



ISLAMIC UNIVERSITY OF TECHNOLOGY (IUT)
ORGANISATION OF ISLAMIC COOPERATION (OIC)
DEPARTMENT OF ELECTRICAL AND ELECTRONIC
ENGINEERING

Project Title :

Comparative Analysis of Machine Learning and Deep Learning Models: Evaluating Regression Based Covid Predictions and CNN Architectures for Image Nearest Neighbours

Project Members: (i) Masrur Ibne Ali (200021131)
(ii) Mustafid Bin Mostofa (200021141)
(iii) Imtiaz Tanweer Rahim (200021150)

Course Title: Artificial Intelligence

Course No. EEE-4709

Instructor Name: Md. Arefin Emon

Submission Date: 09.03.2025

1.INTRODUCTION

1.1 BACKGROUND AND MOTIVATION

The COVID-19 pandemic has highlighted the critical role of data-driven decision-making in healthcare, epidemiology, and public policy. Machine learning (ML) models have been widely employed to predict infection rates, disease severity, and patient outcomes. However, the performance of different ML algorithms varies based on dataset structure, feature engineering, and classification criteria. While KNN, Logistic Regression, and Random Forest Classification are commonly used, their comparative effectiveness in COVID-related predictions remains an open question. Understanding these variations is crucial for selecting the most suitable model for different predictive tasks.

Simultaneously, deep learning has revolutionized image analysis, particularly with the advent of Convolutional Neural Networks (CNNs). In medical imaging and other visual data applications, CNN-based models such as GoogLeNet, ZFNet, and ResNet-101 have demonstrated remarkable feature extraction capabilities. One of their key applications is nearest-neighbor search, which helps in tasks like disease diagnosis, anomaly detection, and image retrieval.

Why is this project important?

COVID-19 datasets exhibit complex patterns, making it essential to determine which ML model best captures these trends. While KNN is commonly used for classification tasks, logistic regression and random forest offer better classification capabilities. Understanding their trade-offs is necessary for improving model selection. A systematic evaluation will help in identifying the

most efficient model based on predictive accuracy, computational cost, and generalizability.

GoogLeNet, ZFNet, and ResNet-101 differ in depth, filter sizes, and computational efficiency, but their performance trade-offs in nearest-neighbor search remain unclear. Understanding these differences will assist in selecting the best architecture for applications such as medical image retrieval, automated diagnosis, and content-based image search.

This study aims to address these gaps by conducting a comparative analysis of both traditional ML models for COVID-19 prediction and CNN architectures for nearest-neighbor search. The findings will provide insights into model selection and optimization, contributing to more effective applications in pandemic response and deep learning-based image processing.

1.2 PROBLEM STATEMENT

The increasing availability of COVID-19 data has led to the widespread application of machine learning techniques for predictive modeling. However, the performance of different machine learning models, such as K-Nearest Neighbours, Logistic Regression, and Random Forest Classification, varies significantly depending on dataset characteristics, feature selection, and preprocessing techniques. Understanding these variations is crucial for optimizing predictive accuracy and reliability in COVID-related studies.

Additionally, deep learning models have gained prominence in image-based tasks, including nearest-neighbor search for medical image retrieval. Convolutional Neural Networks (CNNs) such as GoogLeNet, ZFNet, and

ResNet-101 exhibit differing levels of efficiency and accuracy based on architectural depth and feature extraction capabilities. However, a comprehensive performance comparison of these models in the context of image-based nearest-neighbor search remains limited.

1.3 OBJECTIVES

Comparative Evaluation of Traditional ML Models for COVID-19 Prediction:

- The objective here is to assess the performance of KNN, Random Forest, and Logistic Regression for predicting COVID-19 outcomes based on various features in the dataset.
- **Evaluation Metrics:** The models will be evaluated based on several performance metrics, including accuracy, precision, recall, F1 score, ROC-AUC, MAE, MSE, RMSE, and R2 score. These metrics will provide a comprehensive view of model performance, such as how well the model distinguishes between positive and negative cases (classification metrics), its error rates (regression metrics), and its overall prediction quality.
- **Outcome:** The goal is to identify the most efficient and effective model in terms of both predictive accuracy and generalizability. By evaluating multiple models, we aim to understand their trade-offs and determine the best-suited algorithm for COVID-19 prediction in a real-world scenario.

Comparison of CNN Architectures for Nearest-Neighbor Search:

- The objective is to evaluate the performance of deep learning architectures like GoogLeNet, ResNet, and ZFNet for nearest-neighbor search on a Kaggle image dataset with 10 classes.
- **Evaluation Metrics:** The CNN architectures will be assessed based on training accuracy, training loss, and validation loss. These metrics will provide insight into how well each architecture is able to learn and generalize from the data, as well as how effectively they can minimize error during training and validation.
- **Outcome:** The aim is to identify the most efficient CNN architecture for tasks like medical image retrieval, automated diagnosis, and content-based image search. By comparing these architectures, the study will provide a clearer understanding of the trade-offs between network depth, filter sizes, computational efficiency, and model performance.

