```
[MLF] H. Jiang, "Machine Learning Fundamentals: A Concise Introduction", Cambridge University Press, 2021. (bibtex)
         The purpose of this lab is to apply a simple machine learning method, namely linear regression, to some regression and classification tasks on two popular data
         sets. We will show how linear regression may differ when used to solve a regression or classification problem. As we know, linear regression is simple enough
         so that we can derive the closed-form solution to solve it. In this project, we will use both the closed-form method and an iterative gradient descent method
         (e.g. minibatch SGD) to solve linear regression for these tasks and compare their pros and cons in practice. Moreover, we will use linear regression as a simple
         example to explain some fine-tuning tricks when using any iterative optimization methods (e.g. SGD) in machine learning. As we will see in the up-coming
         projects, these tricks become vital in learning large models in machine learning, such as deep neural networks.
         Prerequisites: N/A
         I. Linear Regression for Regression
         Example 2.1:
         Use linear regression to predict house prices in the popular Boston House data set. Consider to use both the closed-form solution and an iterative method to fit
         to the data and discuss their pros and cons with some experimental results.
In [4]: # download boston house data from Google drive
         !gdown --folder https://drive.google.com/drive/folders/12L9XNwhIH2wQBa4-IdQrhsrtgRFbeIMZ 2> /dev/null
         Processing file 1IZf0tFGW3Zv3Ax5gxKgnklHe1DUI7CAI boston.csv
         Building directory structure completed
In [7]: # load Boston House data set
         import pandas as pd
         import numpy as np
         raw_data = pd.read_csv('Boston/boston.csv', header=None)
         data rows = np.reshape(raw data.to numpy(), (506,14))
         data = data_rows[:,:13]
         target = data_rows[:,13]
         # normalize input features to zero-mean and unit-variance
         data = (data-np.mean(data, axis=0))/np.std(data, axis=0)
         print(data.shape)
         print(target.shape)
         (506, 13)
         (506,)
In [ ]: # use the closed-form solution (Eq(6.9) on page 112)
         # add a constant column of '1' to accomodate the bias (see the margin note on page 107)
         data_wb = np.hstack((data, np.ones((data.shape[0], 1), dtype=data.dtype)))
         print(data_wb.shape)
         # refer to the closed-form solution, i.e. Eq.(6.9) on page 112
         w = np.linalg.inv(data_wb.T @ data_wb) @ data_wb.T @ target
         # calculate the mean square error in the training set
         predict = data wb @ w
         error = np.sum((predict - target)*(predict - target))/data.shape[0]
         print(f'mean square error for the closed-form solution: {error:.5f}')
         (506, 14)
         mean square error for the closed-form solution: 21.89483
         Consider to solve the above linear regression using an iterative optimization, such as gradient descent.
         Refer to eq.(6.8) on page 112, the objective function in linear regression, i.e. the mean square error (MSE), is given as:
                                                         E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i} - y_{i})^{2}
                                                              = \frac{1}{2} \left( \mathbf{w}^{\mathsf{T}} \mathbf{X}^{\mathsf{T}} \mathbf{X} \mathbf{w} - 2 \mathbf{w}^{\mathsf{T}} \mathbf{X}^{\mathsf{T}} \mathbf{y} + \mathbf{y}^{\mathsf{T}} \mathbf{y} \right)
         we can show that its gradient can be computed in several equivalent ways as follows:
                                                     \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^{N} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i} - y_{i}) \mathbf{x}_{i} = \sum_{i=1}^{N} \mathbf{x}_{i} (\mathbf{x}_{i}^{\mathsf{T}} \mathbf{w} - y_{i})
                                                            = \left(\sum_{i=1}^{N} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}}\right) \mathbf{w} - \sum_{i=1}^{N} y_{i} \mathbf{x}_{i}
                                                             = \mathbf{X}^{\mathsf{T}} \mathbf{X} \mathbf{w} - \mathbf{X}^{\mathsf{T}} \mathbf{y}
         where \mathbf{X} and \mathbf{y} are defined in the same way as on page 112.
         In the following, we use the formula from last row to calculate gradients via vectorization. Furthermore, we implement a mini-batch SGD, .i.e. Algorithm 2.3 on
         page 62, to learn linear regression iteratively.
