# NATIONAL INSTITUTE OF TECHNOLOGY SIKKIM

## Machine Learning
(EC16105)
Project Report
Academic Year 2023-2024

**Submitted by**:
Sonu Kumar Bhagat (B210066EC)
Semester VI

**Submitted to**:
Dr. Hemant Kumar Kathania
(Course instructor)
Department of ECE, NIT Sikkim

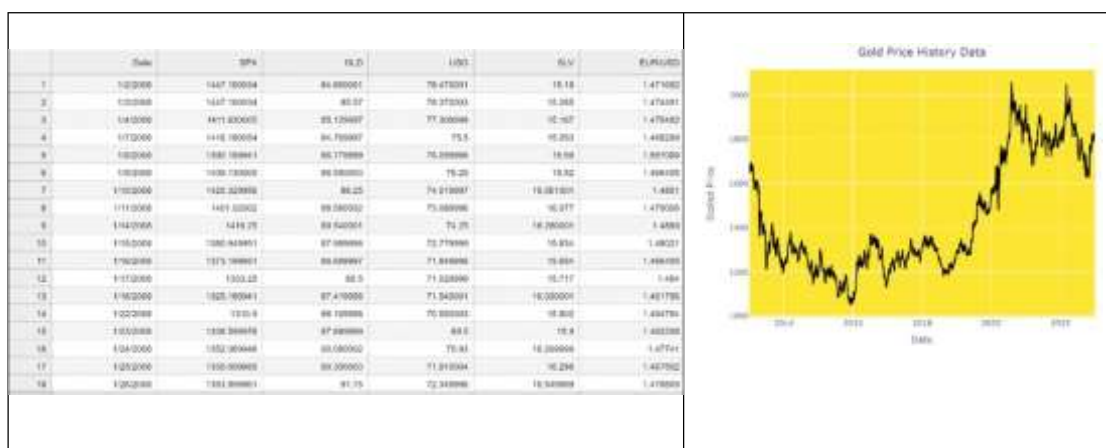## **Contents:**

## Abstract:

The project, titled "GOLD PRICE PREDICTIONS", predicts gold prices based on last year's gold data. The main purpose of the program is to help investors decide when to buy or sell gold by predicting the rise and fall of gold prices every day. Inventory forecasting plays an important role in the financial success of a business. The gold price is calculated by analysing the dataset containing the price of gold from the previous year. The rise in gold prices, along with fluctuations and price declines in other markets such as the capital market and real estate market, has led more and more investors to look at gold as a copy investment tool. There is concern that these high prices will continue and prices will reverse. There are many studies examining the relationship between gold prices and some financial variables. We use machine learning techniques to predict financial fluctuations and focus on predicting the price of gold.

## Objective:

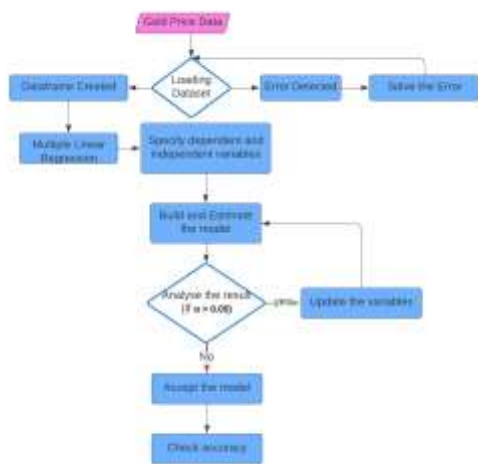The main objectives of the project are:
1. This project is based on the applicability of the proposed machine learning algorithms that had demonstrated their efficiency to predict gold prices with a better predictive rate.
2. To apply the best appropriate Machine Learning procedure.
3. We proposed the development of a prediction model for predicting future gold prices using Multiple Linear Regression, Random Forest, KNN.
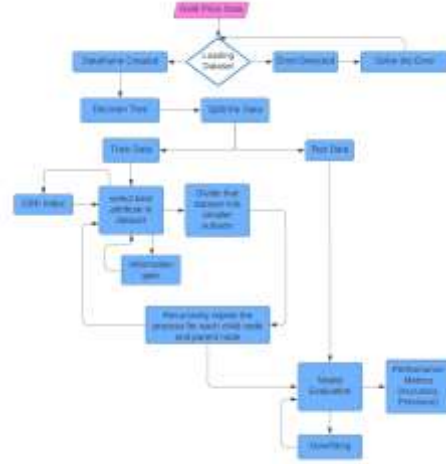
## Dataset Variations:

