at the the end returning the new dataframe comntain the newly calculated column for the specified area during the the specified time duration.

```

1  # calculate the prcentage of changes in a colname value in a predefined
   year, where the base is for every change is the value from the previous
   year.
2  calcPercPreBaseyear = function(ds, areacol, areaname, yearcol, valuecol){
3    for(i in unique(ds[[yearcol]])){
4      base = ds[ds[[yearcol]] == i & ds[[areacol]] %in% areaname, ][[valuecol
       ]]
5      later = ds[ds[[yearcol]] == i+1 & ds[[areacol]] %in% areaname, ][[
       valuecol]]
6      #if(length(later) == 0L) break
7      if(i == 2015L) break
8      sub = later - base
9      if(length(sub) == 0L) next
10     ds[ds[[yearcol]] == i+1 & ds[[areacol]] %in% areaname , paste(valuecol,
       "Percent", sep = "_")] = (sub / later) * 100
11   }
12   ds[ds[[yearcol]] == min(ds[[yearcol]]) & ds[[areacol]] %in% areaname,
       paste(valuecol, "Percent", sep = "_")] = 0
13   return(ds)

```

Listing 5. `SPL_Q7_countries_exploration.R`

Function `calcPercFixBaseyear` also calculate the precentage change, but this time the base year is changing and is not fixed, it takes the value of the pervious year at each calculation for a value in a specific year

```

1 # To plot production data with specific area and itmes
2 prodPlot = function(ds, area, items){
3   p = ggplot(data=ds[ds$Area == area & ds$Item %in% items,], aes(x=Year, y=
4     Percentage, colour=Item)) +
5     geom_line() +
6     geom_point() +
7     ylim(-30, 75) +
8     ggtitle(label=area) +
9     ylab(label="Percentage Production Change") +
10    xlab("Year")
11  return(p)
12 }
```

Listing 6. `SPL_Q7_countries_exploration.R`

We use `prodPlot` function to prудuce the six graphs in **Figure 5. Consumer Price development**, where it show us the change of the products prices in the studied countries in compare the world wide change. The function takes the dataframe along side the area neme and targeted itmes, then it prудuce the plot for the defined parameters.

3.3. Variable Selection

To find out which variables have the biggest impact on food prices, we apply three techniques which were implemented as follows:

3.3.1. VIF based method

Before fitting a linear model, first of all it is necessary to check for correlation among explanatory variables, since high correlation means that variables are not independent from one another. As a way to measure the degree of dependence of variable pairs, Pearson's correlation coefficient was used. Intuitively, explanatory variables should only be correlated to the target variable, not to each other. Otherwise the model will have high standard deviation and therefore high variance and possibly low predictive power. Also significant variables may appear non-significant. However, pair-wise analysis for correlation may not reveal multicollinearity, i.e. a situation where a predictor can be explained through the other predictors. Multicollinearity is excessive correlation among several explanatory variables. A commonly used way to detect that is a metric called variance inflation factor. It helps to quantify the amount of variance that each predictor adds to the model. To calculate the VIF for a variable b_k : Given a linear model :

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \dots + \beta_{p-1} x_{i,p-1} + \epsilon_i$$

By regressing b_k with only one of the other predictors at a time we calculate the variances for b_k . We keep the smallest variance, obtained by:

$$Var(b_k)_{min} = \frac{\delta^2}{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}$$

In order to find out how much b_k increases the overall variance of the model we get the ratio of $Var(b_k)_{min}$ and $Var(b_k)$, the variance of b_k when all remaining variables are regressed on b_k at the same time:

$$\frac{Var(b_k)}{Var(b_k)_{min}} = \frac{\left(\frac{\sum_{i=1}^n \delta^2}{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2} \times \frac{1}{1-R_k^2} \right)}{\left(\frac{\sum_{i=1}^n \delta^2}{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2} \right)} = \frac{1}{1-R_k^2}$$

R_k^2 is the R^2 value when the k^{th} is regressed on all remaining explanatory variables

It follows :

$$VIF_k = \frac{1}{1-R_k^2}$$

The correlation- and VIF based method was implemented as follows: After choosing the explanatory variables and standardizing them we acquire the correlation coefficient for each pair:

```

1 #initial variable selection and normalization
2 colselection_in = c("prod_price", "pr_q1", "pr_q2", "pr_q3", "pr_q4", "tas_q1", "
   tas_q2", "tas_q3", "tas_q4",
3                       "prod_amount", "daily_caloric_supply", "exp_sug", "exp_veg
   ", "exp_cer", "imp_sug", "imp_veg", "imp_cer",
4                       "agri_gdp", "gni_pc", "cp_inflation", "avg_p_barrel", "
   population")
5 target_in = c("prod_price")
6 normalized_in = as.data.frame(scale(india[colselection_in]))
7 feats_in = normalized_in[, !(colnames(normalized_in) %in% target_in)]
8 #Variable selection and modeling
9 #Obtaining pair-wise correlations
10 insign_in = cor(feats_in, method = "pearson", use = "complete.obs")

```

Listing 7. SPL_Q3_india_impvar_vif

After we calculated the correlations for each pair, every variable being part of a pair with a correlation coefficient above 0.70 is removed and a linear model is built.

```

1 # Discovering highly correlated explanatory variables
2 hicorvars_in = findCorrelation(cor(feats_in), cutoff = 0.70)
3 expvarsnohc_in = paste(colnames(feats_in)[-hicorvars_in], collapse = "+")
4 formulanohc_in = paste(target_in, "~", expvarsnohc_in, collapse = "+")
5 mod_varnohc_in = lm(formulanohc_in, data = normalized_in)
6 #For comparison we also apply the VIF-based method to tackle
   multicollinearity:
7 #function for VIF based stepwise removal of multicorrelated variables
8 removeVif = function(explan_vars, cutoffval=10){
9   tempresults = as.data.frame(matrix(ncol = 2, nrow = 0))
10  colnames(tempresults) = c("variable", "vif")
11  #initially calculate VIF for each explanatory variable
12  for (i in 1:NROW(colnames(explan_vars))){
13    temptarget = colnames(explan_vars)[i]
14    tempexpvars = paste(colnames(explan_vars[!(colnames(explan_vars) %in%
       temptarget)]), collapse = "+")
15    tempformula = paste(temptarget, "~", tempexpvars, collapse = " ")
16    tempresults[i,1] = temptarget
17    tempresults[i,2] = VIF(lm(tempformula, data = explan_vars))
18  }
19  print(tempresults[order(tempresults$vif),])
20  #remove variable with highest VIF, calculate new VIF for remaining
   variables until all VIF are below cutoff value
21  while(max(tempresults$vif) >= cutoffval){

```

```

22     tempresults = tempresults[!tempresults$vif == max(tempresults$vif),]
23     tempremvars = tempresults$variable
24     for(j in 1: NROW(tempremvars)){
25         temptarget = tempremvars[j]
26         tempexpvars = paste(tempremvars[!tempremvars %in% temptarget],
27                             collapse = "+")
28         tempformula = paste(temptarget, "~", tempexpvars, collapse = " ")
29         tempresults[j,1] = temptarget
30         tempresults[j,2] = VIF(lm( tempformula, data = explan_vars))
31     }
32     print("Remaining variables:")
33     print(tempresults[order(tempresults$vif),])
34     cat("\n")
35 }
36 return(tempresults)
37 }
38 # for highly correlated variables
39 varslvifhc_in = removeVif(feats_in[, hcorvars_in], 8)
40 # for lower correlated variables
41 varslvifnohc_in = removeVif(feats_in[, -hcorvars_in], 8)
42 #Model without multicollinearity
43 expvars_lovif_in = paste(paste(varslovhc_in$variable, collapse = "+"), "+",
44                          paste(varslvifnohc_in$variable, collapse = "+"), collapse = "+")
45 formula_lovif_in = paste(target_in, "~", expvars_lovif_in, collapse = "+")
46 mod_lovif_in = lm(formula_lovif_in, data = normalized_in)

```

Listing 8. SPL_Q3_india_impvar_vif

The function takes a dataframe containing explanatory variables and a cutoff value that specifies the highest admissible VIF until the function stops removing variables. The r package “fmsb” is required. The function’s algorithm in pseudocode:
insert pseudocode from presentation

Finally, the function returns a dataframe containing the remaining variables and their VIF values

The function is applied to all the highly pair-wise correlated variables identified in line xx as well as to the remaining variables in a separate call. The linear model is then constructed from the remaining predictors and saved as a file.

