# Machine Learning Applications in Healthcare System

S.Vasundhara

Department of Humanities and Mathematics
G.Narayanamma Institute of Technology& Science
Shaikpet Hyderabad

**Abstract:**

Machine learning is the study of computer algorithms that improves automatically through experiences and by the usage of data machine learning algorithms buildamodelbased on sampledataknown as training data in ordertomakepredictions ordecisions without explicitly programmed todosothey are used in widevarietyof applications such as medicine  E-mail filtering speech recognition and computer vision ,where it is difficult o runfeasible to develop, conventional algorithms to perform the needs tasks.Machine learning is a growing technology which enables computers to learn automatically from past data. Machine learning uses various algorithms for constructing mathematical models and making better predictions using the existed data. In this paper we have used kaggledata set of heartdisease whichcontains patients data with different health issues and designs and evaluates incrementallearning solutions by using KNN algorithm to predict the patient who is suffering with heart disease.

Key words: Machine learning,Data Analytics, Algorithms in Machine Learning, literaturereview, applicationsKNN algorithm

**Introduction**:

History of Machine Learning

Before some years (about 40-50 years), machine learning was science fiction, but today it is the part of our daily life. Machine learning is making our day to day life easy from self-driving cars to Amazonvirtualassistant "Alexa". However, the idea behind machine learning is so old and has a long history. Below some milestones are given which have occurred in the history of machine learning:
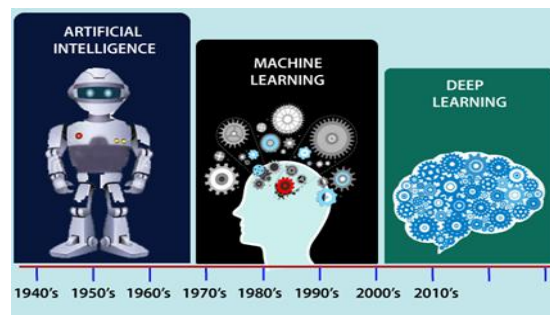


*Figure 1*

Machine Learning:

In the real world, we are surrounded by humans who can learn everything from their experiences with their learning capability, and we have computers or machines which work on our instructions. But can a machine also learn from experiences or past data like a human So here comes the role of Machine Learning.
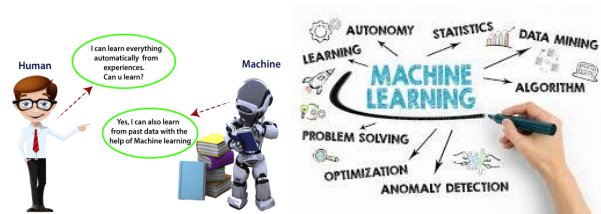


*Figure 2*

Machine Learning Enables a machine to automatically learn from data to improve performance fromexperience and predict things without being explicitly performed

With the help of sample historical data, which is known as training data, machine learning algorithms build a mathematical model that helps in making predictions or decisions without being explicitly programmed. Machine learning brings computer science and statistics together for creating predictive models. Machine learning constructs or uses the algorithms that learn from historical data. The more we will provide the information, the higher will be the performance.A machine   has the ability to learn if it can improve its performance by gaining more data.

A Machine Learning system learns from historical data , builds the prediction models and whenever it receives new data predicts the output for it. e accuracy of predicted output depends upon the amount of data, as the huge amount of data helps to build a better model which predicts the output more accurately.

Suppose we have a complex problem, where we need to perform some predictions, so instead of writing a code for it, we just need to feed the data to generic algorithms, and with the help of these algorithms, machine builds the logic as per the data and predict the output. Machine learning has changed our way of thinking about the problem. The below block diagram explains the working of Machine Learning algorithm:
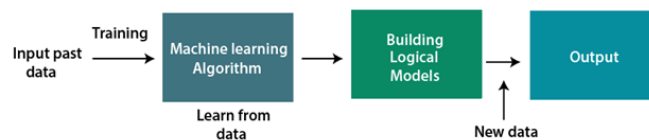


*Figure 3*

Features of Machine Learning:

- o   Machine learning uses data to detect various patterns in a given dataset.

o It can learn from past data and improve automatically.

o It is a data-driven technology.

o Machine learning is much similar to data mining as it also deals with the huge amount of the data.

