

---

# Effective Methods for Capturing Cattle Rustlers

Yung-Hsin Chen<sup>1</sup> and Haoxin Cai<sup>1</sup>

<sup>1</sup> Universität Zürich, Zürich, Switzerland

---

19 December 2022

**T**he goal of the report is to get the most related factors of COVID-19 pandemic cases in countries. 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. With an accuracy of 74.36%, XGBoost model suggests that human development index, life expectancy, population density, hospital beds per thousand and GDP per capita of a country are the leading factors of the severity of the pandemic.

## 1 Introduction

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. Therefore, the report aims to determine the correlation between the severity of the pandemic other information of a country. Notes that the time parameters are taken out of consideration for simplicity.

## 2 Domain Knowledge

In this section, tools and domain knowledges used in the task and reasons of choosing them will be briefly explained including models used (logistic regression, linear perceptron, XGBoost), evaluation metrics (micro-f1, macro-f1) and visualisations (confusion matrix, normalised confusion matrix).

### 2.1 Models

The models used in the task are linear perceptron, logistic regression and XGBoost. The three models are chosen to adapt the small data size of the COVID data. All models have limited parameters and are thus less likely to overfitting, which can happen easily with small datasets. The linear perceptron is the simplest model and will serve as the base-line model.

#### 2.1.1 Linear Perceptron

Linear Perceptron(Sharma, [n.d.](#))(Ghodsi, [n.d.](#)) is a very simple model for binary classification, which can later be adapted to multi-class classification repeating the process with paired up the classes. The architecture of linear perceptron is straightforward. After each data are multiplied by weights and biases, the result (a scalars) will be passed to the activation function. In the case of unit-step function as the activation function, numbers larger than 0 will be assigned to a class, and those smaller than 0 will be assigned to the other class. The whole process can be summarised by [Equation 1](#).

$$y_{\text{predicted}} = \text{unit-step}(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots) \quad (1)$$

The output with a list of predicted classes will be compared with the actual classes by a loss function. The weights and biases will then be updated by applying stochastic gradient descent(Roy, [n.d.](#)) to the loss function. The perceptron is guaranteed to converge if the classes is linearly separable, but has the problem of sticking in the local minima.

Although linear perceptron is fast to implement and is less inclined to overfitting with small datasets, [Equation 1](#) gives away the linearity decision surface nature of logistic regression. It also indicates the no-multicollinearity requirements between independent

variables and the high sensitivity to outliers like linear regression does. Besides, it only takes on yes and no outputs and does not consider calibrated probabilities.

### 2.1.2 Logistic Regression

Logistic regression (Addagatla, n.d.) (Ahmed, n.d.) (Penumudy, n.d.) is a linear model commonly used for classification problems as well. It is very similar to the linear perceptron model except for the activation function. The activation function for logistic regression is Sigmoid function. For logistic regression, regularisation terms are also added to the loss function as a penalty for overfitting.

With the Sigmoid function in mind, it is obvious that logistic regression is mostly used for binary classification and can be extended to multi-class classification like linear perceptron.

Similar to linear perceptron, logistic regression is also sensitive to outliers and does not allow multicollinearity. However, it no longer takes on only yes or no as outputs because the real world is not likely a binary answer world. Instead, it allows continuous probabilities for outputs, which is fitted by the Sigmoid function. With probability over 0.5, the output is assigned to a class, and with probability lower than 0.5, the output will be assigned to the other class. The regularisation terms also prevent the model from overfitting (Geeks, n.d.).

### 2.1.3 XGBoost

XGBoost (Harode, n.d.) (Science, n.d.) is a gradient boosting model with extra structures that has been proved well performed on classification tasks. The boosting model creates weak learner (decision tree) one by one. The previous weak learner set a larger weight for the wrong answers as the input to the next weak learner. By going through the sequence, the classification will eventually be finalised correctly. Each weak learner will have their own predictions, and together they will have to vote to decide what the final answer is. The less accurate weak learners get less votes.

