## PREDICTING THE SPREAD OF CORONAVIRUS BY USING WEKA TOOL

[1]Shruti Tiwari, [2]Dr.Divakar Singh
[1]M.Tech 2[nd] year Student, [2]HOD
Computer Science and Engineering
Barkatullah University Institute of Technology
Bhopal (M.P.)

**ABSTRACT**

2020 was the year of pandemic called Covid-19. The world has never witnessed the lockdown on the global scale that spread like jungle fire but my hunch was whether the predicting of such disease and their spread was possible.

The data mining tool we used for our purpose of predicting is WEKA. The WEKA tool is a collection of machine learning algorithms for data mining tasks. It provides the tools for the data pre-processing, classification, regression, clustering, association rules and visualization. We have used multilayer perception algorithms using the datasets made available by WHO at www.kaggle.com.

We took this task in April 2020 when the datasets related Covid-19 were at earlier stage, but the finding based on our tool's analytics could predict that virus Covid-19 will affect the old age and infants group the most. There two where the easy target due to lack of immunity and our research is proved by the data available  now after 9 months / I year of this pandemic . WEKA tools analyze the medical events expertly.  In Conclusion, we find out that the effect of Covid-19 is less visible in women and children than in men and the comparison between recoverability of men and women.

We find that the possibility of recovery does not depends on on space and genre, it depends only on the patient's immunity systems.

**Keywords:** Using Coronavirus diseases (COVID -19) datasets obtained by WHO[1].Middle East Respiratory Syndrome (MERS-CoV),Severe Acute Respiratory Syndrome (SARS-CoV),apply a multilayer perceptron algorithm on covid19 datasets ,Predict the Spreading of Coronavirus By Using WEKA Tool

## INTRODUCTION

### What is CoronaVirus?

Coronaviruses (CoV) are a large family of viruses that cause illness ranging from the common cold to more severe diseases such as Middle East Respiratory Syndrome (MERS-CoV) and Severe Acute Respiratory Syndrome (SARS-CoV). Coronavirus disease (COVID-19) is a new strain that was discovered in 2019 and has not been previously identified in humans.

Coronaviruses are zoonotic, meaning they are transmitted between animals and people.  Detailed investigations found that SARS-CoV was transmitted from civet cats to humans and MERS-CoV from dromedary camels to humans. Several known coronaviruses are circulating in animals that have not yet infected humans.

### What are the symptoms of coronavirus COVID-19

Common signs of infection include respiratory symptoms, fever, cough, shortness of breath and breathing difficulties.

**Fig. 1.1:A Comparison table between corona, flu and cold respect to Symptoms[7].**



**Fig 1.2 Surface Stability of CoronaVirus HCOV-19[8]**

## SIGNIFICANCE OF THE STUDY

I have used datasets obtained by WHO[1]. And on this We are using the machine learning[4] / data mining tool WEKA[2] Through this, We are trying to understand how the corona virus is affecting the people of which age, how much the infected people are likely to protect. Whether our results are matching with the news that is going on.

## ABOUT THE DATASETS[3]

I took this data from the website www.kaggle.com.[3]

**Fig 2.1**

**Columns-**

**Id,case_in_country,summary,location,country,gender,age .If_onset_approximated ,hosp_visit_date ,exposure_start ,exposure_end,visiting Wuhan from Wuhan ,death ,recovered ,symptomsetc.**

According to the datasets [3], the first case of corona virus infected was found in Wuhan, China on 22 December 2019who died shortly after. He was 61. See below fig.



**Fig 2.2 The first case of corona virus infected.**

Afterthat, the corona virus spread rapidly in China and also took other countries - abroad.Preprocessing is done in which unwanted columns have been removed before using this dataset.

**AFTER PREPROCESSING**

Preprocessing is done in which unwanted columns have been removed before using this dataset.Here we are using 9 attributes which you can see in the picture,here the id attribute will be numeric while all the other attributes will be nominal. Here the main focus is on death and recovery attribute.