In [ ]: # solve linear regression using gradient descent
         import numpy as np
         class Optimizer():
           def init__(self, lr, annealing_rate, batch_size, max_epochs):
              self.lr = lr
              self.annealing_rate = annealing_rate
             self.batch size = batch size
             self.max epochs = max epochs
         # X[N,d]: input features; y[N]: output targets; op: hyper-parameters for optimzer
         def linear_regression_gd(X, y, op):
           n = X.shape[0] # number of samples
           w = np.zeros(X.shape[1]) # initialization
           lr = op.lr
           errors = np.zeros(op.max_epochs)
            for epoch in range(op.max epochs):
             indices = np.random.permutation(n) #randomly shuffle data indices
              for batch_start in range(0, n, op.batch_size):
                X batch = X[indices[batch start:batch start + op.batch size]]
                y batch = y[indices[batch start:batch start + op.batch size]]
                # vectorization to compute gradients for a whole mini-batch (see the above formula)
                w_grad = X_batch.T @ X_batch @ w - X_batch.T @ y_batch
                w -= lr * w grad / X batch.shape[0]
              diff = X @ w - y # prediction difference
              errors[epoch] = np.sum(diff*diff)/n
             lr *= op.annealing rate
              #print(f'epoch={epoch}: the mean square error is {errors[epoch]}')
           return w, errors
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.001, annealing_rate=0.99, batch_size=30, max_epochs=500)
         w, errors = linear regression gd(data wb, target, op)
         plt.title('A larger mini-batch; converging smoothly but still a gap')
         plt.xlabel('number of epoches')
         plt.plot(errors, 'b', 21.89*np.ones(errors.shape[0]), 'c--')
         plt.legend(['gradient descent', 'closed-form solution'])
Out[ ]: <matplotlib.legend.Legend at 0x7ffaed7e1d10>
             A larger mini-batch; converging smoothly but still a gap
                                        gradient descent
                                        --- closed-form solution
          500
          400
          300
          200
          100
                       100
                                               400
                             number of epoches
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.001, annealing_rate=0.99, batch_size=2, max_epochs=100)
         w, errors = linear_regression_gd(data_wb, target, op)
         plt.title('use a small mini-batch; converging quicklly and nicely to the optimum')
         plt.xlabel('number of epoches')
         plt.plot(errors, 'b', 21.89*np.ones(errors.shape[0]), 'c--')
         plt.legend(['gradient descent', 'closed-form solution'])
Out[ ]: <matplotlib.legend.Legend at 0x7ffaed8d6c10>
          use a small mini-batch; converging quicklyy and nicely to the optimum
                                          gradient descent
                                           --- closed-form solution
             300
             250
             200
             150
             100
             50
                                number of epoches
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.01, annealing_rate=0.99, batch_size=20, max_epochs=100)
         w, errors = linear_regression_gd(data_wb, target, op)
         plt.title('use a larger learning rate and larger mini-batch, also converging nicely')
         plt.xlabel('number of epoches')
         plt.plot(errors, 'b', 21.89*np.ones(errors.shape[0]), 'c--')
         plt.legend(['gradient descent', 'closed-form solution'])
Out[ ]: <matplotlib.legend.Legend at 0x7ffaed7e1e50>
          use a larger learning rate and larger mini-batch, also converging nicely
                                           gradient descent
                                           --- closed-form solution
             300
             250
             200
             150
             100
              50
                                number of epoches
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.2, annealing_rate=0.99, batch_size=20, max_epochs=12)
         w, errors = linear_regression_gd(data_wb, target, op)
         plt.title('not converging due to an overly large learning rate is used')
         plt.xlabel('number of epoches')
         plt.plot(errors, 'b', 21.89*np.ones(errors.shape[0]), 'c--')
         plt.legend(['gradient descent', 'closed-form solution'])
Out[ ]: <matplotlib.legend.Legend at 0x7ffaed27b290>
           not converging due to an overly large learning rate is used
                                       gradient descent
                                       --- closed-form solution
          40
          35
          30
                             number of epoches
         Finally, let us show how to solve the above linear regression problem using the scikit-learning implmenetation.
In [ ]: import numpy as np
         from sklearn import linear_model
         from sklearn.metrics import mean squared error
         # Create linear regression object
         l regr = linear model.LinearRegression()
         # Train the model using the training set
         l_regr.fit(data_wb, target)
         # Make predictions using the same training set
         predict = l regr.predict(data wb)
         # The mean squared error
         print("Mean squared error: %.5f" % mean_squared_error(target, predict))
         Mean squared error: 21.89483
         II. Linear Regression for Classification
         Example 2.2:
         Use linear regression to build a binary classifier to classify two digits ('3' and '8') in the MNIST data set. Consider to use both the closed-form solution and an
         iterative method to fit to the data and discuss their pros and cons with experimental results.
