# Predicting Disease Outbreaks Using Time Series Analysis and Classification

This project focuses on predicting disease outbreaks using time series analysis and machine learning techniques. Specifically, it aims to forecast disease spread and classify countries into low, moderate, and high-risk categories based on healthcare infrastructure and socio-economic factors. By leveraging time series forecasting models like ARIMA, Prophet, and LSTM, along with classification algorithms like Random Forests, Decision Trees, and Neural Networks, the project helps in pandemic preparedness and resource allocation. The models can be adapted to predict various infectious diseases, making this tool valuable for global health crisis management.

## 1. Research Data

- **Objective**: Identify reliable datasets related to disease outbreaks, healthcare infrastructure, and socio-economic factors. This data will help in training predictive models.
- **Sources**:
    - **WHO Disease Outbreak News**: Provides up-to-date information on global disease outbreaks.
    - **COVID-19 Dashboards**: Public repositories that track the spread of COVID-19 across countries (e.g., Johns Hopkins University).
    - **Statista**: Provides datasets on healthcare metrics like hospital bed density, ICU capacity, healthcare expenditure, etc.
- **Other Potential Sources**: Government reports, medical journals, and academic research papers.

## 2. Data Retrieval

- **Objective**: Obtain datasets from the sources identified in the research step.
- **Steps**:

- **Web scraping or APIs**: Use tools like requests or BeautifulSoup in Python to scrape data from websites.
    - **Data import**: Download datasets in formats like CSV, JSON, or Excel from online repositories.
    - **Automation**: Set up a data pipeline to regularly fetch updated information using cron jobs or scheduled tasks.
- **Tools**:
    - pandas for data manipulation.
    - requests for API calls.
    - BeautifulSoup for web scraping.

## 3. Data Preparation

- **Objective**: Organize the data so it's ready for analysis.
- **Steps**:
    - **Handling missing values**: Drop or fill missing values (e.g., with the mean, median, or interpolation).
    - **Feature selection**: Identify relevant features (e.g., hospital beds per capita, GDP, testing rates).
    - **Data type conversion**: Convert columns to the appropriate data types (e.g., categorical variables to numeric using encoding).
- **Tools**:
    - pandas for data manipulation.
    - scikit-learn for feature scaling or encoding.

## 4. Data Exploration (EDA)

- **Objective**: Understand the structure and patterns in the data before modeling.
- **Steps**:
    - **Visualizations**: Use histograms, scatter plots, and box plots to explore relationships between variables.
    - **Correlation analysis**: Identify correlations between different features (e.g., healthcare infrastructure and outbreak severity).
    - **Summary statistics**: Check distributions (mean, median, standard deviation) to understand the data's spread.

- **Tools**:
  - matplotlib and seaborn for data visualization.
  - pandas for summary statistics and data inspection.

## 5. Data Modeling

- **Objective**: Build predictive models that classify countries into different outbreak risk categories (low, moderate, high).
- **Steps**:
  - **Time series models** for forecasting outbreaks:
    - **ARIMA**: Useful for univariate time series data with trends and seasonality.
    - **Prophet**: A forecasting tool designed to handle daily data with missing values and outliers.
    - **LSTM**: A type of recurrent neural network that works well for sequence prediction (e.g., predicting future outbreaks).
  - **Classification models** for categorizing risk:
    - **Random Forest**: Robust classifier that handles complex datasets.
    - **Decision Trees**: Useful for understanding feature importance and decision boundaries.
    - **Neural Networks**: Deep learning models that can capture complex patterns in the data.
- **Tools**:
  - statsmodel for ARIMA.
  - prophet for Prophet.
  - tensorflow or keras for LSTM and neural networks.
  - scikit-learn for decision trees and random forests.

## 6. Linear Regression

- **Objective**: Use linear regression to model relationships between healthcare and socio-economic variables, and outbreak severity.
- **Steps**:
  - Fit a linear regression model to predict outcomes like case numbers or mortality based on features (e.g., hospital bed availability).

- Evaluate performance using metrics like R², Mean Squared Error (MSE), or Mean Absolute Error (MAE).
- **Tools**:
  - scikit-learn for implementing linear regression.

## 7. Cleansing and Transforming Data

- **Objective**: Ensure that the data is clean, consistent, and in the right format for modeling.
- **Steps**:
  - **Handling duplicates**: Remove or handle duplicate records.
  - **Normalization**: Scale numerical features to a similar range (e.g., Min-Max scaling).
  - **Encoding categorical variables**: Convert non-numeric features into numerical values using techniques like one-hot encoding.
  - **Outlier detection**: Remove or adjust extreme values that might skew results.
- **Tools**:
  - scikit-learn for scaling and encoding.
  - pandas for data cleansing.

## 8. Exploratory Data Analysis (EDA)

- **Objective**: Perform in-depth analysis to understand the data better.
- **Steps**:
  - Use advanced visualizations like heatmaps to study feature correlations.
  - Investigate relationships between key variables (e.g., GDP vs. COVID-19 cases).
  - Look for any outliers or patterns that could influence model performance.
- **Tools**:
  - matplotlib and seaborn for visualizations.
  - pandas for statistical analysis.

## 9. Building Models

- **Objective**: Train various models and evaluate their performance.
- **Steps**:
    - Split the dataset into training and testing sets (e.g., 80/20 split).
    - Train different models (ARIMA, Prophet, LSTM, Random Forest, etc.).
    - Evaluate models using cross-validation or a separate test set.
    - Tune hyperparameters to improve performance.
    - Select the best-performing model based on accuracy, F1 score, or other evaluation metrics.
- **Tools**:
    - scikit-learn for training and evaluating models.
    - keras or tensorflow for deep learning models (LSTM).

## 10. Presenting and Building the Application

- **Objective**: Build an application or dashboard to present the results.
- **Steps**:
    - Develop a user-friendly interface where users can input data and see the risk classification or forecasted outbreak data.
    - Visualize results using charts and graphs.
    - Provide recommendations for healthcare planning based on the output.
- **Tools**:
    - **Flask** or **Django**: For building web applications.
    - **Streamlit**: For building interactive dashboards quickly.
    - **Plotly** or **Dash**: For interactive visualizations.

In conclusion, this project demonstrates the power of predictive analytics in forecasting disease outbreaks and classifying regions based on risk levels. By utilizing time series forecasting and machine learning techniques, it provides valuable insights for healthcare planning and resource allocation. The models developed can be adapted to various infectious diseases, making them an essential tool for improving global health crisis management and preparedness. Future work will focus on refining these models with more granular data and advanced algorithms.