

Linear Regression from scratch

The goal of this exercise is to implement the linear regression algorithm. The dataset is about predicting salary given gpa and years of experience. The steps to implement are as follows.

1. Read the data from a file (gpa_year_experience.csv)
2. Scale the attributes
3. Compute the error at each iteration and save the error values in vector
4. Plot the error vector as a curve in the end
5. Predict a new instance.
6. Compare with SGDRegressor
7. Create polynomial features and predict new instance

```
# import libraries
import pandas as pd
import numpy as np
import sklearn
from sklearn import *
import matplotlib.pyplot as plt

# ignore warnings
import warnings
warnings.filterwarnings('ignore')

# load data and show first 5 rows
data =
pd.read_csv('https://raw.githubusercontent.com/thomouvic/SENG474/main/
data/gpa_years_experience.csv')
data.head()
```

	gpa	years_of_experience	salary
0	70	1.0	50
1	80	2.0	55
2	65	2.0	45
3	70	2.5	60
4	65	2.7	58

```
# prepare data, split columns into X and y
# X should be a numpy array of shape (m, n), use .values to convert
from dataframe to numpy array
# y should be a numpy array of shape (m,), use .values to convert from
dataframe to numpy array
X = data.drop(['salary'], axis=1).values
Y = data.drop(['gpa', 'years_of_experience'], axis=1).values

# extract m and n from X using X.shape[0] to get m and X.shape[1] to
get n
print('m =', X.shape[0], '\nn =', X.shape[1])
```

```

m = 25
n = 2

# y should be a numpy array of shape (m, 1), use reshape(m, 1) to
# reshape y from (m,) to (m, 1)
Y.reshape(Y.size, 1)
Y.shape

(25, 1)

# normalize X using min-max scaler
(sklearn.preprocessing.MinMaxScaler)
scaler = sklearn.preprocessing.MinMaxScaler()

X = scaler.fit_transform(X)
X

array([[0.3125 , 0.        ],
       [0.625  , 0.125   ],
       [0.15625, 0.125   ],
       [0.3125 , 0.1875  ],
       [0.15625, 0.2125  ],
       [0.625  , 0.25    ],
       [0.9375 , 0.25    ],
       [1.      , 0.275   ],
       [0.      , 0.3125  ],
       [0.3125 , 0.3375  ],
       [0.5     , 0.375   ],
       [0.78125, 0.4375  ],
       [0.625  , 0.5     ],
       [0.      , 0.5625  ],
       [0.125  , 0.6     ],
       [0.      , 0.625   ],
       [0.84375, 0.625   ],
       [0.9375 , 0.6875  ],
       [0.46875, 0.75    ],
       [0.625  , 0.75    ],
       [0.46875, 0.8125  ],
       [0.3125 , 0.875   ],
       [0.625  , 0.9375  ],
       [0.9375 , 0.9625  ],
       [0.78125, 1.      ]])

# add dummy feature to X using scikit-learn dummy feature
(sklearn.preprocessing.add_dummy_feature)
X = sklearn.preprocessing.add_dummy_feature(X)
X

array([[1.      , 0.3125 , 0.        ],
       [1.      , 0.625  , 0.125   ],
       [1.      , 0.15625, 0.125   ],

```

```

[1.      , 0.3125 , 0.1875 ],
[1.      , 0.15625, 0.2125 ],
[1.      , 0.625  , 0.25   ],
[1.      , 0.9375 , 0.25   ],
[1.      , 1.      , 0.275  ],
[1.      , 0.      , 0.3125 ],
[1.      , 0.3125 , 0.3375 ],
[1.      , 0.5     , 0.375  ],
[1.      , 0.78125, 0.4375 ],
[1.      , 0.625  , 0.5     ],
[1.      , 0.      , 0.5625 ],
[1.      , 0.125  , 0.6     ],
[1.      , 0.      , 0.625  ],
[1.      , 0.84375, 0.625  ],
[1.      , 0.9375 , 0.6875 ],
[1.      , 0.46875, 0.75   ],
[1.      , 0.625  , 0.75   ],
[1.      , 0.46875, 0.8125 ],
[1.      , 0.3125 , 0.875  ],
[1.      , 0.625  , 0.9375 ],
[1.      , 0.9375 , 0.9625 ],
[1.      , 0.78125, 1.      ]]

```

```

# print shapes of X and y
# X should be (m, n+1) and y should be (m, 1)
print('Y.shape =', Y.shape)
print('X.shape =', X.shape)

Y.shape = (25, 1)
X.shape = (25, 3)

eta = 0.1 # learning rate
n_epochs = 10
np.random.seed(42) # set random seed to 42 for reproducibility

# create theta, of shape (n+1, 1) and initialize it to random values
using np.random.randn
theta = np.random.rand(3, 1)

E = [] # list to store errors at each epoch
# compute error for initial theta and append to E
error = sklearn.metrics.mean_squared_error(Y, np.dot(X, theta),
squared=False)
E.append(error)

# loop over n_epochs
# for each epoch: compute gradients, update theta, compute error,
append error to E
for epoch in range(n_epochs):

