# INDIVIDUAL ASSIGNMENT

## TECHNOLOGY PARK MALAYSIA

## AICT009-4-2-IDA

## INTRODUCTION TO DATA ANALYTICS

## UCDF2005(1) ICT(DI)

## <u>IDA DOCUMENTATION</u>

**Name:** Chan Ming Li

**TP Number:** TP060774

**Hand out Date:** 09 August 2021

**Hand in Date:** 15 November 2021

**Weightage:** 40%

**Lecturer Name:** Ms. Hema Latha Krishna Nair

## Table of Contents

# 1.0 Introduction

An environment is the surrounding and nature in which the non-living and living things live around. The environment is composed of a variety of essential elements such as air, water, and land. These elements are important to maintain a healthy and balanced nature in an environment. It is undeniable that the environment has been exposed to several threats and causing environmental issues. These environmental issues have enormously increased due to unbalanced nature such as changes of climate, natural disasters, polluted waters, air, land etc. (Vedantu Learn LIVE Online, 2021) . Due to these negative impacts on the environment, organisms have reduced available sources to live such as air, water, and land. In addition, the extinction of biodiversity, lack of minerals, as well as countries would be affected economically. If these problems continue for the following decades, the environment that we live with for years will cease to be exist. Yet we still ignore despite the public awareness, media attention and few of research has done to prevent environmental issues.

As the severity of the environmental issues increases, many problems and challenges have emerged and becoming serious which has a huge impact on our lives and environment. One of the problems caused by environmental issues is human health as it increases the risk of getting diseases such as lung cancer, asthma and many more. Dirty and poisonous gas emissions which causes air pollution have a similar effect as smoking which will destroy our respiratory system (World Health Organization, 2021). Besides, environmental issues threaten the habitat of flora and fauna by habitat fragmentation and habitat degradation. For example, a forest which is mainly the habitat of flora and fauna is destroyed due to deforestation caused by agricultural expansion and so on. Moreover, we found that environmental issues have a huge impact on the economy although some companies have started sustainable development which will mitigate the effects of pollution. As the availability of non-renewable resources is declining, production or business which highly dependent on it is affected (Rinkesh, n.d.).

# 2.0 Business Goal & Objective

## 2.1 Business Goal

The business goal of our group is to conduct analysis on discovering analytic solution that is suitable to solves and handles impacts of global environmental issues effectively. Whether individually or as a group, we designed some solution that should reduce the impacts of global environment by referencing real world solutions that has been implemented. The solution that we identified in reducing impacts of global environment issues is enhancing the health of human by identifying the health issues, controlling the habitat index to be in the normal level and develop plans for sustainable environment and economy growth. We are planning to achieve this solution using different strategies such as predictive modelling to predict and forecast future trends or performance in different aspects. Besides, we will use dashboard solution to visualize significant impacts of environmental issues. By using these strategies, BI report and OLAP is able to carry out activity of discovering patterns and trends in datasets collected.

## 2.2 Objectives

1.  To evaluate the effect of global environment issues on human health based on locations and utilize a clustering model to reduce environmental issue and health disease.
2.  To identify the loss of forest due to environmental issue and implement a statistical model to determine the relationship between the possible drivers of deforestation and loss of forest based on years and locations. **(Chan Ming Li)**
3.  To analyze how economic growth contribute to environmental issues and develop predictive model to support plan for sustainable environment and economy.
4.  To study the loss of population affected by environmental issues and design a predictive model based on different categories in a type of environmental issues based on years.

## 2.3 Scopes

### 2.3.1 Scope 1: OLAP Dashboard on The Change of Forest Loss Impacted by Deforestation

Deliverables:  To analyse and visualize the change of forest loss impacted by deforestation globally based on years by using Analytical Dashboard (OLAP).

- OLAP is known as Online Analytical Processing which is a technology for data discovery includes complex analytical calculations, trend analysis and more. OLAP performs multidimensional analysis which allows extraction of information can be done in a consistence, interactive and fast way. For instance, the severity of the forest loss can be measured and determined by processing the change of forest loss across the years and countries with other factors.

- Layman or other users will have a better understanding on the analysis done as OLAP provides a clear, user-friendly and interactive dashboard which contains information from different aspects to them. Hence, they are able to understand the current situation or the severity of forest loss which leads them to several compliance actions.

- There are several basic operations like roll up, drill down, slice and dice and pivot can be performed by implementing OLAP. All these operations will help users to gain detailed information from different perspectives or dimensions. For example, the regions of forest loss can be drilled down into different countries by fragmenting the data into smaller parts.

## 2.3.2 Scope 2: Predictive Modelling on Forest Loss and Possible Drivers of Deforestation

Deliverables: To develop a linear regression model to determine the relationship between the forest loss and the possible drivers of deforestation in order to implement suitable conservation planning.

- Predictive modelling which is also known as predictive analytics is commonly utilized to predict or forecast the future events or possible outcomes by analysing the patterns from the related input data. For instance, the forest loss can be predicted by analysing the relationship between forest loss and the possible drivers of deforestation.

- The predictive results provided by the model can be used to make fast and critical decisions by users as the results provide an overview on the most likely happen scenario or outcome.

- Based on the results from the prediction model, public or different organisation can provide and implement different solutions to avoid or mitigate the forest loss which is highly affecting our environment.

# 3.0 Data Analytics Life Cycle & Methodologies

## 3.1 Fayyad's KDD Methodology

There are many methodologies often implemented for data analytics and data mining such as Fayyad's KDD, Crisp-DM, SEMMA and many more. Fayyad's knowledge discovery in databases (KDD) is a process or procedure of discovering knowledge that find and transform patterns from a large amount of data by implementing different data mining techniques and algorithms. It is often utilized by researchers from different fields such as machine learning, databases, data visualization, artificial intelligence and many more. KDD methodology plays an important role in data mining field as the main backbone of its process is data mining which allows recognizing and exploring understandable patterns from a vast amount of data (Usama Fayyad, 1996).



*Figure 3(a): Fayyad's KDD Process Model*

## 3.2 Fayyad's KDD Process

Fayyad's KDD is an iterative and interactive process which can be different in a range of 5 to 9 steps based on different situation and the decisions made by users. Firstly, the initial phase in the process is fundamental as problem statement is needed to be identify. The user should have some relevant knowledge and understand the application domain in order to avoid false interpretations and identify their goal (Sharma, 2020). For instance, the understanding of our group domain and identifying the objectives and scopes are conducted in 2.0 part.

After that, data selection takes places as the second phase in the Fayyad's KDD process where includes discovering or finding the relevant, accessible, and critical dataset. As the objectives and goal are identified earlier, users will have a better understanding on which kind

of dataset and variables are needed for the data models (Usama Fayyad, 1996). For example, the resources required for the project such as datasets are determined and selected from different sources which links are provided in 4.0 part. Besides, the usability of the dataset or attributes are determined and selected after exploring and reviewing in 5.1 part to ensure all datasets or attributes are related and helpful based on the objectives and scopes.

Furthermore, data pre-processing and cleaning is carried out as the third phase in the methodology which includes some basic operations to improve the data quality and integrity (Sharma, 2020). Strategies on handling missing or null values as well as noisy data are determined to improve the reliability of the data chosen. For instance, some cleaning operations like removing null values are conducted in 5.2 part based on the necessity of the following analysis and modelling. There is a data cleaning report includes the steps taken in the data cleaning process shown in 5.2 part.

