# Supervised Machine learning to predict COMD-19 Lab result

**DATA602 – Introduction to Machine learning**

Presented by:

Sharath Srinivas

# Why Covid-19 is important?

- With the worst coronavirus (COVID-19) outbreak still ongoing, the necessity for action to prevent ourselves from contracting and transmitting the virus has become an essential aspect of our daily life.

- However, some segments of the population are still not taking the threat as seriously as they should. Whether it's due to denial, living in an echo chamber of misinformation, or simply being unaware of the hazards, significant swaths of the population aren't taking efforts to prevent the spread.

# Why testing is important

- Effective testing is critical in slowing the spread of the virus by identifying persons infected and allowing treatment or isolation.

- Testing is also necessary to learn more about how the virus spreads and how widespread it is in a particular community.

- But, for the same reasons, health officials must ensure that their tests are as efficient as possible.

- In other words, with continued shortages, the imperative of ensuring that those most in need are tested means that not everyone should be checked.

# Introduction

- In this project supervised machine learning techniques are used to develop predictive models for the COVID-19 infection, using an epidemiology labeled dataset for positive and negative COVID-19 cases in Mexico, with supervised learning algorithms such as decision tree, logistic regression, Random Forest and naive Bayes, support vector machine, and k-nearest neighbors

# Problem Statement

- The global pandemic 2019-nCoV or COVID-19 is caused by the latest pandemic outbreak of the novel Severe Acute Respiratory Syndrome-Coronavirus two (SARS-CoV-2). The virus was first discovered in bats in late December 2019 in the Chinese province of Wuhan, with evidence suggesting it spread to people via intermediary hosts such as raccoons, dogs etc.

- Significant COVID-19 symptoms include fever (98%) cough (76%) and diarrhea (3%), which are generally more severe in older persons with chronic conditions , and many patients have reported shortness of breath, which in many cases appears to be a sign of the flu. , 2019 Since its discovery in late 2019, Covid has spread rapidly throughout the world.

# Objective

- To predict covid 19 results based on the symptoms of the patient.

- In this project we will be using the Lab result of covid antigen test as our target variables using symptoms features.