```

```

# compute gradients

gradients = np.dot(X.T, (np.dot(X, theta) - Y)) / len(Y)

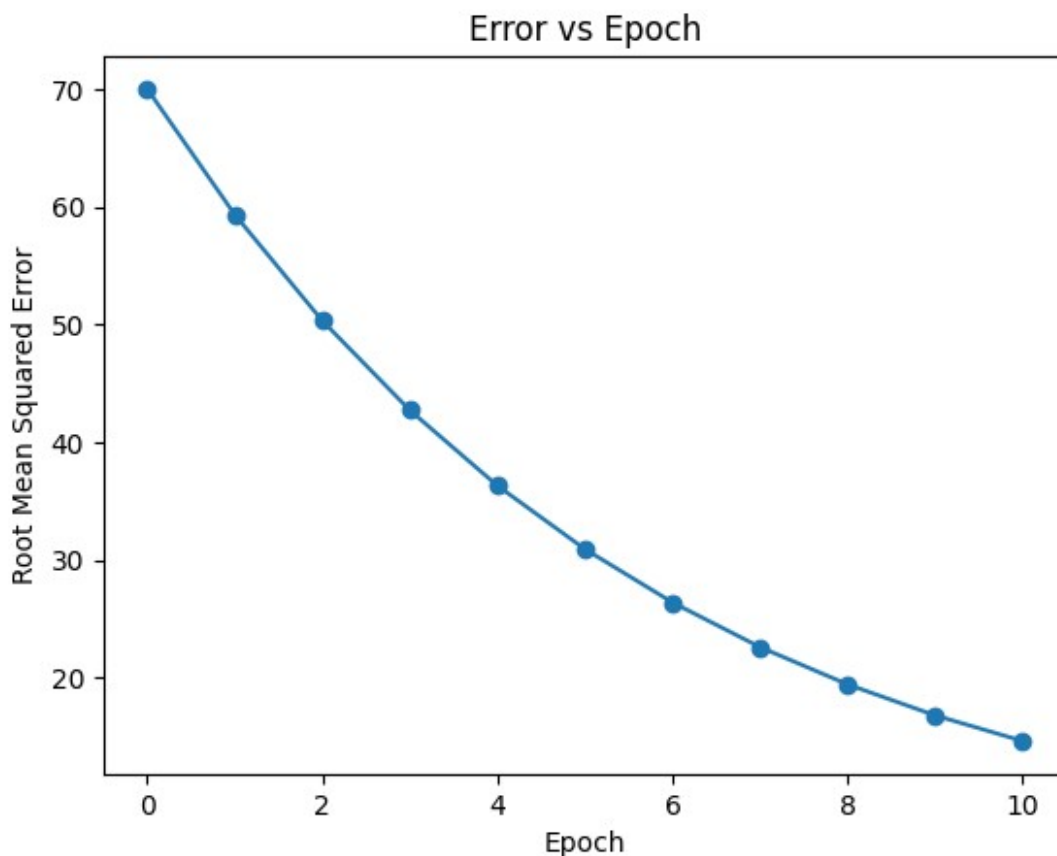
# update theta using the gradients and learning rate
theta = theta - eta * gradients

# compute error and append to E
error = sklearn.metrics.mean_squared_error(Y, np.dot(X, theta),
squared=False)
E.append(error)

# plot error vs epoch
plt.plot(range(n_epochs + 1), E, marker='o')
plt.xlabel('Epoch')
plt.ylabel('Root Mean Squared Error')
plt.title('Error vs Epoch')
plt.show()

# print final theta
print("Final Theta:", theta)

```



```
Final Theta: [[36.85414672]
 [19.29506171]
 [21.06864646]]
```

```
# let's predict the salary for a person who has gpa=70 and
years_of_experience=3.
```

```
# create a numpy array x of shape (1, 2) with these values
```

```
# scale features using the same scaler we used earlier
```

```
# insert dummy feature using dummy feature function
```

```
# Predict salary of x
```

```
x = np.array([[70,3]])
```

```
x = scaler.transform(x)
```

```
x = sklearn.preprocessing.add_dummy_feature(x)
```

```
print('predicted salary using our model =', np.dot(x, theta)[0][0])
```

```
predicted salary using our model = 48.15101511308938
```

```
# Let's compare with scikit-learn's SGDRegressor
```

```
# use SGDRegressor from scikit-learn to fit the data
```

```
# use max_iter=1000, eta0=0.1, random_state=42
```

```
sgd_reg = sklearn.linear_model.SGDRegressor(max_iter=1000, eta0=0.1,
random_state=42)
```

```
sgd_reg.fit(X, Y)
```

```
SGDRegressor(eta0=0.1, random_state=42)
```

```
# predict salary of x using sgd
```

```
print('predicted salary using SGDRegressor =', sgd_reg.predict(x)[0])
```

```
predicted salary using SGDRegressor = 59.444338469177936
```

```
# create polynomial features of degree 2 using scikit-learn
```

```
PolynomialFeatures
```

```
# create X_poly using fit_transform
```

```
# create x_poly using transform
```

```
# fit the data using SGDRegressor
```

```
# predict salary of x using sgd
```

```
poly = sklearn.preprocessing.PolynomialFeatures(2)
```

```
X_poly = poly.fit_transform(X)
```

```
x_poly = poly.fit_transform(x)
```

```
sgd_reg.fit(X_poly, Y)
```

```
print('predicted salary using SGDRegressor and Polynomial Features =',
sgd_reg.predict(x_poly)[0])
```

```
predicted salary using SGDRegressor and Polynomial Features =  
60.056903205155535
```