Moreover, the fourth phase of the process is data transformation which prepares data for data mining and modelling. The data are integrated and consolidated by merging datasets, converting data types and many more according to the usage of following analysis and modelling (Sharma, 2020). For example, the datasets are merged based on the matching columns and several columns are aggregated which helps to build OLAP dashboard and Linear Regression model much easier and effective in 5.3 part.

In addition, data mining is the following phase in the process which is the root phase of the Fayyad's KDD methodology. Different approaches and algorithms are executed in order to extract and reveal information by discovering understandable patterns and trends from the cleaned and transformed data. For instance, a correlation plot is built in order to determine the correlation coefficient of different attributes for linear regression model.

Besides, interpretation or evaluation will be carried out in the sixth phase of the process which interprets the mined patterns or trends obtained from the previous phase. These patterns will be visualized in different forms such as bar graphs, line graphs, pie charts and many more so that the patterns or trends can be clearly viewed by users. For instance, OLAP dashboard that includes different visualization like line graphs, area charts and many more is built and evaluated in 7.0 part.

Lastly, the final phase of Fayyad's KDD process is using the knowledge discovered from the previous phases. The users are able to measures the impact and make decisions based on

the knowledge gained. The knowledge discovered will be documented in a report so that the users can better evaluate the output (Sharma, 2020).

## 3.3 The Advantages & Drawbacks of Fayyad's KDD

There are some advantages and drawbacks when Fayyad's KDD methodology is implemented. Firstly, one of the advantages of KDD is that it could be used for forecasting purposes. For examples, market forecasting which is predicting the customer behaviour and consumer trends by using KDD methodology. Furthermore, the process in KDD is iterative which increases impact of the data by refining the data and gain new knowledge from it. The knowledge that gained from any phase of process can be applied back to the beginning of the process (Data Science Process Alliance, 2021).

However, there are some drawbacks in the process of Fayyad's KDD too. First of all, time consuming is one of the drawbacks of the methodology. As the methodology has an iterative process, the knowledge gained is applied back to the process in the next iteration which leads to time consuming. While returning the process, there is a possible that additional data will be used which needed to run through the whole process again. Hence, the process is considered as one of the time-consuming processes. Besides, the security and privacy of the process is vulnerable which is definitely a drawback of the methodology. When there is a need of customer information for analysis, the data is needed to be collect as much as possible in order to gain more knowledge from it. Therefore, there will be a huge calamity if there is any malicious attacks or data breaches occurs (Data Science Process Alliance, 2021).

# 4.0 Data Understanding

**Dataset Name**: Agricultural Area Dataset

**Data Source**: https://data.worldbank.org/indicator/AG.LND.AGRI.K2

- There are 65 columns and 266 rows about the agricultural land area in square kilometer for different countries from 1960 to 2020.
- The columns that exist in the dataset is "Country Name", "Country Code", "Indicator Name", "Indicator Code", Year columns (1960 – 2020).
- All data in the dataset is in a text data type.
- There are errors like null values, unstandardized, unmatched data types, blank columns, blank rows.

| Column | Problems Identified |
|---|---|
| Year (1960 - 2020) | Null values, unstandardized, unmatched data types, |

**Dataset Name**: Land Area Dataset

**Data Source**: https://data.worldbank.org/indicator/AG.LND.TOTL.K2

- There are 65 columns and 266 rows about the land area of country globally from 1960 to 2020.
- The columns that exist in the dataset is "Country Name", "Country Code", "Indicator Name", "Indicator Code", Year columns (1960 – 2020).
- All data in the dataset is in a text data type.
- There are errors like null values, unstandardized, unmatched data types, duplicate columns, blank rows.

| Column | Problems Identified |
|---|---|
| Year (1960 - 2020) | Null values, unstandardized, unmatched data types, |
| Country Name | Duplicate columns |

**Dataset Name**: Forest Area Dataset

**Data Source**: https://data.worldbank.org/indicator/AG.LND.FRST.K2

- There are 65 columns and 266 rows about the forest area in square kilometers for different countries and years.

- The columns that exist in the dataset is "Country Name", "Country Code", "Indicator Name", "Indicator Code", Year columns (1960 – 2020).

- All data in the dataset is in a text data type.

- There are errors like null values, unstandardized, unmatched data types, blank columns, blank rows.

| Column | Problems Identified |
|---|---|
| Year (1960 - 2020) | Null values, unstandardized, unmatched data types, |

**Dataset Name**: Annual Deforestation Dataset

**Data Source**: https://ourworldindata.org/grapher/annual-deforestation?tab=table

- There are 4 columns and 495 rows about the global deforestation rate annually for certain years which is 1990, 2000, 2020 and 2015.

- The columns that exist in the dataset is "Country", "Code", "Year", "Deforestation", Year columns (1990, 2000, 2020, 2015).

- All data in the dataset is in a text data type.

- There are errors like null values, unstandardized, unmatched data types, unmatched unit of measurement.

| Column | Problems Identified |
|---|---|
| Code | Null values |
| Deforestation | Null values, unstandardized, unmatched data types, unmatched unit of measurement |

**Dataset Name**: Forest Loss Dataset

**Data Source**: https://gfw.global/3kAcKgX

- There are 4 columns and 3968 rows about forest loss for different countries from 2001 to 2020.

- The columns that exist in the dataset is "iso", "umd_tree_cover_loss__year", "umd_tree_cover_loss__ha", "gfw_gross_emissions_co2e_all_gases__Mg".

- All data in the dataset is in a text data type.

- There are errors like null values, unstandardized, unmatched data types, unmatched unit of measurement.

| Column | Problems Identified |
|---|---|
| umd_tree_cover_loss__year | Null values |
| umd_tree_cover_loss__ha | Null values, unmatched data type, unmatched unit of measurement |

**Dataset Name**: Global Forest Loss by Dominant Drivers

**Data Source**: https://gfw.global/3Dr9WtK

- There are 3 columns and 252 rows about the forest loss caused by dominant drivers globally.
- The columns that exist in the dataset is "tsc_tree_cover_loss_drivers__type", "umd_tree_cover_loss__year", "umd_tree_cover_loss__ha", "gfw_gross_emissions_co2e_all_gases__Mg"
- All data in the dataset is in a text data type.
- It will be used to replace the missing countries name from other datasets with the countries code as a matching column.
- The columns in the dataset are unstandardized.

Dataset Name: iso_metadata

Data Source: https://gfw.global/3kAcKgX

- There are 2 columns and 252 rows about countries and their countries codes.
- The columns that exist in the dataset is "name", "iso"
- All data in the dataset is in a text data type.
- It will be used to replace the missing countries name from other datasets with the countries code as a matching column.

Dataset Name: Countries Data

Data Source: https://www.kaggle.com/fernandol/countries-of-the-world

- There are 20 columns and 228 rows about the details of different countries.
- Example columns that exist in the dataset is "Country", "Region", "Population", "Area".
- All data in the dataset is in a text data type.
- The countries and its region is required for the analysis as the other columns are not related to the objectives and scopes.

# 5.0 Data Preparation

Data preparation is a process that includes data cleaning, data transforming and data consolidating which is commonly utilized in business intelligence (BI), machine learning, predictive analytics and more. Data preparation is essential for data analysis or modelling as it helps to ensure data integrity, data quality and the accuracy of the analysis or modelling. The data preparation process is conducted in Power BI which is a user-friendly analytic service or platform provided by Microsoft.

