# Assignment #1

| Name | ID |
|------|-----|
| Essam Alaa Nady | 20216064 |
| Ziyad Alfikie | 20216043 |
| Mariam Behairy | 20216094 |
| Mahitab Yasser | 20216082 |
| Ahmed El Tokhy | 20216014 |

## Introduction:

In this assignment, we are given a company dataset to automate the Loan Eligibility Process in a real time scenario related to customer's detail provided while applying application for home loan forms.

## Loading dataset:

```
1 data = pd.read_csv(r"C:\Users\Essam\Desktop\Assignment 1 ML\loan_old.csv")
2 data.dropna(inplace=True)
3 data = data.drop(["Loan_ID"],axis=1)
```

## Data Analysis:

### 1- Checking whether there are missing values:

```
1 data.isnull().sum()
```

```
Gender              0
Married             0
Dependents          0
Education           0
Income              0
Coapplicant_Income  0
Loan_Tenor          0
Credit_History      0
Property_Area       0
Max_Loan_Amount     0
Loan_Status         0
dtype: int64
```

## 2- Checking the type of each feature (categorical or numerical):

```
1 nf = data.select_dtypes(include=['int', 'float'])
2 cf= data.select_dtypes(include=['object'])
3 print("numerical features: ")
4 print(nf.columns)
5 print("categorical features: ")
6 print(cf.columns)
```

```
numerical features:
Index(['Income', 'Coapplicant_Income', 'Loan_Tenor', 'Credit_History',
       'Max_Loan_Amount'],
      dtype='object')
categorical features:
Index(['Gender', 'Married', 'Dependents', 'Education', 'Property_Area',
       'Loan_Status'],
      dtype='object')
```

## 3- Checking whether numerical features have the same scale:

```
Column: Income
Mean: 5030.730994152047
Standard Deviation: 4469.976642590766
------------
Column: Coapplicant_Income
Mean: 1486.6275243443274
Standard Deviation: 2102.196619596832
------------
Column: Loan_Tenor
Mean: 137.66081871345028
Standard Deviation: 23.139901552505556
------------
Column: Credit_History
Mean: 0.8557504873294347
Standard Deviation: 0.351685495324775
------------
Column: Max_Loan_Amount
Mean: 227.41440545808967
Standard Deviation: 157.63227884192304
------------
```
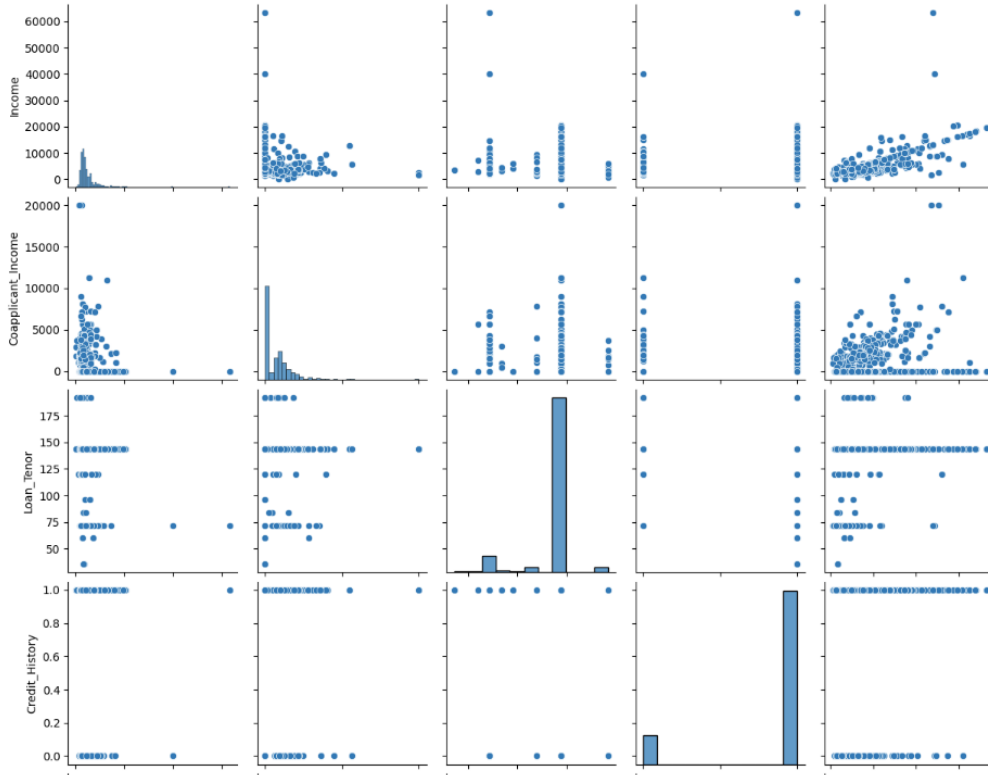
```
[6]: for column in data.select_dtypes(include=['float64', 'int64']).columns:
         mean = data[column].mean()
         std = data[column].std()
         print(f"Column: {column}")
         print(f"Mean: {mean}")
         print(f"Standard Deviation: {std}")
         print("------------")
```

