

K-Nearest Neighbours and K-Fold Cross Validation for Big Data of Covid 19

Kuntoro Kuntoro*

Division of Biostatistics and Population Study, Airlangga University School of Public Health, Surabaya, Indonesia

Abstract

The most popular model in machine learning is K-Nearest Neighbours (KNN). It is used for solving classification. Moreover, K-Fold Cross-validation is an important tool for assessing the performance of machine learning in doing KNN algorithm given available data. Compared to traditional statistical methods, both algorithms are effective to be implemented in big data. A supervised machine learning approach using KNN and K-Fold Cross-Validation algorithms is implemented in this study. For learning process, data of covid 19 is obtained from website. Four predictors such as new case, reproduction rate, new case in ICU, and hospitalized new case are selected to predict the target, new cases will be alive or will die. After cleaning process, 13,223 of 132,645 data sets are selected. This is considered as original data sets. When K-Fold Cross-Validation is executed by Python showing User Warning, the original data sets are replicated to be 264,441 data sets. This is considered as replicated data sets. Performance of KNN algorithm in predicting the target using original data sets shows lower accuracy than that using replicated data sets (75% vs. 92%). The number of members (K) using original data sets is lower than that using replicated data sets (7 vs. 12). Performance of K-Fold Cross-Validation using original data sets shows very small mean accuracy than that using replicated data sets (0.054 vs. 0.998). In using replicated data sets, mean accuracy shows consistent value until 5 splits while in using original data sets mean accuracy only shows in 2 splits. In using big data from various sources, it is recommended to implement appropriate Python libraries which can remove not a number (nan) and messy record effectively. It is also recommended to develop combine and comprehensive algorithm of KNN and K-Fold Cross-Validation.

Keywords: Machine-learning • KNN • K-Fold Cross-Validation • Accuracy

Introduction

The most popular model in machine learning is K-Nearest Neighbours (KNN). It is used for solving classification McCullum [1]. Moreover, KNN model has similarity with regression model in classification. Both models can be used for prediction. Regression model particularly logistic regression model predicts categorical variables via classification. In classification, KNN model splits the data into training, validation, and test sets. This is the basis for prediction in KNN model. Compared to traditional statistical method, KNN model is more effective in big data for predicting dependent variable. The K-Fold Cross validation is an important tool for assessing the performance of machine learning in doing KNN algorithm given available data. According to Lei J [2], for model and tuning parameter selection, K-Fold Cross-Validation is one of the most well-known methods. This method is useful for statistics and machine learning. However, K-Fold Cross-Validation is more effective than traditional statistical methods when big data is implemented in parameter selection. Moreover, the basic idea of K-Fold Cross-Validation is to fit and evaluate each candidate model on separate datasets. Then, performance evaluation is unbiased. Covid 19 cases that threat people globally make countries collect nationwide data daily. Each country has big data of covid 19 cases. The quality of big data in each country depends on the quality of officers who manage the data. Hence, KNN and K-Fold Cross-Validation are expected to help in assessing how good big data has been collected.

**Address for Correspondence:* Kuntoro Kuntoro, Division of Biostatistics and Population Study, Airlangga University School of Public Health, Surabaya, Indonesia, Tel: 62318700289, E-mail: kuntoro@fkm.unair.ac.id

Copyright: © 2022 Kuntoro K. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Date of Submission: 02 December 2022, Manuscript No. jbmbs-22-79491; **Editor assigned:** 05 December, 2022, PreQC No. P-79491; **Reviewed:** 16 December 2022, QC No. Q-79491; **Revised:** 17 November 2022, Manuscript No. R-79491; **Published:** 28 December, 2022, DOI: 10.37421/2155-6180.2022.13.145

Theoretical frame work

Fix E and Hodges JL [3] first developed K-Nearest Neighbours (KNN) algorithm. It is included in a non-parametric classification method in statistics. Then Cover T and Hart P [4] expanded this method. It is used for classification and regression. In both methods, the input consists K-closest training examples in a data set. The output depends on whether KNN is used for classification or regression. Moreover, in KNN classification, the output is a class membership. An object is classified by a plurality vote of its neighbours, with the object being assigned to the class most common among its KNN (K is a positive integer, typically small). If K=1, then the object is simply assigned to the class of that single nearest neighbour. Furthermore, in KNN regression, the output is the property value for the object. This value is the average of the values of KNN.

Concept of classification

A prediction can be made based on current situation and past experienced as well. When a categorical class called a label is to be predicted given a number of variables called features, then we deal with classification. Moreover, an observation is assigned by a classifier with unknown class to a class based on similarity to other observations which are known their classes. Furthermore, observations whose classes are known which are used as basis for prediction are called data training sets. We use this data to train or teach our 3 classifier. Then the classifier is used for prediction of new data that is not known its class [5].

Classification with K-Nearest neighbours (KNN)

A classification algorithm is needed to make prediction for new observation. Then we use KNN classification algorithm. By this algorithm we want to find the K – Nearest or K-Most similar observations in our data training set. The results are used to make prediction for the new observation [5]. Concept of KNN can be explained by (Figures 1 and 2) as follows. Suppose an individual is classified into diseased (red dot) or healthy (green dot) individual based on Euclidean distance of two variables let say var_1 and var_2 . A new observation called unknown diagnosis (blue dot) is decided whether it is classified into diseased or healthy individual depending on the closeness to diseased or healthy one. The closeness is measured by Euclidean distance. The KNN to a