Identification of Optimal Model for COVID-19 Prediction: The project aims to determine the most efficient ML model for predicting COVID-19 outcomes, balancing accuracy, computational cost, and generalizability.

Selection of Best CNN Architecture for Image Classification: By evaluating training accuracy, loss, and validation loss, the study will identify the CNN architecture most suitable for tasks like medical image retrieval and automated diagnosis.

The findings will guide the selection of appropriate models for real-world applications, such as healthcare, where fast and accurate predictions are critical.

1.4 SCOPE AND LIMITATIONS

The models included for Covid detection are KNN, Random Forest, and Logistic Regression. The CNN architectures included for image classification are GoogLeNet, ResNet-101, and ZFNet. Our study will explore how these models and architectures perform across different evaluation criteria, ultimately aiming to identify the most effective approach for both COVID-19 prediction and image classification.

Machine Learning part:

For our project, the dataset for COVID-19 contains entries of 49,173 patients with 21 features. This is a reasonably large dataset that helped us to provide ample information for training and evaluating the machine learning models. However the availability and quality of other COVID-19 dataset may pose challenges. The other datasets may have missing or incomplete records, which could impact model accuracy. Additionally, variations in dataset distribution (e.g., demographic imbalances) may affect generalizability.

There is a risk of overfitting in our case as the data size is huge. Balancing model complexity, regularization, and the size of the training data will be key challenges. We wanted the study's scope to include basic hyperparameter tuning, but unfortunately our laptop resources could not perform exhaustive grid search or cross-validation for such a large dataset, it would just require an unreasonable amount of time.

Deep Learning part:

Deep learning models such as GoogLeNet, ResNet-101, and ZFNet are computationally intensive and require significant processing power, especially for large datasets. Without GPU resources, training times can be unreasonably long and as a result we will try to implement this part of the project on a PC that has a powerful GPU configuration.

2.0 LITERATURE REVIEW/RELATED WORK

2.1 EXISTING STUDIES

For the Machine Learning part, we have mainly gone through the following three articles.

(i) Performance Evaluation of Regression Models for the Prediction of the COVID-19 Reproduction Rate

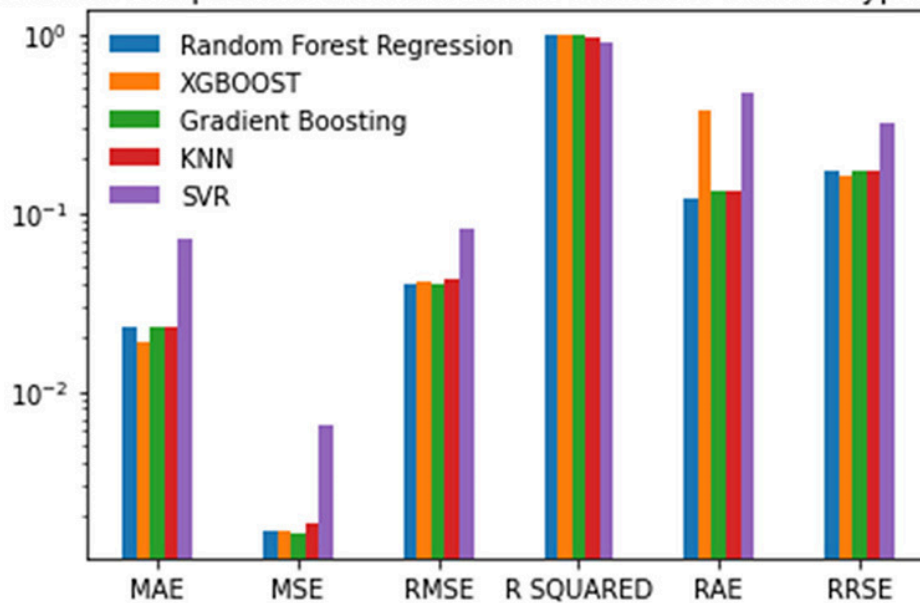
Authors: Jayakumar Kaliappan¹, Kathiravan Srinivasan¹, Saeed Mian Qaisar², Karpagam Sundararajan³, Chuan-Yu Chang^{4} and Suganthan C⁵*

Published: 14 September, 2021

This study evaluated the performance of multiple nonlinear regression techniques—SVR, KNN, Random Forest Regressor, Gradient Boosting, and XGBOOST—for predicting the COVID-19 reproduction rate. Sixteen key features related to testing, deaths, positivity rate, active cases, stringency index, and population density were ranked using Random Forest, Gradient Boosting, and XGBOOST feature selection methods, from which seven were selected. Prediction performance was assessed using MAE, MSE, RMSE, R^2 , RAE, and RRSE.

Since we will not be performing the hyperparameter tuning part for our project, we only mention about the following performance metrics found on this paper:

Performance comparison without Feature Selection without hyperparameter



Sl. No	Performance metrics	Prediction without feature selection and without hyperparameter tuning				
		Random forest regression	XGBOOST	Gradient boosting	KNN	SVR
1	MAE	0.0230122	0.0189412	0.02226608	0.0228918	0.0712651
2	MSE	0.0016347	0.0016482	0.00135535	0.0018072	0.0064267
3	RMSE	0.0404316	0.0405992	0.03681510	0.0425122	0.0801667
4	R-Squared	0.9792338	0.9790759	0.97830657	0.9710729	0.8971356
5	RAE	0.1206129	0.3754830	0.11731593	0.1306129	0.4754830
6	RRSE	0.1700794	0.1605438	0.14728681	0.1700794	0.3207246

(ii) Prediction of COVID-19 Possibilities using KNN Classification Algorithm

Authors: Prasannavenkatesan Theerthagiri, I. Jeena Jacob, A. Usha

Ruby, Vamsidhar Yendapalli

Published: September, 2020

The COVID-19 dataset used for this study contained the patient's details with recovered and deceased status. The vital patient's information was used to diagnose and predict the COVID-19 disease among the infected population. The considered COVID-19 dataset contained 100284 records. The data preprocessing and cleaning process removed the missing and outliers data values from the dataset. The resulting dataset after preprocessing was

reduced to 730 records with three required relevant features of patient details, with 99554 records missing required essential values.

The following accuracy scores were obtained for the different classifiers used:

S. No	Classifier	Accuracy
1.	Logistic Regression (LR)	78.5388
2.	K Neighbors Classifier (KNN)	80.3653
3.	Decision Tree (DT)	75.3425
4.	Support Vector Machines (SVM)	78.9954
5.	Multi-Layer Perceptron (MLP)	77.1689

The following error metrics of the classifiers were found as follows:

S. No	Classifier	MSE	RMSE	Kappa
1.	Logistic Regression (LR)	0.2146	0.4633	0.4109
2.	K Neighbors Classifier (KNN)	0.1963	0.4431	0.469
3.	Decision Tree (DT)	0.2466	0.4966	0.3043
4.	Support Vector Machines (SVM)	0.21	0.4583	0.4266
5.	Multi-Layer Perceptron (MLP)	0.2283	0.4778	0.3411

(iii) Covid-19 detection using Machine Learning

Authors: Bhuvaneswar, Harshitha K, Sahana- 1911T247

This project aimed at comparing different machine learning algorithms like K-nearest neighbors, Random forest and Naive Bayes with respect to their accuracies and then used the best one among them to develop a system which predicts whether a person has COVID or not using the data provided to the model. The data set used was from kaggle.com and it had 5434×21 rows of columns. This dataset contained 20 variables that could be determinants in the prediction of COVID-19. The following performance parameters were found:

	Accuracy	MSE	R2 score	ROC score	Running time
KNN	98.37%	2.57	83.1	98.58	24.252
Logistic Regression	97.03%	3.036	80.086	93.23	0.038
Random Forest	98.39%	2.207	85.51	97.41	213.331

2.2 COMPARISON WITH EXISTING WORK

In our study for the ML part, a significantly larger dataset is utilized, which is expected to yield improved performance metrics compared to findings of the articles that have been shown. A larger dataset typically enhances model performance by providing more diverse and representative samples, reducing variance, and minimizing the risk of overfitting. With more data points available, machine learning models can better capture underlying patterns, leading to higher predictive accuracy and more reliable generalization. The key evaluation metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), R-Squared (R^2), Relative Absolute Error (RAE), and Root Relative Squared Error (RRSE) are anticipated to show noticeable improvements over previously reported results. Our study aims to leverage the advantages of this large dataset to achieve more precise and robust predictions of fatality based on COVID-19 cases.

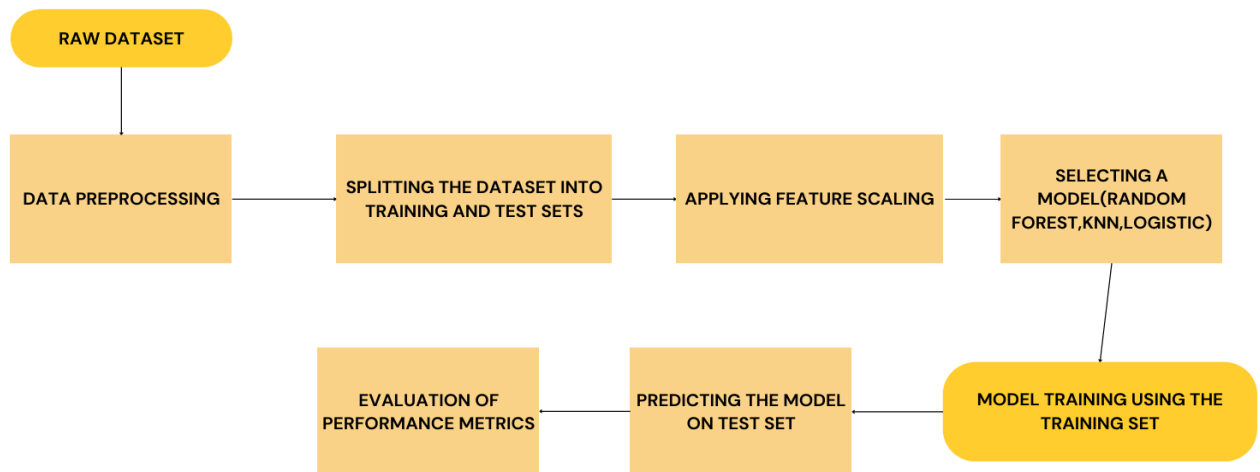
For Deep Learning, we have gone through the following paper:

(iv) Comparative analysis of VGG, ResNet, and GoogLeNet architectures evaluating performance, computational efficiency, and convergence rates
[DOI:10.54254/2755-2721/44/20230676](https://doi.org/10.54254/2755-2721/44/20230676)

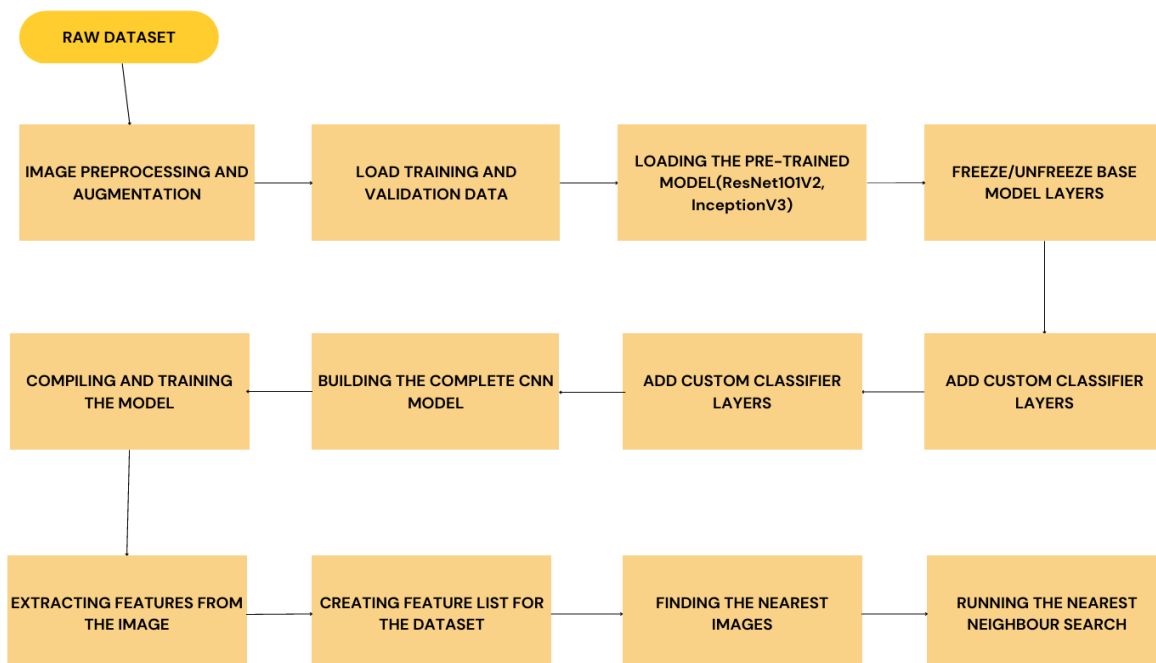
3. SYSTEM ARCHITECTURE / EXPERIMENTAL SETUP

3.1 SYSTEM DESIGN

Workflow Diagram for ML part of project



Workflow Diagram for DL part of project



3.2 SOFTWARE REQUIREMENTS

The project is implemented in VS Code with the Anaconda Python kernel, utilizing Python as the core programming language.

Libraries used for our project:

- ☐ Pandas
- ☐ Numpy
- ☐ Matplotlib and Seaborn
- ☐ Scikit-Learn
- ☐ TensorFlow
- ☐ Keras

3.3 DATA SOURCES AND PREPROCESSING

Data Preprocessing for ML part:

The dataset contains entries of 49173 individuals(rows) and 21 features (these are the columns which include the target variables as well).

In our project, we have performed data preprocessing on the COVID-19 raw dataset from kaggle to present it in such a manner that anyone can understand what each column of this dataset is representing. This is done because the original dataset contains some values that are very ambiguous and can be difficult for people with no knowledge of the raw dataset to capture the essence of what the dataset is really trying to imply.

First, we transformed categorical values into more interpretable labels; for instance, we recoded the 'STATUS' column to indicate whether a patient "Survived" or "Died" and mapped 'SEX' to "male" and "female". Similarly, we modified 'PATIENT_TYPE' to distinguish between "Returned home" and "Hospitalization", while 'CLASIFFICATION_FINAL' was categorized into "Covid positive" and "Covid negative" based on its values. The 'USMER' column, representing healthcare levels, was converted into "Primary care," "Secondary care," and "Tertiary care", and the 'PREGNANT' column was adjusted to indicate "yes" or "no" based on the encoded values. To handle missing or

ambiguous data, we replaced values 97 and 99 with NaN to better account for missing information. Additionally, for multiple binary categorical columns—excluding key attributes like 'SEX', 'AGE', 'PATIENT_TYPE', 'STATUS', 'CLASSIFICATION_FINAL', 'USMER', and 'MEDICAL_UNIT'—we replaced 1 with "yes" and 2 with "no", making the dataset more interpretable. These preprocessing steps ensure that our dataset is clean, structured, and ready for further analysis and model training.