3.3.2. Lasso based method

Lasso (Least Absolute Shrinkage and Selection Operator) is a regression based method for regularization and feature selection .A constraint is put on the absolute sum of model parameters and a penalty is applied to the regression coefficients. The strength of the penalty is defined by the parameter lambda. The bigger lambda gets, the sooner model parameters exceed the threshold value and get dropped from the model. Intuitively, a smaller lambda causes more parameters to stay in the model. Lambda = 0 means there is no penalty at all and the regression becomes an ordinary least square regression Formulation by P. Bühlmann and S Van de Geer(Statistics for High-Dimensional Data: Methods, Theory and Applications. Springer 2011):

Find a solution to the optimization problem minimize

$$\left(\frac{\|Y - X\beta\|_2^2}{n} \right) \text{ subject to } \sum_{j=1}^k \|\beta\|_1 < t$$

t is the upper bound for the sum of the coefficients. Which is equivalent to the parameter estimation

$$\arg \min_{\beta} \left(\frac{\|Y - X\beta\|_2^2}{n} + \lambda \|\beta\|_1 \right)$$

Where

$$\frac{\|Y - X\beta\|_2^2}{n} = \sum_{i=0}^n (Y_i - (X\beta)_i)^2, \|\beta\|_1 = \sum_{j=1}^k |\beta_j| \text{ and } \lambda \geq 0$$

If λ gets 0, t becomes infinite and vice versa. In other words, we find the set of coefficients from the model having the least mean squared error for a sequence of descending lambdas

For application of the lasso method, a function was written: Our implementation of the lasso technique requires the R packages "glmnet" and "plyr"

```

1 #function for LASSO method
2 impVarsLasso = function(ds,targ) {
3   #1. initial variable selection and normalization
4   val = ds[[targ]]
5   x = model.matrix(ds[[targ]]~.-1, ds[!colnames(ds) %in% targ])
6   #2. Applying the Lasso technique
7   lasso = glmnet(x = x, y = val, standardize = TRUE, alpha = 1)
8   fit = cv.glmnet(x = x, y = val, standardize = TRUE, type.measure = "mse"
9     , alpha=1, nfolds=3)
10  #3. Results
11  #with lambda.min
12  lambda_min = which(fit$lambda == fit$lambda.min)
13  #selecting coefficients of variables at lambda where mse is minimal
14  tempmincoefs = as.data.frame(fit$glmnet.fit$beta[, which(fit$
15    lambda == fit$lambda.min)])
16  mincoefs = data.frame(matrix(ncol = 2, nrow = (NROW(
17    tempmincoefs))))
18  mincoefs$variables = as.vector(as.character(labels(tempmincoefs)
19    [[1]]))
20  mincoefs$coefs_minlambda = as.vector(tempmincoefs[[1]])
21  mincoefs$X1 = NULL
22  mincoefs$X2 = NULL
23  #get names in the decreasing order they appear in when lambda is minimal
24  names = names(coef(lasso)[,ncol(coef(lasso))][order(coef(lasso)
25    [,ncol(coef(lasso))],decreasing=TRUE)])
26  names = names[!names %in% c("(Intercept)")]
27  names = as.data.frame(names)
28  colnames(names) = "variables"
29  #add coefficient to names
30  disp_colors = join(names,mincoefs, by = "variables" )
31  disp_colors = disp_colors[!disp_colors$variables %in% c("(Intercept)"),]
32  #set colors for variables when displayed in a graph
33  disp_colors$colors = 0
34  if(NROW(disp_colors[disp_colors$coefs_minlambda >0,])>0) {
35    disp_colors[disp_colors$coefs_minlambda >0,]$colors = c("green")
36  }
37  if(NROW(disp_colors[disp_colors$coefs_minlambda <0,])>0) {
38    disp_colors[disp_colors$coefs_minlambda <0,]$colors = c("red")
39  }
40  #create a list to store the result
41  resultset = vector("list",3)
42  resultset[[1]] = lasso
43  resultset[[2]] = fit
44  resultset[[3]] = disp_colors
45  return(resultset)

```



```

41 | }
42 | #Get most important variables with Lasso function
43 | india_lasso_result = impVarsLasso(india, "prod_price")

```

Listing 9. SPL_Q4_india_impvar_lasso

The function takes a dataframe as input and allows the user to specify the column holding the target variable by passing the name as a string. In line 21 a glmnet model is built, applying the aforementioned lasso method. Then a fitted model is built in the same way. The mean square error of each model for every lambda is calculated using three-fold cross validation for better generalization of the results, i.e. to avoid overfitting. In the following step the coefficients and variable names from the model with the lowest MSE are stored in a dataframe called min-coefs. In order to offer additional information about the correct color when the results are plotted, the variables names are ordered in the decreasing order they appear in when lambda is minimal, i.e. when their curves would cut the right margin of the plot (line 38). If coefficients in mincoeffs are greater 0, they will be displayed in green in the graph, lower 0 in red, all other will be grey. Finally, the model, fitted model and the display information are stored in a list which is the returned

3.3.3. Random Forest based method

The third variable selection technique we applied is based on random forests. Random forests are a non-parametrical statistical method used for regression and classification problems. RF are an ensemble of decision trees built from samples of the whole data. The importance of a variable X^j is calculated as follows: Each tree t runs a prediction on the Out of bag sample OOB_t (the portion of data not included in the sample to create t). The mean square error MSE of this prediction is denoted $errOOB_t$. Now the values for X^j in OOB_t are permuted randomly to get a perturbed sample OOB_t^j . Then the error $errOOB_t^j$ of t on the perturbed sample is calculated. Now the variable importance for X^j is computed as follows:

$$VI(X^j) = \frac{1}{ntree} \sum_t \left(err\widetilde{OOB}_t^j - errOOB_t \right)$$

The function used to obtain the most important variables using RF requires the r package “VSURF” The implementation utilizes a two-step approach to find the most important variables:

- **Threshold step** After computing the VI, each variable is ranked by VI in descending order. Variables with small importance are eliminated. The threshold is derived from the standard deviations of VI. Important variables have a higher variability so only variables with a averaged VI exceeding the threshold are kept in the model
- **Variable Selection Interpretation step:** A nested collection of RF models is built for the first k to m variables. The variables leading to the minimal OOB error are selected, resulting in m' variables being retained. **Prediction step:** The variables from the interpretation step are ordered and sequentially added to an RF model. Variables which lead to an error decrease above a certain threshold are kept.

Note: Since prediction step focuses on predictive power of the model, some variables that in fact have an influence on the target are removed since others are stronger predictors. Our focus lies on generally identifying variables having a significant influence on the target, therefore the prediction step is ignored in our approach.


```

1 #function for finding most important variables based on random forest
2 impVarsRf = function(ds,targ){
3   result_rf = VSURF(ds[[targ]] ~ ., data = ds[!colnames(ds) %in% targ],
4     ntree = 2000,
5     nfor.thres = 50, nmin = 1, nfor.interp = 25, nsd = 1,
6     nfor.pred = 25, nmj = 1, parallel = FALSE, ncores =
7       detectCores() - 1,
8     clusterType = "PSOCK")
9   #create a list to store the result
10  resultset = vector("list",2)
11  resultset[[1]] = result_rf
12  resultset[[2]] = colnames(ds[!colnames(ds) %in% targ])
13  return(resultset)
14 }
15 #apply the function
16 india_v_imp_rf = impVarsRf(india,"prod_price")

```

Listing 10. SPL_Q5_india_impvar_rf

3.4. Result Exploration

Having obtained the most important variables, the results produced in the previous steps are loaded and processed to turn them into text files and images so that our findings are readable and accessible without further programming. The following excerpt gives an impression of how we achieved this

```

1 #loading results for random forest based variable selection
2 india_rf_result = readRDS("..\Qfolder5\Q5_india_rf.rds")
3 #function for plotting VSURF Objects
4 plotVsurf = function(iVsurfOb,iStep,iCountry){
5   header_prefix = "not specified"
6   if(iStep == "thres"){
7     header_prefix = "Thresholding step"
8   }
9   if(iStep == "interp"){
10    header_prefix = "Interpretation step"
11  }
12  plot(iVsurfOb,step = iStep, var.names = FALSE,
13    nvar.interp = length(iVsurfOb$vselect.thres), main = paste(header_
14    prefix,iCountry))
15 }
16 #threshold step
17 #save variables and plot to file
18 sink("..\Qfolder6\Q6_rf_thres_in.txt")
19 print(india_rf_result[[2]][india_rf_result[[1]]$vselect.thres])
20 sink()
21 jpeg("../Qfolder6//Q6_india_rf_thres.jpg", width = 1000, height = 700,
22   units = "px", pointsize = 20,
23   quality = 100)
24 plotVsurf(india_rf_result[[1]],"thres","India")
25 dev.off()
26 #interpretation step
27 sink("..\Qfolder6\Q6_rf_interp_in.txt")
28 print(india_rf_result[[2]][india_rf_result[[1]]$vselect.interp])
29 sink()
30 jpeg("../Qfolder6//Q6_india_rf_interp.jpg", width = 1000, height = 700,
31   units = "px", pointsize = 20,
32   quality = 100)

```

```

31 | plotVsurf(india_rf_result[[1]], "interp", "India")
32 | dev.off()

```

Listing 11. SPL_Q6_result_exploration

In this example the result objects from the random forest based variable selection are loaded. For the creation of the plot a function was written. The function takes a VSURF object and allows the user to specify the corresponding country name and the step for which the results shall be plotted (“thres“ or “interp”). The corresponding variable names are also extracted from the result object and stored in a .txt-file. The results for the other techniques are processed in a similar fashion.