| No. | 1: id Numeric | 2: country Nominal | 3: gender Nominal | 4: age Numeric | 5: If_onset_approximated Nominal | 6: visiting_Wuhan Nominal | 7: from_Wuhan Nominal | 8: death Nominal | 9: recovered Nominal |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1.0 | China | male | 66.0 | 0 | 1 | 0 | 0 | 0 |
| 2 | 2.0 | China | female | 56.0 | 0 | 0 | 1 | 0 | 0 |
| 3 | 3.0 | China | male | 46.0 | 0 | 0 | 1 | 0 | 0 |
| 4 | 4.0 | China | female | 60.0 | NA | 1 | 0 | 0 | 0 |
| 5 | 5.0 | China | male | 58.0 | NA | 0 | 0 | 0 | 0 |
| 6 | 6.0 | China | female | 44.0 | 0 | 0 | 1 | 0 | 0 |
| 7 | 7.0 | China | male | 34.0 | 0 | 0 | 1 | 0 | 0 |
| 8 | 8.0 | China | male | 37.0 | 0 | 1 | 0 | 0 | 0 |
| 9 | 9.0 | China | male | 39.0 | 0 | 1 | 0 | 0 | 0 |
| 10 | 10.0 | China | male | 56.0 | 0 | 1 | 0 | 0 | 0 |
| 11 | 11.0 | China | female | 18.0 | 0 | 1 | 0 | 0 | 0 |
| 12 | 12.0 | China | female | 32.0 | 0 | 1 | 0 | 0 | 0 |
| 13 | 13.0 | China | male | 37.0 | NA | 1 | 0 | 0 | 0 |
| 14 | 14.0 | China | male | 51.0 | 0 | 0 | 1 | 0 | 0 |
| 15 | 15.0 | China | male | 57.0 | 0 | 0 | 1 | 0 | 0 |
| 16 | 16.0 | China | male | 56.0 | 0 | 1 | 0 | 0 | 0 |
| 17 | 17.0 | China | male | 50.0 | 0 | 1 | 0 | 0 | 0 |
| 18 | 18.0 | China | female | 52.0 | 0 | 0 | 1 | 0 | 0 |
| 19 | 19.0 | China | male | 33.0 | 0 | 1 | 0 | 0 | 0 |
| 20 | 20.0 | China | male | 40.0 | 0 | 1 | 0 | 0 | 0 |
| 21 | 21.0 | China | male | 70.0 | NA | 1 | 0 | 0 | 0 |
| 22 | 22.0 | China | male | 51.0 | 0 | 0 | 1 | 0 | 0 |
| 23 | 23.0 | China | male | 0.0 | 0 | 1 | 0 | 0 | 0 |
| 24 | 24.0 | China | female | 28.0 | 0 | 1 | 0 | 0 | 0 |
| 25 | 25.0 | China | male | 37.0 | 0 | 1 | 0 | 0 | 0 |
| 26 | 26.0 | China | male | 19.0 | 0 | 1 | 0 | 0 | 0 |
| 27 | 27.0 | China | male | 29.0 | NA | 1 | 0 | 0 | 0 |
| 28 | 28.0 | China | female | 66.0 | 0 | 0 | 1 | 0 | 0 |
| 29 | 29.0 | China | male | 46.0 | 0 | 0 | 0 | 0 | 0 |
| 30 | 30.0 | China | female | 32.0 | NA | 0 | 1 | 0 | 0 |
| 31 | 31.0 | China | male | 28.0 | NA | 1 | 0 | 0 | 0 |
| 32 | 32.0 | China | male | 55.0 | NA | 0 | 1 | 0 | 0 |
| 33 | 33.0 | China | male | 68.0 | NA | 0 | 1 | 0 | 0 |
| 34 | 34.0 | China | male | 38.0 | 0 | 0 | 0 | 0 | 0 |
| 35 | 35.0 | China | male | 72.0 | 0 | 0 | 1 | 0 | 0 |
| 36 | 36.0 | China | male | 45.0 | 0 | 1 | 0 | 0 | 0 |
| 37 | 37.0 | China | male | 42.0 | 0 | 1 | 0 | 0 | 0 |
| 38 | 38.0 | China | female | 33.0 | NA | 0 | 1 | 0 | 0 |

**Fig 2.3 Screenshot after preprocessing**

## APPLY THE ALGORITHMS ON THE DATASET[3]

After preprocessing, the algorithm will apply to the data. We know that every algorithm has its own characteristics. This paper has used multilayer perceptron algorithm[4] [5] [6] and also put percentage split = 66% for data training and set the **recovered** attribute as a class. Hence 66% of the data will be used for trainingand the remaining data will be used for testing.
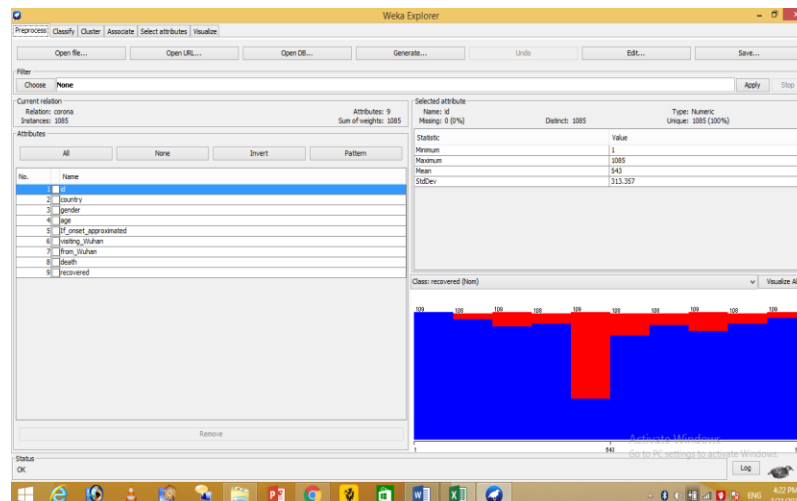