# Data Description

The data set used was provided by Mexico's Epidemiological Surveillance System for Viral Respiratory Diseases and includes information on all COVID19-related cases, including positive and negative tested individuals. This data set is available for downloading.
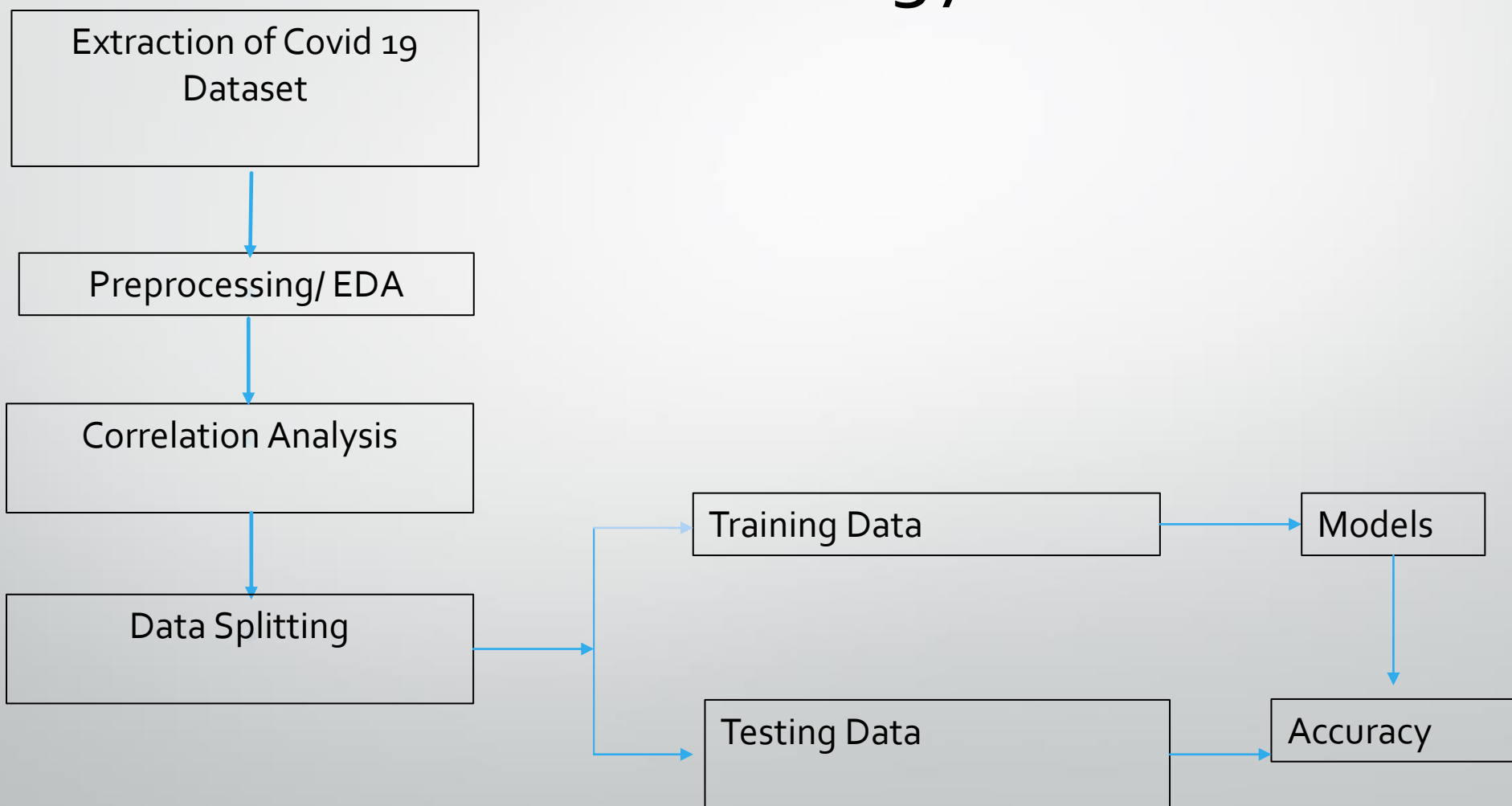
| Sl no | Columns | Description |
|---|---|---|
| 1 | UPDATE DATE | The database is fed daily, this variable allows identifying the date of the last update. |
| 2 | REGISTRATION ID | Case identifier number |
| 3 | ORIGIN | Sentinel surveillance is carried out through the system of health units monitoring respiratory diseases (USMER). The USMER include medical units of the first, second or third level of care, and third level units also participate as USMERs, which due to their characteristics contribute to broadening the epidemiological information panorama, including those with a specialty in pulmonology, infectology or pediatrics. . |
| 4 | SECTOR | Identifies the type of institution of the National Health System that provided the care |
| 5 | ENTITY UM | Identifies the entity where the medical unit that provided the care is located. |
| 6 | SEX | SEX Identifies the sex of the patient. |
| 7 | NAC ENTITY | Identifies the patient's birth entity. |
| 8 | RES ENTITIY | Identifies the patient's residence entity. |
| 9 | MUNICIPALITY RES | Identifies the municipality of residence of the patient. |
| 10 | PATIENT TYPE | Identifies the type of care the patient received in the unit. It is called an outpatient if you returned home or it is called an inpatient if you were admitted to the hospital. |
| 11 | ENTRY DATE | Identifies the date of admission of the patient to the care unit. |
| 12 | SYMPTOMS DATE | Identifies the date on which the patient's symptoms began. |
| 13 | DATE DEF | Identifies the date the patient died. |
| 14 | INTUBATED | Identifies if the patient required intubation. |
| 15 | PNEUMONIA | Identifies if the patient has been diagnosed with pneumonia. |
| 16 | AGE | Identifies the age of the patient. |
| 17 | NATIONALITY | Identifies if the patient is Mexican or foreign. |
| 18 | PREGNANCY | Identifies if the patient is pregnant |
| 19 | INDIGENOUS LANGUAGE SPEAKS | Identifies if the patient speaks an indigenous language |
| 20 | DIABETES | Identifies if the patient has a diagnosis of diabetes. |
| 21 | COPD | Identifies if the patient has a diagnosis of COPD. |
| 22 | ASTHMA | Identifies if the patient has a diagnosis of asthma. |
| 23 | INMUSUPR | Identifies if the patient is immunosuppressed. |
| 24 | HYPERTENSION | Identifies if the patient has a diagnosis of hypertension. |
| 25 | OTHER COM | Identifies if the patient has a diagnosis of other diseases. |
| 26 | CARDIOVASCULAR | Identifies if the patient has a diagnosis of cardiovascular disease. |
| 27 | OBESITY | Identifies if the patient has a diagnosis of obesity. |
| 28 | CHRONIC KIDNEY | Identifies if the patient has a diagnosis of chronic renal failure. |
| 29 | SMOKING | Identifies if the patient has a smoking habit. |
| 30 | OTHER CASE | Identifies if the patient had contact with any other case diagnosed with SARS CoV-2 |
| 31 | RESULT | Identifies if the patient is a case of COVID-19 |
| 32 | MIGRANT | Identifies if the patient is a migrant person |
| 33 | COUNTRY_NATIONALITY | Identifies the nationality of the patient. |
| 34 | COUNTRY_ORIGEN | Identifies the country from which the patient left for Mexico. |
| 35 | ICU | Identifies if the patient required admission to an Intensive Care Unit. |

