
Assignment IV: Linear Regression Models

Introduction to Machine Learning Lab (190.013), SS2023
Björn Ellensohn¹

¹m01435615, bjoern.ellensohn@stud.unileoben.ac.at, Montanuniversität Leoben, Austria
May 12, 2023

This document guides through the process of solving Assignment 3.

1 Introduction

The 4th assignment's task was to implement linear regression models on a provided 'diabetes' dataset. the linear regression models had to be implemented from scratch using only numpy and pandas. the models to create where:

- Least Squares Regression
- Ridge Regression
- Lasso Regression

The laste one was a bonus point task. I will walk through the code in 4 Sections.

2 Part I - Abstract Class

2.1 Preparation

Alongside the 'diabetes.csv' dataset there was also an unfinished Jupyter Notebook provided, which gave an outline on how to solve this task. At the begin, there is an abstract class `Predictor`, that should act as a template for the child classes, implementing the actual regression models. However, I decided to extend the metaclass to include the methods that where the same for all child classe. So I only had to write them once. In my oppinion, this is better for the code's structure and enhances readability.

2.2 Defining the Abstract Class

The class `Predictor` is created that provides the basic outline for the main classes of this exercise. The following functions are included:

- Reading the csv
- Loading and splitting the data into training and test sets
- Normalizing the data

So the dataset is loaded and will be pre-processed accordingly on initiation of the class. The needed data and paramters are stored as attributes.

To mention here are the data loading and pre processing methods. For loading the file, pandas does a good already. The pre-processing was trickier to implement. Datasets tend to have several flaws. Amongst them are outliers and missing values. Missing data can be found with the pandas DataFrame functions `isna` and `isnull`. The outliers of the dataset where caught by boxplotting every column. In this case, the "1.5 x IQR rule" came to shine, as this allows to automate the process of finding the outliers for each of the columns. The last method is the `train_test_split()` method. Here I want to split the dataset into features and targets, as well as dividing the whole dataset into a test set and a train set. This is especially useful for having a refernce for evaluating the performance later.

3 Part II - Least Squares Regression

In the second part of this exercise, a class `LinearRegression(Predictor)` was created to provide the actual implementation of the linear regression.

A formula for Least Squares Regression looks like this:
 $Y = b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + \dots + b_n \cdot X_n$,
 where Y is the response variable, b_0 is the y-intercept, b_i is the coefficient for predictor variable X_i , and X_i is the i th predictor variable.

For evaluation, the predicted data was compared to the test set, so the Mean Squared Error was measured, as well as the R-Squared.

The results are as follows. For comparison, also the unprocessed dataset was evaluated. We can already see, how bad the unprocessed set performs.

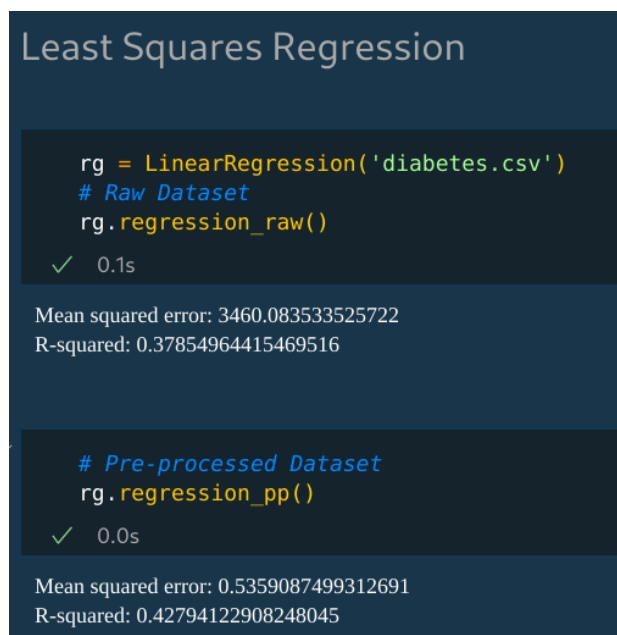


Figure 1: Results of Least Squares Regression

4 Part III - Ridge Regression

Next, a class `RidgeRegression(Predictor)` was made to implement Ridge Regression. Ridge Regression is a regularized linear regression model. The formula looks like this:

Again, the model was evaluated against the test set. The results are as follows:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n |\theta_j|$$