Data Preprocessing for DL part:

Dataset used: CALTECH - 101 from Kaggle; 3000 Images, 10 Classes, Image Size: 224x224 pixel

In our project, we have implemented image preprocessing and data augmentation using Keras' ImageDataGenerator to enhance the performance of our deep learning model. We have applied rescaling, which normalizes pixel values to the range [0,1], along with shearing, zooming, and horizontal flipping to introduce variations in the training data and improve generalization.

Additionally, we have split our dataset, reserving 20% for validation while using the remaining 80% for training. To efficiently load and process the images, we have used `flow_from_directory`, which retrieves images from the "dataset" directory, resizes them to 224x224 pixels, batches them in groups of 32, and organizes them for multi-class classification. This setup ensures that our model learns from diverse augmented images while also being evaluated on separate validation data. By incorporating these techniques, we aim to reduce overfitting, and optimize the performance of CNN architectures such as GoogLeNet, ResNet, and ZFNet, which we are analyzing in our project.

4. METHODOLOGY

4.1 THEORETICAL FOUNDATIONS

With the growth of computer technology, predictive modeling is changing. We are now able to make predictable modeling more efficient, and less expensive than before. In our project, we use various classification algorithms for prediction.

(i) Logistic Regression

Logistic regression is a data categorization technique that uses machine learning. This algorithm models the odds of the potential outcomes of a single experiment using a logistic function. The easiest way to understand the influence of numerous independent factors on a single outcome variable is to use logistic regression, which was designed for this purpose. In general, the algorithm calculates the probability of belonging to a particular class. We have two classes here, $y=0,1$.

(ii) KNN

The oldest supervised machine learning algorithm for classification is KNN, which classifies a given instance according to the majority of categories among its k-nearest neighbours in the dataset. The distance between the item to be categorized and every other item in the data set is calculated by the algorithm.

(iii) Random Forest Regression

This classifier is a meta-estimator that adapts to decision trees on the dataset's different sub-samples and utilizes the average to increase the model's

predicted accuracy and control over-fitting. In most circumstances, this random forest

classifier seems to be more accurate than decision trees, and it also minimizes overfitting. At the Random Forest level, average over all the trees is the final feature importance.

In our model evaluation, we have used several error metrics to determine which model performs the best:

- **Mean Absolute Error (MAE):** It is the average of the absolute differences between predicted and actual values. It gives a simple measure of prediction error in the same units as the data.
- **Mean Squared Error (MSE):** It is the average of the squared differences between predicted and actual values. It penalizes large errors more than small ones, making it sensitive to outliers.
- **Mean Squared Error (RMSE):** It is the square root of MSE. It allows us to understand how large the errors are in the original scale of the data.
- **Cohen's Kappa (κ):** This is a measure of agreement between two raters (or models), correcting for agreement occurring by chance. It ranges from -1 (complete disagreement) to 1 (complete agreement).
- **F1 Score:** This is the harmonic mean of precision and recall, providing a balance between them. This is useful when class distribution is imbalanced and both false positives and false negatives are important.
- **R-Squared (R^2):** It measures how well the model's predictions explain the variance of the true values. A higher value indicates a better model fit.

- **Precision:** This is the proportion of correctly predicted positive instances out of all predicted positives. It shows how many predicted positives are actually true positives.
- **Recall:** This is the proportion of correctly predicted positive instances out of all actual positives. It shows how well the model identifies all positive cases.
- **ROC-AUC Value:** This is the area under the Receiver Operating Characteristic curve. It represents the ability of the model to distinguish between classes; a higher AUC means better performance.
- **ROC-AUC Curve:** This provides a graphical representation showing the trade-off between true positive rate (recall) and false positive rate across different classification thresholds..

In the deep learning task, we have extracted image embeddings from the second-to-last layer of three pre-trained models (GoogLeNet, ResNet, and ZFNet) and then found the nearest 10 neighbors based on those embeddings.

➤ **Image Embedding Extraction:**

In deep learning models, the second-to-last layer (often called the penultimate layer) represents high-level features of an image. The outputs from this layer are used as the image embeddings.

These embeddings are high-dimensional vectors that capture the essence of the image's content in a compressed form.

Mathematically, if an image I is passed through a pre-trained model f_{model} , the embedding e can be represented as:

$$e = f_{\text{model}}(I)$$

where f_{model} represents the pre-trained network (GoogLeNet, ResNet, or ZFNet), and e is the resulting embedding vector.

