

Cmpsc 448

Problem 1)

Original OLS objective: $f_{OLS}(w) = \frac{1}{2} \sum_{i=1}^n (w \cdot x_i - y_i)^2$

To get w take the derivative $\frac{\partial f_{OLS}}{\partial w} = \sum_{i=1}^n (w \cdot x_i - y_i) x_i$

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

Since the means are 0, we can rewrite it as $0 = \sum w x_i^2 - x_i y_i \rightarrow \sum w x_i^2 = \sum x_i y_i$

$$\frac{\sum x_i y_i}{n-1} \cdot \frac{n-1}{\sum x_i^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

Mistake OLS objective: $f_{OLS}(\tilde{w}) = \frac{1}{2} \sum_{i=1}^n (w y_i - x_i)^2$

To get \tilde{w} , take derivative $\frac{\partial f_{OLS}(\tilde{w})}{\partial \tilde{w}} = \sum_{i=1}^n (w y_i - x_i) y_i$

$$\tilde{w} = \frac{\sum x_i y_i}{\sum y_i^2}$$

Since the means are 0, we can rewrite as $0 = \sum \tilde{w} y_i^2 - x_i y_i \rightarrow \sum \tilde{w} y_i^2 = \sum x_i y_i$

$$\frac{\sum x_i y_i}{n-1} \cdot \frac{n-1}{\sum y_i^2} = \frac{\text{Cov}(x, y)}{\text{Var}(y)}$$

Relating w to \tilde{w}

Since $w = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$ and $\tilde{w} = \frac{\text{Cov}(x, y)}{\text{Var}(y)}$

$$\frac{w}{\tilde{w}} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} \cdot \frac{\text{Var}(y)}{\text{Cov}(x, y)} = \frac{\text{Var}(y)}{\text{Var}(x)} \rightarrow w = \tilde{w} \frac{\text{Var}(y)}{\text{Var}(x)}$$

To get $w = \tilde{w}$, we need the variances of x and y to be equal

Problem 2) $\hat{y} = w_1 x_{i1} + w_2 x_{i2}$ $\vec{x}_i = [x_{i1}, x_{i2}]^T$

Objective for OLS is $\frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$ $\vec{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}$ $\vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

$$f(w_1, w_2) = \frac{1}{2} \sum (w_1 x_{i1} + w_2 x_{i2} - y_i)^2$$

To find the optimal solution, take gradient of function

Use given formula for multidimensional OLS solution: $w^* = (\sum \vec{x}_i \vec{x}_i^T)^{-1} \sum y_i \vec{x}_i$

$$\vec{x}_i \vec{x}_i^T = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} \begin{bmatrix} x_{i1} & x_{i2} \end{bmatrix} = \begin{bmatrix} x_{i1}^2 & x_{i1} x_{i2} \\ x_{i1} x_{i2} & x_{i2}^2 \end{bmatrix}$$

$$y_i \vec{x}_i = \begin{bmatrix} x_{i1} y_i \\ x_{i2} y_i \end{bmatrix}$$

$$\sum y_i \vec{x}_i = \begin{bmatrix} \sum x_{i1} y_i \\ \sum x_{i2} y_i \end{bmatrix}$$

$$\sum \vec{x}_i \vec{x}_i^T = \begin{bmatrix} \sum x_{i1}^2 & \sum x_{i1} x_{i2} \\ \sum x_{i2} x_{i1} & \sum x_{i2}^2 \end{bmatrix}$$

$$(\sum \vec{x}_i \vec{x}_i^T)^{-1} = \frac{1}{\sum x_{i1}^2 \sum x_{i2}^2 - \sum x_{i1} x_{i2} \sum x_{i2} x_{i1}} \begin{bmatrix} \sum x_{i1}^2 & -\sum x_{i1} x_{i2} \\ -\sum x_{i2} x_{i1} & \sum x_{i2}^2 \end{bmatrix}$$

$$\vec{w}^* = \begin{bmatrix} \text{Var}(x_2) & -\text{Cov}(x_1, x_2) \\ \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_1, x_2) & \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_2, x_1) \\ -\text{Cov}(x_1, x_2) & \text{Var}(x_1) \\ \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_2, x_1) & \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_1, x_2) \end{bmatrix} \begin{bmatrix} \sum x_{i1} y_i \\ \sum x_{i2} y_i \end{bmatrix}$$