## 5.1 Data Selection

After exploring and discovering the datasets above, there are a total of 228 columns to be selected for the analysis and modelling. The variables which is best related to the objectives and scopes will be selected to remove unused variables and increase the effectiveness of the analysis. The columns that are needed for the further analysis are "countries name", "countries code", "year", "land area", "agricultural land area", "forest area", "annual deforestation", and "forest loss". The other columns that are not selected or unused will be deleted or removed in the data cleaning process. As they are datasets from different sources, the columns from different datasets might have a different unit of measurement. The vary of unit of measurement will causes inaccurate analysis result and reduce the data quality. Therefore, the unit of measurement that is planned to utilize in the analysis and modelling is square kilometer (sq.km). The unit of measurement of certain columns such as "Annual Deforestation" will be changed and converted to square kilometer (sq.km) in the data transformation process. In addition, some aggregated columns like "average agricultural area" are required for the analysis and the columns will be calculated in the data transformation process.

## 5.2 Data Cleaning

The data cleaning and transformation process will be repeated to some selected datasets which contains the same problem. Therefore, the data cleaning process will be documented once so that there will no repetitive and redundant steps in the documentation. Besides, datasets that have involved in certain steps will be recorded in the dataset column and the description will be provided in the table. The illustration from an example dataset will be provided in the table.

## 5.2.1 Delete Unnecessary Rows

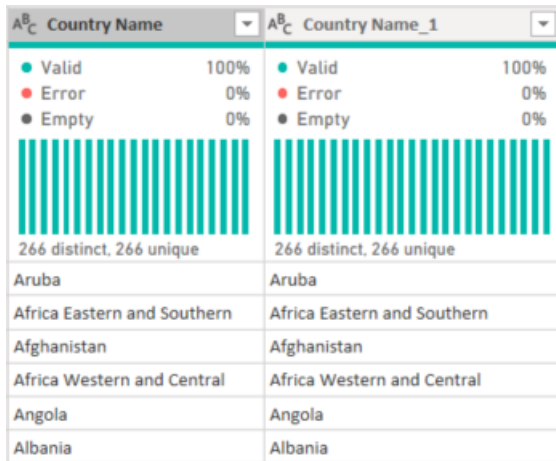| Dataset | Agriculture Land, Land Area, Forest Area |
|---|---|
| Before Cleaning |  |
| Delete unnecessary rows | • Delete the first four rows of the dataset which is unnecessary for our findings by using the remove top rows function in Power BI.<br><br>`= Table.Skip(Source,4)`<br><br> |
| Explanation | The first four rows of the dataset are needed to delete as they are just dataset title and last updated date. Hence, deleting the rows allows the dataset to be |

| | |
|---|---|
| | more standardise and does not cause much significant to the dataset for analysis. |
| Results |  |

## 5.2.2 Use First Row as Header

| | |
|---|---|
| Dataset | Agriculture Land, Land Area, Forest Area, |
| Before Cleaning | Example from Agricultural Land Dataset:<br> |
| Use first row as header | The first row of the dataset is used as the header by using the function provided.<br><br>Example from Agricultural Land Dataset:<br>`= Table.PromoteHeaders(#"Removed Top Rows", [PromoteAllScalars=true])` |

| | |
|---|---|
| |  |
| Explanation | The header for every column is needed as it will provides a brief message on what is the column about. The data is organized and can be read easily when the header is identified and added. |
| Results | Example from Agricultural Land Dataset:  |

## 5.2.3 Remove Duplicate Columns

| | |
|---|---|
| Dataset | Land Area |
| Before Cleaning | Example from Land Area Dataset:  |
| Remove Duplicate Columns | • Remove duplicate columns which is "Country Name_1" |

```
= Table.RemoveColumns(#"Promoted Headers",{"Country Name_1"})
```

| Explanation | The duplicated columns are removed as it will provide redundant data to the following analysis and modelling, and no information can be extracted from the column. |
|---|---|
| Results |  |

## 5.2.4 Delete Null Columns

| Dataset | Agricultural Land, Land Area, Forest Area |
|---|---|
| Before Cleaning | Example from Agricultural Land Dataset:  |

| | |
|---|---|
| | **Column statistics**                    · · · <br><br> Count                                      266 <br> Error                                        0 <br> Empty                                        0 <br> Distinct                                     1 <br> Unique                                       0 <br> Empty string                               266 <br> Min <br> Max <br><br> $^{AB}_C$ Indicator Code ▾   $^{AB}_C$ 1960 ▾   $^{AB}_C$ 1961 ▾ <br> • Valid 100%   • Valid 0%   • Valid 85% <br> • Error 0%   • Error 0%   • Error 0% <br> • Empty 0%   • Empty 100%   • Empty 15% <br><br> 1 distinct, 0 unique   1 distinct, 0 unique   216 distinct, 204 unique <br> AG.LND.AGRI.K2                    20 <br> AG.LND.AGRI.K2                    5326150 <br> AG.LND.AGRI.K2                    377000 <br> AG.LND.AGRI.K2                    3025335.1 <br> AG.LND.AGRI.K2                    571700 |
| Delete Null Columns | • For Agricultural Land dataset, delete the columns in the dataset that are 100% null (1960, 2019, 2020, a blank column). <br><br> • For Land Area dataset, delete the columns in the dataset that are 100 % null (1960, a blank column). <br><br> • For Forest Area dataset, delete the columns in the dataset that are 100 % null (1960 – 1989, a blank column). <br><br> Example from Agricultural Land Dataset: <br><br>  <br><br> `= Table.RemoveColumns(#"Promoted Headers",{"1960", "2019", "2020", ""})` |
| Explanation | The columns that have 100% missing values are deleted as they will adversely affect the performance of the analysis and do not provide any values to the analysis. |
| Results | Example from Agricultural Land Dataset: |

| AᴮC Indicator Code ▼ | AᴮC 1961 ▼ | AᴮC 1962 ▼ |
|---|---|---|
| ● Valid 100% ● Error 0% ● Empty 0% | ● Valid 85% ● Error 0% ● Empty 15% | ● Valid 85% ● Error 0% ● Empty 15% |
| 1 distinct, 0 unique | 216 distinct, 204 unique | 215 distinct, 202 unique |
| AG.LND.AGRI.K2 | 20 | 20 |
| AG.LND.AGRI.K2 | 5326150 | 5322890 |
| AG.LND.AGRI.K2 | 377000 | 377600 |
| AG.LND.AGRI.K2 | 3025335.1 | 3035850 |
| AG.LND.AGRI.K2 | 571700 | 572000 |
| AG.LND.AGRI.K2 | 12320 | 12320 |
| AG.LND.AGRI.K2 | 260 | 260 |

## 5.2.5 Replace Null Values

| Dataset | Agriculture Land, Land Area, Forest Area, Annual Deforestation, Forest Loss |
|---|---|
| Before Cleaning | Example from Agricultural Land Dataset:<br><br>**1961**<br>227 (85%) Valid | 0 (0%) Error | 39 (15%) Empty<br><br>💡 Remove Empty ···<br><br>**Column statistics** ···<br>Count 266<br>Error 0<br>Empty 0<br>Distinct 216<br>Unique 204<br>Empty string 39<br>Min<br>Max 990 |
| Replace null values with the average value across the years for different countries | • For Agricultural Land dataset, replace the missing values which consists of 3% - 15% in each selected column (1961 – 2018) with average value of agricultural land area calculated in 6.2.2.<br><br>• For Land Area dataset, replace the missing values which consists of 1% - 3% in each selected column (1961 – 2020) with average value of land area calculated in 6.2.2.<br><br>• For Forest Area dataset, replace the missing values which consists of 2% - 3% in each selected column (1990 – 2020) with average value of forest area calculated in 6.2.2. |

- For Annual Deforestation dataset, replace the missing values which consists of 3% - 20% in each selected column (1990. 2000, 2010 and 2015) with average value of forest area calculated in 6.2.2 and "Code" column with "N/A".

