# COVID-19 global forecasting

## Problem Description:

The Covid-19 pandemic has affected the entire world in unprecedented ways, with millions of confirmed cases and fatalities. To tackle this problem, this project aims to predict the number of confirmed cases and fatalities worldwide using machine learning models. The dataset used for this project is taken from Kaggle's Covid-19 Global Forecasting Week 5 competition.

### Objectives:

**The objectives of this project are as follows:**

- To predict the number of confirmed cases and fatalities worldwide.

- To identify the most important factors that contribute to the spread of Covid-19.

- To evaluate the performance of different machine learning models for predicting the number of confirmed cases and fatalities.

### Dataset:

The dataset used for this project is taken from Kaggle's Covid-19 Global Forecasting Week 5 competition. The dataset contains two files: train.csv and test.csv. The train.csv file contains information on confirmed cases and fatalities for different countries and regions. The test.csv file contains information on confirmed cases and fatalities for a specific date range. The dataset also contains information on various factors that could contribute to the spread of Covid-19, such as population density, climate, and social distancing measures.

### Deliverables:

**The deliverables for this project are as follows:**

- A machine learning model that can predict the number of confirmed cases and fatalities worldwide.

- A report on the most important factors that contribute to the spread of Covid-19.

- An evaluation of the performance of different machine learning models for predicting the number of confirmed cases and fatalities.

# Suggested Framework to Solve this problem :

## Data Acquisition and Exploration

- Download and load the necessary dataset(s) from the Kaggle link provided

- Explore the data to gain insights into the dataset, including the number of samples, features, missing values, and data types

- Check for outliers and anomalies in the data

- Visualize the data to gain more insights into the data and identify any patterns or trends that may exist

## 2. Data Preprocessing

- Clean the data by addressing missing values, outliers, and anomalies

- Convert the data into a format suitable for training machine learning models

- Perform feature engineering to create new features from existing ones that may improve model performance

- Split the data into training and testing sets

## 3. Model Selection and Training

- Select appropriate machine learning algorithms for the task of predicting the number of confirmed cases and fatalities worldwide

- Train the selected models using the preprocessed training data

- Tune the hyperparameters of the models to improve performance

## 4. Model Evaluation

- Evaluate the performance of the trained models using appropriate evaluation metrics, such as mean squared error or mean absolute error

- Compare the performance of the different models and select the best performing one

## 5. Model Deployment

- Deploy the trained model to predict the number of confirmed cases and fatalities worldwide

- Create a simple and user-friendly interface to allow users to input new data and obtain predictions

- Continuously monitor the model's performance and update the model as needed

## 6. Communication

- Communicate the findings and insights from the project to relevant stakeholders in a clear and concise manner

- Create visualizations and other forms of media to aid in the communication of the findings

- Provide recommendations for future actions based on the findings of the project

# Code Explanation :

Here is the simple explanation for the code you can find at code.py file.

**Section 1:** Importing necessary libraries and loading the dataset In this section, we imported the necessary libraries such as pandas and numpy for data manipulation and sklearn for data preprocessing and modelling. We also loaded the training and testing datasets provided by Kaggle using the pandas library.

**Section 2:** Data preprocessing In this section, we performed data cleaning, feature engineering, and data transformation on the training and testing datasets. We cleaned the data by filling missing values, dropping irrelevant columns, and converting data types. We then created new features such as month, day, and year from the date column. We also encoded categorical variables using one-hot encoding.

**Section 3:** Modelling In this section, we built machine learning models to predict the confirmed cases and fatalities. We used Random Forest Regressor and Gradient Boosting Regressor algorithms to predict confirmed cases and fatalities respectively. We also used GridSearchCV to tune hyperparameters for better model performance.

**Section 4:** Model evaluation In this section, we evaluated the performance of the models using the mean squared error and mean absolute error metrics.

**Section 5:** Prediction In this section, we used the trained models to predict the confirmed cases and fatalities for the testing dataset.

**Section 6:** Submission In this section, we saved the predicted results in the required format and submitted it to Kaggle.

Overall, the code follows a structured approach of data preprocessing, modelling, evaluation and prediction, and submission for this forecasting problem. It also makes use of popular machine learning algorithms and techniques for feature engineering and hyperparameter

tuning. The code is written in a clean and concise manner, making it easy to read and understand.

# Future Work :

**Future Work:** Predicting Covid-19 Confirmed Cases and Fatalities

The following are some potential future work that can be done to improve the predictive models for Covid-19 confirmed cases and fatalities:

1. **Feature Engineering**

- Explore additional data sources: Incorporate other datasets such as population demographics, mobility data, and health system capacity.

- Create new features: Generate new features that might be relevant to predicting Covid-19 spread and impact, such as hospital capacity, government restrictions, and vaccination rates.

- Apply feature selection techniques: Select the most important features using methods such as correlation analysis, principal component analysis (PCA), and recursive feature elimination (RFE).

2. **Model Tuning and Selection**

- Experiment with different models: Try using different machine learning models such as Random Forest, XGBoost, and Neural Networks, and compare their performance against the current models.

- Optimize model hyperparameters: Use grid search or Bayesian optimization techniques to identify the best hyperparameters for the models.

- Evaluate model ensembles: Combine multiple models using techniques such as bagging, boosting, and stacking, and evaluate their performance.

3. **Time-Series Analysis**

- Incorporate time-series analysis: Explore the use of time-series models such as ARIMA, Prophet, and LSTM, which can capture temporal dependencies and patterns in the data.

- Model seasonal effects: Consider the seasonal variations in Covid-19 spread and incorporate seasonal effects in the models.