$$\sum x_{i1}^2 = \text{Var}(x_1)$$

$$\sum x_{i2}^2 = \text{Var}(x_2)$$

$$\sum x_{i1} x_{i2} = \text{Cov}(x_1, x_2)$$

$$\vec{w}^* = \begin{bmatrix} \text{Var}(x_2) \sum x_{i1} y_i - \text{Cov}(x_1, x_2) \sum x_{i2} y_i \\ \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_1, x_2) \\ \text{Var}(x_1) \sum x_{i2} y_i - \text{Cov}(x_1, x_2) \sum x_{i1} y_i \\ \text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_1, x_2) \end{bmatrix} = \begin{bmatrix} w_1^* \\ w_2^* \end{bmatrix}$$

now solve for each feature separately

$$f_{OLS}(\hat{y}) = \frac{1}{2} \sum (\hat{y}_i - y_i)^2 = \frac{1}{2} \sum (w_1 x_{i1} + w_2 x_{i2} - y_i)^2$$

Since we are solving for just w_1 , we set $w_2 = 0$

$$f_{OLS}(w_1) = \frac{1}{2} \sum (w_1 x_{i1} - y_i)^2$$

$$\frac{\partial f}{\partial w_1} = \sum (w_1 x_{i1} - y_i) x_{i1} = 0 \rightarrow \sum w_1 x_{i1}^2 = \sum x_{i1} y_i \rightarrow w_1 = \frac{\sum x_{i1} y_i}{\sum x_{i1}^2} = \frac{\text{Cov}(x_1, y)}{\text{Var}(x_1)}$$

$$\frac{\partial f}{\partial w_2} = \sum (w_1 x_{i1} - y_i) x_{i2} = 0 \rightarrow \sum w_1 x_{i2}^2 = \sum x_{i2} y_i \rightarrow w_2 = \frac{\sum x_{i2} y_i}{\sum x_{i2}^2} = \frac{\text{Cov}(x_2, y)}{\text{Var}(x_2)}$$

$$w_1^{*} = \frac{\text{Cov}(x_1, y)}{\text{Var}(x_1)} \quad w_2^{*} = \frac{\text{Cov}(x_2, y)}{\text{Var}(x_2)}$$

The two values for the 2-dimensional and the 1-dimensional solutions are different. The weights are different to each other. The scenario where

these weights would be the same is if the two variables were uncorrelated.

$$w_1^{*} = \frac{\text{Var}(x_2) \sum x_{i1} y_i - \text{Cov}(x_1, x_2) \sum x_{i2} y_i}{\text{Var}(x_1) \text{Var}(x_2) - \text{Cov}^2(x_1, x_2)} \underset{\text{if uncorrelated}}{\underset{\text{Cov}(x_1, x_2) = 0}{\rightarrow}} \frac{\text{Var}(x_2) \sum x_{i1} y_i - 0}{\text{Var}(x_1) \text{Var}(x_2) - 0} = \frac{\sum x_{i1} y_i}{\text{Var}(x_1)} = w_1'$$

Problem 3)

Objective function for Ridge: $\frac{1}{2} \sum (w_i x_i - y_i)^2 + \lambda \|w\|^2$

$$\text{find } w_1^{*} : \frac{1}{2} \sum (w_1 x_{i1} - y_i)^2 + \lambda w_1^2$$

$$\text{take derivative } \frac{\partial f}{\partial w_1} = \sum (w_1 x_{i1} - y_i) x_{i1} + 2\lambda w_1 = 0 \rightarrow \sum (w_1 x_{i1}^2 - x_{i1} y_i) + 2\lambda w_1 = 0$$

$$w_1 \sum x_{i1}^2 - \sum x_{i1} y_i + w_1 (2\lambda) \rightarrow w_1 = \frac{\sum x_{i1} y_i}{\sum x_{i1}^2 + 2\lambda} \quad w_1^{*} = \frac{\sum x_{i1} y_i}{\sum x_{i1}^2 + 2\lambda}$$

$$\text{now find } [w_1^{*}, w_2^{*}]^T \quad \vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$f(w_1, w_2) = \frac{1}{2} \sum (w_1 x_{i1} + w_2 x_{i2} - y_i)^2 + \lambda (w_1^2 + w_2^2)$$