4. Content related analysis

This section shows the output created in Section 3. The focus is on content-related exploration and interpretation of our findings and on comparison with the discoveries from related papers featured in **Section 1**

4.1. Datasets

The following section will visually highlight some of our findings with the help of graphs. For the sake of comparability, we used percentage changes rather than absolute numbers.



Figure 2. Products Price Change Based on The Cahnge from Previous Year

The first graph **Figure 2** showcases the annual percentage price change of our crop selection. Our selection is based on Harmonized System Codes (HS Code 2017). We focused on the category 07-015 – "Vegetables And Certain Roots And Tubers; Edible".

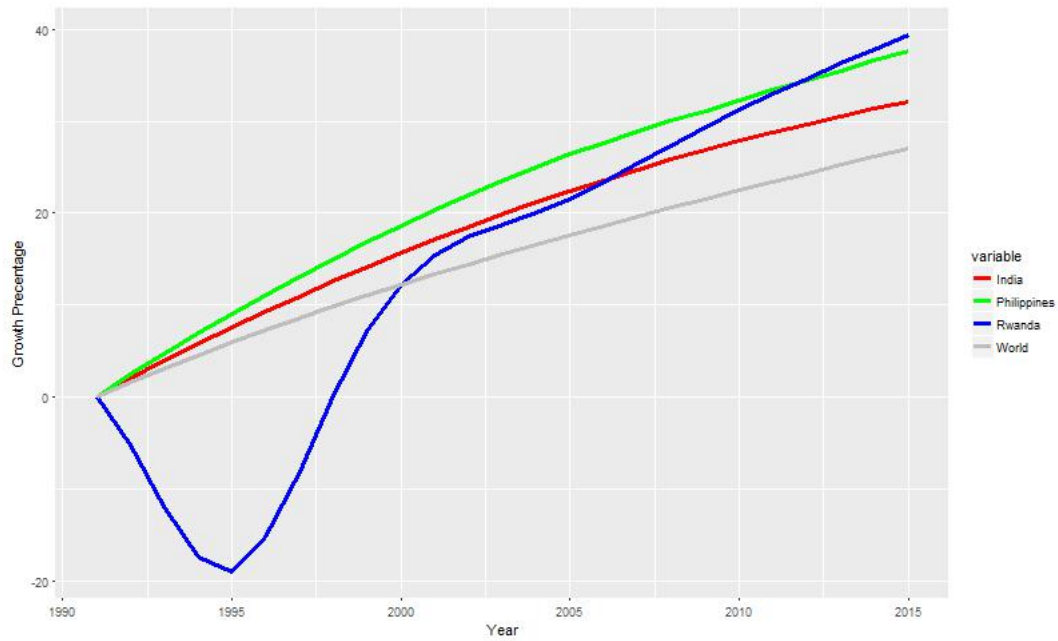


Figure 3. Population Change

The second graph **Figure 3** shows the population percentage population growth against the baseline world population percentage growth from 1990 to 2015. It is interesting to note that our country selection generally had a higher growth percentage than the global percentage trend with the notable outlier of Rwanda in the early 1990s. We assume that this sharp decline was related to the catastrophic civil war that ravaged Rwanda in 1994.

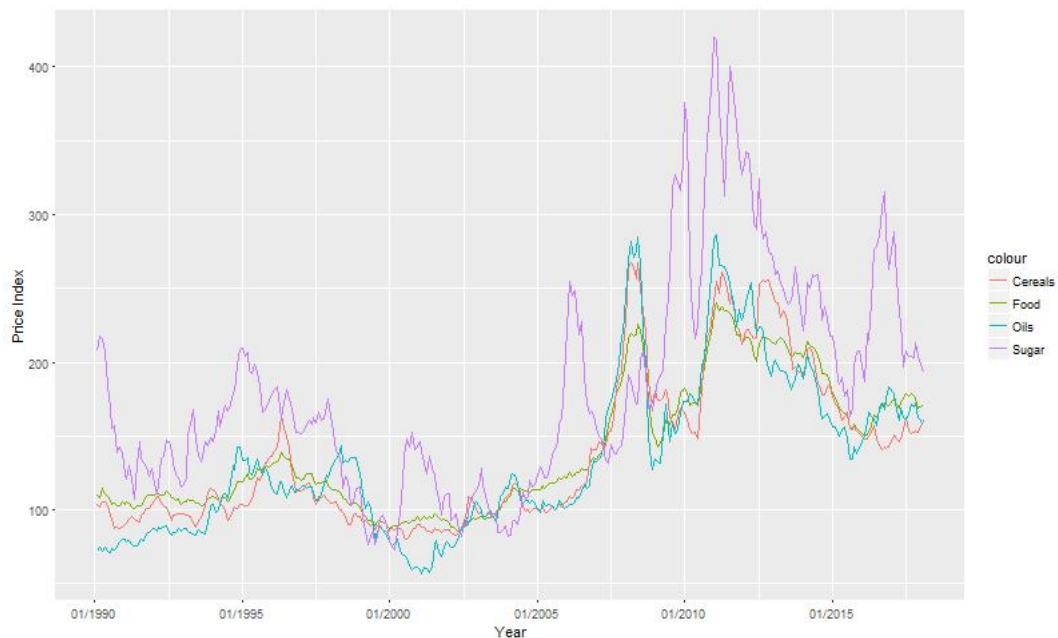


Figure 4. Global Average Food Category Price Change

The next graph **Figure 4** shows the price index percentage change of cereals, vegetable oils and sugar against the global food price index percentage change from 1990 to 2015. A general upward trend can be observed with short term price spikes. Around 2010 a general price index drop can be observed. Sugar price percentage change is higher than the more homogenous

trend development of cereals and vegetable oils which seem to be in line with the global food price index trend. These observations are in line with observations presented by the OECD.



Figure 5. Consumer Price Development

The final set of graphs **Figure 5** shows the consumer price development of our selected goods on a by country level against global consumer price development for the same goods from 1990 - 2015. Starting with India one can note that the general increasing global trend of consumer prices can be observed on the country level as well. Notable exceptions are the substantially higher consumer price development of potatoes as well as sharp price drops of rice and sugar in the early to mid 2000s. The Philippines are similar to India in the sense that the selected product consumer prices seem to follow the general global upward trend. The consumer price increase of Bananas, Coconuts and Rice are above the global trend. The notable exception is sugar which has quite volatile trend with multiple sharp consumer price drops. Rwanda is interesting as its consumer price trend is erratic in comparison to the global generally increasing trend. The global consumer food price for Bananas, Cassava and Maize increase at a constant rate whereas sweet potatoes seem to decrease and stagnate globally around 2005. The Rwandan consumer price for the same goods is in a sharp decline starting in the early 1990s, again our assumption is that this decline is related to the Rwandan civil war. A strong increase in consumer prices can be noted starting in the mid 1990s, spearheaded by Cassava and Maize consumer price development.

4.2. Results

4.2.1. Correlation / VIF based method

4.2.1.1. Correlation overview The first step of this approach was to get an overview of pairwise correlations for each country's dataset, measured by Pearson's correlation coefficient.