# Methodology

```
┌─────────────────────────┐
│  Extraction of Covid 19 │
│         Dataset         │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│     Preprocessing/ EDA  │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│   Correlation Analysis  │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐                ┌─────────────────┐      ┌──────────┐
│      Data Splitting     │─────────────┬─▶│  Training Data  │─────▶│  Models  │
└─────────────────────────┘             │  └─────────────────┘      └──────────┘
                                        │                                 │
                                        │                                 ▼
                                        │  ┌─────────────────┐      ┌──────────┐
                                        └─▶│   Testing Data  │─────▶│ Accuracy │
                                           └─────────────────┘      └──────────┘
```

After importing the data frames dropped the all DATE columns along with nationality and information related to province/state because it will not help for our analysis

The columns "ENTITY RES" and " ENTITY UM" refer to the state of residence and treatment place respectively. These are coded by numbers so that they can be replaced by their respective names.

Symptoms features selected for analysis:

➢ PNEUMONIA
➢ INDIGENOUS_LANG
➢ DIABETES
➢ COPD
➢ ASTHMA
➢ INMUSUPR
➢ HYPERTENSION
➢ OTHERCOM
➢ CARDIOVASCULAR
➢ OBESITY
➢ CHRONIC_KIDNEY
➢ SMOKING
➢ OTHERCASE

Figure 1

# EDA Observation

1. Only Pregnancy, Other case and ICU column has significant amount of Not applicable and Not specified values. In our dataset for symptoms columns, 2 = No, and 1 = Yes. For our analysis we removed all not applicable and not specified and ignored values

2. Result column had three values 2 is negative and 1 is positive and 3 is result pending. We will be dropping the those values which are pending results and convert 2 to 0

```
2      4857 43.87
1       4086 36.9
3      2129 19.23
```

3. Cudad de Mexico state has highest number of positive and negative cases

4. Observed that patient male and female between age 40 to 60 are more likely to get Covid-19

5. top 4 symptoms are Pneumonia, Hypertension, Obesity and diabetes this indicates patients who have history of these decease must be carefully.

# Correlation Matrix



There are no strong positive correlation between the features but there are weak positive correlation between symptoms and result features.

Figure 5

# Model

Data split was done for 60% of train data and 40% of test data

**Logistic Regression** : For Logistic Regression we used C and solver has a hyperparameter with PCA

```
Pipeline(steps=[('pca', PCA(n_components=11)),
                ('logreg', LogisticRegression(C=10, solver='newton-cg'))])
```

```
Validation score: 63.36%
Test score: 64.68%
```

Figure 7

## Decision Tree

```python
param_dt = [
  {'dt__max_depth': [20,25,30,35],
   'dt__min_samples_split':[20, 40,50],
   'dt__min_samples_leaf': [2, 5,10],
   'dt__class_weight':[None, 'balanced']
  }
 ]
```

```
Pipeline(steps=[('dt',
                 DecisionTreeClassifier(max_depth=30, min_samples_leaf=10,
                                        min_samples_split=50))])
```

```
Validation score: 62.07%
Test score: 61.71%
```



Decision Tree Confussion matrix on how the model can predict covid positive

# Random forest

```python
param_rf = [{'rf__max_depth': [5, 8, 10, 12],
             'rf__n_estimators': [10],
             'rf__class_weight': ['balanced', 'balanced_subsample'],
             'rf__max_samples': [1000, 2000, 5000]
            }]
```

```
Pipeline(steps=[('rf',
                 RandomForestClassifier(class_weight='balanced', max_depth=5,
                                        max_samples=1000, n_estimators=10))])
```
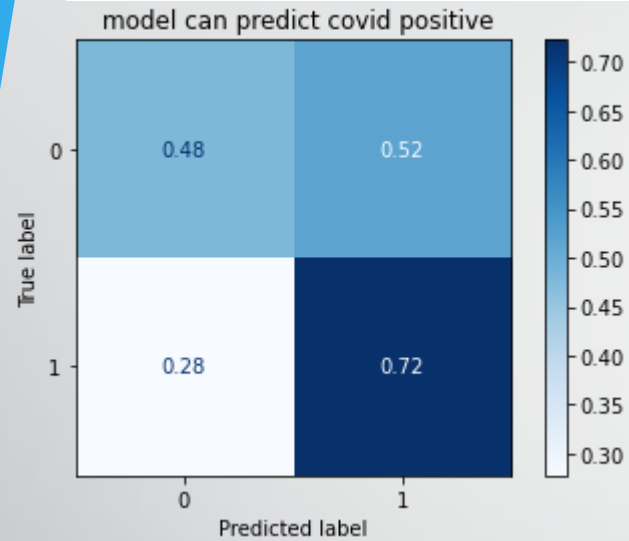
```
Validation score: 78.18%
Test score: 78.80%
```