$$\text{take gradient } \frac{\partial f}{\partial w_1} = \sum (w_1 x_{i1} + w_2 x_{i2} - y_i) x_{i1} + 2\lambda w_1 \rightarrow \sum [(w_1 + w_2) x_{i1}^2 - x_{i1} y_i] + 2\lambda w_1 = 0$$

$$= \sum w_1 x_{i1}^2 + \sum w_2 x_{i1}^2 - \sum x_{i1} y_i + 2\lambda w_1 = 0 \rightarrow \frac{\sum x_{i1} y_i - \sum w_2 x_{i1}^2}{\sum x_{i1}^2 + 2\lambda} = w_1$$

$$\frac{\partial f}{\partial w_2} = \sum (w_1 x_{i1} + w_2 x_{i2} - y_i) x_{i2} + 2\lambda w_2 \rightarrow \frac{\sum x_{i2} y_i - \sum w_1 x_{i2}^2}{\sum x_{i2}^2 + 2\lambda} = w_2$$

the equations for w_1 and w_2 are essentially the same equation

$$\text{we can simplify the equations further } w_1^{*} = \frac{\text{Cov}(x, y) - w_2 \text{Var}(x)}{\text{Var}(x) + 2\lambda}$$

$$w_2^{*} = \frac{\text{Cov}(x, y) - w_1 \text{Var}(x)}{\text{Var}(x) + 2\lambda}$$

from the first part we know $w_1^{*} = \frac{\text{Cov}(x, y)}{\text{Var}(x) + 2\lambda}$

$$w_1^{*} + w_2^{*} = \frac{\text{Cov}(x, y)}{\text{Var}(x) + 2\lambda} = 2w_1^{*} \text{ since } w_1^{*} = w_2^{*} \quad \text{so we get } w_1^{*} = 2w_1'$$

CPMSC 448 Homework 2

Problem 4

Name: Kush Lalwani

PSU id: kgl5163

In this problem will use the Pima Indians Diabetes dataset from the UCI repository to experiment with the k -NN algorithm and find the optimal value for the number of neighbors k . You do not need to implement the algorithm and encouraged to use the implementation in `scikit-learn`.

a)

Download the provided `Pima.csv` data file and load it using `pandas`. As a sanity check, make sure there are 768 rows of data (potential diabetes patients) and 9 columns (8 input features including `Pregnancies`, `Glucose`, `BloodPressure`, `SkinThickness`, `Insulin`, `BMI`, `DiabetesPedigreeFunction`, `Age`, and 1 target output). Note that the data file has no header and you might want to explicitly create the header. The last value in each row contains the target label for that row, and the remaining values are the features. Report the statics of each feature (min, max, average) and the histogram of the labels (target outputs).

```
In [40]: import pandas as pd
import matplotlib.pyplot as plt

varNames = ["Pregnancies", "Glucose", "BloodPressure", "SkinThickness", "Insulin",
            "DiabetesPedigreeFunction", "Age"]

data = pd.read_csv("Pima.csv", names = varNames)

print(f"Dataset shape: {data.shape}\n\nColumn Names: {data.columns}\n")

print(data.drop(columns=["Target"]).agg(["min", "max", "mean"]))

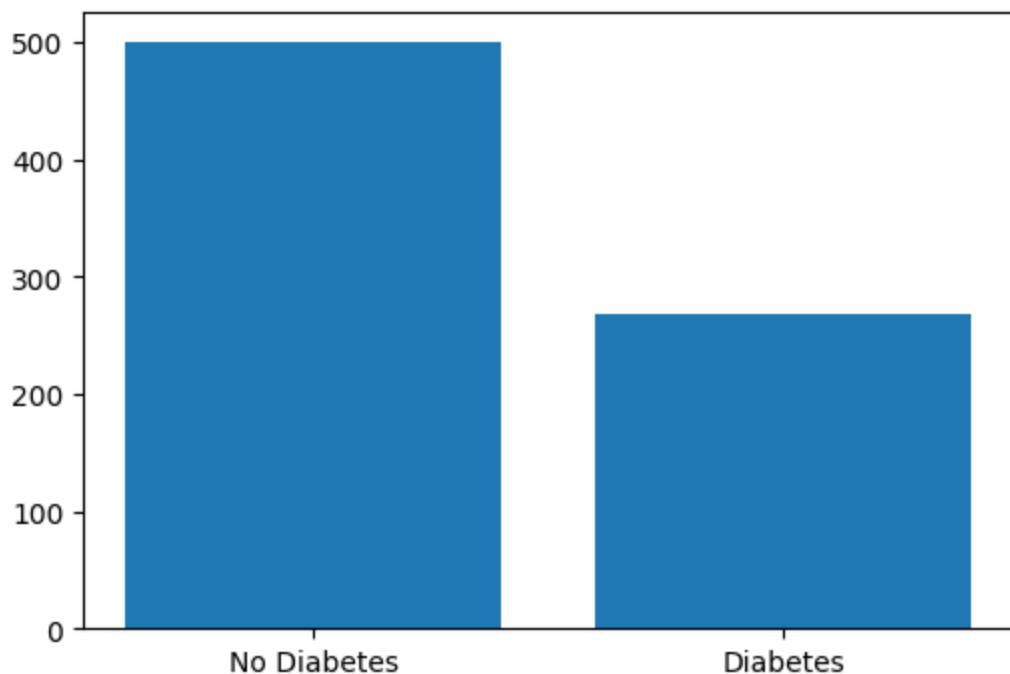
plt.figure(figsize=(6, 4))
data["Target"].hist(bins=2, rwidth=0.8)
plt.xticks([0.25, 0.75], ["No Diabetes", "Diabetes"])
plt.grid(visible = False)
plt.show()
```

```
Dataset shape: (768, 9)
```

```
Column Names: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',  
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Target'],  
      dtype='object')
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	\
min	0.000000	0.000000	0.000000	0.000000	0.000000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	
mean	3.845052	120.894531	69.105469	20.536458	79.799479	

	BMI	DiabetesPedigreeFunction	Age
min	0.000000	0.078000	21.000000
max	67.100000	2.420000	81.000000
mean	31.992578	0.471876	33.240885



b)

Split the data into training and test data with 80% training and 20% test data sizes.

Use 5-fold cross-validation on training data to decide the best number of neighbours k . To this end, you can use the built in functionality in `scikit-learn` such as `cross_val_score`. For $k = 1, 2, 3, \dots, 15$ compute the 5-fold cross validation error and plot the results (with values of k on the x -axis and accuracy on the y -axis). Include the plot in your report and justify your decision for picking a particular number of neighbors k .