Need for Machine Learning:The need for machine learning is increasing day by day. The reason behind the need for machine learning is that it is capable of doing tasks that are too complex for a person to implement directly. As a human, we have some limitations as we cannot access the huge amount of data manually, so for this, we need some computer systems and here comes the machine learning to make things easy for us.

We can train machine learning algorithms by providing them the huge amount of data and let them explore the data, construct the models, and predict the required output automatically. The performance of the machine learning algorithm depends on the amount of data, and it can be determined by the cost function. With the help of machine learning, we can save both time and money.

The importance of machine learning can be easily understood by its uses cases, Currently, machine learning is used in  self-driving cars,cyber fraud detection,face recognition,and friend suggestion by Facebook, etc. Various top companies such as Netflix and Amazon have build machine learning models that are using a vast amount of data to analyze the user interest and recommend product accordingly.

**Following are some key points which show the importance of Machine Learning:**

o Rapid increment in the production of data

o Solving complex problems, which are difficult for a human

o Decision making in various sector including finance

o Finding hidden patterns and extracting useful information from data.

<div align="center"><b>Classification of Machine Learning</b></div>

At a broad level, machine learning can be classified into three types:

1. **Supervised learning**

2. **Unsupervised learning**
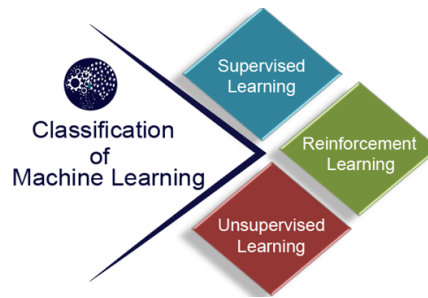
3. **Reinforcement learning**

*Figure 4*

## 1.Supervised Learning

Supervised learning is a type of machine learning method in which we provide sample labeled data to the machine learning system in order to train it, and on that basis, it predicts the output.The system creates a model using labeled data to understand the datasets and learn about each data, once the training and processing are done then we test the model by providing a sample data to check whether it is predicting the exact output or not.

The goal of supervised learning is to map input data with the output data. The supervised learning is based on supervision, and it is the same as when a student learns things in the supervision of the teacher. The example of supervised learning is **spam filtering**.

Supervised learning can be grouped further in two categories of algorithms:

- o **Classification**
- o **Regression**

## 2. Unsupervised Learning

Unsupervised learning is a learning method in which a machine learns without any supervision.

The training is provided to the machine with the set of data that has not been labeled, classified, or categorized, and the algorithm needs to act on that data without any supervision. The goal of unsupervised learning is to restructure the input data into new features or a group of objects with similar patterns.

In unsupervised learning, we don't have a predetermined result. The machine tries to find useful insights from the huge amount of data. It can be further classifieds into two categories of algorithms:

- o **Clustering**
- o **Association**

## 3. Reinforcement Learning

Reinforcement learning is a feedback-based learning method, in which a learning agent gets a reward for each right action and gets a penalty for each wrong action. The agent learns automatically with

these feedbacks and improves its performance. In reinforcement learning, the agent interacts with the environment and explores it. The goal of an agent is to get the most reward points, and hence, it improves its performance.

The robotic dog, which automatically learns the movement of his arms, is an example of Reinforcement learning.

**Classification Algorithm in Machine Learning**

Supervised Machine Learning algorithm can be broadly classified into Regression and Classification Algorithms. In Regression algorithms, we have predicted the output for continuous values, but to predict the categorical values, we need Classification algorithms.

 Classification Algorithm:

The Classification algorithm is a Supervised Learning technique that is used to identify the category of new observations on the basis of training data. In Classification, a program learns from the given dataset or observations and then classifies new observation into a number of classes or groups. Such as, Yes or No, 0 or 1, Spam or Not Spam, cat or dog**,** etc. Classes can be called as targets/labels or categories.

Unlike regression, the output variable of Classification is a category, not a value, such as "Green or Blue", "fruit or animal", etc. Since the Classification algorithm is a Supervised learning technique, hence it takes labeled input data, which means it contains input with the corresponding output.

In classification algorithm, a discrete output function(y) is mapped to input variable(x)

y=f(x), where y = categorical output

The best example of an ML classification algorithm is Email Spam Detector.

