# Effective Methods for Capturing Cattle Rustlers

**Yung-Hsin Chen[1] and Haoxin Cai[1]**

[1]*Universität Zürich, Zürich, Switzerland*

19 December 2022

The COVID-19 pandemic has severe impacts on almost every aspect around the world. It causes not only social and economic disruption, but also drastic death rates. Inevitably, people are curious about the causalities of COVID-19 and how to prevent the pandemic from getting worse. In this report, data are collected from 192 countries, and classification models are applied in order to get a list of feature importance. It is believed that the more-important-features play more crucial role in the severity of the pandemic of a certain country.

## 1 Introduction

The goal of the report is to get the most related factors of COVID-19 pandemic cases in countries. The time parameters are taken out of consideration for simplicity. (write more...)

## 2 Domain Knowledge

## 3 Method

In this section, the methods of this particular task will be explained, including introduction of data, building features for the models, training and predicting classifiers and generating the feature importance table. Classifiers and other used tools are introduced and explained in section 2.

### 3.1 Data

Data used for this task is from covid-19-data from Our World of Data. It is online in a github repository, mak-

**Table 1:** *Data Summary Table*

| Information | Description |
|---|---:|
| number of columns | 67 |
| number of rows | 236386 |
| number of countries | 248 |
| key | {location:date} |
| starting date | 01.Jan.2020 |

ing it easy to access. The data is loaded into the local database by the *request* package of Python.

The data consists of 67 columns and 248 countries with 236386 rows in total. Attribute *location* and *date* together define each unique row of data. The data is still being updated. A summary table of the data is shown in Table 1.

### 3.2 Building Features

The process of building features includes data cleaning and feature selection, and label preparation. Data cleaning and feature selecting is crucial for the accuracy of models. Bad data cleaning can lead to biased results or bad model performances. Relevant attributes will be selected as features to be put into the classifiers, i.e., the models. Finally, since this is a supervised learning task, the label should be prepared for model training.

In the data cleaning phase, the goal is to get a table of one row per country, i.e., each row represents the information of a country. The label is defined as the attribute, *total_cases_per_million*. To achieve this, relavent attributes are first selected for model

**Table 2:** *Data Summary Table*

| Selected Attributes | |
| --- | --- |
| aged_65_older | cardiovasc_death_rate |
| aged_70_older | diabetes_prevalence |
| gdp_per_capita | hospital_beds_per_thousand |
| population_density | human_development_index |
| life_expectancy | median_age |
| people_fully_vaccinated_per_hundred | |

**Table 3:** *Label Interval for Label Preparation*

| Level | Interval |
| --- | --- |
| 0 | 0 - 50'000 |
| 1 | 50'000 - 200'000 |
| 2 | 200'000 - 400'000 |
| 3 | 400'000 - 700'000 |

**Table 4:** *Model Summary*

| Model | Details |
| --- | --- |
| **logistic regression** | |
| penalty | l2 |
| solver | newton-cg |
| **Linear Perceptron** | |
| tolerance | 0.001 |
| penalty | l2 |
| **XGBoost** | |
| learning rate | 0.1 |
| loss | deviance |
| max depth | 3 |
| n_estimators | 100 |
| random state | 21 |

training. For this task, eleven attributes that is speculated to affect the number of COVID cases are selected from the raw data. The selected attributes are listed in **??**. Countries with less than 200 rows of data is then removed. Since each row corresponds to one date, it is not ideal to use the data of a country if less than 200 days of data are recorded. Among the selected attributes, *aged_65_older* and *aged_70_older* are divided by *population* into *aged_65_older_percentage* and *aged_70_older_percentage* respectively so that the models are trained on the percentage of elder people instead of the total number. Except for the attribute *people_fully_vaccinated_per_hundred*, all the other attributes have a single value throughout the dates for each country. However, the attribute *people_fully_vaccinated_per_hundred* and the label attribute *total_cases_per_million* is accumulated day by day. In this case, the data of the latest date is used as the feature value for each country. By doing this, the model will be able to generate the feature importance table according to how all features affect the number of total cases per million in each country. After data cleaning and feature selection, 194 countries/rows and 11 features are left for model training.

After data cleaning and feature selection, label preparation is performed. The attribute *total_cases_per_million* is categorised into four levels of severity. The interval of the categorisation is shown in **??**. Apparently, there is no country with over 700'000 cases per million. [0, 50000, 200000, 400000, 700000]

### 3.3 Classifiers

Classifiers are used for classification tasks. Due to the small dimension of the dataset, several changes are made to adapt the data size. The data is not splitted into training, validation and testing datasets, but only training and testing only. In this task, the training-testing split will be 80% and 20% respectively. Besides, it is not recommended to use models with too many parameters since it might have to higher chance of overfitting. Thus, simple models are picked out for this task including logistic regression, linear perceptron and XGBoost. These models can be easily applied to the dataset with *sklearn* from Python. For hyperparameter selection, grid search cross validation is used. A summary of models used for the task and the best performing hyperparameters chosen by applying grid search cross validation are listed out in **??**.

### 3.4 Feature Importance Table

The feature importance is generated automatically via the trained model. It shows weight of each feature. The weight can be thought of as how significant each feature effects the classification accuracy.

## 4 Results

## 5 Dicussion

## 6 Conclusion