Variables which are part of a correlation pair with coefficient ≥ 0.70 :

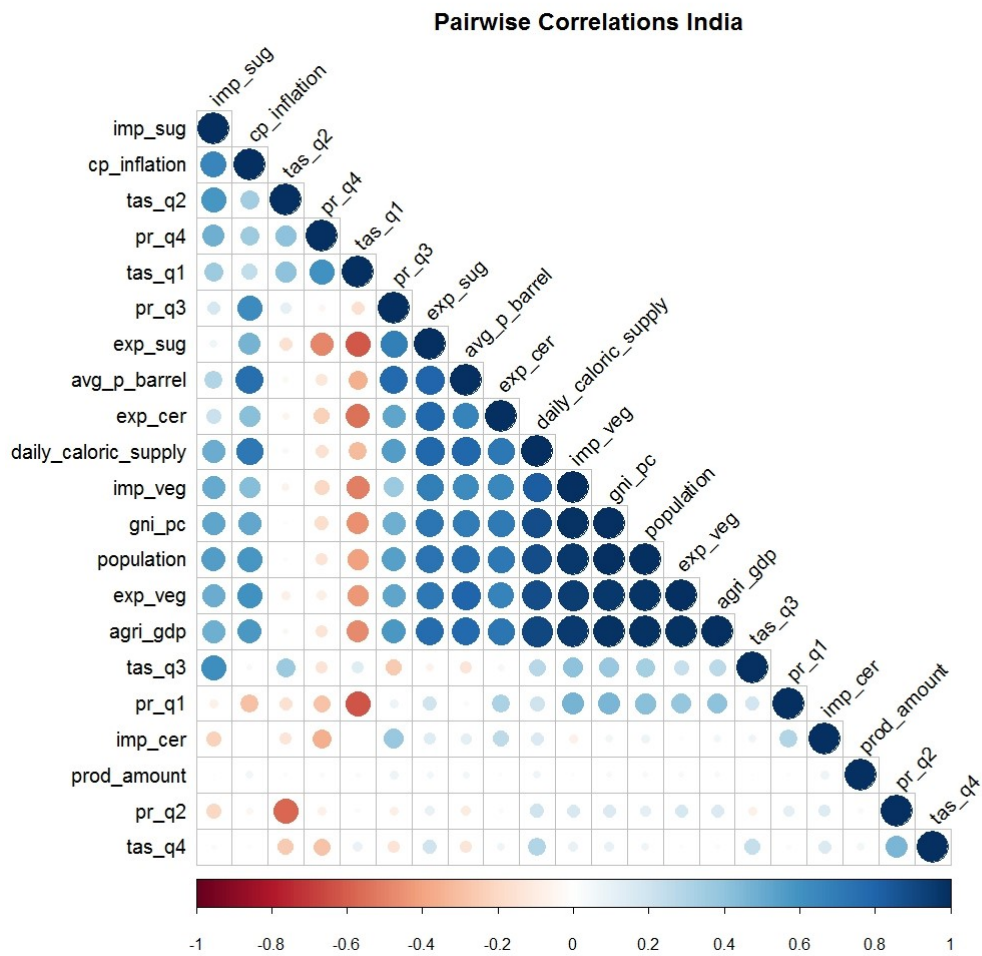


Figure 6. India Correlations

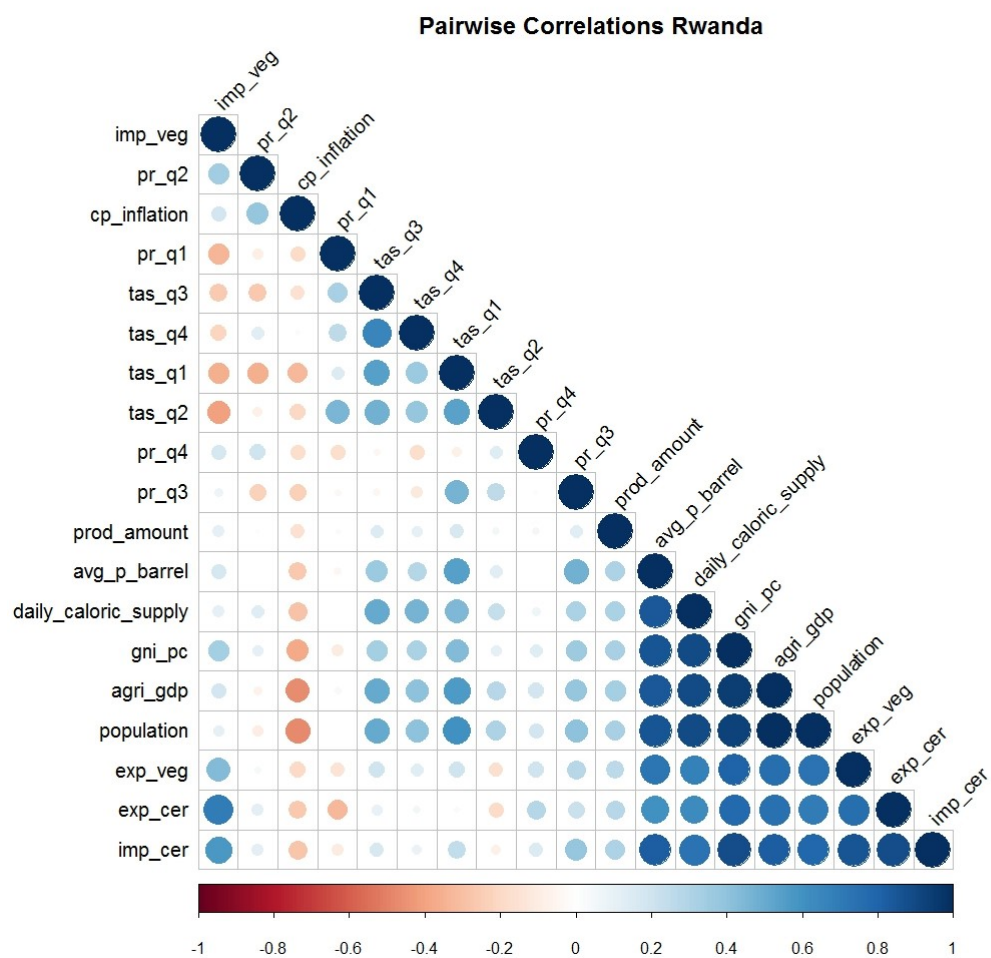


Figure 7. Rwanda Correlations

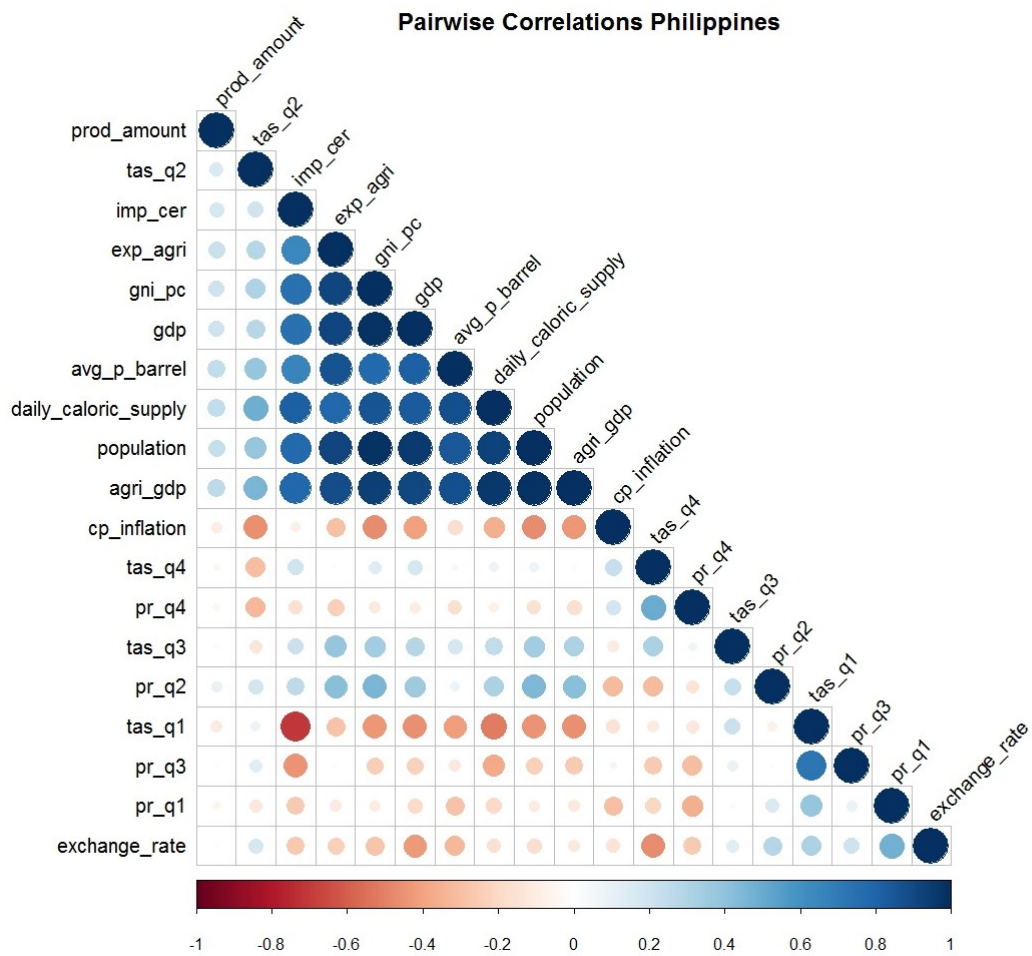


Figure 8. Philippines Correlations

<u>India</u>	<u>Rwanda</u>	<u>Philippines</u>
gni_pc	gni_pc	gni_pc
agri_gdp	agri_gdp	agri_gdp
population	population	population
daily_caloric_supply	daily_caloric_supply	daily_caloric_supply
exp_veg	imp_cer	imp_cer
exp_sug	exp_cer	exp_agri
avg_p_barrel	avg_p_barrel	gdp
		tas_q1

Table 6. My caption

4.2.1.2. Linear Model without highly correlated variables Having removed highly correlated variables, a linear model is built from the remaining predictors:

India

Residuals:

Table 7. My caption

Min	1Q	Median	3Q	Max
-0.68000	-0.23519	0.03321	0.22750	0.65974

Coefficients:

Variable	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	-4.83E-12	4.95E+01	0.000	1.000	
pr_q1	-5.53E+02	3.19E+02	-1.731	0.09243	x
pr_q2	-1.28E+02	1.00E+02	-1.272	0.21217	
pr_q3	-2.23E+02	2.37E+02	-0.944	0.35176	
pr_q4	2.42E+02	2.28E+02	1.063	0.29506	
tas_q1	-5.04E+02	3.24E+02	-1.554	0.12950	
tas_q2	2.77E+02	2.06E+02	1.347	0.18686	
tas_q3	-5.31E+02	5.90E+02	-0.900	0.37424	
tas_q4	1.66E+02	1.22E+02	1.360	0.18283	
prod_amount	6.34E+02	5.05E+01	12.561	2.52e-14	x
exp_cer	-3.41E+02	2.34E+02	-1.455	0.15482	
imp_sug	4.49E+02	6.66E+02	0.673	0.50548	
imp_veg	1.21E+03	3.64E+02	3.316	0.00218	x
imp_cer	5.62E+02	4.02E+02	1.396	0.17184	
cp_inflation	-3.04E+02	3.71E+02	-0.820	0.41774	

Table 8. My caption

Residual standard error: 0.3466 on 34 degrees of freedom
Multiple R-squared: 0.9149, Adjusted R-squared: 0.8799
F-statistic: 26.12 on 14 and 34 DF, p-value: 3.915e-14

Rwanda

Residuals:

Min	1Q	Median	3Q	Max
-1.4427	-0.6418	-0.2439	0.3527	2.8037

Table 9. My caption

Coefficients:

Variable	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	-9.63E-13	6.98E+01	0	1.00000	
pr_q1	1.18E+02	8.92E+01	1.327	0.18638	
pr_q2	9.90E+01	1.01E+02	0.98	0.3287	
pr_q3	3.38E+01	1.01E+02	0.335	0.73769	
pr_q4	1.01E+02	8.75E+01	1.151	0.25153	
tas_q1	2.77E+02	1.20E+02	2.318	0.0217	x
tas_q2	-1.28E+02	1.19E+02	-1.077	0.28327	
tas_q3	-1.79E+01	1.29E+02	-0.139	0.88996	
tas_q4	1.25E+02	1.11E+02	1.125	0.26234	
prod_amount	-2.19E+02	7.39E+01	-2.959	0.00355	x
exp_veg	1.26E+02	1.02E+02	1.234	0.21882	
imp_veg	2.64E+02	1.00E+02	2.643	0.00903	x
cp_inflation	-6.56E+01	8.64E+01	-0.759	0.44901	

