

# exercise\_04\_linear\_regression

November 13, 2021

## 1 Programming Task: Linear Regression

```
[1]: import numpy as np

from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
```

### 1.1 Your task

This notebook provides a code skeleton for performing linear regression. Your task is to complete the functions where required. You are only allowed to use built-in Python functions, as well as any **numpy** functions. No other libraries / imports are allowed.

In the beginning of every function there is docstring which specifies the input and expected output. Write your code in a way that adheres to it. You may only use plain python and anything that we imported for you above such as numpy functions (i.e. no scikit-learn classifiers).

### 1.2 Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is 1. Run all the cells of the notebook (**Kernel -> Restart & Run All**) 2. Export/download the notebook as PDF (**File -> Download as -> PDF via LaTeX (.pdf)**) 3. Concatenate your solutions for other tasks with the output of Step 2. On Linux you can simply use **pdfunite**, there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

**Make sure** you are using **nbconvert Version 5.5 or later** by running **jupyter nbconvert --version**. Older versions clip lines that exceed page width, which makes your code harder to grade.

### 1.3 Load and preprocess the data

In this assignment we will work with the Boston Housing Dataset. The data consists of 506 samples. Each sample represents a district in the city of Boston and has 13 features, such as crime rate or taxation level. The regression target is the median house price in the given district (in \$1000's).

More details can be found here: <http://lib.stat.cmu.edu/datasets/boston>

```
[2]: X , y = load_boston(return_X_y=True)

# Add a vector of ones to the data matrix to absorb the bias term
# (Recall slide #7 from the lecture)
X = np.hstack([np.ones([X.shape[0], 1]), X])
# From now on, D refers to the number of features in the AUGMENTED dataset (i.e.
  ↳ including the dummy '1' feature for the absorbed bias term)

# Split into train and test
test_size = 0.2
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)
```

## 1.4 Task 1: Fit standard linear regression

```
[3]: def fit_least_squares(X, y):
    """Fit ordinary least squares model to the data.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
        Regression targets.

    Returns
    -----
    w : array, shape [D]
        Optimal regression coefficients (w[0] is the bias term).

    """

    ### YOUR CODE HERE ###
    X_T = np.transpose(X)
    X_mul_inv = np.linalg.inv(np.matmul(X_T, X))
    return np.matmul(X_mul_inv, X_T).dot(y)
```

## 1.5 Task 2: Fit ridge regression

```
[4]: def fit_ridge(X, y, reg_strength):
    """Fit ridge regression model to the data.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
```

```

    Regression targets.
    reg_strength : float
        L2 regularization strength (denoted by lambda in the lecture)

    Returns
    -----
    w : array, shape [D]
        Optimal regression coefficients (w[0] is the bias term).

    """
    ### YOUR CODE HERE ###
    _, D = X.shape
    reg_addM = np.zeros((D, D), X.dtype)
    np.fill_diagonal(reg_addM, reg_strength)
    X_T = np.transpose(X)
    mul_inv = np.linalg.inv(np.matmul(X_T, X) + reg_addM)
    return np.matmul(mul_inv, X_T).dot(y)

```

## 1.6 Task 3: Generate predictions for new data

```

[5]: def predict_linear_model(X, w):
    """Generate predictions for the given samples.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    w : array, shape [D]
        Regression coefficients.

    Returns
    -----
    y_pred : array, shape [N]
        Predicted regression targets for the input data.

    """
    return X.dot(w)

```

## 1.7 Task 4: Mean squared error

```

[6]: def mean_squared_error(y_true, y_pred):
    """Compute mean squared error between true and predicted regression targets.

    Reference: `https://en.wikipedia.org/wiki/Mean_squared_error`

    Parameters

```

```

-----
y_true : array
    True regression targets.
y_pred : array
    Predicted regression targets.

Returns
-----
mse : float
    Mean squared error.

"""
### YOUR CODE HERE ###
dif_squared = np.power(y_pred - y_true, 2)
return np.mean(dif_squared)

```

## 1.8 Compare the two models

The reference implementation produces \* MSE for Least squares  $\approx$  **23.96** \* MSE for Ridge regression  $\approx$  **21.03**

Your results might be slightly (i.e.  $\pm 1\%$ ) different from the reference solution due to numerical reasons.

```

[7]: # Load the data
np.random.seed(1234)
X, y = load_boston(return_X_y=True)
X = np.hstack([np.ones([X.shape[0], 1]), X])
test_size = 0.2
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)

# Ordinary least squares regression
w_ls = fit_least_squares(X_train, y_train)
y_pred_ls = predict_linear_model(X_test, w_ls)
mse_ls = mean_squared_error(y_test, y_pred_ls)
print('MSE for Least squares = {0}'.format(mse_ls))

# Ridge regression
reg_strength = 1
w_ridge = fit_ridge(X_train, y_train, reg_strength)
y_pred_ridge = predict_linear_model(X_test, w_ridge)
mse_ridge = mean_squared_error(y_test, y_pred_ridge)
print('MSE for Ridge regression = {0}'.format(mse_ridge))

```

MSE for Least squares = 23.96457138494987

MSE for Ridge regression = 21.034931215918522