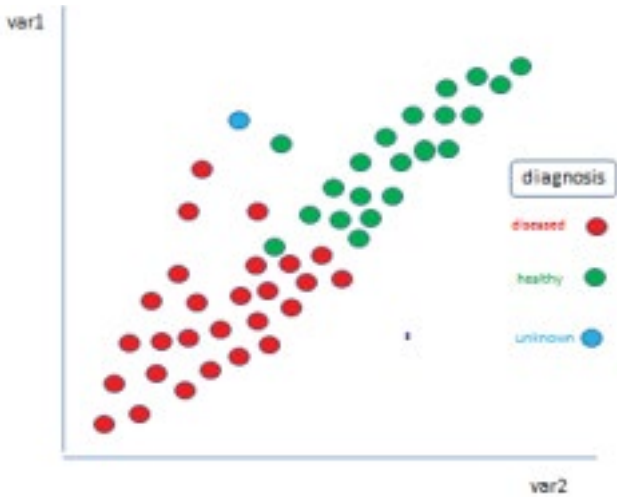


Figure 1. Scatter plot of diseased and healthy individuals based on two variables with a new observation (Modified from Timbers, et al. (2021)).

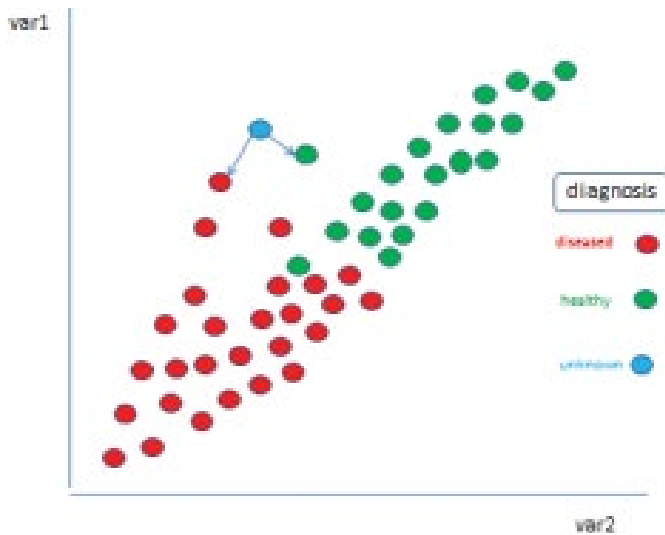


Figure 2. The closeness of a new observation to diseased individual or healthy individual (Modified from Timbers, et al. (2021)).

new observation is obtained by calculating the distance from new observation to each observation in training data. Then select the K – observations which have smallest distance values. Suppose we want to use K=5 neighbours to classify a new observation with certain values of two variables let say var_1 and var_2 shown as blue dot in (Figures 3 and 4). Then we calculate the distance between new observations to each observation in the training set. Moreover, we find K=5 neighbours which are nearest to new observation [5].

This is the equation for calculating distance.

$$Distance = \sqrt{(p\ var_1 - q\ var_1)^2 + (p\ var_2 - q\ var_2)^2} \dots\dots\dots (1)$$

Where p and q are two different observations, each has predictor variables var_1 and var_2 . Moreover, $p\ var_1$ and $q\ var_1$ are the values of var_1 for observations p and q. Furthermore, $p\ var_2$ and $q\ var_2$ are the values of var_2 for observations p and q.

Cross-Validation

Suppose $D = (x_1 - y_1)_{i=1}^n$ is the data drawn independently from a common distribution P on $R^p \times R$ that satisfies (Lei, 2020)

$$y_1 = f(x) + \epsilon_i \dots\dots\dots (2)$$

In which $f : R^p \rightarrow R$ is an unknown function and ϵ satisfies $E(\epsilon|X) = 0$ Moreover, we want to approximate f in such a way that the value of Y can be predicted for future observation of X. Furthermore, \hat{f} is an estimate of f . Then the quality of \hat{f} is evaluated by the predictive risk as follows Lei J [2].

$$Q(\hat{f}) = E[\ell(\hat{f}(X), Y)] \dots\dots\dots (3)$$

In which (X, Y) is a future random drawn from P, while $\ell(\dots)$ is a loss function. Furthermore, the square loss $\ell(y, \hat{y})$ is a typical example of $\ell(\dots)$ Lei J [2].

Validation by sample splitting can be implemented by using a unified notation for tuning parameter selection and model selection as well. Moreover, given a finite set of candidates $M = \{1, 2, \dots, M\}$, in which the meaning of each $m \in M$ can be a model or a tuning parameter value that depends on the context. Furthermore, we can estimate f by \hat{f}_m for each $m \in M$ using an estimating procedure specified by m .

Precision, recall, and F1 score

A classification report is part of classification algorithm like KNN. It is used to measure the quality of prediction concerning the number of predictions which is true or false. Predicting the metrics of a classification report uses true positives, false positives, true negatives and false Negatives. Precision is the ability of a classifier not to label an instance to be positive which is actually negative. For each class it is defined as the ratio of true positives to the sum of true and false positives [4]. The equation of precision is described as follows.

$$Precision = \frac{TP}{TP + FP} \dots\dots\dots (4)$$

In which TP is the number of True Positives and FP is the number of False Positives in a confusion matrix as follows.