Table 10. My caption

Residual standard error: 0.9235 on 162 degrees of freedom

Multiple R-squared: 0.2059, Adjusted R-squared: 0.1471

F-statistic: 3.501 on 12 and 162 DF, p-value: 0.0001285

Philippines:

Residuals

Min	1Q	Median	3Q	Max
-0.9879	-0.4897	-0.2891	0.3079	2.4465

Table 11. My caption

Coefficients:

Variable	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	-1.10E-13	1.02E+02	0	1.00000	
tas_q2	-5.72E+01	1.65E+02	-0.346	0.73059	
tas_q3	-1.15E+01	1.49E+02	-0.077	0.93899	
tas_q4	9.71E+01	1.67E+02	0.583	0.56205	
pr_q1	-1.01E+01	1.59E+02	-0.064	0.94933	
pr_q2	1.95E+02	1.22E+02	1.607	0.11324	
pr_q3	8.93E+00	1.18E+02	0.076	0.93986	
pr_q4	-4.00E+01	1.43E+02	-0.279	0.78124	
avg_p_barrel	4.28E+02	1.44E+02	2.978	0.00419	x
prod_amount	-3.63E+02	1.07E+02	-3.389	0.00124	x
exchange_rate	-1.71E+02	1.69E+02	-1.011	0.31585	
cp_inflation	-1.09E+02	1.45E+02	-0.750	0.45605	

Table 12. My caption

Residual standard error: 0.8693 on 60 degrees of freedom

Multiple R-squared: 0.3614, Adjusted R-squared: 0.2444

F-statistic: 3.087 on 11 and 60 DF, p-value: 0.002419

4.2.1.3. Linear Model with variables left after VIF-removal Since removing highly correlated variables leads to some potentially important variables being dropped, we applied VIF-removal to both the set of highly correlated variables and the remaining set. Creating another linear model with the resulting set of variables we obtained the following results

India:

Residuals

Min	1Q	Median	3Q	Max
-0.66068	-0.22061	0.02824	0.20910	0.67543

Table 13. My caption

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	4.16E-12	5.00E+01	0	1.0000	
population	1.41E+03	6.42E+02	2.188	0.0359	x
daily_caloric_supply	-4.90E+02	1.01E+03	-0.486	0.6301	
exp_sug	2.41E+03	2.54E+03	0.948	0.3501	
avg_p_barrel	-3.99E+02	8.40E+02	-0.475	0.6379	
pr_q1	-2.76E+02	1.69E+02	-1.634	0.1118	
pr_q2	-5.54E+02	5.53E+02	-1.002	0.3236	
pr_q3	-1.18E+03	1.05E+03	-1.130	0.2667	
pr_q4	1.31E+03	1.12E+03	1.172	0.2494	
tas_q2	-2.57E+02	5.92E+02	-0.435	0.6667	
tas_q3	2.03E+01	2.90E+02	0.070	0.9448	
tas_q4	-1.64E+02	2.40E+02	-0.682	0.5001	
prod_amount	6.34E+02	5.10E+01	12.418	5.51e-14	x
exp_cer	-9.57E+02	8.72E+02	-1.097	0.2806	
imp_cer	8.96E+02	8.99E+02	0.996	0.3264	
cp_inflation	-2.06E+02	6.71E+02	-0.307	0.7604	

Table 14. My caption

Residual standard error: 0.35 on 33 degrees of freedom
Multiple R-squared: 0.9158, Adjusted R-squared: 0.8775
F-statistic: 23.92 on 15 and 33 DF, p-value: 1.758e-13

Rwanda

Residuals:

Min	1Q	Median	3Q	Max
-1.5524	-0.6039	-0.2119	0.3888	2.9797

Table 15. My caption

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	-6.61E-13	6.91E+01	0	1.00000	
population	5.28E+02	4.79E+02	1.102	0.27214	
daily_caloric_supply	-4.03E+02	2.85E+02	-1.413	0.15966	
avg_p_barrel	3.42E+02	1.81E+02	1.892	0.06027	x
exp_cer	-1.72E+01	2.91E+02	-0.059	0.9531	
pr_q1	1.35E+02	1.00E+02	1.348	0.17956	
pr_q2	1.21E+02	1.40E+02	0.867	0.38741	
pr_q3	1.44E+01	1.15E+02	0.125	0.90036	
pr_q4	8.60E+01	1.02E+02	0.84	0.40223	
tas_q1	2.73E+01	1.88E+02	0.146	0.88442	
tas_q2	-1.56E+02	1.29E+02	-1210	0.22816	
tas_q3	2.02E+01	1.59E+02	0.127	0.89905	
tas_q4	9.00E+01	1.18E+02	0.761	0.44762	
prod_amount	-2.31E+02	7.38E+01	-3131	0.00208	x
exp_veg	-1.30E+02	1.57E+02	-0.831	0.4075	
imp_veg	2.09E+02	1.77E+02	1.180	0.2397	
cp_inflation	2.12E+01	1.08E+02	0.196	0.8446	

Table 16. My caption

Residual standard error: 0.9138 on 158 degrees of freedom

Multiple R-squared: 0.2417, Adjusted R-squared: 0.1649

F-statistic: 3.148 on 16 and 158 DF, p-value: 0.0001125

Philippines

Residuals:

Min	1Q	Median	3Q	Max
-1.0930	-0.4379	-0.2995	0.2246	2.3262

Table 17. My caption

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	Significant
(Intercept)	-1.14E-13	1.05E+02	0	1.00000	
daily_caloric_supply	4.35E+02	1.26E+03	0.346	0.7308	
exp_agri	3.70E+02	6.26E+02	0.591	0.55687	
imp_cer	-7.12E+01	5.27E+02	-0.135	0.89303	
tas_q1	-1.98E+02	3.95E+02	-0.502	0.61793	
tas_q2	-3.10E+01	2.90E+02	-0.107	0.91522	
tas_q3	-1.15E+01	2.58E+02	-0.044	0.96469	
tas_q4	4.17E+01	1.86E+02	0.225	0.82314	
pr_q1	2.46E+01	2.58E+02	0.095	0.92448	
pr_q2	-4.54E+01	3.47E+02	-0.131	0.89644	
pr_q3	1.92E+02	3.11E+02	0.619	0.53835	
pr_q4	-1.47E+01	2.45E+02	-0.060	0.95243	
avg_p_barrel	-2.64E+02	1.19E+03	-0.222	0.8252	
prod_amount	-3.64E+02	1.10E+02	-3.316	0.00161	x
exchange_rate	-2.01E+02	2.25E+02	-0.894	0.3754	
cp_inflation	-6.22E+01	1.96E+02	-0.317	0.7523	

Table 18. My caption

Residual standard error: 0.8884 on 56 degrees of freedom

Multiple R-squared: 0.3775, Adjusted R-squared: 0.2108

F-statistic: 2.264 on 15 and 56 DF, p-value: 0.01414

4.2.2. Lasso based method

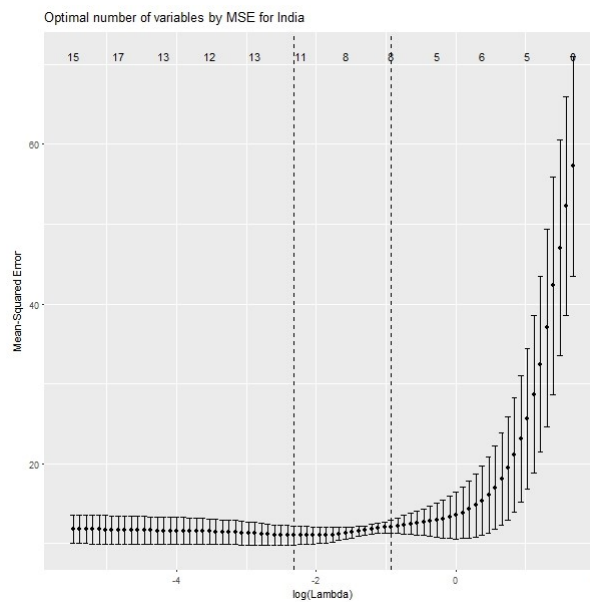


Figure 9. India Optimal Number of Variables

Important variables for India:

variables	coefs_minlambda	colors
tas_q2	8.203321e-01	green
prod_amount	4.289088e-05	green
exp_veg	2.596261e-06	green
agri_gdp	5.487558e-11	green
exp_cer	7.872368e-11	green
imp_sug	7.540223e-07	green
imp_veg	2.962602e-07	green
imp_cer	-2.776436e-06	red
pr_q2	-2.285270e-02	red
pr_q1	-1.124898e-02	red
tas_q3	8.462011e-01	green

Table 19. My caption

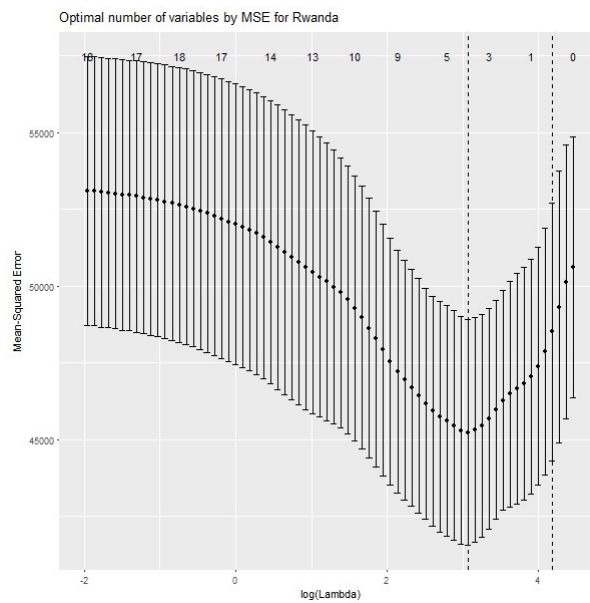


Figure 10. Rwanda Optimal Number of Variables

Important variables for Rwanda:

variables	coefs_minlambda	colors
tas_q4	1.543063e+00	green
avg_p_barrel	1.855989e-01	green
imp_cer	1.654824e-03	green
prod_amount	-2.160013e-05	red

Table 20. My caption

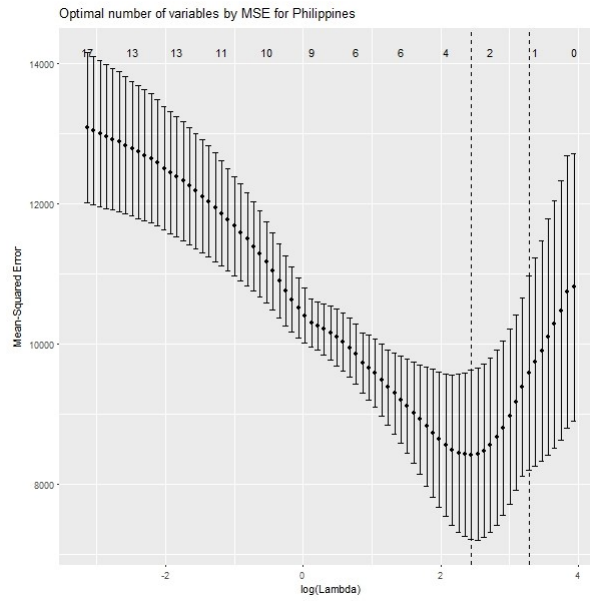


Figure 11. Philippines Optimal Number of Variables

Important variables for Philippines:

<u>variables</u>	<u>coefs_minlambda</u>
gdp	5.429905e-10
prod_amount	-3.595559e-06
avg_p_barrel	8.018443e-02