Some of the numerical features have high standard deviation thus not all at same scale.

## 4- Pairplot:

```
1 sns.pairplot(nf)
2 plt.show()
```

c:\Users\Essam\AppData\Local\Programs\Python\Python310\lib\site-packages\seaborn\axisgrid.py:118: UserWarning: The figure layout has changed to tight
  self._figure.tight_layout(*args, **kwargs)



# Data Processing:

## 1- Removing records containing missing values:

```
[ ]   1 data = data.dropna()
```

```
1 data.isnull().sum()
```

```
Gender                0
Married               0
Dependents            0
Education             0
Income                0
Coapplicant_Income    0
Loan_Tenor            0
Credit_History        0
Property_Area         0
Max_Loan_Amount       0
Loan_Status           0
dtype: int64
```

## 2- Separating features and targets:

```
[ ]    1 X = data.drop(columns=["Max_Loan_Amount","Loan_Status"], axis=1)
       2 y = data["Max_Loan_Amount"]
```

## 3- Data is shuffled and split into training and testing sets:

```
[ ]   1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42,shuffle=True)
```

## 4- Categorical attributes are encoded:

```
[ ]   1 columns_to_encode = ["Gender", "Married", "Dependents", "Education", "Loan_Tenor", "Credit_History", "Property_Area"]
      2
      3 for column in columns_to_encode:
      4     encoder = LabelEncoder()
      5     X_train[column] = encoder.fit_transform(X_train[column])
      6     X_test[column] = encoder.transform(X_test[column])
      7
      8
```

```
[ ]   1 X_train = X_train.to_numpy()
      2 X_test = X_test.to_numpy()
```

## 5- Numerical features are standardized:

```
[ ]    1 scaler = StandardScaler()
       2 X_train = scaler.fit_transform(X_train)
       3 X_test = scaler.transform(X_test)
       4
       5 # y_train = scaler.fit_transform(y_train)
       6 # y_test = scaler.transform(y_test)
       7
       8 y_train = np.log1p(y_train)
       9 y_test = np.log1p(y_test)
```

# Linear Regression Model:

## 1- Predicting loan account:

```python
1  # Initialize the Linear Regression model
2  linear_model = LinearRegression()
3
4  # Train the model
5  linear_model.fit(X_train, y_train)
6
7  # Make predictions
8  y_pred = linear_model.predict(X_test)
9
10 # Calculate Mean Squared Error and R-squared for the test set
11 mse = mean_squared_error(y_test, y_pred)
12 r2 = r2_score(y_test, y_pred)
13
14
15 y_train_pred = linear_model.predict(X_train)
16 y_test_pred = linear_model.predict(X_test)
```

