# Predicting Percentage Increase in Covid-19 Cases Worldwide using Regression Decision Trees

Utkarsh Patel
18EC30048

Akhil
18CS10070

## Abstract

This report is submitted as part of an assignment of the course CS60050 – Machine Learning, IIT Kharagpur. The dataset contained the time-series percentage increase in COVID-19 cases worldwide. The attributes were date, confirmed cases, recovered cases, number of deaths and percentage increase in number of cases. The target attribute was the percentage increase in cases. It was required to build a regression decision tree for this task. This report briefly explains the procedure we followed to achieve the goal and the results obtained.

## 1 Generating Containers

We first created a class `Record` to stores the records of the dataset in Python virtual environment. The attributes of the class `Record` are: `date`, `confirmed`, `recovered`, `deaths` and `increaseRate`. While other attributes represented the same thing as in dataset, we modified the interpretation of `date` attribute. Instead of representing a string literal, we converted it to number of days since `EPOCH`, where `EPOCH` is a reference date. We used `datetime` library to do the same. Later, the records (which are now `Record` instances) are fused to form a `pandas.DataFrame` object. This is done because it is quite simple to handle data using such objects.



**Fig 1.** Dataset before and after generating containers

## 2 Generating Decision Tree for given `max_depth`

The user is asked to enter a `max_depth`. Taking this as input, we build a decision tree that doesn't exceed this maximum depth. While building the trees of a given maximum depth, the entire dataset is divided into training set, validation set and test set in the ratio 60:20:20. This is done 10 times with independent random splitting. Each tree remembers what was the partition used for its training and testing (and validation, however validation set is not used in this part). For choosing the best attribute for given split, *residual squared* sums of the target concept is employed. The average test accuracy is reported as the test accuracy of these 10 trees, while the best one is chosen for further tasks. To define the accuracy of a decision tree, we used *mean squared error* over the predicted values and actual values.
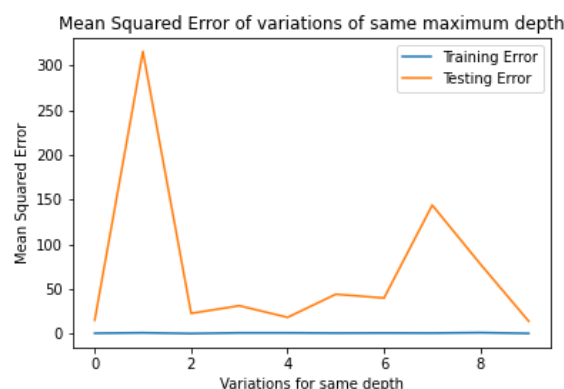


**Fig 2.** Training and testing error observed over different split of the dataset for same maximum depth (`max_depth` = 6 in this case)

## 3 Analysing the Model for Different Values of `max_depth`

In this part, the algorithm used in previous part is used for different values of `max_depth`. The range chosen was 1 to 20. The training and test error for each value of `max_depth` is plotted for analysis.
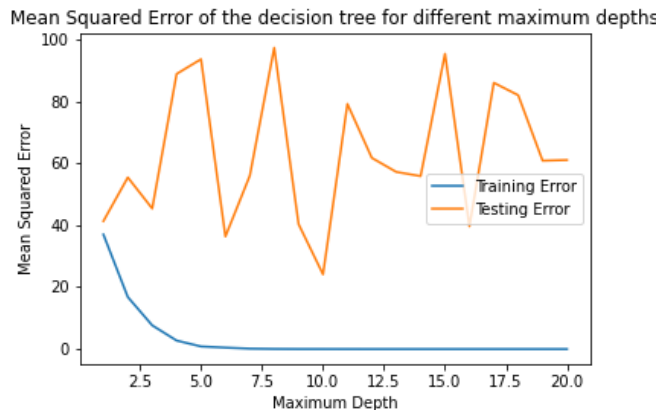


**Fig 3.** Training and testing error for different values of maximum depth

As it can be observed in the plot, training error decreases monotonically as maximum depth of the tree is increased. However, this is not the case with testing error. Hence, we need to restrict what should be the maximum allowed depth of the decision tree. As per the plot, the optimal maximum depth is 10. Beyond that, the testing error is always greater. (This may change if you run the program again as the splitting is done randomly and independently.) Hence, in our case, best possible depth limit is 10.

## 4 Pruning the Best Tree over Validation Set

We use the tree which gave the best results (i.e. its depth is 10). It can be further pruned so as to improve its performance. For this case, the validation set is used. Note that it is not necessary that the pruned tree will perform better than the unpruned one over the test set (as pruning is done only on the validation set). For pruning, we are using *reduced-error pruning* algorithm. It is done in top-down manner recursively. For a given node, it checks whether after pruning this node, the error on validation set reduces or not. If it does, the node is pruned by making it a leaf node and its prediction is set to the mean of training instances which were incident on this node. In our case, the unpruned tree performed better on test set as compared to the pruned one. (But the scenario may be different if the program is executed again.)

**Table 1.** Error of the unpruned and pruned decision tree (*depth* = 10)

| *Pruned* | **Training Error** | **Validation Error** | **Testing Error** |
|---|---|---|---|
| *False* | 0.0065786405928254426 | 31.491522818820886 | 15.459176608103375 |
| *True* | 14.15782870734966 | 20.653508126041466 | 32.293096850435674 |

## 5 Printing the Decision Tree

A recursive approach is used to print the decision tree. Some part of the trees is presented below.



**Fig 4.** Part of unpruned (left) and pruned (right) decision tree (*depth* = 10)