- For Forest Loss dataset, replace the missing values which consists of 3% - 10% in each selected column (2001 - 2020) with average value of forest loss calculated in 6.2.2.

Example from Agricultural Land Dataset:



```
= Table.ReplaceValue(#"Inserted Average",null,each [AVERAGE],Replacer.ReplaceValue,
{"1961", "1962", "1963", "1964", "1965", "1966", "1967", "1968", "1969", "1970", "1971",
"1972", "1973", "1974", "1975", "1976", "1977", "1978", "1979", "1980", "1981", "1982", "1983", "1984", "1985", "1986", "1987", "1988", "1989", "1990", "1991", "1992",
"1993", "1994", "1995", "1996", "1997", "1998", "1999", "2000", "2001", "2002", "2003", "2004", "2005", "2006", "2007", "2008", "2009", "2010", "2011", "2012", "2013",
"2014", "2015", "2016", "2017", "2018"})
```

| | |
|---|---|
| Explanation | The datasets consist of 3% - 15 % missing values in some columns which is a big number of missing values. If removing them, it will cause the sample data size smaller and information loss which is unbeneficial to the analysis. Therefore, the average of the values across the years for different countries is used to replace the null values in the selected columns. However, there are rows that are 100 % null which causes average cannot be calculated. Hence, the remaining null rows will be deleted in the following data cleaning process. |
| Results | Example from Agricultural Land Dataset: |

## 5.3 Data Transformation

### 5.3.1 Change Data Type

| Dataset | Agricultural Land, Land Area, Forest Area, Annual Deforestation, |
|---|---|
| Before Transformation | Example from Agricultural Land Dataset:  |
| Change Data Type | <ul><li>For Agricultural Land dataset, the data type of columns (1961 – 2018) are changed from "text" to "whole number".</li><li>For Land Area dataset, the data type of columns (1961 – 2020) are changed from "text" to "whole number".</li><li>For Forest Area dataset, the data type of columns (1990 – 2020) are changed from "text" to "whole number".</li><li>For Annual Deforestation dataset, the data type of columns (1990, 2000, 2010, 2015) are changed from "text" to "whole number".</li></ul> Example from Agricultural Land Dataset: |

| | |
|---|---|
| | Data Type: Text ▾    ¹₂ Replac<br><br>Decimal Number<br>Fixed decimal number<br>Whole Number<br>Percentage<br>Date/Time<br>Date<br>Time<br>Date/Time/Timezone<br>Duration<br>Text<br>True/False<br>Binary<br><br>⊢ Table.TransformColumnTypes(#"Removed Columns",{{"1961", Int64.Type}, {"1962", Int64.Type}, {"1963", Int64.Type}, {"1964", Int64.Type}, {"1965", Int64.Type}, {"1966", Int64.Type}, {"1967", Int64.Type}, {"1968", Int64.Type}, {"1969", Int64.Type}, {"1970", Int64.Type}, {"1971", Int64.Type}, {"1972", Int64.Type}, {"1973", Int64.Type}, {"1974", Int64.Type}, {"1975", Int64.Type}, {"1976", Int64.Type}, {"1977", Int64.Type}, {"1978", Int64.Type}, {"1979", Int64.Type}, {"1980", Int64.Type}, {"1981", Int64.Type}, {"1982", Int64.Type}, {"1983", Int64.Type}, {"1984", Int64.Type}, {"1985", Int64.Type}, {"1986", Int64.Type}, {"1987", Int64.Type}, {"1988", Int64.Type}, {"1989", Int64.Type}, {"1990", Int64.Type}, {"1991", Int64.Type}, {"1992", Int64.Type}, {"1993", Int64.Type}, {"1994", Int64.Type}, {"1995", Int64.Type}, {"1996", |
| Explanation | The data type is changed because the default data type which is text does not match the value which is a whole number. This may cause errors and the columns cannot be aggregated for the following cleaning and transformation. Besides, the data is rounded off automatically when the data type is changed to "whole number". |
| Results | Example from Agricultural Land Dataset:<br><br>¹²₃ 1961          ▾    ¹²₃ 1962          ▾<br><br>Data Type: Whole Number ▾ |

## 5.3.2 Create New Column with Average Value by Year and Country
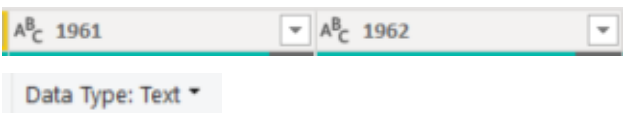
| | |
|---|---|
| Dataset | Agricultural Land, Land Area, Forest Area, Annual Deforestation, Forest Loss |
| Before Transformation | Example from Agricultural Land Dataset:<br><br>ᴬᵇᴄ 2017          ▾   ᴬᵇᴄ 2018          ▾<br><br>● Valid          97%    ● Valid          97%<br>● Error          0%     ● Error          0%<br>● Empty          3%     ● Empty          3%<br><br>252 distinct, 245 unique    252 distinct, 245 unique<br>20                         20<br>6521411.4                  6532568.32<br>379100                     379190<br>3612043.2                  3614109.8<br>563974.3                   569524.9 |

| Create New Column | • For Agricultural Land dataset, create a new column by calculating the average value from 1961 to 2018 for different countries. |
|---|---|
| | • For Land Area dataset, create a new column by calculating the average value from 1961 to 2020 for different countries. |
| | • For Forest Area dataset, create a new column by calculating the average value from 1990 to 2020 for different countries. |
| | • For Annual Deforestation dataset, create a new column by calculating the average value from 1990, 2000, 2010 and 2015 for different countries. |
| | • For Forest Loss dataset, create a new column by calculating the average value from 2001 to 2020 for different countries. |
| | Example from Agricultural Land Dataset:<br> |
| Explanation | The average of values across the year for different countries will be used to fill in the null or missing values respectively in different datasets. Besides, the average values for each country can also be used in the following analysis. |
| Results | Example from Agricultural Land Dataset: |

### 5.3.3 Unpivot Columns

| Dataset | Agricultural Land, Land Area, Forest Area, Annual Deforestation, Forest Loss |
|---|---|
| Before Transformation | Example from Agricultural Land Dataset:  |
| Unpivot Columns, Rename the columns | • For Agricultural Land dataset, there are 58 columns (1961 – 2018) are selected to unpivot. <br><br> • For Land Area dataset, there are 60 columns (1961 – 2020) are selected to unpivot. <br><br> • For Forest Area dataset, there are 31 columns (1990 – 2020) are selected to unpivot. <br><br> • For Annual Deforestation dataset, there are 4 columns (1990, 2000, 2010 and 2015) are selected to unpivot. |

|  | |
|---|---|
| | • For Forest Loss dataset, there are 20 columns (2001 - 2020) are selected to unpivot.
• Rename the "Attribute" as "Year", "Value" as "Agricultural Land Size (sq.km)" and "Average" as "Average Agricultural Land Size (sq.km)".
• Rename the "Attribute" as "Year", "Value" as "Land Area (sq.km)" and "Average" as "Average Land Area (sq.km)".
• Rename the "Attribute" as "Year", "Value" as "Forest Area (sq.km)" and "Average" as "Average Forest Area (sq.km)".
• Rename the "Attribute" as "Year", "Value" as "Annual Deforestation (ha)" and "Average" as "Average Annual Deforestation (ha)".
• Rename the "Attribute" as "Year", "Value" as "Forest Loss (ha)" and "Average" as "Average Forest Loss (ha)".