XGBoost has the advantage of tackling non-linear problems and does not require non-multicollinearity on data. It is also very efficient due to the parallel tree boosting and less parameters to be trained comparing to deep neural networks. Although the boosting models are greedy, resulting in higher chance of overfitting, setting the correct maximum depth of the trees can effectively avoid this problem.

## 2.2 Evaluation

Since it is a classification task, f1-score (Leung, n.d.) is commonly used for model evaluation. However, f1-score is only calculated per class. The multi-class case will require the aggregation of f1-scores of each class to determine the overall performance of the models on all classes. Micro-f1 and macro-f1 are two different ways of aggregating the f1-scores of each class and will be used as the evaluation metrics of this task. In addition, the confusion matrix and the normalised confusion matrix will be used as the visualisation of the evaluation.

### 2.2.1 Metrics

After training, the model will be tested on the testing dataset. The performance of the models will then be tested by the evaluation metrics. Metrics used for this task are micro-f1 and macro-f1.

F1 score is a metric of taking precision<sup>1</sup> and recall<sup>2</sup> into consideration at the same time per class. F1-score is defined as Equation 2.

$$\begin{aligned} \text{f1-score} &= 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \\ &= \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FP} + \text{FN}} \end{aligned} \quad (2)$$

Since f1-score is calculated per class, to calculate the aggregation of multi-class will become tricky. This is where micro-f1 and macro-f1 come into play. They are two different ways of aggregating multi-class f1-scores. Macro-f1 calculates the average of f1-scores of all classes as Equation 3. The equation shows that macro-f1 gives same weights to each class no matter how large the class is.

$$\text{macro-f1} = \frac{\text{sum(f1-score)}}{\text{number of classes}} \quad (3)$$

On the other hand, micro-f1 is defined by Equation 4. This equation is the same as the one for f1-score (Equation 2). However, the TP, FP and FN stands for the sum of the metrics for all classes.

$$\text{micro-f1} = \frac{\text{TP}}{\text{TP} + \frac{1}{2}(\text{FP} + \text{FN})} \quad (4)$$

This shows that f1-score views each observation points equally important. This might cause a bias of measurement with imbalanced datasets. With f1-score, larger classes with more observation points will have a larger impact on the f1-score. A comparison and discussion of the two metrics will be explained more in section 5.

<sup>1</sup>According to Layman definition, precision means, of all the positive predictions I made, how many of them are truly positive?

<sup>2</sup>According to Layman definition, recall means, of all the actual positive examples out there, how many of them did I correctly predict to be positive?

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

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	

### 2.2.2 Visualisation

To visualise how well the model performs, a confusion matrix will be the most suitable tool. The confusion matrix aims to check if the model is classifying correctly or confusing different classes. Each number in the row represents the number of instances of the actual class, while each number in the column is the number of instances of the predicted class. The brighter the squares indicates the more number of instances. Ideally, the diagonal squares should be the brightest. However, some class has fewer data, resulting in darker colour of squares even the data are all classified correctly. Therefore, in the report, both confusion matrix and normalised confusion matrix will be shown. The normalised confusion matrix helps eliminate the unbalanced dataset issue.

## 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<sup>3</sup>, making it easy to access. The data is loaded into the local database (Menestrel, n.d.) by the *request* package of Python. Each time the data is downloaded, the newest data (up-to-date) will be fetched from the github repository. Thus, the result might be slightly different every time. The result of this report is executed with data up to 30. Nov. 2022.

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.

<sup>3</sup>The data can be found here: <https://github.com/owid/covid-19-data/tree/master/public/data>

### 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, relevant attributes are first selected for model training. For this task, eleven attributes that are speculated to affect the number of COVID cases are selected from the raw data. The selected attributes are listed in Table 2. Countries with less than 200 rows of data are 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 Table 3. Apparently, there is no country with over 700'000 cases per million.