## 2- Calculate Mean Squared Error and R-Squared Score:

```python
10 # Calculate Mean Squared Error and R-squared for the test set
11 mse = mean_squared_error(y_test, y_pred)
12 r2 = r2_score(y_test, y_pred)
13
14
15 y_train_pred = linear_model.predict(X_train)
16 y_test_pred = linear_model.predict(X_test)
17
18 # Calculate Mean Squared Error for training and testing sets
19 train_mse = mean_squared_error(y_train, y_train_pred)
20 test_mse = mean_squared_error(y_test, y_test_pred)
21
22 print("R-squared score:", r2)
23 print("Train Mean Squared Error:", train_mse)
24 print("Test Mean Squared Error:", test_mse)
```

```
R-squared score: 0.04212693602124051
Train Mean Squared Error: 0.12564368348868454
Test Mean Squared Error: 0.3702009888223079
```

MSE on the training set is significantly lower than the MSE on the test set, it could be an indication that the model is overfitting.

## 3- Try lasso model to reduce the overfitting:

```
[ ]    1 # Initialize Lasso Regression model
       2 lasso_model = Lasso(alpha=0.1)  # Adjust the alpha value for regularization
       3
       4 # Train the Lasso model
       5 lasso_model.fit(X_train, y_train)
       6
       7 # Make predictions
       8 y_pred_lasso = lasso_model.predict(X_test)
       9
      10
      11 # Assess the Lasso model
      12 mse_lasso = mean_squared_error(y_test, y_pred_lasso)
      13 r2_lasso = r2_score(y_test, y_pred_lasso)
      14
      15 print("R-squared score (Lasso):", r2_lasso)
      16
      17 # Calculate Mean Squared Error for training and testing sets with Lasso Regression
      18 train_mse_lasso = mean_squared_error(y_train, lasso_model.predict(X_train))
      19 test_mse_lasso = mean_squared_error(y_test, y_pred_lasso)
      20
      21 print("Train Mean Squared Error (Lasso):", train_mse_lasso)
      22 print("Test Mean Squared Error (Lasso):", test_mse_lasso)
      23
      24 joblib.dump(lasso_model, 'lasso_model_weights.pkl')
      25
```

```
R-squared score (Lasso): 0.3219890671473171
Train Mean Squared Error (Lasso): 0.16842204623089424
Test Mean Squared Error (Lasso): 0.26203922754838466
['lasso_model_weights.pkl']
```

Better results!!

# Logistic Regression Model:

## 1- Sigmoid Function:

```
1 def sigmoid(z):
2     return 1 / (1 + np.exp(-z))
3
```

## 2- Cost Function:

```
5 def cost_f(X, t, weights):
6     m = X.shape[0]
7     prop = sigmoid(np.dot(X, weights))
8     positive = -t * np.log(prop)
9     negative = - (1 - t) * np.log(1 - prop)
0     cost = 1 / m * np.sum(positive + negative)
1     return cost
```

## 3- Derivative Function:

```
14 def f_dervative(X, t, weights):
15     m = X.shape[0]
16     prop = sigmoid(np.dot(X, weights))
17     error = prop - t
18     gradient = X.T @ error / m
19     return gradient
```

## 4- Gradient Descent Function:

```
22 def gradient_descent(X, t, step_size=0.1, precision=0.0001, max_iter=7000):
23     examples, features = X.shape
24     iter = 0
25     costs = []   # Move this line outside the loop
26     cur_weights = np.random.rand(features)  # random starting point
27     last_weights = cur_weights + 100 * precision
28
29     # print(f'Initial Random Cost: {cost_f(X, t, cur_weights)}')
30
31     while norm(cur_weights - last_weights) > precision and iter < max_iter:
32         last_weights = cur_weights.copy()  #  copy
33         gradient = f_dervative(X, t, cur_weights)
34         cur_weights -= gradient * step_size
35         #print(cost_f(X, cur_weights))
36         iter += 1
37         current_cost = cost_f(X, t, cur_weights)
38         costs.append(current_cost)
39
40     # Plot the cost over iterations
41     plt.plot(range(1, iter + 1), costs, linestyle='-', color='b')
42     plt.xlabel('Iterations')
43     plt.ylabel('Cost')
44     plt.title('Cost Function vs iterations')
45     plt.show()
46     print(f'Total Iterations {iter}')
47     print(f'Optimal Cost: {cost_f(X, t, cur_weights)}')
48     return cur_weights
```