**Fig 2.4 Dataset is ready for MLP algorithm**

Run the Multilayer Perceptron (MLP) algorithm against the data. MLP uses backpropagation[6] to classify instances.The sigmoid nodes are the nodes used in backpropagation[6] and the associated data. This is the network itself (its weights and attributes). The nodes in the hidden layer of this network are all sigmoid but the output nodes are linear units (eg. linear node 0 is your output unit and sigmoid nodes 1-6 are your six hidden units). All the values given are your interconnection weights. You can

use them to manually calculate your results (which is done for you below the network).The bottom part is the final results calculated from the network.
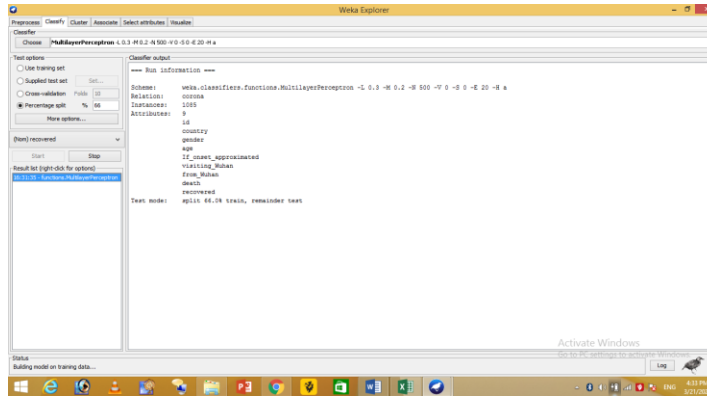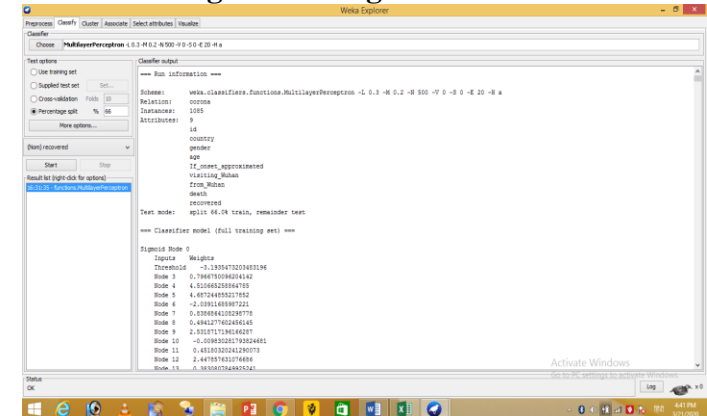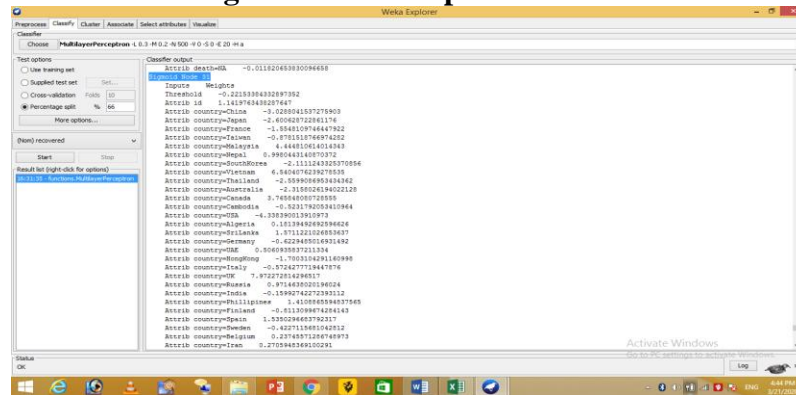


**Fig. 2.5: At Sigmoid node 0**



**Fig.2.6 :At sigmoid mode 31**

And end with sigmoid node 31. Hence 32 sigmoid nodes are used.