**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
<b>Logistic Regression</b>	
penalty	l2
solver	newton-cg
<b>Linear Perceptron</b>	
tolerance	0.001
random state	0
<b>XGBoost</b>	
learning rate	0.1
loss	deviance
max depth	3
n_estimators	100
random state	21

**Table 5:** Model Accuracy Summary

Model	Accuracy	Micro-f1	Macro-f1
Logistic Regression	64.10%	0.6410	0.5863
Linear Perceptron	48.72%	0.4872	0.2976
XGBoost	76.92%	0.7692	0.7749

**Table 6:** Feature Importance from XGBoost

	Feature	Importance
1	human_development_index	0.462551
2	life_expectancy	0.462551
3	population_density	0.138737
4	hospital_beds_per_thousand	0.082238
5	gdp_per_capita	0.065485
6	people_fully_vaccinated_per_hundred	0.059118
7	cardiovasc_death_rate	0.045620
8	aged_70_older_percentage	0.033860
9	median_age	0.029610
12	diabetes_prevalence	0.026468
11	aged_65_older_percentage	0.025800

[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 [Table 4](#).

### 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. The more positive the feature importance score, the more the feature helps reduce the loss while training. However,

if the feature importance is negative, the feature increases the loss while training.

## 4 Results

By applying the test dataset on the model, the performance of the models can be evaluated. The accuracy summary of each model and the feature importance table from the best performing model is shown in [Table 5](#) and [Table 6](#)<sup>4</sup>. Among the models, XGBoost has the best performance. Among the features, the model suggests that human development index<sup>5</sup>, life expectancy<sup>6</sup>, population density, hospital beds per thousand and gdp per capita are the top five items that affect the covid cases per milliom in a country.

## 5 Dicussion & Conclusion

Besides the accuracies and f1-scores of the models, visualisations also help with the understanding of how the models perform. Since XGBoost is the best performing model, the visualisation of it will be shown.

<sup>4</sup>The result can be replicated by taking the data up to 30. Nov. 2022.

<sup>5</sup>According to [Our World In Data](#), human development index is: A composite index measuring average achievement in three basic dimensions of human development—a long and healthy life, knowledge and a decent standard of living. Values for 2019, imported from <http://hdr.undp.org/en/indicators/137506>

<sup>6</sup>According to [Our World In Data](#), Life expectancy is: Life expectancy at birth in 2019