## 6- Accuracy Function:

```
42 def accuracy(X, t, weights, threshold = 0.5):
43     m = X.shape[0]
44     prop = sigmoid(np.dot(X, weights))
45     labels = (prop >= threshold).astype(int)
46     correct = np.sum((t == labels))
47     return correct / m * 100.0
```

## 7- Separating features and targets:

```
[ ]   1 X_trainl, X_testl, y_trainl, y_testl = train_test_split(Xl, yl, test_size=0.2, random_state=42,shuffle=True)
```

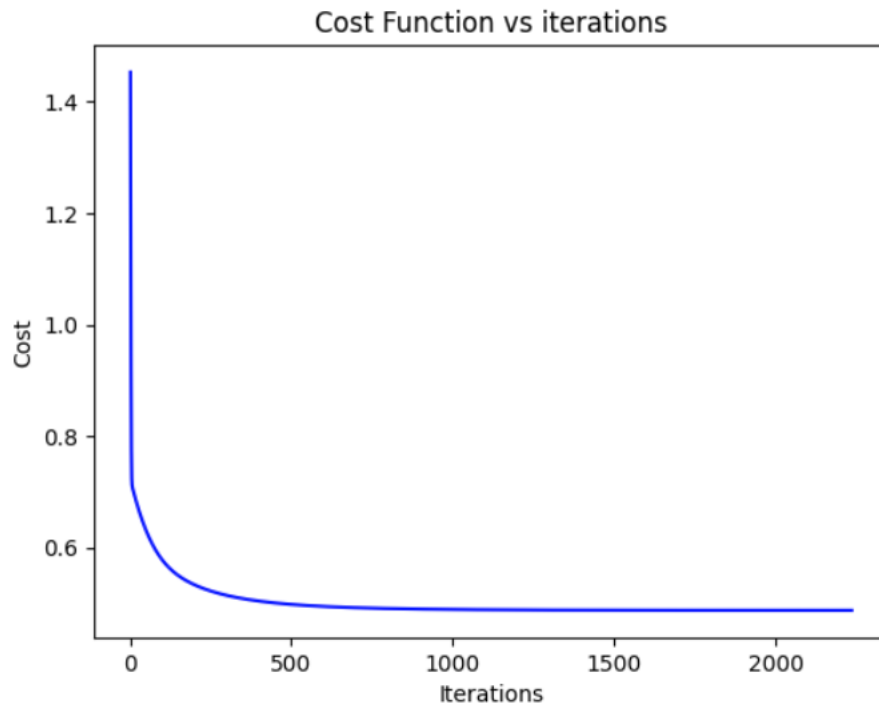## 8- Categorical attributes are encoded:

```
[ ]   1 columns_to_encode = ["Gender", "Married", "Dependents", "Education", "Loan_Tenor", "Credit_History", "Property_Area"]
      2
      3 for column in columns_to_encode:
      4     encoder = LabelEncoder()
      5     X_trainl[column] = encoder.fit_transform(X_trainl[column])
      6     X_testl[column] = encoder.transform(X_testl[column])
      7
      8 encod = LabelEncoder()
      9 y_trainl = encod.fit_transform(y_trainl)
     10
     11 # Transform the 'Loan_Status' column in y_test
     12 y_testl = encod.transform(y_testl)
```