Figure 2: Formula of Ridge Regression

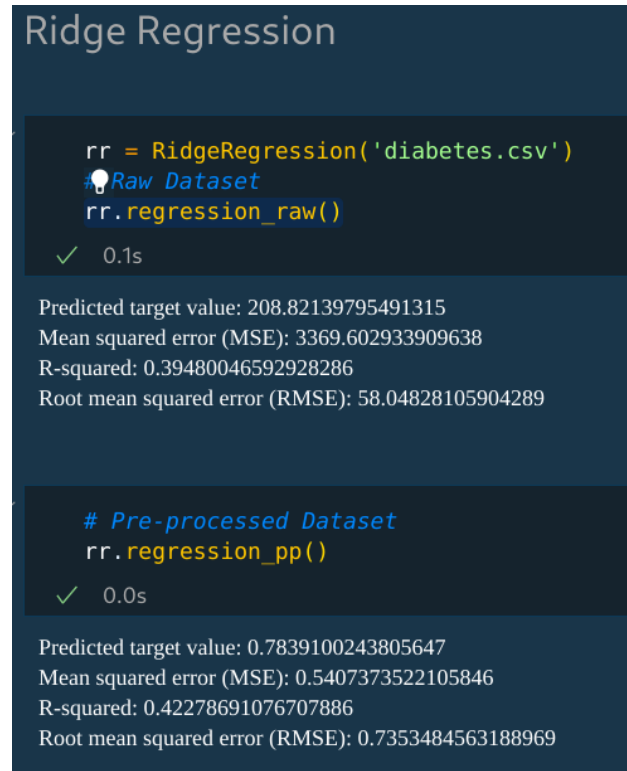


Figure 3: Results of Ridge Regression

5 Part IV - Lasso Regression (Bonus)

The Lasso Regression is also a regularized linear regression model. The formula is like this:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Figure 4: Formula of Lasso Regression

Also, the resulting test performance was evaluated.

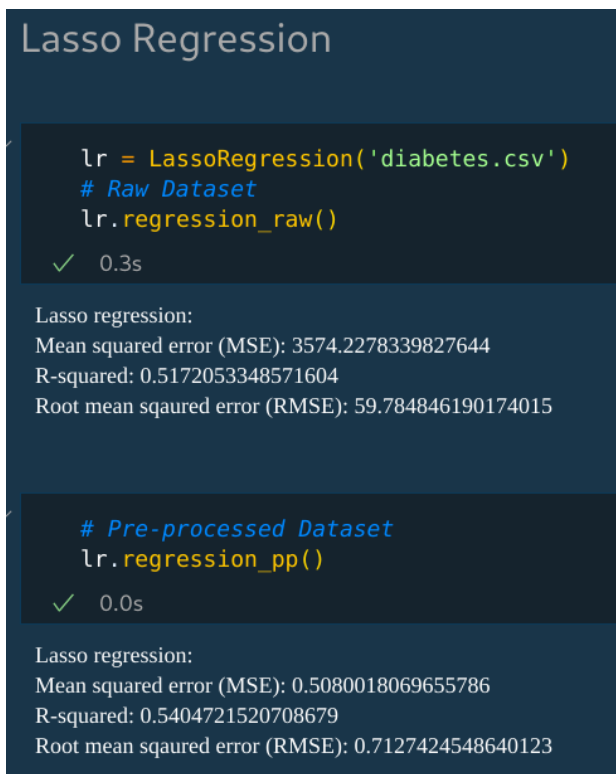


Figure 5: Results of Lasso Regression

APPENDIX

This section houses the code.