Example from Agricultural Land Dataset:

 |
| **Explanation** | When the columns are all unpivot, the values are now transformed from columns into rows which allows us to evaluate and generate charts or graphs more easily. When unpivot columns is done, the rows that have 100 % null values will not include in the new column. Therefore, all the columns have no null values now as others null values is replaced by the average value across the years for different countries. Besides, the two columns generated is renamed as the column name provides us a better understanding on the data. |
| **Results** | Example from Agricultural Land Dataset:

 |

## 5.3.4 Pivot Columns

| | |
|---|---|
| Dataset | Annual Deforestation, Forest Loss |
| Before Transformation | Example from Annual Deforestation Dataset:<br><br> |
| Pivot Columns | • For Annual Deforestation Dataset, the values in "Year" column are used to create new columns and the values of "Deforestation" are used as values in the columns.<br><br>Example from Annual Deforestation Dataset:<br><br> |
| Explanation | Pivot columns is required to calculate the average of the value across the years for different countries in order to replace the null values in the datasets. |
| Results | Example from Annual Deforestation Dataset:<br><br> |

## 5.3.5 Convert Unit of Measurement

| Dataset | Annual Deforestation, Forest Loss, Global Forest Loss by Dominant Drivers |
|---|---|
| Before Transformation | Example from Annual Deforestation Dataset:<br><br> |
| Convert Unit of Measurement | <ul><li>For Annual Deforestation dataset, convert the values from hectares (ha) to square kilometre (sq.km) in selected columns (Annual Deforestation (ha), Average Annual Deforestation (ha)).</li><li>For Forest Loss dataset, convert the values from hectares (ha) to square kilometre (sq.km) in selected columns (Forest Loss (ha), Average Forest Loss (ha)).</li><li>For Global Forest Loss by Dominant Drivers dataset, convert the values from hectares (ha) to square kilometre (sq.km) in selected columns (Forest Loss (ha)).</li><li>The formula used to convert hectares (ha) to square kilometre (sq.km) is [sq.km = ha/100].</li></ul><br>Example from Annual Deforestation Dataset:<br><br> |

| | |
|---|---|
| |  |
| Explanation | The unit of measurement is converted from hectares to square kilometre to standardize the dataset and avoid bias in the analysis and modelling. The result of analysis or modelling will inaccurate if the unit of measurement of the values in dataset is different. |
| Results | Example from Annual Deforestation Dataset:  |

### 5.3.6 Merge Datasets

| | |
|---|---|
| Dataset | Forest Loss, iso_metadata, |
| Before Transformation | Example from Forest Loss Dataset:  |

| | |
|---|---|
| | Example from iso_metadata dataset:  |
| Merge Dataset | • Inner join the Forest Loss dataset and iso_metadata dataset with a matching column, "iso". <br><br> Example from Forest Loss & iso_metadata Dataset:  |
| Explanation | Both datasets have a matching column which allows merging together so that the following analysis would be easier as the variables are now stored in the same table. |
| Results | Example from Forest Loss Dataset: |

The five cleaned and transformed datasets which are Land Area, Agricultural Land, Forest Area, Annual Deforestation and Forest Loss datasets are ready to merge and integrate in order to provide an overview on all variable in one table. This allows us to perform more efficient searches and the following analysis would be much easier to carry out. Firstly, Land Area dataset is merged with Agricultural Land dataset based on "Country Name" and "Year" column as illustrated in Figure 5.3.6(a) and Figure 5.3.6(c). After that, the columns like "Agricultural Land (sq.km)" and "Average Agricultural Land (sq.km)" from Agricultural Land (sq.km) is expanded in the table.



*Figure 5.3.6(a): GUI for Merging Land Area Dataset and Agricultural Land Dataset*



```
= Table.NestedJoin(#"Land Area", {"Country Name", "Year"}, #"Agricultural Land", {"Country Name", "Year"}, "Agricultural Land", JoinKind.LeftOuter)
```

*Figure 5.3.6(b): Code for Merging Land Area Dataset and Agricultural Land Dataset*

After that, other datasets are merged by using the method shown above. As we are using left join for merging datasets, there are some null values or missing values exists in numerical columns. Therefore, replacing the null values with "N/A" would be good approach to avoid null values in the analysis as average and median could not be aggregated from the columns.

## 5.4 Data Schema

# 6.0 Prediction & Modelling Techniques

## 6.1 OLAP Report

**OLAP Report**: Forest Loss and Deforestation against Agricultural, Forest and Land Area by Regions, Countries and Year.

**Attributes**: Land Area (sq.km), Forest Area (sq.km), Agricultural Land Area (sq.km), Average Agricultural Land Area (sq.km), Average Forest Area (sq.km), Average Land Area (sq.km), Forest Loss, Dominant Drivers, Annual Deforestation (sq.km), Regions, Countries, Year.

**Purpose**: OLAP Dashboard allows users to visualize on the summary of forest loss and deforestation across the years for different countries. They are able to compare the forest area, land area and agricultural land area from the visualisation done in order to determine the severity of the forest loss and deforestation. The basic operations such as drill down can be done to view the forest loss, agricultural land area, forest area, annual deforestation by regions, countries and years. The users can have a better understanding and make decisions on mitigating the forest loss and deforestation by viewing the relationship between different attributes based on the visualization in the dashboard.

**Advantages:**

There are several advantages when OLAP is chosen to implement in data analysis to explore and discover understandable trends and patterns. Firstly, OLAP provides multidimensional data representation includes drill down, roll up, slice and dice and pivot which allows users to have a better understanding on the relationships between different attributes. For instance, users can view the more detailed forest loss by regions and also countries by drilling down. Furthermore, OLAP allows users to create measures and filters in order to review the information from a specific aspect. Besides, the processing speed of the information is another advantage of OLAP report. The execution time of a single query is short so that users does not have to spend much time and extract the information in a shorter time.

**Disadvantages:**

OLAP report provides many advantages to users but there are some disadvantages too. First of all, one of the disadvantages of OLAP is high dependency on IT as it requires users to have at least basic IT skills in order to create the report. For some cases, SQL scripts or any programming codes are needed to create OLAP report. Furthermore, it requires a vast amount of data to create an informative and insightful report as there are operations like drill down can be done.

## 6.2 Predictive Modelling

### 6.2.1 Linear Regression Model

In this project, the prediction technique used is Linear Regression model which is a model that shows the relationship between an independent variable and dependent variable and predict the future outcome.