## 9- Numerical features are standardized:

```
1 numeric_cols = ["Income", "Coapplicant_Income"]
2 scaler = StandardScaler()
3
4 # Fit and transform on training set
5 X_trainl[numeric_cols] = scaler.fit_transform(X_trainl[numeric_cols])
6
7 # Use the same scaler to transform the test set
8 X_testl[numeric_cols] = scaler.transform(X_testl[numeric_cols])
```

## 10- Training the model using gradient descent and evaluating its accuracy:

```
1 weights = gradient_descent(X_trainl,y_trainl)
2 print(f'Accuracy: {accuracy(X_testl,y_testl,weights)}')
3 np.save('Logistic_model_weights.npy', weights)
```

### Cost Function vs iterations



```
Total Iterations 2236
Optimal Cost: 0.48674880863286263
Accuracy: 80.58252427184466
```

## Loading the Testing Set:

```
1 data_new = pd.read_csv(r"C:\Users\Essam\Desktop\Assignment 1 ML\loan_new.csv")
2 # Removing spaces in column name
3 data_new = data_new.dropna()
4 # Write the cleaned data to a new CSV file
5 data_new.to_csv(r"C:\Users\Essam\Desktop\Assignment 1 ML\cleaned_data.csv", index=False)
```

```
1 data_new = data_new.drop(["Loan_ID"],axis=1)
```

```
1 for column in columns_to_encode:
2     encoder = LabelEncoder()
3     data_new[column] = encoder.fit_transform(data_new[column])
```

```
1
2 numeric_cols = data_new.select_dtypes(include=['int64', 'float64']).columns
3
4 # Subsetting the DataFrame to only include these columns
5 numeric_data = data_new[numeric_cols]
6
7 # Initialize the StandardScaler
8 scaler = StandardScaler()
9
10 # Fit and transform only the selected columns
11 scaled_data = scaler.fit_transform(numeric_data)
12
13 # Replace the original values with the scaled values
14 data_new[numeric_cols] = scaled_data
15 data_new
```

|  | Gender | Married | Dependents | Education | Income | Coapplicant_Income | Loan_Tenor | Credit_History | Property_Area |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 0.208582 | -0.656381 | 0.314250 | 0.46082 | 2 |
| 1 | 1 | 1 | 1 | 0 | -0.349612 | -0.013323 | 0.314250 | 0.46082 | 2 |
| 2 | 1 | 1 | 2 | 0 | 0.056577 | 0.115289 | 0.314250 | 0.46082 | 2 |
| 4 | 1 | 0 | 0 | 1 | -0.307389 | -0.656381 | 0.314250 | 0.46082 | 2 |
| 5 | 1 | 1 | 0 | 1 | -0.541940 | 0.810649 | 0.314250 | 0.46082 | 2 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 361 | 1 | 1 | 1 | 0 | -0.519984 | 0.272624 | 0.314250 | 0.46082 | 1 |
| 362 | 1 | 1 | 3 | 1 | -0.152640 | 0.105429 | 0.314250 | 0.46082 | 2 |
| 363 | 1 | 1 | 0 | 0 | -0.121183 | -0.352429 | 0.314250 | 0.46082 | 2 |
| 365 | 1 | 1 | 0 | 0 | 0.056577 | 0.369511 | 0.314250 | 0.46082 | 0 |
| 366 | 1 | 0 | 0 | 0 | 0.943270 | -0.656381 | -3.088317 | 0.46082 | 0 |