Code for Part I

```

1 class Predictor:
2     def __init__(self, dataset):
3         self.coefficients = None
4         self.df = pd.read_csv(dataset)
5         self.df_pp = self.preprocess(self.df)
6         self.X_train_raw, self.X_test_raw, self.y_train_raw, self.y_test_raw =
7             ↪ self.train_test_split(self.df)
8         self.X_train_pp, self.X_test_pp, self.y_train_pp, self.y_test_pp =
9             ↪ self.train_test_split(self.df_pp)
10
11     def preprocess(self, df):
12         # It is better to do this after splitting the dataset. --> need to change that
13         # Handle missing values
14         df.replace(0, np.nan, inplace=True)
15
16         # Remove outliers using iqr rule:
17         df_cleaned = df.copy()
18         for column in df:
19             q1 = df[column].quantile(q=0.25)
20             q3 = df[column].quantile(q=0.75)
21             med = df[column].median()
22
23             iqr = q3 - q1
24             upper_bound = q3+(1.5*iqr)
25             lower_bound = q1-(1.5*iqr)
26
27             df_cleaned[column][(df[column] <= lower_bound) | (df[column] >= upper_bound)] =
28                 ↪ df_cleaned[column].median()
29
30         # Normalize data:
31         df_normalized = df_cleaned.copy()
32         for column in df_cleaned:
33
34             mean = df_cleaned[column].mean()
35             std = df_cleaned[column].std()
36
37             df_normalized[column] = (df_cleaned[column] - mean) / std
38         df_processed = df_normalized
39         return df_processed
40
41     def train_test_split(self, df, test_size=0.2):
42         # Shuffle the rows of the dataset randomly
43         df_randomized = df.sample(frac=1, random_state=42).reset_index(drop=True)
44
45         # Extract the features and target variable
46         X = df_randomized.drop('target', axis=1)
47         y = df_randomized['target']
48
49         # Split the dataset into training and testing sets
50         split_ratio = 1 - test_size
51         split_index = int(split_ratio * len(df_randomized))
52
53         X_train = X[:split_index]
54         X_test = X[split_index:]
55         y_train = y[:split_index]
56         y_test = y[split_index:]

```

```

55         return X_train, X_test, y_train, y_test
56
57     def fit(self, X, y):
58         pass
59
60     def predict(self, X):
61         pass

```

Code for Part II

```

1  class LinearRegression(Predictor):
2      def __init__(self, dataset):
3          super().__init__(dataset)
4
5      def fit(self, X, y):
6          X = np.insert(X, 0, 1, axis=1)
7          y = y.values.reshape(-1, 1)
8          self.coefficients = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(y)
9
10     def predict(self, X):
11         X = np.insert(X, 0, 1, axis=1)
12         return X.dot(self.coefficients)
13
14     def regression_raw(self):
15         # copy the values
16         X_train = self.X_train_raw
17         y_train = self.y_train_raw
18         X_test = self.X_test_raw
19         y_test = self.y_test_raw
20
21         # implement linear regression
22
23         # Implement the formula for the least-squares regression line
24         X_train_T = np.transpose(X_train)
25         beta = np.linalg.inv(X_train_T.dot(X_train)).dot(X_train_T).dot(y_train) # these are the
26         ↪ weights
27
28         # Train the model on the training set using the least-squares regression line
29         y_pred_train = X_train.dot(beta) # the prediction on the train set
30
31         # Evaluate the performance of the model on the testing set using metrics such as mean
32         ↪ squared error and R-squared
33         y_pred_test = X_test.dot(beta) # the prediction on the test set
34
35         # calculate the mean squared error
36         mse = np.mean((y_test - y_pred_test)**2)
37         r_squared = 1 - (np.sum((y_test - y_pred_test)**2) / np.sum((y_test - np.mean(y_test))**2))
38
39         print('Mean squared error:', mse)
40         print('R-squared:', r_squared)
41
42         # that should do
43
44     def regression_pp(self):
45         # copy the values
46         X_train = self.X_train_pp
47         y_train = self.y_train_pp
48         X_test = self.X_test_pp
49         y_test = self.y_test_pp
50
51         # implement linear regression

```

Assignment IV: Linear Regression Models

```
51 # Implement the formula for the least-squares regression line
52 X_train_T = np.transpose(X_train)
53 beta = np.linalg.inv(X_train_T.dot(X_train)).dot(X_train_T).dot(y_train) # these are the
    ↳ weights
54
55 # Train the model on the training set using the least-squares regression line
56 y_pred_train = X_train.dot(beta) # the prediction on the train set
57
58 # Evaluate the performance of the model on the testing set using metrics such as mean
    ↳ squared error and R-squared
59 y_pred_test = X_test.dot(beta) # the prediction on the test set
60
61 # calculate the mean squared error
62 mse = np.mean((y_test - y_pred_test)**2)
63 r_squared = 1 - (np.sum((y_test - y_pred_test)**2) / np.sum((y_test - np.mean(y_test))**2))
64
65 print('Mean squared error:', mse)
66 print('R-squared:', r_squared)
```