Recall is the ability of a classifier to find all positive instances. For each class it is defined as the ratio of true positives to the sum of true positives and false negatives [6]. The equation of recall is described as follows.

$$Recall = \frac{TP}{TP + FN} \dots\dots\dots (5)$$



Figure 3. Confusion matrix (Source: Data to Fish, 2021).

```

IDLE Shell 3.9.4 - C:/Users/kuntoro001/Documents/MACHINE LEARNING 041121...
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:14:02) [MSC v.1920 64 bit (AMD64)] on win32
Type "help()" "copyright()" "credits()" or "license()" for more information.
>>> # Importing libraries needed
>>> # Import data set into Python script
>>> # Import pandas module
>>> df = pd.read_csv('D:/cov19/covid19.csv')
>>> df
   newscases  reproduce   ICU  hosp  class
0           631         0.94    215   856     0
1           410         0.64    219   823     0
2           395         0.78    245   829     0
3           257         0.71    245   826     0
4           270         0.69    244   712     0
...
13217      84369         0.98   11754  41936     1
13218      81666         0.98   11513  41309     1
13219      81218         0.98   11311  40814     1
13220      29656         0.98   11667  40036     1
13221      21888         0.98   11177  40392     1
[13222 rows x 5 columns]
>>> # splitting data set into training data and test data
>>> from sklearn.model_selection import train_test_split
>>> train, test = train_test_split(df, test_size = 0.3)
>>> x_train = train.drop('class', axis=1)
>>> y_train = train['class']
>>> x_test = test.drop('class', axis=1)
>>> y_test = test['class']
>>> # standardizing data set to be similar scale
>>> from sklearn.preprocessing import MinMaxScaler
>>> scaler = MinMaxScaler(feature_range=(0, 1))
>>> x_train_scaled = scaler.fit_transform(x_train)
>>> x_test_scaled = scaler.fit_transform(x_test)
>>> # preparing model and training
>>> from sklearn import neighbors
    
```

Figure 4. Preparing libraries, importing, standardizing, training - testing data sets, and classifying.

In which TP is the number of True Positives and FN is the number of False Negatives in a confusion matrix mentioned above. The F1 score is a weighted harmonic mean of precision and recall such that the best score is 1.0 and the worst score is 0.0. Generally speaking, F1 scores are lower than accuracy measures as they embed precision and recall into their calculation. As a rule of thumb, the weighted average of F1 should be used to compare classifier models, not global accuracy [6]. The equation of recall is described as follows.

F1Score = (2 * (Recall + Precision)) / (Recall + Precision) (6)

Using figure 3 we can calculate precision, recall, and F1 score

Precision = TP / (TP + FP) = 4 / (4 + 2) = 4 / 6 = 0.667

Recall = TP / (TP + FN) = 4 / (4 + 1) = 4 / 5 = 0.800

F1Score = (2 * (Recall + Precision)) / (Recall + Precision) = (2 * (0.667 * 0.800)) / (0.667 + 0.800) = 1.0672 / 1.4670 = 0.727

Methodology

A supervised machine learning approach using KNN and K-Fold Cross-Validation is implemented in this study. For learning process, data sets of covid 19 is obtained by downloading in website as follows. https://github.com/owid/covid-19-data/blob/master/public/data/latest/owid-covid-latest.csv Size of data sets is 132,645. Four variables, X are selected for KNN modelling, these are new case, reproduction rate (reprate), new case in ICU (Intensive Care Unit) or icu, and hospitalized new case (hosp). These variables are considered as predictors (X) used for predicting dependent variable (Y) as target. After data cleaning process, the complete data sets are to be 13,223. This data is stored in folder D of computer in csv format with name covid251121ver13.csv. For KFold Cross-Validation, variable names (new case, reprate, icu, and hosp) and class column are deleted and the file is stored as covid251121cvver13.csv. Then data sets reduce to be 13,222.

The steps for implementing KNN are as follows Zach [7]. First, downloading data sets from website mentioned above. Second, loading Python libraries such as Numpy, Pandas, Scikitlearn. Third, importing data sets into Python script using Pandas library with code df=pd.read_csv('D:\covid251121ver13.csv'). Fourth, standardizing data sets to adjust all X values in order to have the same scale to make better prediction. Fifth splitting data sets into training data and test data. Sixth, predicting the target (Y) with a KNN algorithm. Seventh, measuring accuracy of the model. The steps for implementing K-Fold Cross-Validation are as follows McCullem N [1]. First, loading Python libraries such as Numpy, Pandas, Scikit-learn. Second, importing data sets into Python script using Pandas library with code df=pd.read_csv('D:\covid251121cvver13.csv'). Third, performing K-Fold Cross-Validation in K-Fold Cross-Validation, when User Warning appears, which means the least populated class in y has only 1 member, which is less than n_splits=2, we replicate data sets nearly 20 times to be 264,441 data sets. Then we use n_splits=2, 3, 4, and 5.

Results and Discussion

K-Nearest neighbours for data sets with size 13,223

In implementing Python, preparing libraries is essential, these are Pandas, Numpy, and Scikit learn. Pandas is needed for importing data sets which is placed in folder of computer to Python script. Numpy is needed for numerical operation. Scikit-learn are needed when we use supervised machine learning. All data sets should be numbers not character symbols. It should be converted into number 1 or 0. All data sets should have the same scale. Hence, standardization is a step of implementing KNN to overcome these problems. Scikit-learn library provides standardization of data sets. Data sets which meet requirements enter further step, data training, data validating, and data testing.

We split data sets into two parts, training data dan test data. We choose 30% of data sets as training data, the rest as test data. Scikit-learn library provides splitting data sets into both data [7,8]. Further step we classify data sets into class with, one member, 2 members, 3 members, 4 members etc. based on distance closeness. The number of member is called K, in which K=1, 2, 3, 4, This is a basic concept of KNN. Scikit-learn library provides this process.