314 rows × 9 columns

```python
1 loaded_weights = np.load('Logistic_model_weights.npy')
2 def predict_loan_status(X, weights):
3     probabilities = sigmoid(np.dot(X, weights))
4     return (probabilities >= 0.5).astype(int)
5
6 # Predict loan status for data_new
7 predicted_loan_status = predict_loan_status(data_new, loaded_weights)
8 predicted_loan_status
```

```
array([1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
       1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1,
       0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
       1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1,
       1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
       1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
       1, 1, 1, 1, 1, 1])
```

```python
1 cleaned_data = pd.read_csv(r"C:\Users\Essam\Desktop\Assignment 1 ML\cleaned_data.csv")
2
3 # Adding the predicted values as a new column 'Loan_Status'
4 cleaned_data['Loan_Status'] = predicted_loan_status
5 cleaned_data['Loan_Status'] = cleaned_data['Loan_Status'].apply(lambda x: 'Y' if x == 1 else 'N')
6
7 # Save the updated data to a new CSV file
8 cleaned_data.to_csv(r"C:\Users\Essam\Desktop\Assignment 1 ML\cleaned_data.csv", index=False)
```

```python
1 loaded_model = joblib.load('lasso_model_weights.pkl')
2 predictions_lasso = loaded_model.predict(data_new)
3 predictions_lasso
```

```
c:\Users\Essam\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\base.py:457: UserWarning: X has feature names, but Lasso was fitted without feature names
  warnings.warn(
array([ 5.2356461 ,  5.13765852,  5.31898798,  5.03543999,  5.21500048,
        4.94942673,  5.08500001,  5.45879464,  5.15344357,  5.02028527,
        5.06593521,  5.76425902,  5.14930517,  5.23130448,  5.37306832,
        5.10354751,  6.14875079,  4.43949758,  5.04982882,  4.22927292,
        5.0754976 ,  5.55575039,  6.66471071,  5.67796407,  4.48797995,
        4.97006991,  5.34389808,  5.1678561 ,  5.20917087,  5.11337606,
        5.0251184 ,  4.97929681,  5.17354321,  5.24490276,  5.22071417,
        5.29286924,  5.06558559,  5.07046769,  5.50511847,  5.06738777,
        5.11046333,  5.44218475,  5.10020587,  5.29708174,  4.47487898,
        5.20066737,  4.99235144,  5.12767897,  4.6344282 ,  5.28443221,
        4.44528908,  5.10098515,  5.36489889,  5.17404422,  5.05778906,
        5.13488771,  5.28933231,  5.15209036,  5.04928403,  5.39505734,
        5.37207373,  5.13267594,  4.39176549,  5.25212484,  5.43948688,
        5.33218722,  5.21546714,  5.28774556,  5.36712529,  5.41832188,
        5.16795232, 10.70846511,  5.24757687,  5.45919866,  4.42366289,
        5.07087014,  5.21117899,  5.01823734,  5.19140341,  5.18141678,
        5.87443792,  5.32223612,  5.28764284,  5.32910153,  5.27988149,
        5.3721949 ,  5.3107605 ,  5.56323375,  5.13177484,  5.09406368,
        5.08663836,  4.27875103,  5.14875087,  5.13677181,  5.13570688,
        5.24433278,  5.15553089,  5.04699034,  5.27495618,  5.48431768,
        4.9062383 ,  5.12221236,  5.1039155 ,  5.02307047,  5.25888751,
        5.16484796,  5.47901427,  5.79104601,  5.14938708,  5.29954166,
        5.38145905,  4.31037114,  5.14052565,  5.04486049,  5.16087462,
        5.06361957,  4.52799649,  5.11588287,  5.13570688,  5.16333415,
        5.15430213,  4.83218261,  5.36537328,  4.52565562,  5.54379573,
        5.03543999,  5.48000568,  5.21222265,  5.25039123,  5.35935096,
        5.11017773,  5.17666558,  5.15544164,  5.16478755,  4.81030292,
        5.45062261,  5.16148355,  5.44404657,  5.41381642,  5.44354922,
        5.00988177,  5.22442341,  4.98891091,  4.95286726,  5.02002652,
        5.26667969,  5.14249864,  5.0590322 ,  5.09130765,  5.13431428,
```