➤ **Nearest Neighbor Search:**

In order to measure how similar the embeddings of different images are, the two common methods used are:

(i) Cosine Similarity: Cosine similarity measures the similarity between two vectors (in this case, image embeddings) by calculating the cosine of the angle between them. The closer the cosine similarity is to 1, the more similar the two embeddings are, meaning the images are likely to be similar. A cosine similarity of 0 would indicate no similarity, and -1 indicates completely opposite directions.

(ii) Euclidean Distance: Euclidean distance is a metric that measures the straight-line distance between two points (in this case, the embeddings of two images) in the high-dimensional space. A smaller Euclidean distance means that the images are more similar, as the distance between their embeddings is smaller. Larger distances imply more dissimilarity.

The dimensionality of an embedding refers to the number of features or components in the vector representing an image. For example, an embedding may have hundreds or thousands of dimensions, each capturing different aspects of the image. The dimensionality affects how detailed or compressed the representation is.

Nearest neighbors are the images whose embeddings are most similar to a given image based on either cosine similarity or Euclidean distance. The algorithm identifies the top 10 images whose embeddings are closest to the target image, helping in finding visually similar images.

4.2 EXPERIMENTAL SETUP

For ML part,the following steps were taken to perform the study:

➤ Loading and Preprocessing the Dataset

The dataset UPDATED_COVID.csv was loaded using `pandas.read_csv()`, with column names manually defined and the first row skipped.

The 'NaN', 'ICU', and 'INTUBED' columns were removed as these columns had many missing values and hence were not needed for analysis. The target column 'STATUS' was moved to the extreme right of the dataset to make it the target feature. Missing values in the 'PNEUMONIA' column were dropped, and missing values in the 'AGE' column were imputed using the mean value. Categorical columns such as 'USMER', 'SEX', 'PATIENT_TYPE', and 'CLASIFFICATION_FINAL' were encoded using `LabelEncoder` to convert them into numeric values. The 'STATUS' column was converted to binary values where 'Died' was mapped to 1 and 'Survived' to 0, and other categorical values like 'yes' and 'no' were replaced with 1 and 0.

➤ Splitting Data into Features and Target Variables

The feature set X (independent variables) and target variable y (dependent variable) were separated. The dataset was split into training and testing sets using `train_test_split()` with a test size of 20%.

➤ Feature Scaling

The 'AGE' feature was scaled using `StandardScaler`. The scaling was performed separately for the training and testing datasets. Similarly, the 'MEDICAL_UNIT' column was also scaled.

➤ **Model Training and Evaluation**

KNN, Logistic Regression and Random Forest Regression models were used.

➤ **Model Performance Evaluation**

The confusion matrix was computed to evaluate the performance of the 3 models. Several performance metrics were computed for the KNN model, including: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), R^2 Score, F1 Score, Precision, Recall, Cohen's Kappa, and ROC-AUC.

➤ **Visualization of Results**

Confusion Matrix Visualization: The confusion matrix was plotted using `ConfusionMatrixDisplay()` to visually assess the classifier's performance.

Bar Chart Comparison: A bar chart comparing the error metrics (MSE, RMSE, Kappa) across Logistic Regression, Random Forest, and KNN models was generated.

ROC Curve Comparison: The ROC curve for KNN, Random Forest, and Logistic Regression models was plotted and compared, showing the AUC values for each model.

Bar Plot for Model Performance Metrics: A bar plot comparing Accuracy, Precision, Recall, and F1-Score across KNN, Random Forest, and Logistic Regression models was generated.

Below are some code snippets of how we have processed the raw dataset, splitted it into training and test sets and also performed feature scaling.

MAKING SENSE OUT OF THE RAW DATASET

```
dataset['STATUS'] = dataset['STATUS'].apply(lambda x: 'Survived' if x == '9999-99-99' else 'Died')
dataset['SEX'] = dataset['SEX'].replace({1: 'female', 2: 'male'})
dataset['PATIENT_TYPE'] = dataset['PATIENT_TYPE'].replace({1: 'Returned home', 2: 'Hospitalization'})
dataset['CLASIFFICATION_FINAL'] = dataset['CLASIFFICATION_FINAL'].apply(lambda x: "Covid positive" if x in [1, 2, 3] else "Covid negative")
dataset['USMER'] = dataset['USMER'].replace({1: 'Primary care', 2: 'Secondary care', 3: 'Tertiary care'})
dataset['PREGNANT'] = dataset['PREGNANT'].apply(lambda x: 'yes' if x == 98 else 'no')
```

Fig: A code snippet of how the raw dataset is transformed to a meaningful representation. (We have in fact designed a separate .ipynb file just to make the raw dataset seem unambiguous and understandable)

dataset.head()

	USMER	MEDICAL_UNIT	SEX	PATIENT_TYPE	STATUS	INTUBED	PNEUMONIA	AGE	PREGNANT	DIABETES	...	ASTHMA	INM...
0	Secondary care	1	female	Returned home	Died	NaN	yes	65.0	no	no	...	no	
1	Secondary care	1	male	Returned home	Died	NaN	yes	72.0	no	no	...	no	
2	Secondary care	1	male	Hospitalization	Died	yes	no	55.0	no	yes	...	no	
3	Secondary care	1	female	Returned home	Died	NaN	no	53.0	no	no	...	no	
4	Secondary care	1	male	Returned home	Died	NaN	no	68.0	no	yes	...	no	

5 rows × 21 columns

Fig: The dataset contains many columns which are categorical in type and needs to be converted to numerical representations.

Label Encoding the categorical features

```
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
columns_to_encode = ['USMER', 'SEX', 'PATIENT_TYPE', 'CLASIFFICATION_FINAL']
for column in columns_to_encode:
    data[column] = label_encoder.fit_transform(data[column])
```

```
Splitting the Dataset into Predictors(Independent variables) and Target Column(Dependent variable)
```

```
X = data.iloc[:, 0:-1].values
y = data.iloc[:, -1].values
```

Python

Fig: Label Encoding categorical features and splitting the dataset into independent variables and dependent variable

```
Splitting the Dataset into Train and Test sets
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 100)
```

Python

```
Feature Scaling
```

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()

age_column_index = 5
# Fitting the scaler object on the AGE column in the training set and transforming it
X_train[:, age_column_index] = scaler.fit_transform(X_train[:, age_column_index].reshape(-1, 1)).flatten()

# Transforming the AGE column in the test set using the same scaler object
X_test[:, age_column_index] = scaler.transform(X_test[:, age_column_index].reshape(-1, 1)).flatten()
```

Python

Fig: Splitting the dataset using help from a module of the sklearn library and also performing feature scaling

For DL part,the following steps were taken to perform the study:

➤ **Data Preprocessing:**

Image Resizing: The image data was scaled using ImageDataGenerator to normalize the pixel values by dividing by 255, ensuring that all images are in the range of 0 to 1.

Augmentation: Data augmentation was applied (shear, zoom, and horizontal flip) to introduce variability in the training data, which helps improve the generalization of the model. A validation split of 20% was also specified.

Image Size: Images were resized to (299, 299) pixels to match the input size required by the pretrained models.

➤ **Loading the Dataset:**

The dataset was loaded from a directory ('dataset') using `flow_from_directory`. This method automatically labels images based on their directory structure and applies data augmentation for training and validation sets. The images were loaded in batches of 32, and the labels were one-hot encoded using `class_mode='categorical'`.

➤ **Loading Pretrained Model:**

A base model like (InceptionV3) was loaded with pretrained weights from ImageNet (`weights='imagenet'`) without the top layers (`include_top=False`). This allows for feature extraction without interference from the model's original classification layers. The model's layers were frozen until layer 40 to prevent overfitting and speed up training. Layers from layer 40 onward were left trainable, enabling fine-tuning.

➤ **Creating a New Classifier:**

A new classifier was added on top of the base model:

GlobalAveragePooling2D: It reduces the spatial dimensions of the feature map.

Dense Layer (1024 units): This is a fully connected layer with ReLU activation function.

Dropout: Regularization technique (50%) was used to prevent overfitting.

Output Layer (Dense with 10 units): Output layer with softmax activation was used for multi-class classification.

➤ **Model Compilation:**

The model was compiled using the Adam optimizer and categorical cross-entropy loss for multi-class classification, with accuracy as the evaluation metric.

➤ **Model Training:**

The model is trained for 10 epochs using the fit method on the training data (training_set) and validated using the validation data (validation_set).

Performance Metrics: Training accuracy, training loss, validation accuracy, and validation loss were evaluated during each epoch.

➤ **Feature Extraction:**

After training, the model was used to extract features from images.

extract_features_fixed: This function loads an image, resizes it to the required input size (299x299), normalizes it, and extracts the features using the base model (InceptionV3).

These features were then stored in a list for all images in the dataset.

➤ **Nearest Neighbor Search:**

Using the NearestNeighbors algorithm (with cosine distance), the features of a query image were compared with the stored features of all images in the dataset. For each query image, the nearest n_neighbors (10 in our case) were identified, and the images were displayed along with their distance values.

This was done for multiple query images (e.g., one from the 'Motorbikes' class and one from a different class, 'plane').

➤ **Comparison Between Models:**

GoogLeNet, ResNet, and ZFNet all were used in our study for finding nearest neighbors, following similar steps as with InceptionV3.

For each model, the following metrics were recorded: Training Accuracy, Training Loss, Validation Accuracy, and Validation Loss.

➤ **Results Analysis and Comparison:**

The performance of each model (InceptionV3, GoogLeNet, ResNet, ZFNet) was compared using the following:

Training Accuracy and Loss: To assess how well the model fits the training data.

Validation Accuracy and Loss: To evaluate the model's generalization capabilities.

Nearest Neighbor Results: To compare the ability of each model to find similar images in terms of feature space.

➤ Visualization:

For each query image, the nearest images were displayed alongside their distances. In addition plots of performance metrics (accuracy, loss, etc.) for each model were presented for visual comparison.

Below are some code snippets of the libraries that have been used and the process of data preprocessing.