The last step is to predict target, Y (1=covid 19 cases who will be alive or 0=covid 19 cases who will die) based on a number of predictors X, new case, reprate, icu, and hosp using KNN model. Accuracy in prediction is measured by root mean squared error (rmse) for every number of class membership in K=1, 2, 3, 20. The results are plotted in a graph of rmse on different K. The highest rmse occurs in K=1 that means each class has one member. The value of rmse is 0.49499639066447704. As K increases to K=7 then rmse decreases. The lowest rmse occurs in K=7, in which each class has seven members. The value of rmse is 0.4035386704306. When K=8 to K=20 rmse increases, except in K=11 rmse is lower than rmse in K=10. We select maximum number of K is equal to 20 in which rmse is 0.4115685931061294. We omit K=1 with the highest value of rmse because we do not want each class only has one member. In this study, it is better to select K=7 in order to obtain the lowest value of rmse. In this case each class has seven members, the optimum number of members that gives the lowest value of rmse or the lowest error in prediction for KNN model. In (Predicting KNN model 7) shows plot of rmse on various values of K. In the beginning when K=1 rmse is high then rmse decreases sharply until the lowest value of rmse that occurs when K is equal to 7. Then rmse increases gradually until K is equal to 20 (Figures 5-7).

Classification report and confusion matrix

In KNN model classification report, the results display accuracy, precision, recall, F1 score according to target is equal to 0 that means prediction of covid 19 cases will die and target is equal to 1 that means prediction of covid 19

Figure 5. Predicting KNN model.

Figure 6. Predicting KNN model (continued).

cases will be alive given a number of predictors new case, replate, icu, and hosp.(Figure 9) shows that in predicting covid 19 cases will die have precision, recall, and F1 score respectively 0.63, 0.63, and 0.63 while in predicting covid 19 cases will be alive have precision, recall, and F1 score respectively 0.81, 0.81, and 0.81. Calculating the average of target, each of precision, recall, and F1 score has value 0.72, while in calculating the weighted 10 average of target, each of precision, recall, and F1 score has value 0.75. Weighted average is calculated by considering the proportion of sample size in each class target (0,1). The average accuracy of KNN algorithm is 0.75 or 75%. It is obtained in K=7 and the smallest rmse 0.4035386704306. In predicting target or variable Y (covid 19 will be alive or will die) given predictors new case, replate, icu, and hosp, KNN algorithm shows good performance with average accuracy 75%.

K-Fold Cross-Validation for data sets with size 13,222

(Figure 10) explains that we need to prepare Python libraries such as Pandas, Numpy, and Skicit-learn. The next step we import csv formatted file, covid261121cvver13, csv from folder D of computer. Then we clean datasets from non a number (nan) as requirement for further analysis. The following step is to perform K-Fold Cross-Validation. In this step, first of all we define predictors, X and response variable, Y. Second, we build model K Neighbours Classifier. Third we define cross-validation method. (Figure 11) shows mean and standard deviation of accuracy are respectively 0.054 and 0.001. Mean of accuracy is too small although its standard deviation is too small. Small of standard deviation means accuracy across testing is stable. Mean of accuracy too small because the least populated class in response variable, Y has only one members, which is less than n_splits=2. We already have 13,223 records in original data sets. We try to increase the number of records almost 20 times to be 264,441 as replicated data sets. Then we cross-validate again. These are the results.

K-Fold Cross -Validation for data sets with size 264,441

(Figures 12 and 13) show that by increasing data set almost 20 times, the mean and standard deviation of accuracy respectively 0.998 and 0.000 by splitting data set into 2 parts. Then we split data into 3, 4, and 5 parts. The mean and standard deviation of accuracy are the same. By splitting data set into 2, 3, 4, and 5 parts, each accuracy is high with no variation.

K-Nearest neighbours for data sets with size 264,441

The highest rmse occurs in K=1 that means each class has one member. The value of rmse is 0.28039074233295097. As K increases to K=12 then rmse decreases. The lowest rmse occurs in K=12, in which each class has twelve members. The value of rmse is 0.20697044925406405. Then K=13 to K=20 rmse increases. We select maximum number of K is equal to 20 in which rmse is 0.23044759082966193. We omit K=1 with the highest value of rmse because we do not want each class only has one member. In this study, it is better to select K=12 in order to obtain the lowest value of rmse. In this case each class has twelve members, the optimum number of members that gives the lowest value of rmse or the lowest error in prediction for KNN model