Code for Part III

```
1 class RidgeRegression(Predictor):
2     def __init__(self, dataset, alpha=1):
3         super().__init__(dataset) # this command initializes the parent class (Predictor) and
    ↳ passes the dataset.
4         self.alpha = alpha
5
6     def regression_raw(self):
7         """
8         Fits a ridge regression model on the training data using the specified regularization
    ↳ parameter alpha.
9         Using raw dataset
10        """
11        X_train = self.X_train_raw
12        X_test = self.X_test_raw
13        y_train = self.y_train_raw
14        y_test = self.y_test_raw
15        alpha = self.alpha
16
17        # Add a column of 1s to the training data to have the correct dimension.
18        X_train = np.hstack([np.ones((X_train.shape[0], 1)), X_train])
19
20        n_features = X_train.shape[1]
21        I = np.eye(n_features)
22        w = np.linalg.inv(X_train.T.dot(X_train) + alpha * I).dot(X_train.T).dot(y_train)
23
24        #X_test = (X_test - X.mean()) / X.std() #this should not be needed, since the
    ↳ normalization is happening in the constructor
25        X_test = np.hstack([np.ones((X_test.shape[0], 1)), X_test])
26        y_pred = X_test.dot(w)
27
28        # calculate the mean squared error
29        mse = np.mean((y_test - y_pred)**2)
30        r_squared = 1 - (np.sum((y_test - y_pred)**2) / np.sum((y_test - np.mean(y_test))**2))
31
32        # calculate root mean squared error (RMSE)
33        rmse = np.sqrt(mse)
34
35        print("Predicted target value:", y_pred[0])
36        print("Mean squared error (MSE):", mse)
37        print("R-squared:", r_squared)
38        print("Root mean squared error (RMSE):", rmse)
```

```

39
40 def regression_pp(self):
41     """
42     Fits a ridge regression model on the training data using the specified regularization
    ↪ parameter alpha.
43     Using processed dataset.
44     """
45     X_train = self.X_train_pp
46     X_test = self.X_test_pp
47     y_train = self.y_train_pp
48     y_test = self.y_test_pp
49     alpha = self.alpha
50
51     # Add a column of 1s to the training data to have the correct dimension.
52     X_train = np.hstack([np.ones((X_train.shape[0], 1)), X_train])
53
54     n_features = X_train.shape[1]
55     I = np.eye(n_features)
56     w = np.linalg.inv(X_train.T.dot(X_train) + alpha * I).dot(X_train.T).dot(y_train)
57
58     X_test = np.hstack([np.ones((X_test.shape[0], 1)), X_test])
59     y_pred = X_test.dot(w)
60
61     # calculate the mean squared error
62     mse = np.mean((y_test - y_pred)**2)
63     r_squared = 1 - (np.sum((y_test - y_pred)**2) / np.sum((y_test - np.mean(y_test))**2))
64
65     # calculate root mean squared error (RMSE)
66     rmse = np.sqrt(mse)
67
68     print("Predicted target value:", y_pred[0])
69     print("Mean squared error (MSE):", mse)
70     print("R-squared:", r_squared)
71     print("Root mean squared error (RMSE):", rmse)

```

Code for Part IV

```

1 class LassoRegression(Predictor):
2     def __init__(self, dataset, alpha=1, max_iter=1000, tol=0.0001):
3         super().__init__(dataset)
4         self.alpha = alpha
5         self.max_iter = max_iter
6         self.tol = tol
7
8     def regression_raw(self):
9         """
10        Fits a lasso regression model on the training data using the specified regularization
    ↪ parameter alpha, iterations, and tolerance.
11        Using raw dataset.
12        """
13        # Define hyperparameters
14        alpha = self.alpha # regularization strength
15        max_iterations = self.max_iter # number of gradient descent iterations
16        tolerance = self.tol
17
18        # Load the data
19        diabetes = pd.read_csv("diabetes.csv")
20
21        diabetes.insert(0, "Intercept", 1)
22
23        train_size = int(0.8 * len(diabetes))
24