The main goal of the Classification algorithm is to identify the category of a given dataset, and these algorithms are mainly used to predict the output for the categorical data.

Classification algorithms can be better understood using the below diagram. In the below diagram, there are two classes, class A and Class B. These classes have features that are similar to each other and dissimilar to other classes.
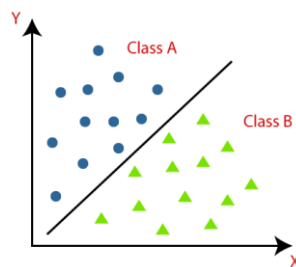


*Figure 5*

The algorithm which implements the classification on a dataset is known as a classifier. There are two types of Classifications:

- o **Binary Classifier:** If the classification problem has only two possible outcomes, then it is called asBinary Classifier.
  **Examples:** YES or NO, MALE or FEMALE, SPAM or NOT SPAM, CAT or DOG, etc.

- o **Multi-class Classifier:** If a classification problem has more than two outcomes, then it is called as Multi-class Classifier.
  **Example:** Classifications of types of crops, Classification of types of music.

Types of ML Classification Algorithms:

Classification Algorithms can be further divided into the Mainly two categories:

- o **Linear Models**

    - o Logistic Regression

    - o Support Vector Machines

- o **Non-linear Models**

    - o K-Nearest Neighbors

    - o Kernel SVM

    - o Naïve Bayes

    - o Decision Tree Classification

    - o Random Forest Classification

**K-Nearest Neighbor(KNN) Algorithm for Machine Learning**

- o K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.

- o K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

- o K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.

- o K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.

- o K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data.

o It is also called a **lazy learner algorithm** because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.

o KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

o **Example:** Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.

o category.



KNN Classifier

Input value          Predicted Output

*Figure 6*

Why do we need a K-NN Algorithm?

Suppose there are two categories, i.e., Category A and Category B, and we have a new data point x1, so this data point will lie in which of these categories. To solve this type of problem, we need a K-NN algorithm. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:
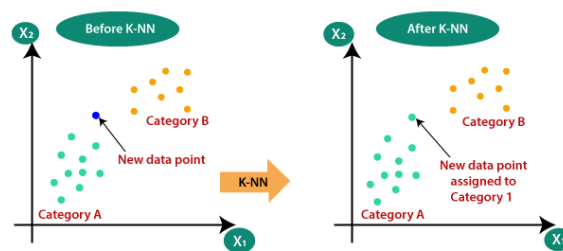


*Figure 7*

How does K-NN work?

The K-NN working can be explained on the basis of the below algorithm:

o **Step-1:** Select the number K of the neighbors

o **Step-2:** Calculate the Euclidean distance of **K number of neighbors**

o **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.

o **Step-4:** Among these k neighbors, count the number of the data points in each category.

o **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.

o **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:
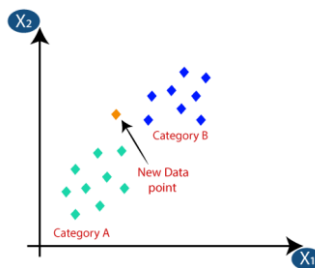


*Figure 8*

o Firstly, we will choose the number of neighbors, so we will choose the k=5.

o Next, we will calculate the **Euclidean distance** between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:



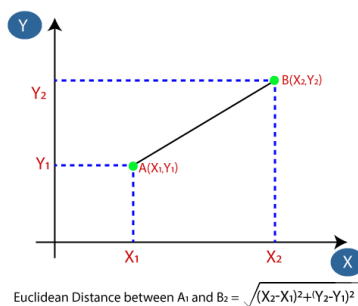Euclidean Distance between $A_1$ and $B_2 = \sqrt{(X_2-X_1)^2+(Y_2-Y_1)^2}$

*Figure 9*

o By calculating the Euclidean distance, we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:
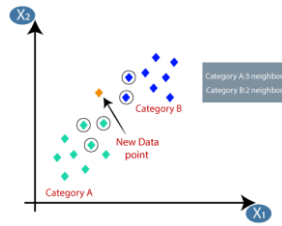
*Figure 10*

- o As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.

How to select the value of K in the K-NN Algorithm?

Below are some points to remember while selecting the value of K in the K-NN algorithm:

- o There is no particular way to determine the best value for "K", so we need to try some values to find the best out of them. The most preferred value for K is 5.

- o A very low value for K such as K=1 or K=2, can be noisy and lead to the effects of outliers in the model.

- o Large values for K are good, but it may find some difficulties.

Advantages of KNN Algorithm:

- o It is simple to implement.

- o It is robust to the noisy training data

- o It can be more effective if the training data is large.

Disadvantages of KNN Algorithm:

- o Always needs to determine the value of K which may be complex some time.

- o The computation cost is high because of calculating the distance between the data points for all the training samples.

**Method**

Python implementation of the KNN algorithm

To do the Python implementation of the K-NN algorithm, we will use the problem  and dataset which related to heart disease to  improve the performance of the model. Below is the problem description:

**Problem for K-NN Algorithm:** The data base contains of 5patients with 14 attributesof health conditions like age,sex,cp,trtbs,fbs,etc….The goalfield is here toverify the person who is suffering

with heart disease.So for this problem, we have a dataset that contains patient's information through the Cleveland database. The dataset contains lots of information using the information to build a model. Below is the dataset:

| Age | Sex | Cp | Trtbs | Chol | Fbs | Resteg | Thalach | Exng | Oldpeak | Slp | Caa | Thall | output |
|-----|-----|-----|-------|------|-----|--------|---------|------|---------|-----|-----|-------|--------|
| 63 | 1 | 3 | 145 | 233 | 1 | 0 | 150 | 0 | 2.3 | 0 | 0 | 1 | 1 |
| 37 | 1 | 2 | 130 | 250 | 0 | 1 | 187 | 0 | 3.5 | 0 | 0 | 2 | 1 |
| 41 | 1 | 1 | 130 | 204 | 0 | 0 | 172 | 0 | 1.4 | 2 | 0 | 2 | 1 |
| 56 | 1 | 1 | 120 | 236 | 0 | 1 | 178 | 0 | 0.8 | 2 | 0 | 2 | 1 |
| 57 | 1 | 0 | 120 | 354 | 0 | 1 | 163 | 1 | 0.6 | 2 | 0 | 2 | 1 |

Terminology:

Age,Sex,

Cp:Chestpain

Trestbps: Resting Blood pressure

Chol: Cholesterol

Fbs:Fasting Blood Sugar

Resteg: Resting Cardio electric

Thalach:Maximumhert rate achieved

Exang:Exercise induged angina

Oldpeak:depression induced by exercise relative to rest.

Slp:the slopeofthe peakexercise

Caa:number ofmajor vessels colored by fluorosopy

In[ ]: import pandas as pd

importnumpyasnp

```
fromsklearn.mDdel_selectionimporttrain_test_splitfromsklearn.preprocessingimportStandardScaler

#from sklearn.preprocessing import MinMaxScaler
```

```
Iron sklearn .ne1ghbor s import KNe1ghbor sClass1f1er from
sk1ear n. netr1c s Import confus1on_matr 1x

from sklearn .netr1cs import f 1_score

from sklearn .metr1cs Import accuracy_score
```

In[  ] :dataset=pd . read_csv( 'heat fi , c sv' )

|› in t ( ten( dataset ) )

print(dataset.head())

303

| age | | sex | cp | trtbps | cho1lbs | | thall | output | exng | oldpeak | slp | caa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 63 | 1 | 3 | 145 | 233 | 1 | ... | | 0 | 2.3 | 0 | 0 |
| 1 | 1 | | | | | | | | | | |
| 1 37 | 1 | 2 | 130 | 250 | 0 | | | 0 | 3.5 | 0 | 0 |
| 2 | 1 | | | | | | | | | | |
| 2 41 | 0 | 1 | 130 | 204 | D | .. | | 0 | 1. 4 | 2 | 0 |
| 2 | 1 | | | | | | | | | | |
| 3 5d | 1 | 1 | 1 20 | 23g | 0 | | | 0 | 0.8 | 2 | 0 |
| 2 | 1 | | | | | | | | | | |
| 4 57 | 0 | 0 | 1 20 | 354 | 0 | | | 1 | 0.6 | 2 | O |
| 2 | 1 | | | | | | | | | | |

[5 rows x 14columns]