```
In [85]: from sklearn.model_selection import train_test_split, cross_val_score  
from sklearn.neighbors import KNeighborsClassifier  
import numpy as np
```

```

np.random.seed(448)

X = data.drop(columns=["Target"])
y = data["Target"]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

accuracies = []
klist = range(1,16)

for k in klist:
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn, X_train, y_train, cv=5, scoring="accuracy")
    accuracies.append(scores.mean())

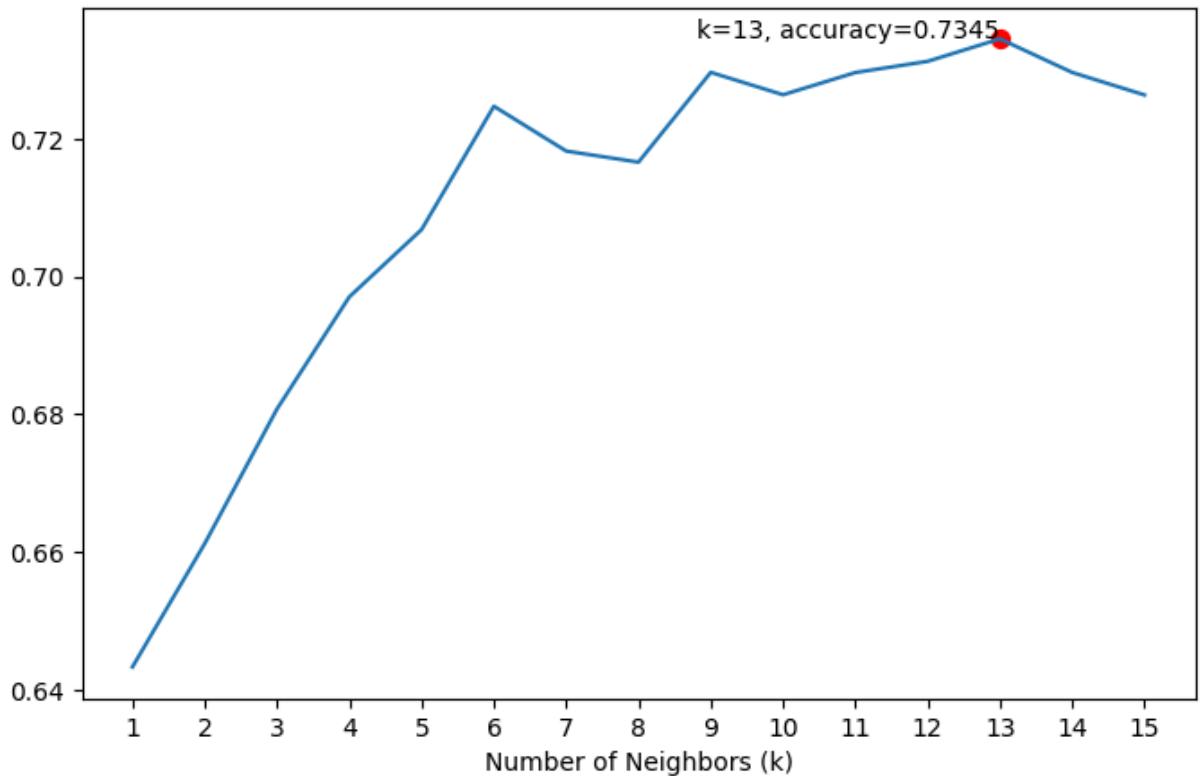
best_accuracy = max(accuracies)
best_k = accuracies.index(best_accuracy) + 1