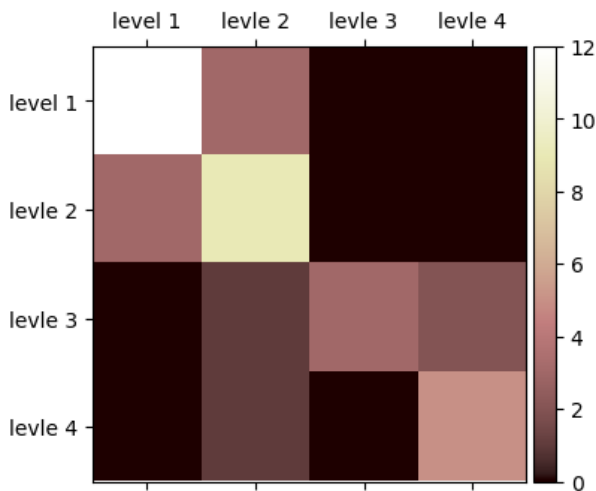


Figure 1: Confusion Matrix of XGBoost Result

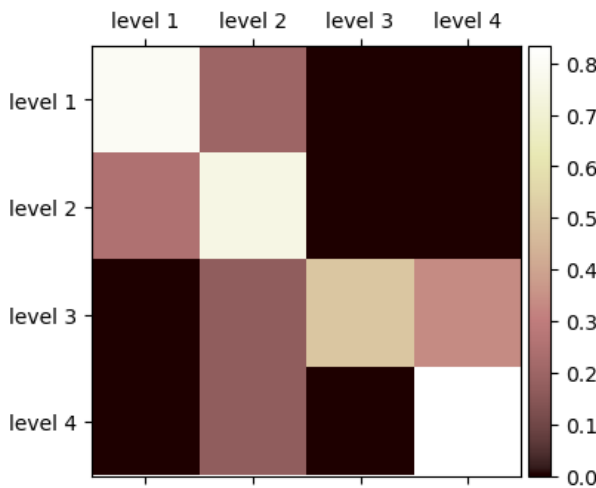


Figure 2: Confusion Matrix of XGBoost Result

The confusion matrix (Figure 1) shows that most classes are classified correctly with the diagonal squares brighter than the others. However, due to the fact that level 4 has fewer data size than the others, it is hard to tell if it actually performs worse than the other levels.

In the normalised confusion matrix (Figure 2), it is clear that the darker square of level 4 is caused by the smaller datasize of the class. The model is actually performing well even on smaller data size classes. That is why the macro-f1 is larger than micro-f1.

From the feature importance table, it can be concluded that the severity of covid is more related to the human development index. The human development index is a score on the quality of living in the country including average life expectanvy, GDP, education opportunity, etc. The number of people fully vaccinated, the population composition or the diabetes prevalence

do not play the most important role in the total cases as expected.

## Bibliography

- Addagatla, Arun (n.d.). *Understanding Logistic Regression*. URL: <https://medium.com/nerd-for-tech/understanding-logistic-regression-782baa868a54>. (accessed: 11.12.2022).
- Ahmed, Mazen (n.d.). *Logistic Regression Explained*. URL: <https://linguisticmaz.medium.com/logistic-regression-explained-1849a0f5ce54>. (accessed: 11.12.2022).
- Geeks, Geek for (n.d.). *Advantages and Disadvantages of Logistic Regression*. URL: <https://www.geeksforgeeks.org/advantages-and-disadvantages-of-logistic-regression/>. (accessed: 11.12.2022).
- Ghods, Ali (n.d.). *Statistical Learning- Classification*. URL: <https://sas.uwaterloo.ca/~aghodsib/courses/f07stat841/notes/lecture10.pdf>. (accessed: 11.12.2022).
- Harode, Rohan (n.d.). *XGBoost: A Deep Dive into Boosting*. URL: <https://medium.com/sfu-csmp/xgboost-a-deep-dive-into-boosting-f06c9c41349>. (accessed: 11.12.2022).
- Leung, Kenneth (n.d.). *Micro, Macro & Weighted Averages of F1 Score, Clearly Explained*. URL: <https://towardsdatascience.com/micro-macro-weighted-averages-of-f1-score-clearly-explained-b603420b292f>. (accessed: 11.12.2022).
- Menestrel, Thomas Le (n.d.). *Loading a csv file from GitHub in Python*. URL: <https://medium.com/towards-entrepreneurship/importing-a-csv-file-from-github-in-a-jupyter-notebook-e2c28e7e74a5>. (accessed: 11.12.2022).
- Penumudy, Tanvi (n.d.). *A Comprehensive Guide to Logistic Regression*. URL: <https://medium.com/analytics-vidhya/a-comprehensive-guide-to-logistic-regression-e0cf04fe738c>. (accessed: 11.12.2022).
- Roy, Abhijit (n.d.). *An Introduction to Gradient Descent and Backpropagation*. URL: <https://towardsdatascience.com/an-introduction-to-gradient-descent-and-backpropagation-81648bdb19b2>. (accessed: 11.12.2022).
- Science, ODSC - Open Data (n.d.). *XGBoost: Enhancement Over Gradient Boosting Machines*. URL: <https://odsc.medium.com/xgboost-enhancement-over-gradient-boosting-machines-73abafa49b14#:~:text=XGBoost%20provides%20a%20parallel%20tree,problems%20beyond%20billions%20of%20examples..> (accessed: 11.12.2022).
- Sharma, Sagar (n.d.). *What the Hell is Perceptron?* URL: <https://towardsdatascience.com/what-the-hell-is-perceptron-626217814f53>. (accessed: 11.12.2022).