In [1]: # download MNIST data
         !gdown --folder https://drive.google.com/drive/folders/1r20aRjc2iu903kN3Xj9jNYY2uMgcERY1 2> /dev/null
         # install python_mnist
         !pip install python_mnist
         Processing file 1Jf2XqGR7y1fzOZNKLJiom7GmZZUzXhfs t10k-images-idx3-ubyte
         Processing file 1qiYu9dW3ZNrlvTFO5fI4qf8Wtr8K-pCu t10k-labels-idx1-ubyte
         Processing file 1SnWvBcUETRJ53rEJozFUUo-hOQFPKxjp train-images-idx3-ubyte
         Processing file 1kKEIi_pwVHmabByAnwZQsaMgro9XiBFE train-labels-idx1-ubyte
         Building directory structure completed
         Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
         Collecting python_mnist
           Downloading python_mnist-0.7-py2.py3-none-any.whl (9.6 kB)
         Installing collected packages: python mnist
         Successfully installed python_mnist-0.7
In [2]: #load MINST images
         from mnist import MNIST
         import numpy as np
         mnist loader = MNIST('MNIST')
         train_data, train_label = mnist_loader.load_training()
         test_data, test_label = mnist_loader.load_testing()
         train_data = np.array(train_data, dtype='float')/255 # norm to [0,1]
         train_label = np.array(train_label, dtype='short')
         test_data = np.array(test_data, dtype='float')/255 # norm to [0,1]
         test_label = np.array(test_label, dtype='short')
         #add small random noise to avoid matrix singularity
         train_data += np.random.normal(0,0.0001,train_data.shape)
         print(train_data.shape, train_label.shape, test_data.shape, test_label.shape)
         (60000, 784) (60000,) (10000, 784) (10000,)
In [3]: # prepare digits '3' and '8' for linear regression
         digit_train_index = np.logical_or(train_label == 3, train_label == 8)
         X_train = train_data[digit_train_index]
         y_train = train_label[digit_train_index]
         digit_test_index = np.logical_or(test_label == 3, test_label == 8)
         X_test = test_data[digit_test_index]
         y_test = test_label[digit_test_index]
         # add a constant column of '1' to accomodate the bias (see the margin note on page 107)
         X_train = np.hstack((X_train, np.ones((X_train.shape[0], 1), dtype=X_train.dtype)))
         X_test = np.hstack((X_test, np.ones((X_test.shape[0], 1), dtype=X_test.dtype)))
         # convert labels: '3' => -1, '8' => +1
         CUTOFF = 5 # any number between '3' and '8'
         y_train = np.sign(y_train-CUTOFF)
         y_test = np.sign(y_test-CUTOFF)
         print(X_train.shape)
         print(y_train)
         print(X_test.shape)
         print(y_test)
         (11982, 785)
         [-1 -1 -1 \dots 1 -1 1]
         (1984, 785)
         [-1 -1 -1 ... -1 1 -1]
In [ ]: # use the closed-form solution
         # refer to the closed-form solution, i.e. Eq.(6.9) on page 112
         w = np.linalg.inv(X train.T @ X train) @ X train.T @ y train
         # calculate the mean square error and classification accuracy on the training set
         predict = X_train @ w
         error = np.sum((predict - y_train)*(predict - y_train))/X_train.shape[0]
         print(f'mean square error on training data for the closed-form solution: {error:.5f}')
         accuracy = np.count_nonzero(np.equal(np.sign(predict),y_train))/y_train.size*100.0
         print(f'classification accuracy on training data for the closed-form solution: {accuracy:.2f}%')
         # calculate the mean square error and classification accuracy on the test set
         predict = X test @ w
         error = np.sum((predict - y_test)*(predict - y_test))/X_test.shape[0]
         print(f'mean square error on test data for the closed-form solution: {error:.5f}')
         accuracy = np.count_nonzero(np.equal(np.sign(predict),y_test))/y_test.size*100.0
         print(f'classification accuracy on test data for the closed-form solution: {accuracy:.2f}%')
         mean square error on training data for the closed-form solution: 0.19629
         classification accuracy on training data for the closed-form solution: 96.99%
         mean square error on training data for the closed-form solution: 1.30940
         classification accuracy on test data for the closed-form solution: 95.92%
In [ ]: # use linear regression from sklearn
         import numpy as np
         from sklearn import linear model
         from sklearn.metrics import mean squared error
         # Create linear regression object
         l_regr = linear_model.LinearRegression()
         # Train the model using the training set
         l_regr.fit(X_train, y_train)
         # Make predictions using the same training set
         predict = l_regr.predict(X_train)
         print("Mean squared error on training data: %.5f" % mean squared error(y train, predict))
         # Make predictions using the test set
         predict = 1 regr.predict(X test)
         print("Mean squared error on test data: %.5f" % mean_squared_error(y_test, predict))
         Mean squared error on training data: 0.19629
         Mean squared error on test data: 1.30940
         Next, let us consider to use mini-batch stochastic gradient descent (SGD) to learn linear regression models for this binary classification problem. When we fine-
         tune any SGD method for a classification problem in machline learning, it is very important to monitor the following three learning curves:
          1. Classification Accuracy on the training set (curve A): this is the goal of the empirical risk mininization (ERM) of the zero-one loss for classification (see Eq.