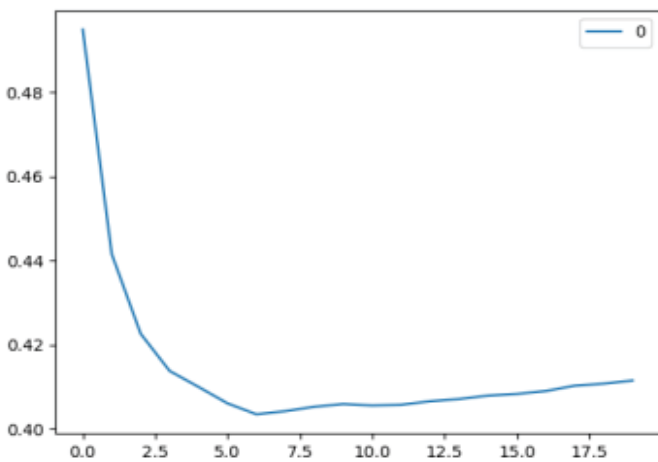


Figure 7. Plotting rsme on various K values.

```

IDLE Shell 3.9.4 C:/Users/kuntoro01/Documents/MACHINE LEARNING 041121...
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
D64)] on win32
Type "help()", "copyright()", "credits()" or "license()" for more information.
>>> # K NEAREST NEIGHBORS 13223 covid 19 301121
>>> # Python libraries needed
>>> # import data set into Python script
>>> import numpy as np
>>> import pandas as pd
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
>>> from pandas import read_csv
>>> df=pd.read_csv('D:\covid261121cvver13.csv')
>>> df
   newcase  replate   icu  hosp  class
0         0         0     215   856     0
1         1         0     219   823     0
2         1         0     245   829     0
3         3         0     245   826     0
4         4         0     244   712     0
...
13217  84369   0.98  11754  41936     1
13218  81566   0.98  11167  40392     1
13219  96665   0.98  11319  40981     1
13220  29886   0.98  11167  40396     1
13221  21968   0.98  11177  40392     1

[[13222 rows x 5 columns]]
>>> # standardizing data set to be similar scale
>>> from sklearn.preprocessing import StandardScaler
>>> scaler=StandardScaler()
>>> scaler.fit(df.drop('class',axis=1))
StandardScaler()
>>> scaled_features=scaler.transform(df.drop('class',axis=1))
>>> scaled_data=pd.DataFrame(scaled_features,columns=df.drop('class',axis=1).col
umns)
>>> # splitting data set into training data and test data
>>> from sklearn.model_selection import train_test_split
>>> x=scaled_data
>>> y=df['class']
>>> x_training,x_test,y_training,y_test=train_test_split(x,y,test_size=0.3)
Ln: 26 Col: 0
  
```

Figure 8. Process of classification report.

```

IDLE Shell 3.9.4 C:/Users/kuntoro01/Documents/MACHINE LEARNING 041121...
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
D64)] on win32
Type "help()", "copyright()", "credits()" or "license()" for more information.
>>> # K NEAREST NEIGHBORS 13223 covid 19 301121
>>> # Python libraries needed
>>> # import data set into Python script
>>> import numpy as np
>>> import pandas as pd
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
>>> from pandas import read_csv
>>> df=pd.read_csv('D:\covid261121cvver13.csv')
>>> df
   newcase  replate   icu  hosp  class
0         0         0     215   856     0
1         1         0     219   823     0
2         1         0     245   829     0
3         3         0     245   826     0
4         4         0     244   712     0
...
13217  84369   0.98  11754  41936     1
13218  81566   0.98  11167  40392     1
13219  96665   0.98  11319  40981     1
13220  29886   0.98  11167  40396     1
13221  21968   0.98  11177  40392     1

[[13222 rows x 5 columns]]
>>> # standardizing data set to be similar scale
>>> from sklearn.preprocessing import StandardScaler
>>> scaler=StandardScaler()
>>> scaler.fit(df.drop('class',axis=1))
StandardScaler()
>>> scaled_features=scaler.transform(df.drop('class',axis=1))
>>> scaled_data=pd.DataFrame(scaled_features,columns=df.drop('class',axis=1).col
umns)
>>> # splitting data set into training data and test data
>>> from sklearn.model_selection import train_test_split
>>> x=scaled_data
>>> y=df['class']
>>> x_training,x_test,y_training,y_test=train_test_split(x,y,test_size=0.3)
>>> # processing nearest neighbors
>>> from sklearn.neighbors import KNeighborsClassifier
>>> model=KNeighborsClassifier(n_neighbors=1)
>>> model.fit(x_training, y_training)
>>> predictions=model.predict(x_test)
>>> from sklearn.metrics import classification_report
>>> from sklearn.metrics import confusion_matrix
>>> print(classification_report(y_test,predictions))
 precision    recall  f1-score   support

 0.         0.63         0.63         0.63         1350
 1.         0.81         0.81         0.81         2609

 accuracy         0.72         0.72         0.75         3967
 macro avg         0.72         0.75         0.75         3967
 weighted avg         0.75         0.75         0.75         3967

[[ 661 407]
 [ 806 2104]]
Ln: 26 Col: 0
  
```

Figure 9. Accuracy, precision, recall, and F1 score.

```

cross validation 13223 covid 19 301121.py - C:/Users/kuntoro01/Documents/M...
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
D64)] on win32
Type "help()", "copyright()", "credits()" or "license()" for more information.
>>> # CROSS VALIDATION 13223 COVID 19 301121
>>> # loading Python libraries needed
>>> from numpy import mean
>>> from sklearn.model_selection import cross_val_score
>>> from sklearn.model_selection import RepeatedStratifiedKFold
>>> from sklearn.neighbors import KNeighborsClassifier
>>> from sklearn.preprocessing import LabelEncoder
>>> from sklearn.preprocessing import StandardScaler
>>> # import data set into Python script
>>> from pandas import read_csv
>>> import numpy import nan
>>> df= read_csv('D:\covid261121cvver13.csv', header=None, na_values='?')
>>> df.dropna(inplace=True)
>>> print(df.shape)
(13222, 4)
>>> # performing k fold cross validation
>>> # define predictor and response variables
>>> data = df.values
>>> X, y = data[:, :-1], data[:, -1]
>>> X = X.astype('float32')
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> # building model
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> frame = StandardScaler()
>>> model = KNeighborsClassifier()
>>> pipeline = Pipeline(steps=[('s', frame), ('m', model)])
>>> # define cross validation method to use
>>> cv = RepeatedStratifiedKFold(n_splits=2, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy',cv=cv, n_jobs=
Warning (from warnings module):
  File "C:\Users\kuntoro01\AppData\Local\Programs\Python\Python39\lib\site-packa
warnings.warn(
UserWarning: The least populated class in y has only 1 members, which is less th
Warning (from warnings module):
Ln: 1 Col: 0
  
```

Figure 10. Preparing libraries, importing, cleaning non a number (nan), performing cross validation.

(Figures 14 -16). Shows plot of rmse on various values of K. In the beginning when K=1 rmse is high then rmse decreases sharply until the lowest value of rmse that occurs when K is equal to 12. Then rmse increases gradually until K is equal to 20.