# Flowcharts:

| Multiple Linear Regression: | Decision Tree: |
|---|---|
|  |  |

| Random Forest: | KNN: |
|---|---|
|  |  |

# Codes:

| Multiple Linear Regression: | Decision Tree: |
|---|---|

**Multiple Linear Regression:**

```python
class MultipleLinearRegression:

    def __init__(self, learning_rate=0.01, iterations = 2000):
        self.learning_rate=learning_rate
        self.iterations=iterations
        self.cofficients = None
        self.intercept = None

    def predict(self, X):
        return self.intercept + self.cofficients * X

    def fit(self,X,y):
        self.cofficients = 0
        self.intercept = 0
        n=len(X)

        for i in range(self.iterations):
            y_predicted = self.predict(X)

            d_coefficient = (-2/n)*sum(X*(y-y_predicted))
            d_intercept =  (-2/n) *sum((y-y_predicted))

            self.cofficients -= self.learning_rate * d_coefficient
            self.intercept -= self.learning_rate * d_intercept
```

**Decision Tree:**

```python
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error
from sklearn.metrics import r2_score
import pandas as pd

data = pd.read_csv('gld_price_data.csv')

X = data.drop(['Date','GLD'], axis = 1)
y = data['GLD']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

regressor = DecisionTreeRegressor()

# Train the model
regressor.fit(X_train, y_train)

# Make predictions
predictions = regressor.predict(X_test)
```

**Random Forest:**

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn import metrics

gold_data = pd.read_csv('gld_price_data.csv')
gold_data.head()

#to check any missing values
gold_data.isnull().sum()

#splitting the features and target
X = gold_data.drop(['Date','GLD'], axis = 1)
y = gold_data['GLD']

X_train , X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 2)

regressor.fit(X_train, y_train)
regressor = RandomForestRegressor(n_estimators = 100)

#prediction on test data
test_data_prediction = regressor.predict(X_test)
```

**KNN:**

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsRegressor
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error

gold_data = pd.read_csv('gld_price_data.csv')

X = gold_data.drop(['Date','GLD'], axis = 1)
y = gold_data['GLD']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

k = 3
regressor = KNeighborsRegressor(n_neighbors=k)

regressor.fit(X_train_scaled, y_train)

predictions = regressor.predict(X_test_scaled)
```

# Accuracy:

| Algorithms | R2_Score | Mean Square Error |
|---|---|---|
| Multiple Linear Regression | 0.8975640982991402 | 56.16559421500603 |
| Decision Tree | 0.9831109316901961 | 9.260274392156793 |
| Random Forest | 0.9892766469309671 | 5.655962881475242 |
| KNN | 0.9941851307615867 | 3.188292196731859 |

## Outputs:

| | |
|---|---|
| Year: 1986 <br> Month: 12 <br> Day: 6 <br> Predict <br> Predicted gold price for 1986-12-6: $798.79 | Year: 2021 <br> Month: 1 <br> Day: 1 <br> Predict <br> Predicted gold price for 2021-1-1: $2405.91 |
| Year: 2024 <br> Month: 4 <br> Day: 22 <br> Predict <br> Predicted gold price for 2024-4-22: $2561.77 | Year: 2028 <br> Month: 8 <br> Day: 2 <br> Predict <br> Predicted gold price for 2028-8-2: $2763.61 |

## Conclusion:

This project aimed at understanding the relationship between gold prices and options that affect gold prices such as stocks, oil prices, rupee dollar exchange rate, silver price. This study uses monthly price data from January 2000 to December 2018. The data is divided into two periods. The first period was from January 2000 to October 2011, when gold prices increased; Four machine learning algorithms were used to analyse this data: linear regression, decision tree, random forest regression and KNN. The relationship between the variables is strong in stage I, and strong in stage II. It was found to be weak at this stage. While this model fit the first-level data well, the second-level fit was poor. KNN has a better prediction accuracy for the entire time- period, while multiple linear regression shows less accuracy as it best suits for the data having linear relationship between its features and target.

## Kernel Functions:

### 1. Polynomial Kernel:
It represents the similarity of vectors in the training set of data in a feature space over polynomials of the original variables used in the kernel.

$$f(x1, x2) = (x1^T . x2 + 1)^d$$

## 2. Sigmoid Kernel:

It is used for taking input, mapping them to a value of 0 and 1 so that they can be separated by a simple straight line.

$$f(x1, x2) = tanh(\alpha x^T. y + x)$$

## 3. RBF Kernel:

**R**adial **B**asis **F**unction (RBF) is used to create non-linear combinations of our features to lift our samples onto a higher-dimensional feature space where we can use a linear decision boundary to separate our classes. It is the most used kernel in SVM classifications.

Mathematically it can be represented as:

$$f(x1, x2) = e^{-||x1-x2||^2/2\sigma^2}$$

| *Polynomial Kernel* | *Sigmoid Kernel* |
|---|---|
|  |  |

| *RBF Kernel* |
|---|
|  |