             (5.6) on page 99).
          2. Classification Accuracy on an unseen test/development set (curve B): we need to compare the curves A and B over the learning course to monitor whether
             overfitting or underfit occurs. Overfitting happens when the gap between A and B is overly big while underfitting happens when A and B get very close and
            both of them yield fairly poor performance. Moreover, we can also monitor the curves A and B to determine when to terminate the learning proces for the
             best possible performance on the test/devopement set.
          3. The value of the learning objective function (curve C): because the zero-one loss is not directly minimizable, we will have to establish a proxy objective
             function according to some criteria (see Tabe 7.1 on page 135). These objective functions are closelly related to the zero-one loss but they are NOT the
             same. When we fine-tune an iterative optimization method, the first thing is to ensure that the value of the chosen objective function descreases over the
             entire learning course. If we cannot reduce the objective function (even when a sufficiently small learning rate is used), it is very likely that the
             implementation or code is buggy. Furthermore, if curve C is going down while curve A is not going up, this is another indicator that someting is still wrong
             in the implementation.
In [ ]: # solve linear regression using gradient descent
         import numpy as np
         class Optimizer():
           def __init__(self, lr, annealing_rate, batch_size, max_epochs):
              self.lr = lr
              self.annealing_rate = annealing_rate
             self.batch_size = batch_size
             self.max epochs = max epochs
         # X[N,d]: training features; y[N]: training targets;
         # X2[N,d]: test features; y2[N]: test targets;
         # op: hyper-parameters for optimzer
         # Note: X2 and y2 are not used in training
                  but only for computting the learning curve B
         def linear_regression_gd2(X, y, X2, y2, op):
           n = X.shape[0] # number of samples
           w = np.zeros(X.shape[1]) # initialization
           lr = op.lr
           errorsA = np.zeros(op.max_epochs)
           errorsB = np.zeros(op.max_epochs)
           errorsC = np.zeros(op.max epochs)
            for epoch in range(op.max epochs):
             indices = np.random.permutation(n) #randomly shuffle data indices
              for batch_start in range(0, n, op.batch_size):
                X_batch = X[indices[batch_start:batch_start + op.batch_size]]
                y batch = y[indices[batch start:batch start + op.batch size]]
                # vectorization to compute gradients for a whole mini-batch (see the above formula)
                w_grad = X_batch.T @ X_batch @ w - X_batch.T @ y_batch
                w -= lr * w_grad / X_batch.shape[0]
              # for learning curve C
              diff = X @ w - y # prediction difference
              errorsC[epoch] = np.sum(diff*diff)/n
              # for learning curve A
              predict = np.sign(X @ w)
              errorsA[epoch] = np.count nonzero(np.equal(predict,y))/y.size
              # for learning curve B
              predict2 = np.sign(X2 @ w)
             errorsB[epoch] = np.count_nonzero(np.equal(predict2,y2))/y2.size
             lr *= op.annealing rate
             print(f'epoch={epoch}: the mean square error is {errorsC[epoch]:.3f} ({errorsA[epoch]:.3f}, {errorsB[epoch]:.3f})')
           return w, errorsA, errorsB, errorsC
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.001, annealing_rate=0.99, batch_size=50, max_epochs=20)
         w, A, B, C = linear_regression_gd2(X_train, y_train, X_test, y_test, op)
         fig, ax = plt.subplots(2)
         fig.suptitle('monitoring three learning curves (A, B, C)')
         ax[0].plot(C, 'b', 0.196*np.ones(C.shape[0]), 'c--')
         ax[0].legend(['curve C', 'closed-form solution'])
         ax[1].plot(A, 'b', B, 'r')
         ax[1].legend(['curve A', 'curve B'])
         epoch=0: the mean square error is 0.398 (0.911,0.920)
         epoch=1: the mean square error is 0.309 (0.929,0.942)
         epoch=2: the mean square error is 0.279 (0.938,0.950)
         epoch=3: the mean square error is 0.264 (0.945,0.957)
         epoch=4: the mean square error is 0.254 (0.949,0.959)
         epoch=5: the mean square error is 0.248 (0.951,0.962)