Table 21. My caption

4.2.3. Random Forest based method

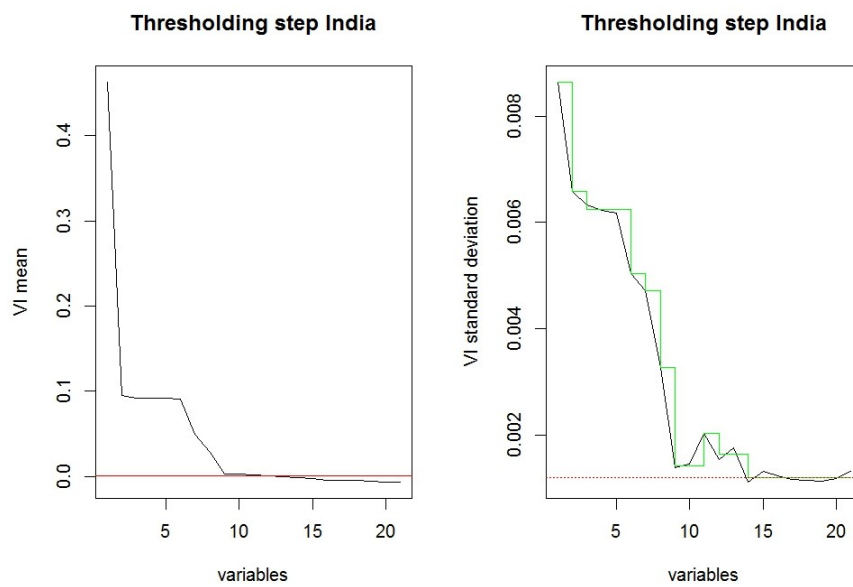


Figure 12. India Thresholding Step

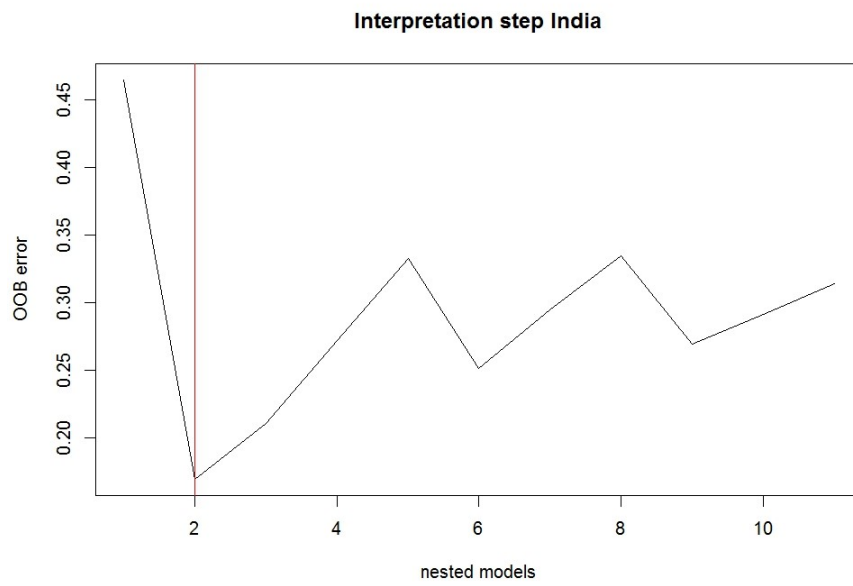


Figure 13. India interpretation Step

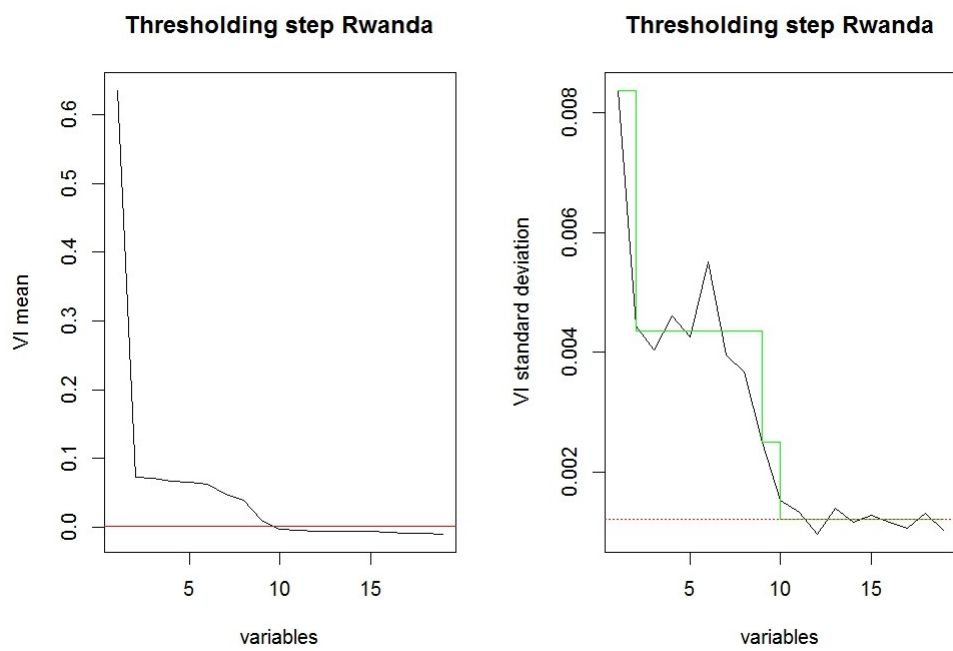


Figure 14. Rwanda Thresholding Step

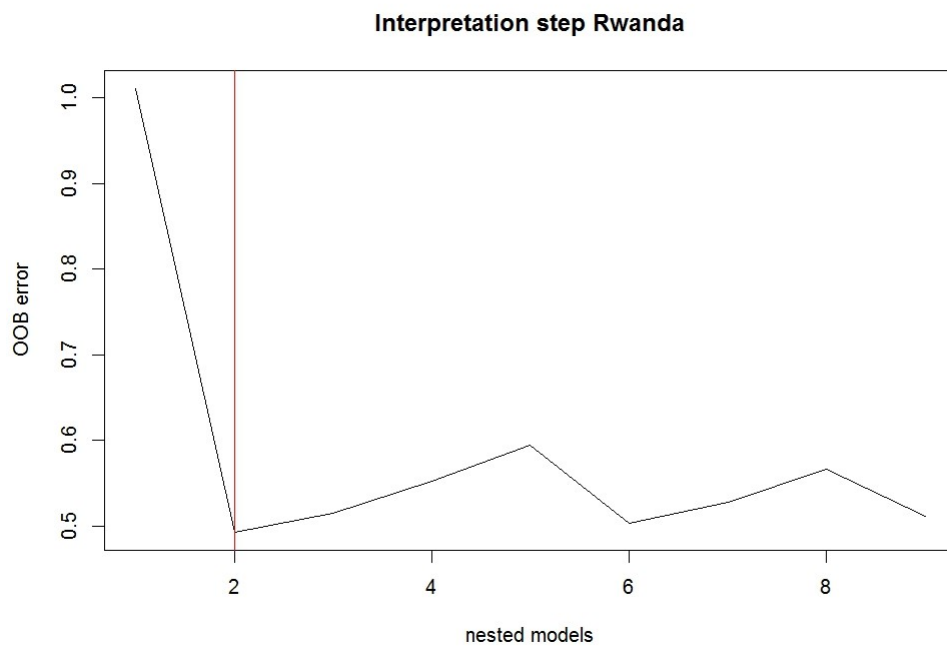


Figure 15. Rwanda interpretation Step

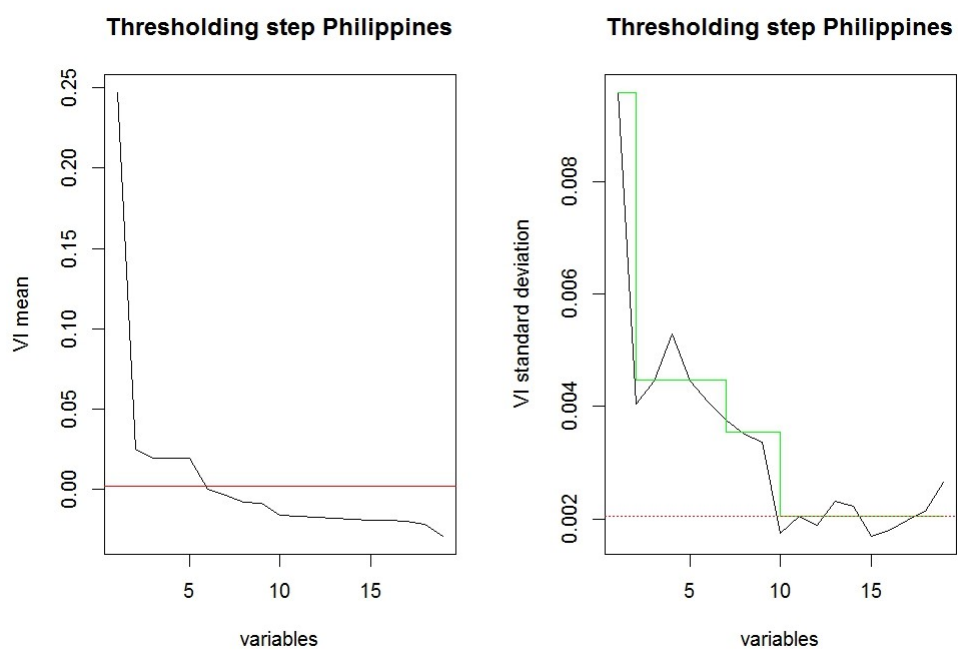


Figure 16. Philippines Thresholding Step

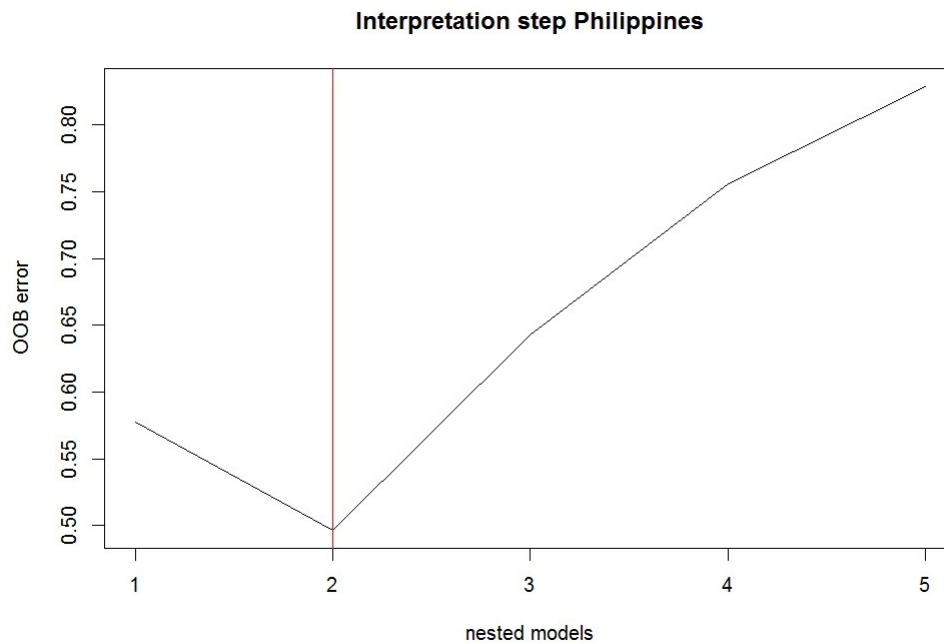


Figure 17. Philippines interpretation Step

India		Rwanda		Philippines	
threshold step	interpretation step	threshold step	interpretation step	threshold step	interpretation step
prod_amount	prod_amount	prod_amount	prod_amount	prod_amount	prod_amount
agri_gdp	agri_gdp	agri_gdp		agri_gdp	agri_gdp
gni_pc		gni_pc		gni_pc	
population		population		population	
imp_veg					
exp_veg		exp_veg			
daily_caloric_supply		daily_caloric_supply			
imp_sug					
cp_inflation					
imp_cer		imp_cer	imp_cer		
exp_sug					
		exp_cer			
		avg_p_barrel			
				gdp	

Table 22. My caption

4.2.4. Summary

India:

The linear model with non highly correlated variables reveals that the amount of rain in the first quarter, the produced amount and the amount of imported vegetables significantly affect the price. The model made after VIF reduction has only two significant variables, population and produced amount. It appears that population, the only significant demand related factor , outweighs climatic and macroeconomic supply related factors in terms of predictive power in the linear models. However produced amount seems to be the most important factor since it is always included.

The lasso method gives a more differentiated impression. Again, produced amount seems to be important. Imports of sugar, vegetables and cereals as well as demand related factors ex-

ports of vegetables and cereals also have an impact on price. Climatic factors play a bigger role here since rain in the first half of the year and temperatures from April to September also appear to be a driving factor, probably because of their impact on the agricultural output. The fact that agricultural GDP is also on the list supports this.

The findings produced by the threshold step of the RF method generally go in the same direction as produced amount and agri_gdp seem to be important again as well as imports of sugar and vegetables and export of vegetables. In addition to that population size, gross national income per capita, daily caloric supply and inflation have an influence on price too. After the interpretation step only produced amount and agri_gdp are left.

The importance of the produced amount is evident as it is included in every result. On the demand side, population size and food exports appear to be the most influential factors.

Rwanda:

According to the findings produced by the linear model with no highly correlated variables, only supply related factors are to be considered. Produced amount appears to be very important, temperature in the first quarter and the amount of imported vegetables are also influential. The VIF based model reveals that apart from produced amount, food prices in Rwanda are also dependent on the petrol oil price. These findings are supported by the results of the lasso model, as in addition to produced amount and oil price, rain in the fourth quarter of the year and the import of cereals are also affecting food prices.

The RF model indicates, that in addition to agricultural output, demand related variables such as population size, GNI per capita, daily caloric supply and the amount of exported vegetables contribute to food price development. After the interpretation step only produced amount is left.

Again, the single most important factor appears to be produced amount. Combined with demographic factors and the dependency of Rwanda's agriculture on petroleum most of the food prices in this country can be explained.

Philippines:

The linear model with low correlated variables for the Philippines gives only two significant factors : oil price per barrel and produced amount. The VIF based model only has produced amount. While not adding any further information these findings are in line with the results for the other two countries. This is also supported by the findings obtained from the lasso model, as only gdp, produced amount and oil price per barrel appear to be important. The interpretation step of the RF model adds population and gni_pc to the list of relevant variables, produced_amount and agri_gdp are considered important as well. After the interpretation step only produced amount and agri_gdp remain.