Classification report and confusion matrix

In KNN model classification report, the results display accuracy, precision, recall, F1 score according to target is equal to 0 that means prediction of covid 19 cases will die and target is equal to 1 that means prediction of covid 19 cases will be alive given a number of predictors new case, replate, icu, and hosp. (Figures 17 and 18) show process of classification report (Figure 8). (Figure 18) shows that in predicting covid 19 cases will die have precision,

```

C:\Users\kuntoro1\Documents\M...
File Edit Format Run Options Window Help
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
Type "help", "copyright", "credits" or "license()" for more information.
>>> # import data set into Python script
>>> from pandas import read_csv
>>> from numpy import nan
>>> df = read_csv('D:\covid271121ovbigdata.csv', header=None, na_values='?')
>>> df.dropna(inplace=True)
>>> print(df.shape)
(13222, 4)
>>> # performing k fold cross validation
>>> # define predictor and response variables
>>> data = df.values
>>> X, y = data[:, 1:-1], data[:, -1]
>>> X = X.astype('float32')
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> # building model
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> trans = StandardScaler()
>>> model = KNeighborsClassifier()
>>> pipeline = Pipeline(steps=[('t', trans), ('m', model)])
>>> # define cross validation method to use
>>> cv = RepeatedStratifiedKFold(n_splits=2, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
>>> # print mean and standard deviation of accuracy
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.954 (0.001)
>>>
Ln 1 Col 0

```

Figure 11. Mean and standard deviation of accuracy.

```

IDLE Shell 3.9.4 C:\Users\kuntoro1\Documents\MACHINE LEARNING 041121...
File Edit Shell Debug Options Window Help
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
Type "help", "copyright", "credits" or "license()" for more information.
>>> # loading Python libraries needed
>>> from sklearn.model_selection import cross_val_score
>>> from sklearn.model_selection import RepeatedStratifiedKFold
>>> from sklearn.neighbors import KNeighborsClassifier
>>> from sklearn.preprocessing import StandardScaler
>>> from sklearn.pipeline import Pipeline
>>> # import data set into Python script
>>> from pandas import read_csv
>>> from numpy import nan
>>> df = read_csv('D:\covid271121ovbigdata.csv', header=None, na_values='?')
>>> df.dropna(inplace=True)
>>> print(df.shape)
(26440, 4)
>>> # performing k fold cross validation
>>> # define predictor and response variables
>>> data = df.values
>>> X, y = data[:, 1:-1], data[:, -1]
>>> X = X.astype('float32')
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> # building model
>>> trans = StandardScaler()
>>> model = KNeighborsClassifier()
>>> pipeline = Pipeline(steps=[('t', trans), ('m', model)])
>>> # define cross validation method to use
>>> cv = RepeatedStratifiedKFold(n_splits=2, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
>>> # print mean and standard deviation of accuracy
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.9 (0.000)
>>> # splitting n into 3 to 5
>>> cv1 = RepeatedStratifiedKFold(n_splits=3, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv1, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.998 (0.000)
>>> cv2 = RepeatedStratifiedKFold(n_splits=4, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv2, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.999 (0.000)
>>> cv3 = RepeatedStratifiedKFold(n_splits=5, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv3, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.998 (0.000)
>>>
Ln 32 Col 66

```

Figure 12. Mean and standard deviation of accuracy for n_splits=2.

```

IDLE Shell 3.9.4 C:\Users\kuntoro1\Documents\MACHINE LEARNING 041121...
File Edit Shell Debug Options Window Help
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
Type "help", "copyright", "credits" or "license()" for more information.
>>> # loading Python libraries needed
>>> from sklearn.model_selection import cross_val_score
>>> from sklearn.model_selection import RepeatedStratifiedKFold
>>> from sklearn.neighbors import KNeighborsClassifier
>>> from sklearn.preprocessing import StandardScaler
>>> from sklearn.pipeline import Pipeline
>>> # import data set into Python script
>>> from pandas import read_csv
>>> from numpy import nan
>>> df = read_csv('D:\covid271121ovbigdata.csv', header=None, na_values='?')
>>> df.dropna(inplace=True)
>>> print(df.shape)
(26440, 4)
>>> # performing k fold cross validation
>>> # define predictor and response variables
>>> data = df.values
>>> X, y = data[:, 1:-1], data[:, -1]
>>> X = X.astype('float32')
>>> y = LabelEncoder().fit_transform(y.astype('str'))
>>> # building model
>>> trans = StandardScaler()
>>> model = KNeighborsClassifier()
>>> pipeline = Pipeline(steps=[('t', trans), ('m', model)])
>>> # define cross validation method to use
>>> cv = RepeatedStratifiedKFold(n_splits=2, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
>>> # print mean and standard deviation of accuracy
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.998 (0.000)
>>> # splitting n into 3 to 5
>>> cv1 = RepeatedStratifiedKFold(n_splits=3, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv1, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.998 (0.000)
>>> cv2 = RepeatedStratifiedKFold(n_splits=4, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv2, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.999 (0.000)
>>> cv3 = RepeatedStratifiedKFold(n_splits=5, n_repeats=3, random_state=1)
>>> n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv3, n_jobs=-1)
>>> print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.998 (0.000)
>>>
Ln 32 Col 66

```

Figure 13. Mean and standard deviation of accuracy for n_splits=3, 4 and 5.

recall, and F1 score respectively 0.89, 0.89, 0.89. These values are higher than the values in (Figure 9). These values increase when size of data sets increases. Moreover, predicting covid 19 cases will be alive has precision, recall, and F1 score respectively 0.94, 0.94, 0.94. Again, these values are higher than the values in (Figure 19). These values increase when size of data sets increases. Calculating the average of target, each of precision, recall, and F1 score has value 0.91, while in calculating the weighted average of target, each of precision, recall, and F1 score has value 0.92. Weighted average is calculated by considering the proportion of sample size in each class target (0, 1). Compared to original data sets, replicated data sets gives higher average accuracy (92% vs. 75%) more members of class (12 vs. 7), and less rmse (0.20697044925406405 vs. 0.4035386704306). Hence, increasing size of data sets will increase the performance of KNN algorithm in predicting target or variable Y.