- Forecast long-term trends: Use time-series models to forecast Covid-19 cases and fatalities beyond the current time period.

4. **Visualization and Interpretation**

- Develop interactive dashboards: Build interactive dashboards that enable users to explore the data and model predictions visually.

- Explain model predictions: Use techniques such as SHAP values, partial dependence plots (PDP), and LIME to explain the model predictions and identify the key features that contribute to the predictions.

- Communicate results effectively: Use clear and concise language to communicate the findings and insights to stakeholders and the general public.

**Step-by-Step Guide:**

1. Collect additional data sources such as population demographics, mobility data, and health system capacity.

2. Create new features such as hospital capacity, government restrictions, and vaccination rates.

3. Use feature selection techniques such as correlation analysis, principal component analysis (PCA), and recursive feature elimination (RFE) to select the most important features.

4. Experiment with different machine learning models such as Random Forest, XGBoost, and Neural Networks, and optimize their hyperparameters using grid search or Bayesian optimization techniques.

5. Combine multiple models using techniques such as bagging, boosting, and stacking, and evaluate their performance.

6. Use time-series models such as ARIMA, Prophet, and LSTM, to capture temporal dependencies and patterns in the data.

7. Build interactive dashboards that enable users to explore the data and model predictions visually.

8. Use techniques such as SHAP values, partial dependence plots (PDP), and LIME to explain the model predictions and identify the key features that contribute to the predictions.

9. Communicate the findings and insights to stakeholders and the general public using clear and concise language.

# Exercise Questions :

1. **What are the different models that can be used for time-series forecasting? How do they differ?**

- There are various models for time-series forecasting such as ARIMA, SARIMA, LSTM, and Prophet. The difference between them lies in their ability to capture seasonality, trend, and noise in the data. ARIMA models capture linear trends and stationary data, while SARIMA models can handle seasonality in addition to linear trends. LSTM models use recurrent neural networks to capture non-linear dependencies in the data, and Prophet models are designed specifically for time-series forecasting with the ability to capture multiple seasonalities and trend changepoints.

2. **Can you explain the concept of feature engineering in time-series forecasting?**

- Feature engineering involves transforming raw data into useful features that can be used to build predictive models. In the context of time-series forecasting, this involves creating lagged variables or time-based features such as rolling averages, moving averages, and exponential smoothing. Feature engineering is an important step in time-series forecasting because it can help capture seasonality and trend in the data, as well as other factors that might impact the outcome variable.

3. **What are some methods for handling missing data in time-series forecasting?**

- There are several methods for handling missing data in time-series forecasting, such as linear interpolation, last observation carried forward (LOCF), forward or backward filling, and seasonal decomposition. Linear interpolation is a simple method that fills in missing values with a linearly interpolated value based on the neighboring data points. LOCF fills in missing values with the most recent non-missing value, while forward or backward filling fills in missing values with the nearest non-missing value. Seasonal decomposition involves decomposing the time-series into its seasonal, trend, and residual components, and then forecasting each component separately.

**4. Can you explain the difference between point forecasting and probabilistic forecasting?**

• Point forecasting involves predicting a single value for the outcome variable, while probabilistic forecasting involves predicting a probability distribution over the possible values of the outcome variable. Point forecasting is useful when we are interested in a specific value, while probabilistic forecasting is useful when we want to understand the range of possible outcomes and their likelihoods.

**5. How can we evaluate the performance of a time-series forecasting model?**

• There are various metrics that can be used to evaluate the performance of a time-series forecasting model, such as mean absolute error (MAE), mean squared error (MSE), root mean squared error (RMSE), and mean absolute percentage error (MAPE). These metrics can be calculated on the training set and the test set to evaluate the model's ability to fit the data and generalize to new data. Additionally, visual inspection of the predicted values and their confidence intervals can be useful for understanding the model's performance.

# Concept Explanation :

The algorithm we used in this project is the Random Forest algorithm. Now, let's take a closer look at what it is and how it works!

Random Forest is an ensemble learning algorithm used for both classification and regression tasks. Ensemble learning refers to the technique of combining multiple models to improve the overall performance of the model. The idea behind the Random Forest algorithm is to build many decision trees and average their predictions to make more accurate predictions.

Imagine you are trying to decide whether or not to go to the beach. There are several factors that could influence your decision, such as the temperature, wind speed, humidity, and cloud cover. In a decision tree model, you would start with the root node (the top of the tree) and ask a question about one of the factors, such as "Is the temperature above 80 degrees Fahrenheit?" Depending on the answer, you would move down the tree to a new node and ask another question until you reach a leaf node (the bottom of the tree), which represents the final decision. For example, if you answered "yes" to the temperature question, you might move down the tree to a node that asks "Is the humidity below 60%?" and so on until you reach a leaf node that says "Go to the beach!" or "Stay home!"

In a Random Forest model, you would build many decision trees (hence the name "forest") and average their predictions to get a more accurate result. Each tree in the forest is trained on a randomly selected subset of the data and a randomly selected subset of the features (in our case, the weather factors). This randomness helps to prevent overfitting, which is when a model performs well on the training data but poorly on new, unseen data.

Let's say we built 100 decision trees in our Random Forest model. To make a prediction, we would feed the test data into each of the 100 trees and get a prediction from each tree. We would then average the predictions to get the final result.

So, to sum it up: Random Forest is an ensemble learning algorithm that builds many decision trees and averages their predictions to improve accuracy. Each tree is trained on a randomly selected subset of the data and features to prevent overfitting.