In the Philippines, food prices appear to be mostly driven by the produced amount and the population, while other factors only play a minor role.

5. Outlook

Some other techniques or combinations of techniques, more data, different sources, etc, think of something

India: produced amount, population growth, imports, exports, rain January – June, temperatures April - September
Rwanda: produced amount, population growth and income, vegetable and cereal imports, oil price, rain October – December, temperature January- March.
Philippines: produced amount, population growth

6. Conclusion

Erokhin highlighted the fact that one important contributing factor to food security are food prices in his work. The OECD identified an understanding on food price development as a paramount basis for sound food security policy decision making on the national level. With these observations in mind we started this work with the goal to see if we could replicate the findings of Erokhin and the OECD on a country specific level. The initial step was the country selection. We decided on India, Rwanda and the Philippines as this selection of countries provided us with a varied economic, political and climatic basis of research, which in turn allowed us to form a comprehensive conclusion if the global OECD findings were applicable on a country specific level. We will first contrast the global OECD findings on consumer food price development with the country specific selection of consumer goods food price development. All goods selected were chosen based on the Harmonized System Codes (HS Code 2017). We focused on the category 07-015. The OECD analyzed annual consumer food price development of wheat, coarse grains rice and oil seeds from 1971 to 2007. The conclusion is that consumer food prices are volatile and price spikes are a common occurrence. The OECD has identified unfavorable climatic conditions in 2005 and 2007 in major crop producing regions as one factor contributing to an increase in consumer food prices (OECD, 2008). Further strong demand growth for food and animal feed led to increase in price development. Other factors contributing to consumer price trend development as noted by the OECD are the following: Macro-economic conditions such as GDP growth in developing countries and oil price development. The OECD predicted based on the above factors that global consumer food prices would increase in comparison to historic trends but also that consumer prices would decrease from the 2007-2008 spike price level. Our findings support this prediction as one can note a general upward trend in the global as well as national consumer food prices for our selection of goods as shown in section 4.1. Figure X shows that consumer prices of our selection of cereals, vegetables oils and sugar are increasing in accordance with global food consumer prices. One can observe that while this price development is above the historic trend and generally increasing, the price spike in 2007-2008 has normalized as predicted by the OECD. Our findings show that global perspective and prediction of the OECD in regards to the consumer food price index hold true on a by country level. OECD stance on an increasing trend of consumer food prices in regards to food security is that while the overall impact on developing countries is modest, urban poor populations will be strongly negatively impacted. This discrepancy of modest national impact to strong negative impact for the poor urban population is based on the share of disposable income available and its use to purchase food at higher consumer food prices. According to Erokhin another negative aspect of an increase in consumer food prices is that exporting goods becomes more profitable to producers from countries with a weak local currency (Erokhin, 2017). This can further exacerbate food security for the poorest populations of a country.

References

- [1] Michel Goossens, Frank Mittelbach, and Alexander Samarin. *The L^AT_EX Companion*. Addison-Wesley, Reading, Massachusetts, 1993.
- [2] Erokhin. V. *Establishing Food Security and Alternatives to International Trade in Emerging Economies*. Case of Adana. Pak. J. Nutr. 2013, 13, 1–6