**Attributes**: Agricultural Land (sq.km), Annual Deforestation (sq.km), Regions

By implementing linear regression model, the relationship between agricultural land area (sq.km) and annual deforestation (sq.km) by regions can be determined and the future annual deforestation in different regions will be predicted. The annual deforestation will be predicted by entering a value of agricultural land area. Hence, the users will have a better understanding on the severity of deforestation and how agricultural land area will impact the deforestation in different regions. After that, users are able to make decisions on different solutions to mitigate the deforestation rate with the predicted results.

**Advantages:**

Linear regression model is widely used by different analyst or users as it is very helpful in data analysis or predictive modelling. Firstly, linear regression is user-friendly, and it only requires simple implementation for analysis and modelling. Users are able to train models efficiently and effectively by using a lower computational power compared to other modelling algorithms. Furthermore, overfitting can be avoided by implementing regularization which can be implement easily.

**Disadvantages:**

One of the disadvantages of linear regression model is that the model is sensitive to noisy data or outliers as they might affect the performance or accuracy of the model. Besides, the model will assume that there is a relationship between the independent and dependent variables which is not accurate and not reliable.

## 6.2.2 Correlation Plot

In this project, a correlation plot is utilized to determine and investigate the strength of relationship between variables. There is a numerical measure between two variables known as correlation coefficient shown in the correlation plot. The correlation plot will help us to determine which variable have relation with another so that the two variables can be used to fit into the linear regression model. The correlations of the variables are mapped in the plot as illustrated in Figure 6.2.2(a). The positive and negative correlations can be determined by the colours which is blue to red or the numerical values which is range between -1 to 1. In our correlation plot, the blue colour indicates positive correlation while red colour indicates the negative correlation. Based on the correlation plot, there are several strong correlations includes agricultural land area and annual deforestation, land area and annual deforestation and more. Hence, agricultural land area and annual deforestation is selected for the linear regression model as they have strong correlation and agricultural land area is one of the main possible drivers to deforestation.

## Correlation Plot



*Figure 6.2.2(a): Correlation Plot*

# 7.0 Analysis and Recommendations

## 7.1 OLAP Report

The data that was selected in data selection process related to forest loss and deforestation will be used to create a OLAP dashboard. Hence, the dashboard can be viewed by different users and understand the severity of deforestation.

### 7.1.1 Building OLAP Dashboard

In this project, OLAP Dashboard is created by using Power BI which is an analytical service provided by Microsoft. The dashboard provides several visualizations such as graphs and charts to analyse and extract information from the data. There are some aggregated data are used in the dashboard to allows users to visualize in another aspect.



*Figure7.1.1(a): Main Dashboard on Forest Loss and Deforestation*

The main objective of the main dashboard as illustrated in Figure 7.1.1(a) is to provides insights and descriptive analysis of several variables includes forest loss, annual deforestation, land area, agricultural land area, forest area and the dominant drivers of deforestation. The dashboard consists of 2 cards, a map, an area chart, a line chart, a clustered bar chart, a waterfall chart and a stacked column chart.

*Figure 7.1.1(b): Filter by Year and Country is Provided in the Main Dashboard*



*Figure 7.1.1(c): Filter Result on the Dashboard*

Besides, the main dashboard provides a filter function shown in Figure 7.1.1(b) where users can filter the graphs and charts by year and countries. After choosing a specific year or country, the dashboard will be sliced into the selected year or country. For instance, "Argentina" is selected as the country on the filter panel and the dashboard will show the related information of Argentina on the main dashboard as illustrated in Figure 7.1.1(c).

*Figure7.1.1(d): Example of Drilling Down in the Main Dashboard*

Furthermore, the main dashboard can be drilled down into different layers directly from the graphs or charts. For instance, users can drill down into a specific region which is "Western Europe" to visualize different information in that specific region as shown in Figure 7.1.1(d). The others graphs and charts will also show the information of the specific region automatically.



*Figure 7.1.1(e): 2 Cards Showing Land Area and Forest Area in sq.km*

The cards illustrated in Figure 7.1.1(e) is showing the total land area and forest area in square kilometer (sq.km). When there is any operation done such as drilling down or slicing, the value of the cards will turn accordingly. Users can view the value of the cards immediately as the land area and forest area will provide insights in determining the difference between the land area and the forest area in the region or country.

*Figure 7.1.1(f): Map that shows the Average Land Area, Agricultural Land Area & Forest Area (sq.km) with the locations.*

As shown in Figure 7.1.1(f), the regions exist in our dataset are plotted in the map to provide an overview on where is the regions located for users. The country will be plotted when the map is drilled down into a specific region. For instance, countries like Sweden, Norway, Belgium and many more will be shown if the map is drilled down into "Western Europe". The values of average land area, agricultural land area and forest area will appear when cursor is pointing at the region or country.



*Figure 7.1.1(g): Area Chart Showing Land Area and Forest Area by Country*



*Figure7.1.1(h): Area Chart After Drilled Down*

The area chart in Figure 7.1.1(g) shows the land area (sq.km) which is in brown colour and forest area (sq.km) which is in blue colour of different countries. This chart allows users to compare the land area and forest area of different countries visually so that the relationship between them can be determined easily. The values can be viewed by users when the cursor is pointing in the area chart. The chart can be drilled down into different years of different countries. For example, users can view the land area and forest area of a country such as "China" from year 1990 to 2020 as illustrated in Figure 7.1.1(h).



*Figure 7.1.1(i): Line Graph Showing Agricultural Land Area and Forest Area (sq.km) by Year*

The line graph illustrated in Figure 7.1.1(i) demonstrates the relationship between agricultural land area and forest area (sq.km) across the years. The yellow line indicates the agricultural land area whereas the blue line indicates the forest area. For instance, we can see that the forest area is declining when the agricultural land area is increasing between the year 2000 and 2020. Hence, users can analyse some meaningful patterns and extract information from the graph above.



*Figure 7.1.1(j): Clustered Bar Chart Showing Annual Deforestation by Region*

*Figure 7.1.1(k): Clustered Bar Graph After Drilled Down*



*Figure7.1.1(l): Clustered Bar Graph After Second Drilled Down*

As illustrated in Figure 7.1.1(j), the clustered bar chart shows the annual deforestation for different regions. Hence, users can identify which region has the highest annual deforestation by looking at the chart. In this chart, users can also drill down the chart the show a more detailed information by countries. For example, when users drill down into a region such as "Latin America", the information of the countries in Latin America such as "Brazil", "Mexico" and many more will show up as illustrated in Figure 7.1.1(k). After that, users can drill down one more layer which is year in order to visualize the deforestation across years for the country as shown in Figure 7.1.1(l). Therefore, a more detailed analysis can be carried out by the users.



*Figure 7.1.1(m): Waterfall Chart Showing Forest Loss (sq.km) by region*

*Figure 7.1.1(n): Waterfall Chart After Drilling Down*

The waterfall chart as shown in Figure 7.1.1(m) demonstrates the forest loss (sq,km) in different regions. The green bar indicates that the forest loss has increased while the red bar indicates that the forest loss has decreased, and the total amount of forest loss is illustrated with a blue bar. The chart can also be drilled down into different countries within a particular region in order to analyse the forest loss across the countries which is located in the same region. For example, the forest loss of countries such as "Brazil", "Argentina" and more will be shown when "Latin America" has been drilled down.