**Fig. 2.7:After complete the run**



## EXPLANATION OF THE RESULT

Since we had initially put percentage split = 66% for data training. Hence 66% of the data will be used for trainingand the remaining data will be used for testing.Therefore, after the process of perceptron algorithm, the result obtained in the confusion matrix is the result of testing the remaining data.A confusion matrix[4]is a summary of prediction results on a classification problem. The number of correct and incorrect predictions are summarized with count values and broken down by each class. This is the key to the confusion matrix. The confusion matrix shows the ways in which your classification model is confused when it makes predictions.
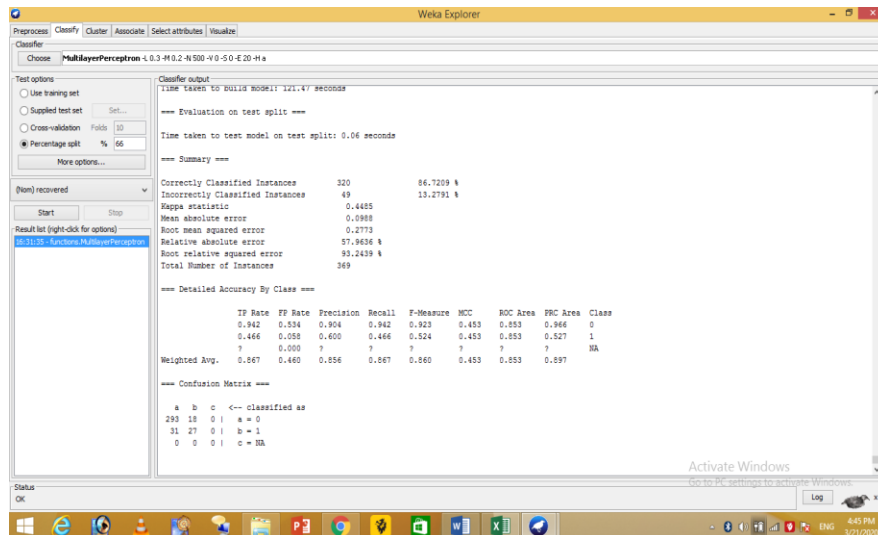
**Fig. 3.1: Summary and confusion matrix**

The data of 369 instance, ie 369 infected people has been tested. And as a result this matrix is formed. Let us try to understand this matrix. Out of 369 people, 293 people have not been recovered yet, all of them are undergoing treatment. Out of them, there is every hope of the rescue of 18 people. While 58 people have recovered, out of which 31 people seemed less likely to survive.

Those who have recovered successfully, the age is the most important. According to science and physicians, as a person's age increases, so does his / her immunity system week, only few people are able to maintain the immunity system even in their old age, or say that even after infection their immunity system quickly becomes strong.
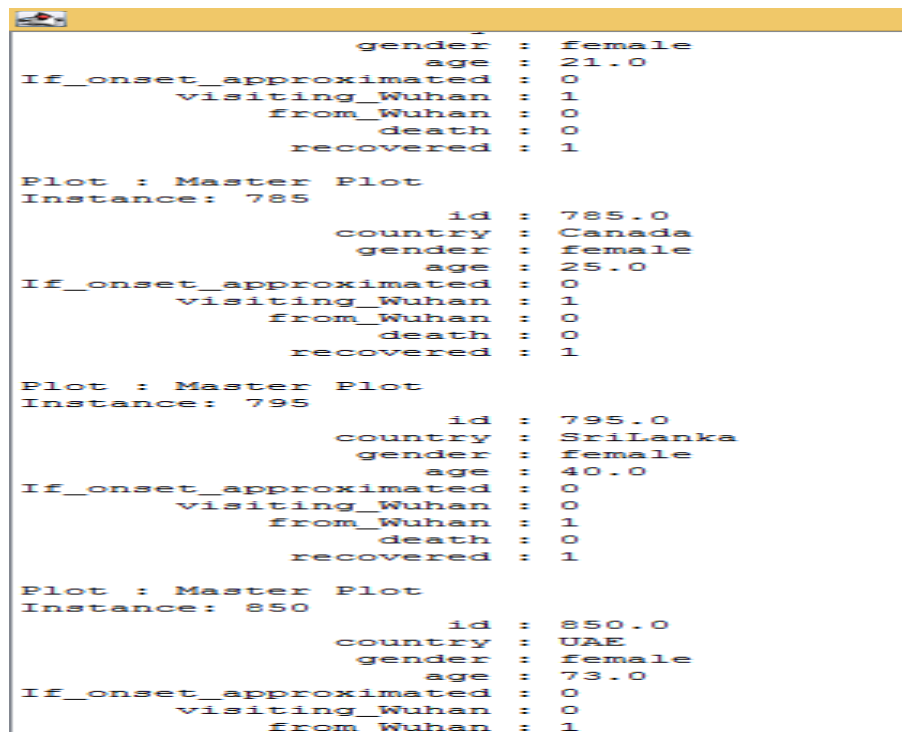


**Fig. 3.2 Women's recoverability**
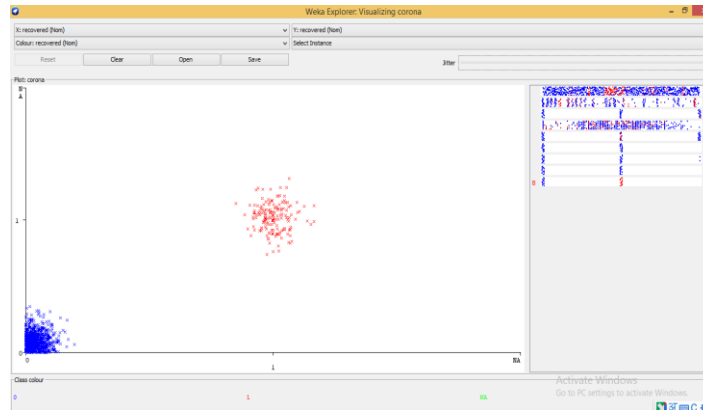
**Fig. 3.3: Men's recoverability**



Fig. **3.4:Total recoverability**

0 = no recover, is present by blue dots. And 1= recover, is present by red dots.

