

Optional Lab: Linear Regression using Scikit-Learn

There is an open-source, commercially usable machine learning toolkit called [scikit-learn \(https://scikit-learn.org/stable/index.html\)](https://scikit-learn.org/stable/index.html). This toolkit contains implementations of many of the algorithms that you will work with in this course.

Goals

In this lab you will:

- Utilize scikit-learn to implement linear regression using Gradient Descent

Tools

You will utilize functions from scikit-learn as well as matplotlib and NumPy.

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 from sklearn.linear_model import SGDRegressor
        4 from sklearn.preprocessing import StandardScaler
        5 from lab_utils_multi import load_house_data
        6 from lab_utils_common import dlc
        7 np.set_printoptions(precision=2)
        8 plt.style.use('./deeplearning.mplstyle')
```

Gradient Descent

Scikit-learn has a gradient descent regression model [sklearn.linear_model.SGDRegressor \(https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDRegressor.html#examples-using-sklearn-linear-model-sgdregressor\)](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDRegressor.html#examples-using-sklearn-linear-model-sgdregressor). Like your previous implementation of gradient descent, this model performs best with normalized inputs. [sklearn.preprocessing.StandardScaler \(https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html#sklearn.preprocessing.StandardScaler\)](https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html#sklearn.preprocessing.StandardScaler) will perform z-score normalization as in a previous lab. Here it is referred to as 'standard score'.

Load the data set

```
In [2]: 1 X_train, y_train = load_house_data()
        2 X_features = ['size(sqft)', 'bedrooms', 'floors', 'age']
```

Scale/normalize the training data

```
In [3]: 1 scaler = StandardScaler()
        2 X_norm = scaler.fit_transform(X_train)
        3 print(f"Peak to Peak range by column in Raw X:{np.ptp(X_train,axis=0)}")
        4 print(f"Peak to Peak range by column in Normalized X:{np.ptp(X_norm,axis=0)}")
```

```
Peak to Peak range by column in Raw X:[2.41e+03 4.00e+00 1.00e+00 9.50e+01]
Peak to Peak range by column in Normalized X:[5.85 6.14 2.06 3.69]
```

Create and fit the regression model

```
In [4]: 1 sgdr = SGDRegressor(max_iter=1000)
        2 sgdr.fit(X_norm, y_train)
        3 print(sgdr)
        4 print(f"number of iterations completed: {sgdr.n_iter_}, number of weight updates:
        ~
```

```
SGDRegressor()
number of iterations completed: 128, number of weight updates: 12673.0
```

View parameters

Note, the parameters are associated with the *normalized* input data. The fit parameters are very close to those found in the previous lab with this data.

```
In [5]: 1 b_norm = sgdr.intercept_
        2 w_norm = sgdr.coef_
        3 print(f"model parameters:                w: {w_norm}, b:{b_norm}")
        4 print( "model parameters from previous lab: w: [110.56 -21.27 -32.71 -37.97], b:
        ~
```

```
model parameters:                w: [110.16 -21.11 -32.53 -38.06], b:[363.1]
model parameters from previous lab: w: [110.56 -21.27 -32.71 -37.97], b: 363.16
```

Make predictions

Predict the targets of the training data. Use both the `predict` routine and compute using w and b .

```
In [6]: 1 # make a prediction using sgdr.predict()
        2 y_pred_sgd = sgdr.predict(X_norm)
        3 # make a prediction using w,b.
        4 y_pred = np.dot(X_norm, w_norm) + b_norm
        5 print(f"prediction using np.dot() and sgdr.predict match: {(y_pred == y_pred_sgd)}")
        6
        7 print(f"Prediction on training set:\n{y_pred[:4]}")
        8 print(f"Target values \n{y_train[:4]}")
        ~
```

```
prediction using np.dot() and sgdr.predict match: True
Prediction on training set:
[295.13 485.76 389.44 491.93]
Target values
[300.  509.8 394.  540. ]
```

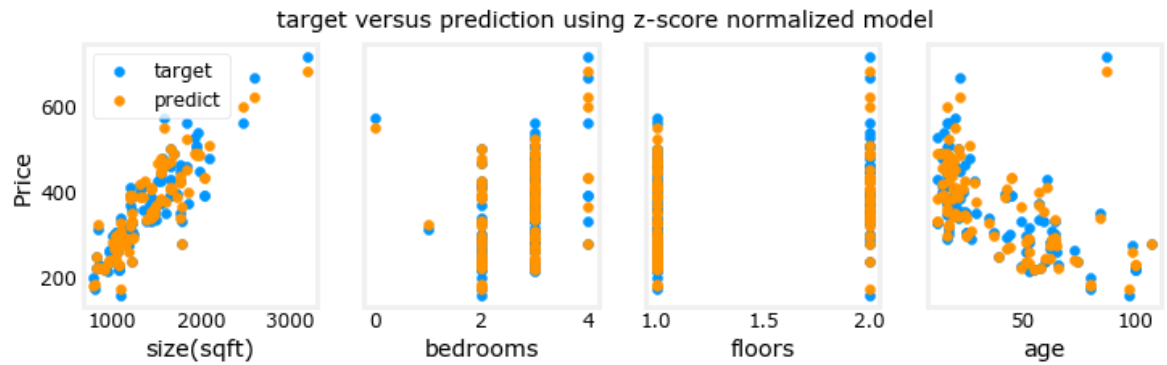
Plot Results

Let's plot the predictions versus the target values.

```

In [7]: 1 # plot predictions and targets vs original features
2 fig,ax=plt.subplots(1,4,figsize=(12,3),sharey=True)
3 for i in range(len(ax)):
4     ax[i].scatter(X_train[:,i],y_train, label = 'target')
5     ax[i].set_xlabel(X_features[i])
6     ax[i].scatter(X_train[:,i],y_pred,color=dlc["dlorange"], label = 'predict')
7 ax[0].set_ylabel("Price"); ax[0].legend();
8 fig.suptitle("target versus prediction using z-score normalized model")
9 plt.show()

```



Congratulations!

In this lab you:

- utilized an open-source machine learning toolkit, scikit-learn
- implemented linear regression using gradient descent and feature normalization from that toolkit

In []: