Browse ⌄    My Settings ⌄    Help ⌄          Institutional Sign In

Institutional Sign In

All    ⌄                                                    🔍

ADVANCED SEARCH

Conferences  >  2020 3rd International Confer...  ⍰

# Implementation of Stacking Ensemble Learning for Classification of COVID-19 using Image Dataset CT Scan and Lung X-Ray

**Publisher: IEEE**    | Cite This |    📄 PDF

Annisa Utama Berliana ;  Alhadi Bustamam    **All Authors** •••

**13**          **1627**
Cites in          Full
Papers          Text Views

🔓   ®   🔗   ©   📁   🔔

**Alerts**

Manage Content Alerts

Add to Citation Alerts

🔓 Free

---

**Abstract**

Document Sections

I.  Introduction

II.  Related Work

III.  Materials and Methods

IV.  Results

V.  Conclusion

Authors

Figures

References

Citations

Keywords

Metrics

More Like This

📄

Downl

PDF

**Abstract:**
Novel Coronavirus Disease (COVID-19) is a disease caused by SARS-CoV-2, which has become a global pandemic. COVID-19 was first discovered in Wuhan, China, and has already... **View more**

⌄ **Metadata**
**Abstract:**
Novel Coronavirus Disease (COVID-19) is a disease caused by SARS-CoV-2, which has become a global pandemic. COVID-19 was first discovered in Wuhan, China, and has already spread to various countries, which, until now, still haven't found a proper way to deal with it. Various studies related to COVID-19 have been carried out, including initial screening to control the disease's spread. X-ray images and Computed Tomography (CT) can be utilized for initial screening in diagnosing lung conditions for patients with COVID-19 symptoms. Machine learning has been at the forefront of many fields, such as analyzing X-Ray and CT Images. Machine learning shows an outstanding performance compared to other methods. In this paper, we present an ensemble learning with stacking to analyze X-Ray and CT in calcifying COVID-19, which was previously pre-documented using the Gabor feature. The ensemble learning model is built with two levels of learning, namely the base-learners and the meta-learner. The base-learners we use to build the model are Support Vector Classification (SVC), Random Forest (RF), and K-Nearest Neighbors (KNN), and the meta-learner we use is Support Vector Classification (SVC). The proposed method's performance is implemented on a publicly available COVID-19 data set, including 1140 chest X-Ray images and 2400 CT Images. The proposed method shows that the stacking ensemble learning of Support Vector Classification (SVC), Random Forest (RF), and K-Nearest Neighbors (KNN) can provide accuracy above 97% for CT Images and 99% for chest X-Ray images.

▶ **ISBN Information:**           **Conference Location:** Yogyakarta, Indonesia

≡ **Contents**

## SECTION I.
# Introduction

Wuhan City, Hubei Province, China, is where the first report on Novel Coronavirus Disease (COVID-19) was reported in November 2019. The World organization (WHO) announced that this virus could cause respiratory disease with fever, cough, and pneumonia [1], [2]. The virus is spread from person to person by droplets when the patient coughs, talks, or sneezes. Close contact is also the transmission source, such as hands contaminated with the virus and direct contact with the mouth, nose, or conjunctiva of the eyes [2]. Therefore, COVID-19 can quickly spread to various countries, not only in China, so the WHO declared this pandemic as an emergency health problem. On September 12, 2020, 216 countries were recorded, 911,877 patients died, and 28,329,790 cases were confirmed worldwide [3].

The SARS-CoV-2 virus causes novel Coronavirus Disease (COVID-19). The SARS-CoV-2 virus is a family of *corona Miridae*, the same family as viruses that caused the Middle East Respiratory Syndrome (MERS) and Severe Acute Respiratory Syndrome (SARS) [3]. MERS was first discovered in Saudi Arabia. While SARS was first reported in Southern China in 2003 and spread to various countries, SARS was caused by the SARS-CoV virus [4].

CT Images and X-rays can be used as the most effective and functional method for screening patients' initial screening in the severity and degree of lung inflammation in COVID19 [5]. According to the Chinese National Health Commission's confirmation, China's Hubei Province presented a radiographic presentation of pneumonia for a standard clinical diagnosis [6]. However, analyzing CT and X-ray images requires a radiologist and takes a long time. The rapid spread of the virus and an increase in the number of patients who exceed hospital capacity have resulted in very long initial skinning queues for COVID-19. This has a significant impact on handling the spread of Covid-19. Therefore, a machine learning model is required to analyze CT and X-ray images quickly and with high accuracy [7].

The ensemble learning method is the most frequently used machine learning technique. The ensemble learning method shows that combining several learners produces a strong learner, reducing bias and variance. The ensemble learning method principle is that each learning model covers other learning models' weaknesses so that better results are obtained [8]. So that many researchers are interested in using the ensemble learning method in processing CT images and X-Ray images. Machine learning models that are usually used with the ensemble learning approach are CNN, Naive Bayes (NB), Artificial Neural Network (ANN), SVM (Support Vector Machine), Random Forest (PF), Decision Tree (DT), K-Nearest Neighbors (KNN) and etc. However, several papers that use the ensemble learning approach. Therefore, in this paper, we apply the multilevel ensemble stacking model. The first level consists of three base learners, namely Random Forest (PJ), Support Vector Classification (SVC), and KNearest Neighbors (KNN), which are then stacked. After that, we will obtain temporary results that will be used at the second level. Our second tier uses meta-learners to generate final predictions. The meta-student we use is Support Vector Classification (SVC). The purpose of using the Stacking ensemble machine learning model is to improve the machine learning model's accuracy. We applied it to the 2400 CT Image dataset, including 1200 COVID-19 and 1200 Normal patients, and 1140 X-Ray Image dataset, including 570 COVID-19 and 570 Normal patients. We extract the features using the Gabor feature to convert the image into a vector form to be processed better. So, it is hoped that the model we propose can help doctors predict COVID-19 through CT Images and chest X-Ray Images data more quickly and accurately. It can also ease the doctor's work to focus swiftly and immediately handle patients and reduce long patient queues.

## SECTION II.
# Related Work

Similar research underlies us in conducting this research, starting from a study conducted by Imad et al. [9]. He classified the chest X-Ray Images of the COVID-19 with feature extraction using the Histogram of Oriented Gradients (HOG). He compared the accuracy of several machine learning classification models, namely SVM, KNN, PF, NB algorithm, and DT. The results showed that the SVM has the highest accuracy

than other machine learning classification models. In the same year, Arora et al. [22] apply deep learning (DenseNet and GoogleNet) for feature extraction and stacking ensemble learning to X-Ray image data processing COVID-19 prediction. The result is that the final accuracy of the proposed model is better than other models.

In the previous year, Boukellouz et al. [8] applied three machine learning models, namely ANN, PF, and KNN, as base-learners at the first level and a meta-learner at the second level of the ensemble stacking learning the proposed he model. However, this model applies to a dataset of pseudo-CT Images driven by magnetic resonance where feature extraction has previously been performed using several feature extraction algorithms, namely Gabor feature, Histogram of oriented gradients (HOG), Discrete wavelet transform (DWT), and Histogram. of local binary patterns (LBPH). In the experiment, the mean absolute error (MAE) and the average mean error (ME) his proposed better than the PF models. In his research, he also increased the value increments and calculated the bias in electron counting. Therefore, we are interested in applying the stacking ensemble learning model to machine learning, which consists of SVC, KNN, and PJ as base learner and meta-learner using SVC with feature extraction Gabor Feature on data chest XRay Images and CT Images of COVID-19 and non-COVID19 patients (Normal).

## SECTION III.
# Materials and Methods

### A. Data description and pre-processing

This study used two image data from the online platform Kaggle, namely CT Images [10] and chest X-Ray images [11]. We use 2,400 data images in CT images, including images of 1,200 COVID-19 patients and 1,200 non-COVID19 patients (Normal). We use 1, 140 data images covering 570 COVID-19 patients in the chest X-Ray image and 570 nonCOVID-19 patients (Normal). In both image data, we label a folder using the name COVID-19 and Normal, which is done manually. Then we divided it into training data and test data with a ratio of 8: 2 (80% for training data and 20% for test data).

We first resize CT Images and chest X-ray Images to 64 x64 pixels in the pre-processing stage. We did feature extraction on both image data using Gabor Feature, which produced 80 x80 feature vectors. The feature extraction results are normalized first before inputting them into the multilevel stacking ensemble learning models. The process stages that we will carry out in this study can be seen briefly in Figure 1. Our computation process uses the python programming language.



**Fig. 1.**
Flowchart of the proposed workflow

### B. Feature Extraction

In extracting features from the image, the Gabor filter can provide good perception. By constructing a different scale and orientation from $G_k(z)$, the Gabor feature in the $Z$ position is formed, and the Gabor kernel convolution is defined as follows [12],

$$G_k(z) = I(z)^*\psi(k, z) \tag{1}$$

View Source

I(z)is the original image, and $\psi(k, z)$ is the Gabor features in a position of z.

The initial stage of feature extraction is converting the RGB image value into a grayscale value, and then a convolution process is carried out on the Gabor filter. In this study, we use 40 filters, size =10x10, f=1, $\theta$=0, and $o$=3. Also, we only use the odd-numbered Gabor Kernel $(G_i = img(\phi_{2_i,0}))$, which is edge and location sensitive. After that, we look for the mean value, the mean value will be used as a reference for the classification process. In the final stage for each pixel (x,y) in the normalized key point region, four Gabor filter fuse expression followed.

$$f_i(x, y) = G_i(x, y)^*I(x, y) \ \ i = 0, 1, 2, 3 \tag{2}$$

## C. Normalization

We use the zero-mean normalization technique and the unit-variance normalization, which normalizes the extracted vector features which can be defined as follows: [13]

$$X_i^* = \frac{x_i - \mu}{o}, i = 1, 2, 3, \ldots, n \tag{3}$$

where $\mu$ is the mean value of the feature vector x, and $\sigma$ is the standard deviation. The feature vector x is transformed into a random variable with a mean of zero and a variance of one.

## D. Stacking Ensemble Learning Model

Ensemble learning is usually built by *Base-learners* using several homogeneous or heterogeneous learners whose prediction results are combined. The combination method types are averaging, voting, or learning. [8] Stacking is joining several base learners at the first level to meta-learner at the second level [14].

In this study, we use three *Base-learner*, namely Random Forest (RF), Support Vector Classification (SVC), and KNearest Neighbors (KNN), to take advantage of their diversity, which is the main contributor to building a performing ensemble model [18]. The following is a description of the three *Base-learner* that we use [8]:

- In 1992 at the Annual Workshop on Computational Learning Theory, SVM was first introduced by Boser, Guyon, and Vapnik. The concept of SVM is a hyperplane design that can classify all training data into two classes. The distance between the hyperplane and the closest pattern (support vector) of each class is called the margin. There are two main ways to define this hyperplane: linear separation and nonlinear separation [15]. In linear separation data, the hyperplane is defined as:

$$< w, x > +b = 0, \; \forall x \in H_n \tag{4}$$

w is the normal vector, b is the bias, and $< w, x >$ is the inner product between w and x. The optimal hyperplane is one that has the largest margin $\frac{2}{\|w\|}$.

The final optimization of the margins on the nonlinear data separation is:

$$\min \left[ \frac{\|w\|^2}{2} + c \sum_{i=1}^{n} \xi_i \right] \tag{5}$$

Subject to $y_i(< w, x > +b) \geq 1 - \xi_i, \xi_i \geq 0, \forall_i \in [-1, 1]$. $xi_i$ is the slack variable for querying a misclassified sample, and the constant C controls the number of penalties. The variables introduced are by Fauvel et al. [20]. In this study, we use with kernel Radial Basis Function (RBF).

The random forest is a classification and regression algorithm that is part of the ensemble learning group. The model base classifier used is a decision tree to form a forest. The random forest was made popular by Leo Breiman. The classification process is based on the most votes from the returned decision tree. The process for building a Random Forest is the bagging, the bootstrap aggregating, and the random subspace method. Random Forest can handle data with high dimensions better than other classifier models. Unlike the usual Decision Tree, overfitting is handled by maintaining the forest's model tree's variance. Tree formation in the Random Forest algorithm is done by conducting training on the sample. The variables used for the split are taken randomly, and classification is carried out after all trees are formed. The classification determination in Random Forest is taken based on each tree's votes, and the most votes become the winner. In determining the division that will be the root/node, you can use the Gini Index value, which is defined as follows [16]:

$$Gini_t = 1 - \sum_{j=1}^{v} P^2(Y_h) \tag{6}$$

View Source ⊙

Following Eq. (6) in the case of binary classification the splitting criterion is evaluated as:

$$\text{Gini}_t^{\text{split}} = \frac{D_1}{D}\text{Gini}_{t_1} + \frac{D_2}{D}\text{Gini}_{t_2} \rightarrow \min \qquad (7)$$

View Source ⊙

where D is the Number of objects in the tree node at time t (parent node). So if $D_1$ means at time $t_1$ and if $D_2$ means at time $t_2$.

- KNN is the simplest algorithm among other algorithms in predicting a class in a sample. To classify a class, KNN works based on the closest distance between objects in the following way:**Step 1**, Calculate the distance of all training vectors $(x_i)$ to test vectors $(x_j)$ is possible as follows [17]:

$$d(ij) = f(x_i, x_j) \qquad (8)$$

View Source ⊙

where $f(x_i, x_j)$ is a scalar distance function. The following equation determines the distance vector

$$D(i) = \left\{ \begin{array}{l} d(i,j) \mid i = 1, 2, \ldots, N_{\text{test}}, \\ j = 1, 2, \ldots, N_{\text{train}} \end{array} \right\} \qquad (9)$$

View Source ⊙

**Step 2**, K value that is closest to the vector value (the so-called K-nearest neighbor) as follow:

$$D_N(i) = \text{sort} Ascending(D(i))$$
$$V = \{\delta(D_N(i)(1)), \ldots, \delta(D_N(i)(K))\} \qquad (10)$$

View Source ⊙

**Step 3**, Calculate the average value and select, which yields the best score for all test features, can result in *K* optimal according to the best accuracy in the KNN algorithm.