## POSSIBILITY OF RECOVERY IS DEPENDING ON COUNTRY?

There is also a question that the possibility of recovery of an infected person is also based on the country and abroad.If we talk about Singapore then you can seehere. From the data tested, 23 infected persons of Singapore were found, out of which 6 people have been cured and the remaining 17 are undergoing treatment (not recovered).
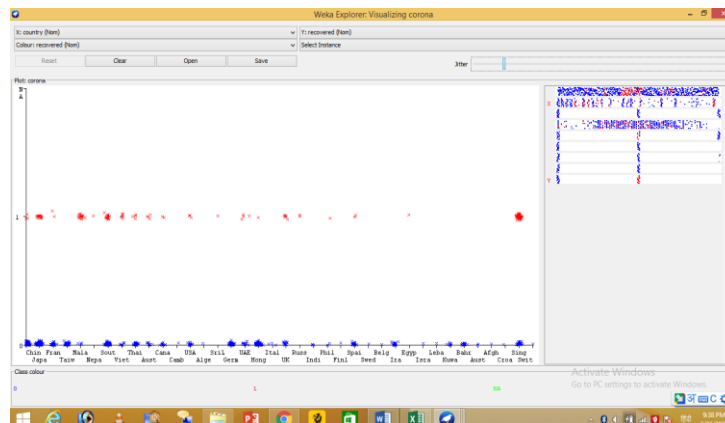


**Fig 3.1.1: Visualization of recoverability**

Where darker colors (red and blue dots), there are more cases of infected cases, and accordingly the rates of Recover and No Recover are seen to be increasing. It is clear that the names of the foreign countries do not matter, how the patients are being treated there, and what are the preventive methods to prevent the virus from spreading. For example, WHO[1] has released a list of several methods of rescue so that it is necessary to be adopted by the individual, the country, that is, all of them. Since no vaccine has yet been made, only rescue is the only solution. The visualization of the data tested above is clear that if the infected cases increase then there may be a problem in recovery and the rate of getting infected by it will keep increasing more rapidly. And that's why it has been declared a global epidemic.

### How likely is death from corona virus?

Now, discuss about **that how likely is death from corona virus?**

Here obtained a new confusion matrix by applying the same procedure by applying multilayer perceptron algorithm and classifying the death attribute.Since we had initially put percentage split = 66% for data training. Hence 66% of the data will be used for trainingand the remaining data will be used for testing.Therefore, after the process of perceptron algorithm, the result obtained in the confusion matrix is the result of testing the remaining data.The data of 369 instance, ie 369 infected people has been tested. And as a result this matrix is formed. Let us try to understand this matrix. Out of 369 people, 15 have died, while 9of them had little chance of survival. In addition, there is a possibility that 3 other people may die. While some of these 351 people have recovered and some are undergoing treatment.
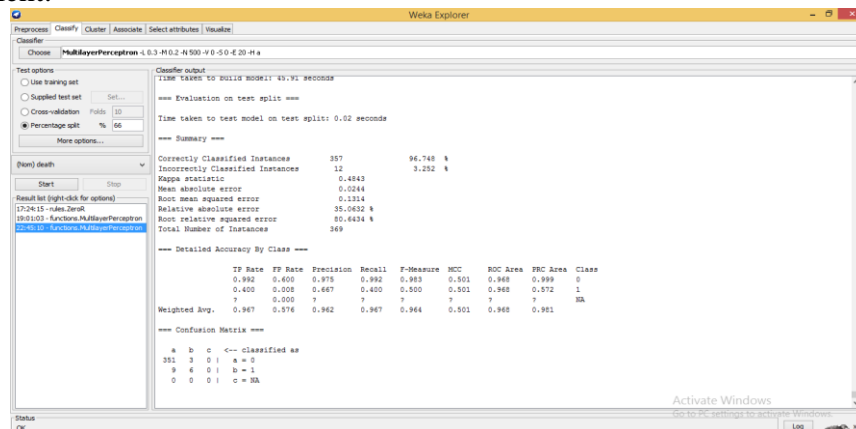