```

```

25 X_train = diabetes.iloc[:train_size, :-1].values
26 y_train = diabetes.iloc[:train_size, -1].values
27 X_test = diabetes.iloc[train_size:, :-1].values
28 y_test = diabetes.iloc[train_size:, -1].values
29
30 theta_lasso = np.zeros(X_train.shape[1])
31 for i in range(max_iterations):
32     theta_prev = theta_lasso.copy()
33     for j in range(X_train.shape[1]):
34         if j == 0:
35             theta_lasso[j] = np.mean(y_train)
36         else:
37             xj = X_train[:, j]
38             rj = y_train - X_train @ theta_lasso + xj * theta_lasso[j]
39             zj = xj @ xj
40             if zj == 0:
41                 theta_lasso[j] = 0
42             else:
43                 if np.sum(xj * rj) > alpha / 2:
44                     theta_lasso[j] = (np.sum(xj * rj) - alpha / 2) / zj
45                 elif np.sum(xj * rj) < - alpha / 2:
46                     theta_lasso[j] = (np.sum(xj * rj) + alpha / 2) / zj
47                 else:
48                     theta_lasso[j] = 0
49             if np.sum((theta_lasso - theta_prev) ** 2) < tolerance:
50                 break
51
52 sst = np.sum((y_test - np.mean(y_test)) ** 2)
53
54 y_pred_lasso = X_test @ theta_lasso
55 mse_lasso = np.mean((y_test - y_pred_lasso) ** 2)
56 ssr_lasso = np.sum((y_pred_lasso - np.mean(y_test)) ** 2)
57 r_squared_lasso = 1 - (ssr_lasso / sst)
58
59 rmse_lasso = np.sqrt(mse_lasso)
60
61 print("Lasso regression:")
62 print("Mean squared error (MSE):", mse_lasso)
63 print("R-squared:", r_squared_lasso)
64 print("Root mean squared error (RMSE):", rmse_lasso)
65
66 def regression_pp(self):
67     """
68     Fits a lasso regression model on the training data using the specified regularization
69     → parameter alpha, iterations, and tolerance.
70     Using processed dataset.
71     """
72     # Define hyperparameters
73     alpha = self.alpha # regularization strength
74     max_iterations = self.max_iter # number of gradient descent iterations
75     tolerance = self.tol
76
77     # Load the data
78     diabetes_norm = pd.read_csv("diabetes_norm.csv")
79
80     diabetes_norm.insert(0, "Intercept", 1)
81
82     train_size = int(0.8 * len(diabetes_norm))
83
84     X_train = diabetes_norm.iloc[:train_size, :-1].values
85     y_train = diabetes_norm.iloc[:train_size, -1].values
86     X_test = diabetes_norm.iloc[train_size:, :-1].values
87     y_test = diabetes_norm.iloc[train_size:, -1].values

```



```

88 theta_lasso = np.zeros(X_train.shape[1])
89 for i in range(max_iterations):
90     theta_prev = theta_lasso.copy()
91     for j in range(X_train.shape[1]):
92         if j == 0:
93             theta_lasso[j] = np.mean(y_train)
94         else:
95             xj = X_train[:, j]
96             rj = y_train - X_train @ theta_lasso + xj * theta_lasso[j]
97             zj = xj @ xj
98             if zj == 0:
99                 theta_lasso[j] = 0
100             else:
101                 if np.sum(xj * rj) > alpha / 2:
102                     theta_lasso[j] = (np.sum(xj * rj) - alpha / 2) / zj
103                 elif np.sum(xj * rj) < - alpha / 2:
104                     theta_lasso[j] = (np.sum(xj * rj) + alpha / 2) / zj
105                 else:
106                     theta_lasso[j] = 0
107             if np.sum((theta_lasso - theta_prev) ** 2) < tolerance:
108                 break
109
110 sst = np.sum((y_test - np.mean(y_test)) ** 2)
111
112 y_pred_lasso = X_test @ theta_lasso
113 mse_lasso = np.mean((y_test - y_pred_lasso) ** 2)
114 ssr_lasso = np.sum((y_pred_lasso - np.mean(y_test)) ** 2)
115 r_squared_lasso = 1 - (ssr_lasso / sst)
116
117 rmse_lasso = np.sqrt(mse_lasso)
118
119 print("Lasso regression:")
120 print("Mean squared error (MSE):", mse_lasso)
121 print("R-squared:", r_squared_lasso)
122 print("Root mean squared error (RMSE):", rmse_lasso)

```