To validate the effectiveness of our proposed method, we use the values, sensitivity, and specificity measured by the following equation [19]:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \times 100\% \qquad (11)$$

View Source ⊙

$$\text{precision} = \frac{TP}{TP + FP} \times 100\% \qquad (12)$$

View Source ⊙

$$\text{recall} = \frac{TP}{TP + FN} \times 100\% \qquad (13)$$

View Source ⊙

Where, TP (True Positive), TN (True Negative), FN (False Negative), and FP (False Positive).

## SECTION IV.
# Results

Figure 1 shows that we carried out the pre-processing stage on the two image data (CT images and chest X-rays) that we used. In our CT images, 2,400 data images include images of 1,200 COVID-19 patients and 1,200 non-COVID19 patients (Normal). In the chest X-ray, we used 1,140 data images covering 570 COVID-19 patients and 570 nonCOVID-19 patients (Normal). The pre-processing stage includes labeling a folder named COVID-19 and Normal, resizing 64 x64 pixels, dividing the data into train data and data sets with an 8: 2 grayscale ratio. After that, feature extraction with collaborative features is carried out to convert the image data into vectors using the Gabor feature. The feature extraction results are 80 x80 feature vectors for CT images and 80 x80 feature vectors for chest X-ray images. We use the feature extraction results as input into the multilevel stacking ensemble learning models that we have previously normalized.

In this study, we use multilevel stacking ensemble learning, which consists of first and second levels. At the first level, we used three base learners, namely SVC, KNN, and PJ. The base learner's selection is based on several trials we have conducted on several ensemble learning models, including ANN, NB, and DT. From several ensemble learning models that we stacked, it was found that the best AUC (Area Under the Curve) was stacking three ensemble learning models that we made as base learners, namely SVC, KNN, and PH. We use the output from the first level as input for the second level. At the second level, we use SVC as a meta-learner because SVC has a good AUC (Area Under the Curve) at the first level.

In the final stage, 97% accuracy is generated for CT Images and 99% for chest X-Ray images. As an evaluation material for the results of the multilevel stacking ensemble learning models that we use, we process individual models for each of the models we use. The results we get from the multilevel stacking ensemble learning model on CT image and chest X-Ray images can be seen in Table II.

**Table I** The Classification Performance of Stacking Ensemble Learning on Ct Scan and X-Ray Images


Table I- The Classification Performance of Stacking Ensemble Learning on Ct Scan and X-Ray Images

From the table, it can be seen that the stacking ensemble learning model we are proposing gave us 100% AUC, 99% ACC, 98% precision, and 100% recall for X-Ray images. It gave us 99.2% AUC, 97% ACC, 97% precision, and 97% recall for CT Images. So, it can be said that the stacking ensemble learning model proposed gives more accurate classification than the other model on both datasets.

## SECTION V.
# Conclusion

In this study, we proposed an efficient machine learning classifier for the diagnosis of COVID-19 disease from chest X-Ray images and CT Images. We use a 2400 CT Image dataset and 1140 chest X-ray datasets. We processed the dataset using a multilevel stacking ensemble learning model. At the first level, we used three base learners, namely SVC, KNN, and PF. At the second level, we use SVC as a metalearner. This model gave us 100% AUC, 99% ACC, 98% precision, and 100% recall for X-Ray images. and gave us 99.2% AUC, 97% ACC, 97% precision, and 97% withdrawal for CT images. The Histogram of the results of the multilevel stacking ensemble learning models that we use can be seen in Figure 2. The experiment result shows model stacking ensemble learning gives more accurate classification than the individual model on both datasets.


Fig. 2. - The comparison of the results in individual models with the multilevel stacking ensemble learning model on chest x-ray images and CT Images.

**Fig. 2.**
The comparison of the results in individual models with the multilevel stacking ensemble learning model on chest x-ray images and CT Images.

**ACKNOWLEDGMENT**

| Authors | ⌄ |
|---|---|
| Figures | ⌄ |
| References | ⌄ |
| Citations | ⌄ |
| Keywords | ⌄ |
| Metrics | ⌄ |

**More Like This**

Fracture detection in x-ray images through stacked random forests feature fusion

2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI)

Published: 2015

Multi-Modal Fusion of Deep Learning with CNN based COVID-19 Detection and Classification Combining Chest X-ray Images

2023 7th International Conference on Intelligent Computing and Control Systems (ICICCS)

Published: 2023

**Show More**

» Contact & Support

About IEEE *Xplore* | Contact Us | Help | Accessibility | Terms of Use | Nondiscrimination Policy | Sitemap | Privacy & Opting Out of Cookies