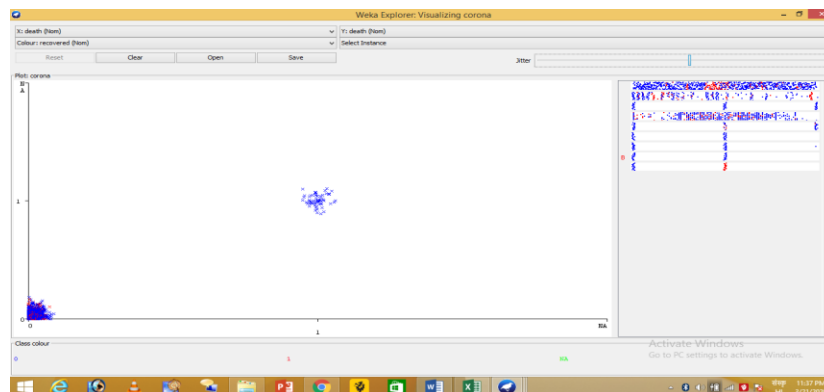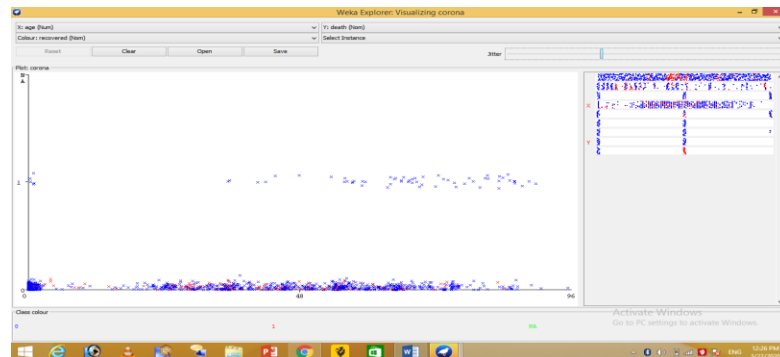**Fig. 3.2.1: Summary and confusion matrix**



**Fig. 3.2.2: Visualization after apply the MLP Algorithm.**

At 0, some peoples are recovered and some peoples are undergoing treatment, while at 1, peoples are died , which is the maximum in the number , but is the minimum rate than recover rate.Like I already saidthat the visualization of the data tested above is clear that if the infected cases increase then there may be a problem in recovery and the rate of getting infected by it will keep increasing more rapidly. If the number of infected persons increases as fast as possible, then the time taken to recover may result in a higher death situation.One reason for this, transition area will also have to increase. However, which infected people are more likely to die?  Let's have a look here



**Visualization diagram - death, recovered and currently in treatment people**

In the visualization picture, the darkest blue dots are higher at the bottom. These dark blue dots are the number of infected individuals who are currently being treated. Red dots are also visible among them, it is showing the recovered people. While there are light blue dots on the top, they show dead people.This plot is showing the relation between age and death numbers Most of the death numbers are seen above 50 people. But if you look down, you get confirmation of recovery of more than 50 years. This means that the whole game is the strength of the immune system.

But the thing to note is that the effect of coronavirus is less visible in women and children than men.If we look at the figures of the dead, the number of women and children in it is less.



**Fig. 3.2.4: Visualization of recoverability of various gender.**

In this diagram you can also see that most of the dots upwards are coming from men, rather than women. A study by the Chinese Centers of Disease Control revealed that of the 44,000 people infected with coronaviruses, 2.8% were killed and 1.7% were women. Talking about age, while 0.2% of the children and adolescents infected with the virus have died, 15% of people over the age of 80 have died.

**SHOULD IT BE UNDERSTOOD FROM THESE FIGURES THAT WOMEN AND CHILDREN HAVE LESS FEAR OF BEING CORONAVIRUS**

One reason may be that infections in women and children are less or that their body can fight the virus better. Doctor Bharat Pankhania from the University of Exeter says, "Everyone is infected with the new virus that comes in, it's the most important thing. The reason is that no one has the immunity to fight that virus.".[7][8]
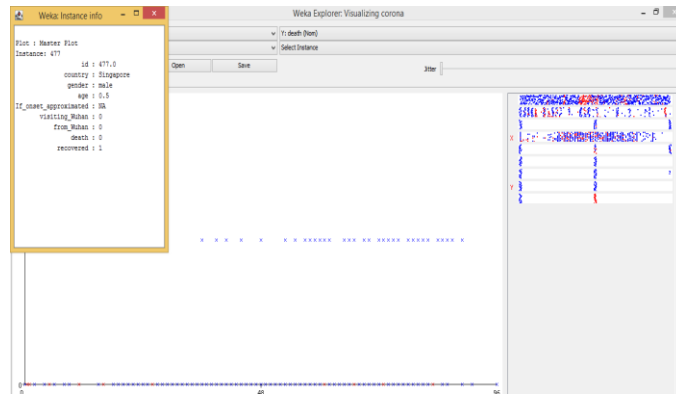


**Fig.3.3.1: Visualization of Recovery in 5 month old child patient**

Doctor Nathalie McDermitt[7] of King's College London explains, "One of the reasons behind decreasing infections in children may be that parents protect children more. Protects them from the risk of infection.

You will be surprised to know that the number of women died in coronavirus is less than men. However, scientists are not surprised by this at all. The same is seen in other infections including flu. The reason for this is that due to their lifestyle, the health of men is worse than that of women. Their lifestyle includes smoking and alcohol more than women. However, the difference in how the immune system of men and women reacts to infection also matters.

Professor Paul Hunter[7] at the University of East Anglia says, "Women have intrinsically different immune responses than men; women are at greater risk of auto-immune diseases (diseases caused by over-activation of the immune system). And there is also considerable evidence that women produce better antibodies to flu vaccines.