*Figure 7.1.1(o): Forest Loss by Dominant Drivers*



*Figure 7.11(p): Stacked Column Chart After Drilled Down*

Based on the stacked column chart shown in Figure 7.1.1(o), there are 6 dominant drivers includes forestry, commodity driven deforestation, wildfire, shifting agriculture, unknown and

urbanization are used in the chart. This allows users to have an overview and determine which dominant drivers causes the most forest loss globally. Besides, the users can drill down into each dominant drivers to view the forest loss caused by the specific dominant driver across the years. For instance, Figure 7.1.1(p) demonstrates the forest loss caused by 'Forestry" globally across the years.

## 7.1.2 Analysis on OLAP Dashboard

After building OLAP dashboard, it provides a ton of insights and information on deforestation and forest loss which matches the scope of the project. The visualizations gained from the OLAP dashboard will be analysed to extract meaningful information which will help in decision making and achieve the scope of the project. Firstly, the countries are having a larger land area compared to the forest area in an overall basis as illustrated in Figure 7.1.1(g). After drilling down as illustrated in 7.1.1(h), the forest area in most of the countries increases between the year 1990 to 2000 whereas it decreases after 2000 year. This indicates that the forest area in most of the countries are reducing when the forestry in 1990 to 2000 is not as developed as now. Before 2000s, the global development is slow compared to the current situation, there is still a large amount of land that has not been developed. Hence, there is no need to remove or develop the forest area. However, the global development after 2000s is growing rapidly and the land area that can be developed in some countries are getting scarce and the forestry has a rising demand. Therefore, the forest area in some countries is declining after 2000s as the deforestation has increases as illustrated in Figure 7.1.1(l).

Besides, the agricultural land area increases while the forest area decreases from year 2000 to 2020 as shown in the Figure 7.1.1(i). We can know that agricultural land may be one of the possible drivers of deforestation based on the graph. This is because when the demand of agricultural industry increases, the use of agricultural land will expand to the forest area but not land use area as the land use area are usually utilized for urbanization purpose. Furthermore, we can see that the forest area in Brazil is larger compared to other countries in Latin America and the Brazil has the highest annual deforestation in Latin America by slicing down to Brazil. Based on the scenario, we can assume that countries that have larger forest area may have a higher deforestation rate. When the country has a large forest area as their natural resources, the government or organisation of the country might want to fully utilize the natural resources

to increase the economy of the country. However, there might be other factors or possible drivers such as different policy of countries will affect the annual deforestation of the country.

Moreover, the forest loss caused by commodity driven deforestation globally is lower before year 2010 in an overall basis. This indicates that the global demand of commodity that drives deforestation from human increases as the time goes on. Therefore, there are more commodity that requires the resources from forest which will increase the forest loss globally. With this information, the public or users will know that the deforestation will increases time by time as the global development is still growing rapidly, Besides, they will also understand that the expansion of agricultural land and commodity driven deforestation might have a large impact on the forest loss. The countries that are having large forest area should be alert on the deforestation rate of their countries. Lastly, the users or public should take action on solutions based on forest loss and deforestation that has been analysed and visualized according to locations and time.

## 7.2 Linear Regression Modelling

The data that was selected from the data selection process related to the annual deforestation and its possible drivers will be used for linear regression model. Based on the correlation plot done previously, agricultural land area (sq.km) and annual deforestation (sq.km) will be utilized in the linear regression model.

### 7.2.1 Building Linear Regression Model

In this project, linear regression model is created and developed by using Power BI. The relationship between agricultural land area (sq.km) and annual deforestation (sq.km) will be modelled by fitting into linear equation which is y = mx + b. In order to fit the linear equation, some values should be calculated by using different formula using DAX in Power BI.

```
. xsq = 'Final Dataset'[Agricultural Land.Agricultural Land Area (sq.km)]^2
```

*Figure 7.2.1(a): Calculate x square*

In our case, x represents the agricultural land area whereas y represents the annual deforestation. Agricultural land area is selected as it is one of the possible drivers of deforestation. Besides, the annual deforestation is the variable chosen to be predicted as it provides us the future deforestation which meets the scope of this project. Firstly, x square is calculated as it is required for the following calculations to determine the scope and b-intercept.

```
xy = 'Final Dataset'[Agricultural Land.Agricultural Land Area (sq.km)] * 'Final Dataset'[Annual Deforestation.Annual Deforestation (sq.km)]
```

*Figure 7.2.1(b): Calculate x * y*

```
n = COUNTROWS('Final Dataset')
```

*Figure 7.2.1(c): Count Rows of the Dataset Used*

```
xysum = SUM('Final Dataset'[xy])
```

*Figure 7.2.1(d): Sum of xy*

```
xsum =
SUM('Final Dataset'[Agricultural Land.Agricultural Land Area (sq.km)])
```

*Figure 7.2.1(e): Sum of x*

```
ysum =
SUM('Final Dataset'[Annual Deforestation.Annual Deforestation (sq.km)])
```

*Figure 7.2.1(f): Sum of y*

```
xsqrsum =
SUM('Final Dataset'[xsq])
```

*Figure 7.2.1(g): Sum of xsq*

After that, x times y which represents agricultural land area times annual deforestation is also required for the following calculations as illustrated in Figure 7.2.1(b). Besides, the total number of rows are needed to calculate as it will be used in calculating slope and b-intercept formula as illustrated in Figure 7.2.1(c). The total of xy, total of x and total of y is also be calculated by using the in-built formula provided by Power BI.

```
1  m (Slope) =
2  DIVIDE(
3      [n]*[xysum]-[Xsum]*[Ysum],
4      [n]*[xsqrsum]-[Xsum]^2,
5      0
6  )
```

*Figure 7.2.1(h): Calculate slope*

Then, the slope is calculated as illustrated in Figure 7.2.1(h) by using the formula and the measures that are calculated previously as they are needed to fit into the linear equation which is y = mx + b.

```
1  b (Intercept) =
2  DIVIDE([ysum]*[xsqrsum]-[Xsum]*[xysum], [n]*[xsqrsum]-[Xsum]^2,
3      0
4  )
```

*Figure 7.2.1(i): Calculate b-intercept*

Next, b-intercept is calculated too by using the formula as shown in Figure 7.2.1(i).

```
1  Predicted Annual Deforestation =
2  ([m (Slope)]*
3  'Agricultural Land Area (sq.km)'[Agricultural Land Area (sq.km) Value]+
4  [b (Intercept)]
5  )
```

*Figure 7.2.1 (j): Calculate the Predicted Annual Deforestation*

*Figure 7.2.1(k): Card Visualization for Predicted Annual Deforestation*

Lastly, the predicted annual deforestation will be calculated by using the forumla as shown in Figure 7.2.1(j) and the slope and b-intercept calculated previously. The predidcted value will be calculated based on the agriculturcal land area value entered by users as the formula consists of the value of agricultural land area from the what-if parameter. The value of predicted annual deforestation will be shown in a card visualization as illustrated as Figure 7.2.1(k) which helps users to determine the value more easier.



*Figure 7.2.1(l): Visualization of Linear Regression Model*

Linear regression model as shown in Figure 7.2.1(l) is developed to shows the relationship between agricultural land area and annual deforestation and the future value of annual deforestation will be predicted based on the value of agricultural land. This model gives an overview on the future annual deforestation value to the users so that they could have a better understanding on the probable future situation and make decision on different solutions.