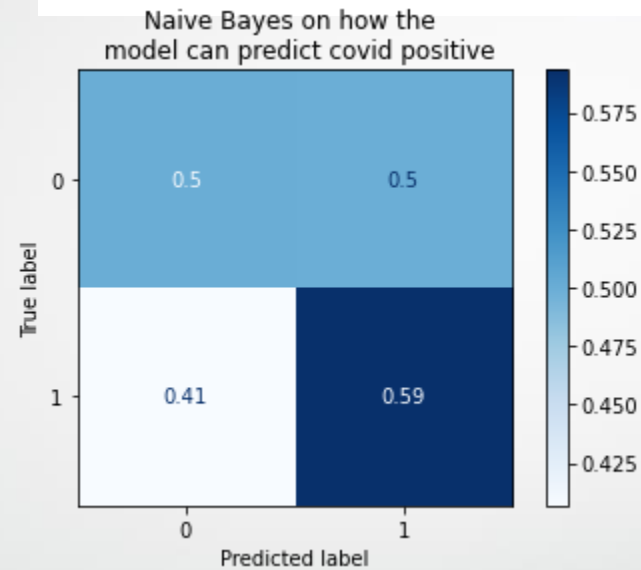Random forest Confussion matrix on how the model can predict covid positive

## KNN

Validation score: 59.31%
Test score: 61.24%

Knn on how the
model can predict covid positive

## ADA Boost

AdaBoost Validation Score: 64.82%
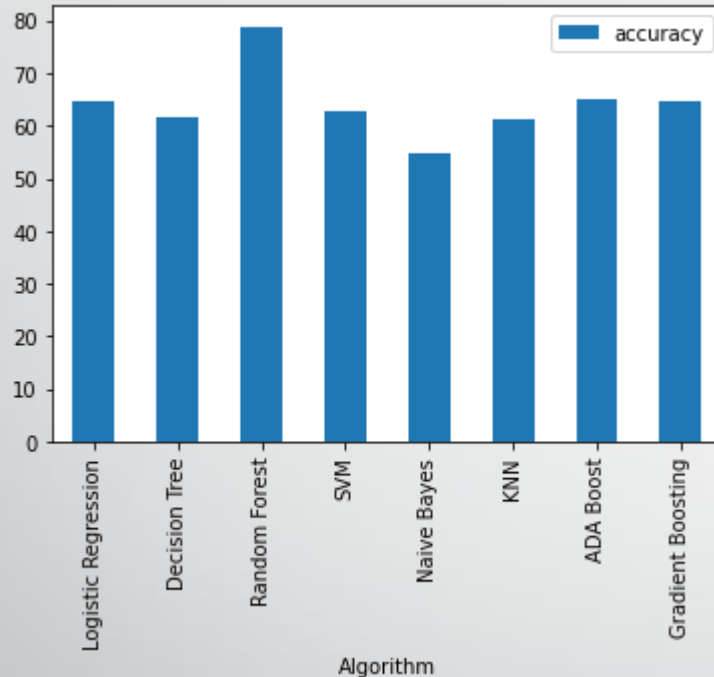AdaBoost Test score : 64.93%

## Gradient Boosting

Validation score: 65.57%
Test score: 64.84%

# Conclusion



Early COVID-19 prediction can aid in minimizing the undue burden on healthcare systems by assisting in the diagnosis of COVID-19 patients. In this project prior to creating the models, the correlation coefficient analysis between various dependent and independent features was performed to establish the strength of the association between each dependent and independent feature of the dataset. 60% of the training dataset was utilized to train the models, while the remaining 40% was used to test the models. The Random forest model has the highest accuracy of 78.80% percent followed by logistic regression, ADA boost and Gradient Boosting which lies around 65%.The performance of all models was evaluated based on accuracy parameters

Thank you !!