         epoch=6: the mean square error is 0.244 (0.953,0.961)
         epoch=7: the mean square error is 0.240 (0.954,0.962)
         epoch=8: the mean square error is 0.237 (0.955,0.964)
         epoch=9: the mean square error is 0.235 (0.956,0.963)
         epoch=10: the mean square error is 0.234 (0.957,0.963)
         epoch=11: the mean square error is 0.232 (0.957,0.964)
         epoch=12: the mean square error is 0.231 (0.958,0.964)
         epoch=13: the mean square error is 0.230 (0.958,0.964)
         epoch=14: the mean square error is 0.229 (0.958,0.964)
         epoch=15: the mean square error is 0.228 (0.959,0.964)
         epoch=16: the mean square error is 0.227 (0.959,0.962)
         epoch=17: the mean square error is 0.226 (0.959,0.963)
         epoch=18: the mean square error is 0.226 (0.959,0.965)
         epoch=19: the mean square error is 0.225 (0.960,0.965)
Out[ ]: <matplotlib.legend.Legend at 0x7ffadb5ad890>
                   monitoring three learning curves (A, B, C)
           0.4
                                        — curve C
                                        --- closed-form solution
           0.3
               0.0 2.5 5.0 7.5 10.0 12.5 15.0 17.5
          0.96
          0.94
                                                 curve A
          0.92
                                                  curve B
                    2.5 5.0 7.5 10.0 12.5 15.0 17.5
         In the above setting, we use a large mini-batch (50), which leads to fairly smooth convergence. As we can see, even though there is a big gap in the objective
         function between SGD (curve C) and the closed-form solution, classification accuracy of SGD exceeds that of the closed-form solution on either the training or
         testing set. This indicates that MSE used in linear regression is NOT a good learning criterion for classification (see why in section 7.1.1 on page 136).
In [ ]: import matplotlib.pyplot as plt
         op = Optimizer(lr=0.001, annealing_rate=0.99, batch_size=5, max_epochs=10)
         w, A, B, C = linear_regression_gd2(X_train, y_train, X_test, y_test, op)
         fig, ax = plt.subplots(2)
         fig.suptitle('monitoring three learning curves (A, B, C)')
         ax[0].plot(C, 'b', 0.196*np.ones(C.shape[0]), 'c--')
         ax[0].legend(['curve C', 'closed-form solution'])
         ax[1].plot(A, 'b', B, 'r')
         ax[1].legend(['curve A', 'curve B'])
         epoch=0: the mean square error is 0.236 (0.956,0.963)
         epoch=1: the mean square error is 0.224 (0.960,0.964)
         epoch=2: the mean square error is 0.220 (0.961,0.963)
         epoch=3: the mean square error is 0.217 (0.962,0.964)
         epoch=4: the mean square error is 0.215 (0.962,0.963)
         epoch=5: the mean square error is 0.214 (0.963,0.963)
         epoch=6: the mean square error is 0.213 (0.964,0.965)
         epoch=7: the mean square error is 0.215 (0.963,0.962)
         epoch=8: the mean square error is 0.212 (0.963,0.964)
         epoch=9: the mean square error is 0.214 (0.963,0.963)
Out[]: <matplotlib.legend.Legend at 0x7ffadb49df50>
                     monitoring three learning curves (A, B, C)
                                          curve C
                                          closed-form solution
            0.22
            0.20
          0.9650
          0.9625
          0.9600
                                                   curve A
         In this setting, we use the same learning rate but a much smaller mini-batch size. A smaller mini-batch means more model updates in each epoch. As a result,
         the classification accuracy on the training set (curve A) is improved over the previous setting. However, we can see that the classification accuracy on the
         unseen set (curve B) starts to go down after epoch 6, which indicates that we should terminate the learning at epoch 6.
         Exercises
         Problem 2.1:
```

Use Ridge regression to solve the regression problem in Example 2.1 as well as the classification problem in Example 2.2, also implement both closed-form

Use LASSO to solve the regression problem in Example 2.1 as well as the classification problem in Example 2.2, compare the results of LASSO with those of

and iterative approachs, compare the results of Ridge regression with those of linear regression.

Problem 2.2:

linear regression and Ridge regression.

Lab 2: Linear Regression

NOTE: This is a lab project accompanying the following book [MLF] and it should be used together with the book.