In     :X=dataset.iloc[:.0:13] y=dataset.iloc[:,13J

X_train,X_test,y_train,y_test=tra1n_test_sp11t(X,y,random_state=          0, test_s1ze=0.2)

In        ]:scalingX=StandardScaler()

Xtrain=scalingX.fittransform(X_train) X_test=scalingX.transform(X_test) X_train

Out

```
array [[-1.32773282. - 1. 43641607.   0 .98584243 ,                    ]
         -0 .70710678,  -0 . 46472917 ] ,
       [ 1.24903178,  - 1. 43641607,   0. 98584243 , ...
         0 . 2651 d504,  -0 . 46472917 ] ,
       [ 0 . 3527g583 ,   0.6961771 2,  0.98584243,
         - 0 .70710678,  -0 . 46472917 ] .

       [ 0. 12869935 .   0.6961771 2,  1.94013791,
         - 0.70710678.   1. 14190596 J .
       [ - 0. 87959984,   0. 6961771 2.  0.98984243,
         - 0.707t0678.  -0. 46472917] .
       [ 0. 35276583 .   0. 6961771 2.  0. 0315469g ,
         -0.70710678.  -0. 46472917] } )
```

3636

-0 .6d1g9316.

.  -0 .66169316,

0.95577901,

-0.66169316,

**0.9557790],**

-0.66169316,

In [ j: classify=KNeighborsClassifier(n_neighbors:10.p:2,metric= e«clrle

Classify.fit(X_train,y_train)

Out[ ] KNeighborsClassifier(algorithm='auto',leaf_size=30.Metric='euclidean' ,

Let s=10, p=2,

metr1c_params=None ,n_jobs=None , n_ne1ghbor

weights='uniform')

In[ ]: y_predicted=classify.predict(X_test) yredicted

'i‹I      array ([0,          1,  0,  0,  0,  1,  0,  0,  0,  0,  1,  1,  0,  1,  1,  1,  0, 1, 0,

1, 1. 0,                      1,  0,  0 . 1,  1,  0,  0,  1,  1,  1,  0,  0,  1,  0, 0, 1,
                             0.  1,

0 .                           1 , 1,  0,  0,  1 ,  1.  0.  1.  1 ,  1.  0,  1.  1.  1 ,  1,  1 ])}

, 1, 0,

0,

print(cm)

In I'cm=confus1on_matr1x(y_test, y_pred1cted )

[[243]

[ 4 30]]

In[  ]: print( f 1_score( y_test,y_ predicted) )

0.8955223880597014

I n[ ]:   print(accuracy_ score(y_test, y_predicted))

0. 885245901 6393442

**Conclusion:**

This research focused on the study of art of AI and ML applications, selecting literature on what has now become a particularly hot topic in scientific research. This document describes the systematic selection of the most relevant literature was implemented. It provides the review of applications in various scientific fields using ML techniques.For the selection ofdocuments,and objective and clearmethodsof investigationsused. Implementing Python language to predict the best model for heart disease is used.

**References:**

1.  Theory: Devroye, Gyorfi and Lugosi, *A Probabilistic Theory of Pattern Recognition*, 1996
2.  Bishop, *Pattern Recognition and Machine Learning*, 2006
3.  Tan, Steinbach, and Kumar, Introduction to Data Mining, Addison-Wesley, 2005.
4.  Sutton and Barto, Reinforcement Learning: An Introduction, MIT Press, 1998.
5.  Bertsekas and Tsitsiklis, Neuro-Dynamic Programming, Athena Scientific, 1996
6.  Hastie et al, Bishop, and Duda et al. all have chapters on LDA, logistic regression, and other linear classifiers.
7.  On Over-fitting in Model Selection and Subsequent Selection Bias in Performance Evaluation, Gavin C. Cawley, Nicola L. C. Talbot; JMLR 11(Jul):2079-2107, 2010.
8.  Discussion of how certain model selection strategies are more biased than others; essential reading if you are doing comparative studies of different machine learning methods.
9.  Descent Methods for Tuning Parameter Refinement, Alexander Lorbert, Peter Ramadge ; AISTATS 2010. A natural idea.
10. An entry level discussion of the bootstrap, cross-validation, and other error estimates is given in Efron and Tibshirani, An Introduction to the Bootstrap, 1993
11. Data repositories: kaggle.com