```
import tensorflow as tf
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import GlobalAveragePooling2D, Conv2D, MaxPooling2D, Flatten, Dense, Dropout
from keras.preprocessing import image
from tensorflow.keras.applications import InceptionV3
import numpy as np
from sklearn.neighbors import NearestNeighbors
import matplotlib.pyplot as plt
import os
```

```
train_datagen = ImageDataGenerator(rescale=1./255,
                                   shear_range=0.2,
                                   zoom_range=0.2,
                                   horizontal_flip=True,
                                   validation_split=0.2)

training_set = train_datagen.flow_from_directory('dataset',
                                                  target_size=(299, 299),
                                                  batch_size=32,
                                                  class_mode='categorical',
                                                  subset='training')

validation_set = train_datagen.flow_from_directory('dataset',
                                                    target_size=(299, 299),
                                                    batch_size=32,
                                                    class_mode='categorical',
                                                    subset='validation')
```


4.3 ASSUMPTIONS AND CONSTRAINTS

Assumptions for the ML study:

- It was assumed that the dataset provided was clean and representative of the problem being studied, except for missing values, which were handled through imputation or removal.
- It was assumed that imputing missing values (e.g., replacing missing 'AGE' values with the mean) was an appropriate approach, assuming the missing data did not introduce bias.
- The dataset was assumed to be sufficiently large for training and validation of machine learning models, ensuring reliable results.

Limitations in the methodology for ML:

- Columns like 'ICU' and 'INTUBED' were dropped. The decision to exclude these features could lead to loss of important information that might have improved the model's predictive power.
- There is a lack of cross validation since the study used a single train-test split (80% train, 20% test). While this is common, it may lead to variability in the results depending on how the data is split.
- The presence of correlated features could lead to redundancy in the models, potentially reducing their performance.
- Since models like Random Forest and KNN were used without hyperparameter tuning, there remains a risk of overfitting.

Assumptions for the DL study:

- It is assumed that the images in the dataset are of good quality and that preprocessing steps (rescaling, augmentation, etc.) will improve the model's ability to generalize without introducing significant noise.
- The study assumes that pretrained models like InceptionV3, GoogLeNet, ResNet, and ZFNet can effectively extract relevant features from the images and that these features can be used for nearest neighbor search, irrespective of the image class or dataset.

- It is assumed that cosine distance is an appropriate metric for finding the nearest neighbors in the feature space.
- The study assumes that fine-tuning the last layers of the pretrained models will allow for better performance on the new dataset while avoiding overfitting. Only the latter layers are made trainable to prevent overfitting on small datasets.

Limitations in the methodology for DL:

- Fine-tuning the models could potentially lead to overfitting on the training data, especially if the dataset is small. If the model is too complex for the available data, it may memorize the data instead of learning generalized patterns.
- The study only compares the performance of InceptionV3, GoogLeNet, ResNet, and ZFNet. There are many other models that could potentially outperform these, such as EfficientNet, DenseNet, or other newer architectures.
- Deep learning models like InceptionV3, GoogLeNet, ResNet, and ZFNet are computationally expensive. Running these models on a CPU may lead to longer training times and even inability to train the models effectively.

8. ETHICAL CONSIDERATIONS AND SUSTAINABILITY

8.1 ETHICAL ISSUES

Both of our studies handle sensitive data. The COVID-19 study must ensure health data is anonymized and secure. In image classification, the use of personal images could raise privacy concerns. The compliance with data protection regulations like GDPR is essential. Biases in data can lead to unfair

predictions or classifications. In the COVID-19 project, non-representative data could affect healthcare decisions. In image classification, biased datasets may result in poor performance for certain groups. There could be ethical concerns over using models for decision-making in healthcare, where human oversight may be necessary. Similarly, image classification could raise issues around surveillance and privacy if misused.

8.2 SUSTAINABILITY

The COVID-19 model could improve healthcare efficiency and reduce costs. Image classification models could automate processes and reduce operational costs, but could also lead to job displacement in certain sectors. Both projects could benefit society by improving healthcare and efficiency. However, there are concerns about AI-driven surveillance and job displacement. Sustainable and fair AI deployment is essential for positive societal impact.

10. REFERENCES

- *Performance Evaluation of Regression Models for the Prediction of the COVID-19 Reproduction Rate*
Authors: Jayakumar Kaliappan¹, Kathiravan Srinivasan¹, Saeed Mian Qaisar², Karpagam Sundararajan³, Chuan-Yu Chang^{4*} and Suganthan C⁵
- *Prediction of COVID-19 Possibilities using KNN Classification Algorithm*
Authors: Prasannavenkatesan Theerthagiri, I. Jeena Jacob, A. Usha Ruby, Vamsidhar Yendapalli
- *Covid-19 detection using Machine Learning*
Authors: Bhuvaneswar, Harshitha K, Sahana- 191IT247
- *Comparative analysis of VGG, ResNet, and GoogLeNet architectures evaluating performance, computational efficiency, and convergence rates*
DOI:10.54254/2755-2721/44/20230676