Children may develop corona virus infection. The youngest case yet, is of a day-old child. Very little is known about the symptoms of Covid-19 in children, but the symptoms are mild such as fever, runny nose and cough. Even young children can become ill from it. The same happens in the case of flu in which children below five years of age (especially less than two years) are at greater risk.

"People become more ill as they age, because their immunity decreases," says Doctor Pankhania[8]. More infections have been found in older people or in people already suffering from diseases like weak immunity and severe asthma. They will be at greater risk, but the effect of the virus in children has been found to be mild.

**Children's immune system is better?**

There are important differences in the immune system of a child and adult. In childhood our immune system is immature and can react excessively, so it is normal for children to have fever. Hyper activation of the immune system is also not good because it can damage the rest of the body. This is also one reason for coronavirus to be fatal.

**CONCLUSION**

As the result, we find out that the effect of Covid-19 is less visible in women and children than in men and the comparison between recoverability of men and women.

We find that the possibility of recovery does not depends on space and genre, it depends only on the patient's immunity systems.

**Table-1: Countries, areas or territories with cases[9].**

| | | | | |
|---|---|---|---|---|
| China :81416 cases | Italy :47021 cases | Spain :19980 cases | Iran (Islamic Republic of) :19644 cases | Germany :18323 cases |
| United States of America :15219 cases | France :12475 cases | Republic of Korea :8799 cases | Switzerland :4840 cases | The United Kingdom :3983 cases |
| Netherlands :2994 cases | Austria :2649 cases | Belgium :2257 cases | Norway :1742 cases | Sweden :1623 cases |
| Denmark :1255 cases | Australia :1081 cases | Malaysia :1030 cases | Portugal :1020 cases | Japan :1007 cases |
| Czechia :904 cases | Canada :846 cases | Israel :712 cases | International conveyance (Diamond Princess) :712 cases | Ireland :683 cases |
| Turkey :670 cases | Brazil :621 cases | Greece :495 cases | Pakistan :495 cases | Luxembourg :484 cases |
| Qatar :470 cases | Finland :450 cases | Indonesia :450 cases | Chile :434 cases | Poland :425 cases |
| Thailand :411 cases | Iceland :409 cases | Singapore :385 cases | Ecuador :367 cases | Saudi Arabia :344 cases |
| Slovenia :341 cases | Romania :308 cases | Philippines :307 cases | Bahrain :297 cases | Egypt :285 cases |
| Estonia :283 cases | India :258 cases | Russian Federation :253 cases | South Africa :240 cases | Peru :234 cases |
| Iraq :193 cases | Lebanon :187 cases | Kuwait :176 cases | Mexico :164 cases | Serbia :159 cases |
| San Marino :151 cases | Colombia :145 cases | United Arab Emirates :140 cases | Slovakia :137 cases | Panama :137 cases |
| Armenia :136 cases | Argentina :128 cases | Bulgaria :127 cases | Croatia :126 cases | Costa Rica :113 cases |
| Latvia :111 cases | Uruguay :94 cases | Algeria :94 cases | Viet Nam :91 cases | Morocco :86 cases |
| Hungary :85 cases | Jordan :84 cases | Faroe Islands :80 cases | Brunei Darussalam :78 cases | Andorra :75 cases |
| Sri Lanka :72 cases | Dominican Republic :72 cases | Albania :70 cases | North Macedonia :70 cases | Lithuania :69 cases |
| Cyprus :67 cases | Republic of | Malta :64 cases | Burkina Faso :64 | Tunisia :60casesBelarus :5 |