The K Fold Cross-Validation using replicated data sets compared to original data sets gives higher mean accuracy and lower standard deviation (0.998 and 0.000 vs. 0.054 and 0.001). The number of class to be splitted in replicated data sets is more than original data sets (2, 3, 4, 5 vs. only 2). Original data sets can only be splitted not more than 2. Original data sets will give smaller members in each class. When in a class has only one member which is less than minimum n_split=2, it will give less performance of K-Fold Cross-Validation indicated by smaller mean accuracy. However, K-Fold Cross-Validation as supervised machine learning still gives good performance in variability perspective indicated by small standar deviation. The KNN and K-Fold Cross-Validation algorithms do not give similar results although their performances are evaluated in the same data sets. However, the results of both algorithms are in line. In higher size of data sets both give in line performance in accuracy perspective. In KNN algorithm variable names such as new case, reparate, icu and hosp to be included in data sets and we should add class column for placing target to be predicted. In K-Fold Cross Validation algorithm variable names should be deleted so that size of data sets is reduced one and no additional class column [9,10].

```

IDLE Shell 3.9.4 C:\Users\kuntoro1\Documents\MACHINE LEARNING 041121...
File Edit Shell Debug Options Window Help
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
Type "help", "copyright", "credits" or "license()" for more information.
>>> # Python libraries needed
>>> # K NEAREST NEIGHBORS SIZE 264441 COVID 19 01 DECEMBER 2021
>>> # import data set into Python script
>>> import pandas as pd
>>> from pandas import read_csv
>>> df = pd.read_csv('D:\covid271121ovbigdata.csv')
>>> df

```

	newcase	reparate	icu	hosp	class
0	531	0.98	218	186	0
1	410	0.84	239	823	0
2	395	0.75	245	829	0
3	257	0.71	245	826	0
4	270	0.69	244	712	0
...
264435	84369	0.98	11784	41936	1
264436	81566	0.98	11533	41309	1
264437	96665	0.98	11319	40581	1
264438	29866	0.98	11167	40336	1
264439	21968	0.98	11177	40392	1

```

>>> # splitting data set into training data and test data
>>> from sklearn.model_selection import train_test_split
>>> train, test = train_test_split(df, test_size = 0.3)
>>> X_train = train.drop('class', axis=1)
>>> y_train = train['class']
>>> X_test = test.drop('class', axis = 1)
>>> y_test = test['class']
>>> # standardizing data set to be similar scale
>>> from sklearn.preprocessing import MinMaxScaler
>>> scaler = MinMaxScaler(feature_range=(0, 1))
>>> train_scaled = scaler.fit_transform(X_train)
>>> X_train = pd.DataFrame(X_train_scaled)
>>> X_test_scaled = scaler.fit_transform(X_test)
>>> X_test = pd.DataFrame(X_test_scaled)
>>> # processing nearest neighborhood
>>> from sklearn import neighbors
>>>
Ln 94 Col 0

```

Figure 14. Preparing libraries, importing, standardizing, training - testing data set, and classifying.

```

IDLE Shell 3.9.4 C:\Users\kuntoro1\Documents\MACHINE LEARNING 041121...
File Edit Shell Debug Options Window Help
Python 3.9.4 (tags/v3.9.4:1f2e308, Apr 6 2021, 13:40:21) [MSC v.1928 64 bit (AM
Type "help", "copyright", "credits" or "license()" for more information.
>>> # preparing plot of rmse on different k
>>> import matplotlib.pyplot as plt
>>> # plotting the rmse values against k value
>>> rmse = pd.DataFrame(rmse_val) # show curve
>>> rmse.plot()
>>> plt.show()
>>>
Ln 94 Col 0

```

Figure 15. Predicting KNN model.

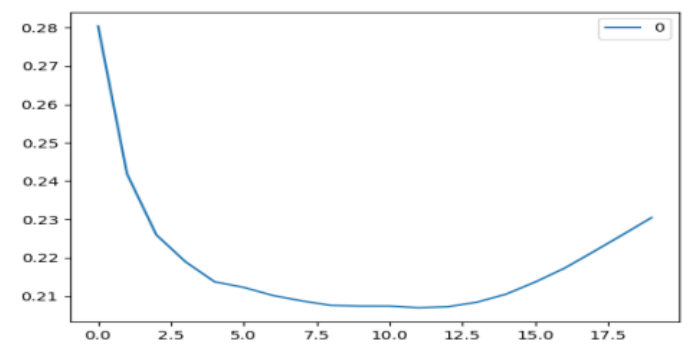


Figure 16. Plotting rmse on various k values.

```

IDLE Shell 3.9.4 - C:/Users/kuntoro01/Documents/MACHINE LEARNING 041121...
File Edit Shell Debug Options Window Help
Python 3.9.4 (tags/v3.9.4:11f2e306, Apr 6 2021, 13:40:21) [MSC v.1926 64 bit (AMD64)] on win32
>>> # K NEAREST NEIGHBORS 264441 covid 19 011221
>>> # Python libraries needed
>>> # import data set into python script
>>> import numpy as np
>>> import pandas as pd
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
>>> from pandas import read_csv
>>> df=pd.read_csv('D:\covid271121bigdata.csv')
>>> df

```

	newcase	retrate	icu	hoop	class
0	531	0.96	215	856	0
1	418	0.84	219	823	0
2	395	0.75	245	829	0
3	257	0.71	245	826	0
4	270	0.69	244	712	0
...
264435	84369	0.98	11754	41936	1
264436	81566	0.98	11533	41309	1
264437	96665	0.98	11319	40881	1
264438	29856	0.98	11167	40392	1
264439	21960	0.98	11177	40392	1