*Figure 7.2.1(m): What-If Parameter and Slicer*

What-if parameter as shown in Figure 7.2.1(m) is utilized as the assumption or condition of the model in order to determine and view the predicted value. When the value of agricultural land area is entered into the what-if parameter, the predicted annual deforestation will appear accordingly in the card visualization. Besides, regions for the linear regression model can be selected by using the slicer. This is because different regions might have different annual deforestation and agricultural land area which will affect the predicted future value of annual deforestation. Therefore, users could view the predicted value of annual deforestation for different regions in order to extract the more accurate results so that different approaches can be made for different regions.

## 7.2.2 Analysis of Linear Regression Model

From the linear regression model above, there are several visualization and prediction results based on the regions. The visualization and predicted results will be analysed in order to meet the scope and objectives of this project. Firstly, the linear regression model for "Latin America & Carib" as illustrated in Figure 7.2.1(l) is showing that it is positive correlation as it has a positive gradient. This indicates that when the agricultural land area increases, the annual deforestation also increases. Hence, we can know that expansion of agricultural land area most probably is one of main drivers of deforestation based on the plotted graph. Besides, we can know that if the agricultural land area is around 0.5 million, the predicted annual deforestation

value will be 5.32k. Based on the prediction result, users are able to know the probable value of annual deforestation so that actions can be made to mitigate or reduce the annual deforestation. Users can also have a better understanding on the severity on annual deforestation which may caused by the expansion of agricultural land area according to the prediction results. Based on the linear regression model, when the value is entered into the what-if parameter which is the agricultural land area increases, the value of annual deforestation will also increase. This indicates that the expansion of agricultural land area is highly affecting the annual deforestation rate. In addition, the agricultural is expanding as the need of global food and energy is increasing rapidly due to population growth of the world. As some of the countries have reserved land area for future development, the agricultural land can only be expanded to the forest area which leads to deforestation.



*Figure 7.2.2(a): Visualization on Northern America*

For Northern America linear regression model as illustrated in Figure 7.2.2(a), the relationship between agricultural land area and annual deforestation is positive correlation where the gradient is positive. Hence, we know that the expansion of agricultural land area is also affecting the annual deforestation in Northern America based on the graph. Both Northern America and Latin America might have the same issues as the gradient showed in the model is almost the same which is also increasing linearly. In addition, if the agricultural land area is around 600k, the predicted annual deforestation will be around 462 sq.km. Therefore, dedicated

organisation or department can limit the agricultural land area to be expanded to mitigate the annual deforestation by using the predicted results.

## 7.3 Recommendations

### 7.3.1 Social Issues and Ethical Issues

There are a possibility that social and ethical issues will be raised when data mining is carried out. During the process of the project, we utilized the datasets that engage with the public includes the agricultural land area datasets and so on which may raise some social and ethical issues. Firstly, privacy is one of the main social and ethical issues arises in research or data mining process. For instance, the data collected should always stored in a protected and secure location to avoid data breach and malicious attacks which loses the confidentiality of the data. In our case, all data includes agricultural land area, annual deforestation and many more should protect and handle with care in the data preparation process like data cleaning and data transformation.

Furthermore, another social and ethical issue is that getting permission from the owner or sources before processing the data. In our project, all datasets have the permission to be collected for analysis and modelling as the data collected from the sources are open -source or available for public download. Hence, the possible social and ethical issue can be avoided if actions are taken carefully before conducting the analysis and modelling.

# 8.0 Conclusion

As deforestation is one of the most concerning global environmental issues, the analysis and modelling in the project will help to determine the relationship between the possible driver and deforestation. Based on the project, agricultural land area is considered as one of the possible drivers that have a large impact on deforestation and forest loss. The objectives and scopes identified earlier are achieved successfully. The public or dedicated department could propose solutions and make decisions in order to mitigate the forest loss and deforestation.

## 8.1 Workload Matrix Report

| Name & TP Number | Tasks |
|---|---|
| Chan Ming Li – TP060774 | • Linear Regression Predictive Model<br>• OLAP |
| Cheong Jia Yao – TP060827 | • Time-series Predictive Model<br>• OLAP |
| Emily Kuan Yu Xuan – TP060808 | • Clustering Model<br>• OLAP |
| Praveenaraaj Sangar – TP061061 | • Distribution Graph<br>• OLAP |

## 8.2 Personal Reflection Report

In a nutshell, I have achieved some learning outcomes by working on this project such as data cleaning, data mining techniques, knowledge discovery and more. Firstly, we have to identify the objectives and scopes based on the domain chosen so that analysis can be conducted well. After identifying the objectives and scopes, data analytics methodologies are researched, and we have selected Fayyad's KDD as our group's methodology. I have learned what is Fayyad's KDD methodology and its phases includes data selection, data prepressing, data transformation, data mining, interpretation or evaluation and applying knowledge gained. The methodology phases are implemented in our group's project as well. Data understanding, selection, cleaning and transformation has done by using the techniques in Power BI taught in class and also online resources. I have used power query featured by Power BI to clean and transform the null and noisy data in order to improve the data integrity and improve the performance of the analysis.

Moreover, I have learned and implemented OLAP Dashboard and Linear Regression Model in Power BI to achieve my objectives and scopes. After the models are developed, analysis is carried out in order to extract different information which will help users to have a better understanding on forest loss and deforestation and make decisions on the solutions which will mitigate the deforestation rate. In addition, I have researched the social or ethical issues that may faced in our case. In short, the objectives, scopes, learning outcomes are achieved by implementing data analytical skills that I have learned.

# 9.0 References

Arciniega, D. (17 9, 2021). *How to do Simple Linear Regression in Power BI*. Retrieved from
        IterationInsights: https://iterationinsights.com/article/linear-regression-in-power-bi/

Brownlee, J. (25 3, 2016). *Linear Regression for Machine Learning*. Retrieved from Machine Learning
        Mastery: https://machinelearningmastery.com/linear-regression-for-machine-learning/

Data Science Process Alliance. (2021). *KDD and Data Mining*. Retrieved from Data Science Process
        Alliance: https://www.datascience-pm.com/kdd-and-data-mining/

O'Connor, E. (26 5, 2021). *CAN I DO DATA CLEANING IN POWER BI?* Retrieved from EPCGroup:
        https://www.epcgroup.net/can-i-do-data-cleaning-in-power-bi/

Rinkesh. (n.d.). *What is Habitat Loss and Destruction*. Retrieved from Conserve Energy Future:
        https://www.conserve-energy-future.com/causes-effects-solutions-for-habitat-loss-and-
        destruction.php

Sharma, R. (23 11, 2020). *KDD Process in Data Mining: What You Need To Know?* Retrieved from
        upGrad: https://www.upgrad.com/blog/kdd-process-data-
        mining/#What_is_KDD_in_Data_Mining

Usama Fayyad, G. P.-S. (1996). From Data Mining to Knowledge Discovery in Databases. *American
        Association for Artificial Intelligence*, 37 - 47.

Vedantu Learn LIVE Online. (8 September, 2021).
        *(https://www.vedantu.com/english/environmental-issue-essay)*. Retrieved from
        (https://www.vedantu.com: (https://www.vedantu.com/english/environmental-issue-essay)

World Health Organization. (2021). *How air pollution is destroying our health*. Retrieved from World
        Health Organization: https://www.who.int/news-room/spotlight/how-air-pollution-is-
        destroying-our-health