| | Moldova :66 cases | | cases | 7 cases |
|---|---|---|---|---|
| Senegal :56 cases | New Zealand :53 cases | Kazakhstan :53 cases | occupied Palestinian territory :52 cases | Oman :52 cases |
| Cambodia :51 cases | Guadeloupe :51 cases | Azerbaijan :44 cases | Bosnia and Herzegovina :44 cases | Georgia :43 cases |
| Venezuela (Bolivarian Republic of) :36 cases | Liechtenstein :34 cases | Uzbekistan :33 cases | Martinique :32 cases | Réunion :28 cases |
| Cameroon :27 cases | Ukraine :26 cases | Bangladesh :24 cases | Afghanistan :24 cases | Honduras :24 cases |
| Democratic Republic of the Congo :23 cases | Nigeria :22 cases | Ghana :19 cases | Rwanda :17 cases | Bolivia (Plurinational State of) :16 cases |
| Cuba :16 cases | Jamaica :16 cases | French Guiana :15 cases | Guam :14 cases | Montenegro :14 cases |
| Puerto Rico :14 cases | Maldives :13 cases | Paraguay :13 cases | Jersey :12 cases | Monaco :12 cases |
| Guatemala :12 cases | Mauritius :12 cases | French Polynesia :11 cases | Mongolia :10 cases | Gibraltar :10 cases |
| Trinidad and Tobago :9 cases | Côte d'Ivoire :9 cases | Ethiopia :9 cases | Togo :9 cases | Kenya :7 cases |
| Seychelles :7 cases | Kyrgyzstan :6 cases | Equatorial Guinea :6 cases | United Republic of Tanzania :6 cases | Aruba :5 cases |
| Barbados :5 cases | Guyana :5 cases | Bahamas :4 cases | Saint Martin :4 cases | Mayotte :4 cases |
| Cayman Islands :3 cases | Curacao :3 cases | Saint Barthelemy :3 cases | United States Virgin Islands :3 cases | Cabo Verde :3 cases |
| Central African Republic :3 cases | Congo :3 cases | Gabon :3 cases | Liberia :3 cases | Madagascar :3 cases |
| Namibia :3 cases | Fiji :2 cases | New Caledonia :2 cases | Greenland :2 cases | Bhutan :2 cases |
| Sudan :2 cases | Bermuda :2 cases | Haiti :2 cases | Saint Lucia :2 cases | Suriname :2 cases |
| Angola :2 case | Benin :2 cases | Guinea :2 cases | Mauritania :2 cases | Zambia :2 cases |
| Zimbabwe :2 cases | Papua New Guinea :1 cases | Guernsey :1 cases | Holy See :1 cases | Isle of Man :1 cases |
| Nepal :1 cases | Timor-Leste :1 | Djibouti :1 | Somalia :1 cases | Antigua and Barbuda :1 |

|  | cases | cases |  | cases |
|---|---|---|---|---|
| El Salvador :1 cases | Montserrat :1 cases | Nicaragua :1 cases | Saint Vincent and the Grenadines :1 cases | Sint Maarten :1 cases |
| Chad :1 cases | Eswatini :1 cases | Gambia :1 cases | Niger :1 cases |  |

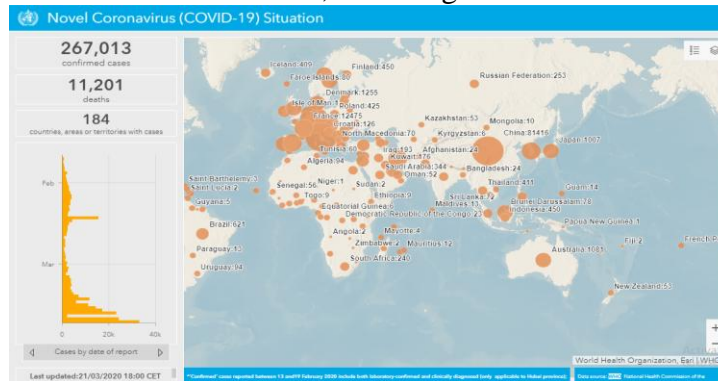So far, the following infected cases have arrived, including 184 countries and more .



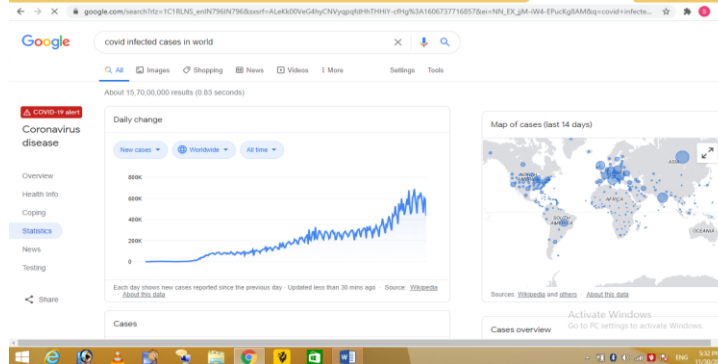**Fig. 4.2: According to the data received by 21 March 2020**[9]**.**



**Fig. 4.3: Coronavirus Disease (Covid-19) infected cases in worldby WHO**

## REFERENCES

1.  About WHO https://www.who.int/
2.  About WEKA visit here https://sourceforge.net/projects/weka
3.  Data was taken From- https://www.kaggle.com/
4.  Hastie, Trevor. Tibshirani, Robert. Friedman, Jerome. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer, New York, NY, 2009.
5.  Rosenblatt, Frank. x. Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms. Spartan Books, Washington DC, 1961
6.  Rumelhart, David E., Geoffrey E. Hinton, and R. J. Williams. "Learning Internal Representations by Error Propagation". David E. Rumelhart, James L. McClelland, and the PDP research group. (editors), Parallel distributed processing: Explorations in the microstructure of cognition, Volume 1: Foundation. MIT Press, 1986.
7.  Fig. 1.1:A Comparison table between corona, flu and cold respect toSymptoms. https://experience.arcgis.com
8.  Fig.1.2:Surface Stability of CoronaVirus-HCOV-19https://www.bhaskar.com
9.  Fig. 4.2 : According to the data received by 21 March 2020 **,**
    Fig.4.1: Table-1: Countries, areas or territories with cases https://www.who.int/