Figure 17. Process of classification report.

```

IDLE Shell 3.9.4 - C:/Users/kuntoro01/Documents/MACHINE LEARNING 041121...
264437 96665 0.98 11319 40881 1
264438 29856 0.98 11167 40392 1
264439 21960 0.98 11177 40392 1
[264440 rows x 5 columns]
>>> # standardizing data set to be similar scale
>>> from sklearn.preprocessing import StandardScaler
>>> scaler=StandardScaler()
>>> scaler.fit(df.drop('class',axis=1))
StandardScaler()
>>> scaled_features=scaler.transform(df.drop('class',axis=1))
>>> scaled_data=pd.DataFrame(scaled_features,columns=df.drop('class',axis=1).columns)
>>> # splitting data set into training data and test data
>>> from sklearn.model_selection import train_test_split
>>> x_scaled_data
>>> y=df['class']
>>> x_training,x_test,y_training,y_test=train_test_split(x_scaled_data,y,train_size=0.3)
>>> # processing nearest neighbors
>>> from sklearn.neighbors import KNeighborsClassifier
>>> model=KNeighborsClassifier(n_neighbors=1)
>>> model.fit(x_training,y_training)
KNeighborsClassifier(n_neighbors=1)
>>> predictions=model.predict(x_test)
>>> from sklearn.metrics import classification_report
>>> from sklearn.metrics import confusion_matrix
>>> print(classification_report(y_test,predictions))

```

	precision	recall	f1-score	support
0	0.89	0.89	0.89	27333
1	0.94	0.94	0.94	81999
accuracy			0.92	79332
macro avg	0.91	0.91	0.91	79332
weighted avg	0.92	0.92	0.92	79332

```

>>> print(confusion_matrix(y_test,predictions))
[[2487  216]
 [ 3036 48963]]
>>>

```

Figure 18. Accuracy, precision, recall, and F1 score.

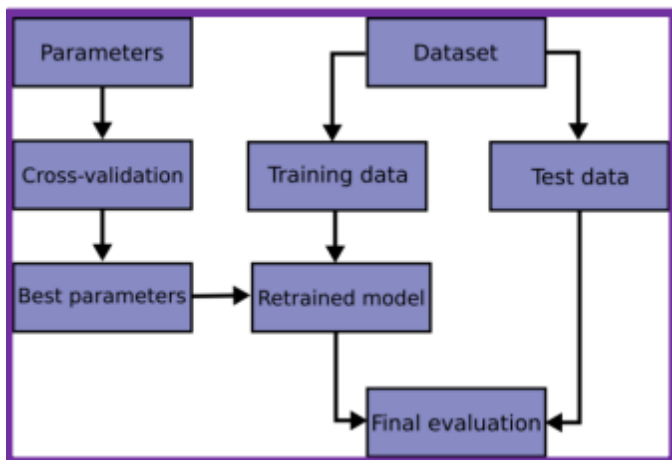


Figure 19. KNN and K-Fold cross-validation (Source: Scikit learn, 2021).

Figure 19 shows the relationship between KNN and K-Fold Cross-Validation algorithms. As supervised machine learning, both algorithms work

in their own ways that need partly the same Python libraries and partly the different Python libraries depending on the requirements. Functionally both algorithms have relationship in which K-Fold Cross-Validation reinforces KNN performance. Moreover, a best parameter obtained from K-Fold Cross-Validation is used to retrain the model. That means performance of both algorithms are in line. Better performance of K-Fold Cross-Validation will result in better performance of KNN. A challenge for algorithm developers for merging both algorithms in one entity algorithm by considering strength and weakness of each algorithm.

Conclusion and Recommendation

Performance of KNN algorithm in original data sets shows good average accuracy 75%, while performance of KNN in replicated data sets shows better average accuracy 92%. The optimum class members using original data sets are K=7, while the optimum class members using replicated data sets are higher, K=12. Performance of K- Cross-Validation algorithm in original data sets shows very small mean accuracy 0.054, while performance of K-Cross-Validation in replicated data sets shows very high mean accuracy 0.998 and very small standard deviation accuracy 0.000. When using big data from various sources it is recommended to apply supervised machine learning with appropriate Python libraries which have high performance in eliminating non a number (nan) and messy data sets. It is recommended to develop combined and integrated KNN and K-Fold Cross-Validation algorithm which can reinforce the weakness of each algorithm.

References

- Mccullum, Nick. "K nearest neighbors in python - A step-by-step guide." (2021).
- Lei, Jing. "Cross-validation with confidence." *J Am Stat* 115 (2020): 1978-1997.
- Fix, Evelyn and Joseph Lawson Hodges. "Discriminatory analysis, nonparametric discrimination: Consistency properties." *Int Stat Rev* 57 (1989): 238-247.
- Cover, Thomas and Peter Hart. "Nearest neighbor pattern classification." *IEEE Trans Inform Theory* 13 (1957): 21 - 27.
- Timbers, Tiffany, Trevor Campbell and Melissa Lee. " Data science: A first introduction." *CRC Press* (2022).
- Muthukrishnan. "Understanding the classification report through sklearn." (2021).
- Zach. " K-Fold Cross Validation in Python." (2020).
- <https://www.nickmccullum.com/python-machine-learning/k-nearest-neighbors-python/>
- <https://github.com/owid/covid-19-data/blob/master/public/data/latest/owid-covid-latest.csv>
- Learn, Scikit. " Cross-validation: Evaluating estimator performance." (2022).

How to cite this article: Kuntoro, Kuntoro. "K-Nearest Neighbours and K-Fold Cross Validation for Big Data of Covid 19." *J Biom Biosta* 13 (2022): 145.