print(f"Best k: {best_k} with accuracy: {best_accuracy:.4f}")

plt.figure(figsize=(8, 5))
plt.plot(klist, cv_scores)
plt.xlabel("Number of Neighbors (k)")
plt.xticks(klist)
plt.scatter(best_k, best_accuracy, color="red", s=50)
plt.text(best_k, best_accuracy, f"k={best_k}, accuracy={best_accuracy:.4f}",
plt.show()

```

Best k: 13 with accuracy: 0.7345



c)

Evaluate the k -NN algorithm on test data with the optimal number of neighbours you obtained in previous step and report the test error.

```
In [94]: from sklearn.metrics import accuracy_score

#train model with ideal k
knn_best = KNeighborsClassifier(n_neighbors=best_k)
knn_best.fit(X_train, y_train)

y_hat = knn_best.predict(X_test)

test_accuracy = accuracy_score(y_test, y_hat)

test_error = 1 - test_accuracy

print(f"Test Error: {test_error:.4f}")
```

Test Error: 0.2338

d)

Process the input data by subtracting the mean (a.k.a. centralization) and dividing by the standard deviation (a.k.a. standardization) over each dimension (feature), repeat the previous part and report the accuracy. Do centralization and standardization affect the accuracy? Why?

```
In [100...]: X_mean = X_train.mean(axis=0) #find the means and standard deviations by the
X_std = X_train.std(axis=0)

X_train_scaled = (X_train - X_mean) / X_std
X_test_scaled = (X_test - X_mean) / X_std

knn_standard = KNeighborsClassifier(n_neighbors=best_k)
knn_standard.fit(X_train_scaled, y_train)

y_hat_scaled = knn_standard.predict(X_test_scaled)

test_accuracy_scaled = accuracy_score(y_test, y_hat_scaled)

test_error_scaled = 1 - test_accuracy_scaled

print(f"Test Error: {test_error_scaled:.4f}\n\nTest Accuracy: {test_accuracy_scaled:.4f}")
```

Test Error: 0.2532

Test Accuracy: 0.7468

As we can see, the test error and the accuracy slightly increased for the scaled knn regression model. This could be because standardizing the features in the dataset could potentially decrease the euclidean distance between points. Since,

the scaling is different for each feature, some points may have gotten closer to others. This could lead to certain points being classified differently as they may now be closer to other points.

This notebook was converted with convert.ploomber.io

```

1 import random
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from Problem5 import bgd_l2, sgd_l2
5
6 def plot_loss_history(history_fw, params, type = "Regular"):
7
8     plt.figure(figsize=(8, 5))
9     plt.plot(range(len(history_fw)), history_fw, color="r")
10    plt.xlabel("Iteration")
11    plt.ylabel("Objective Function Value")
12    plt.title(rf"\{type} Gradient Descent Progress, \eta: {params[0]}, \delta: {params[1]}, \lambda: {params[2]}, iterations: {params[3]}")
13    plt.grid()
14    plt.show()
15
16 def load_data():
17     data = np.load("c:/Users/giris/OneDrive - The Pennsylvania State University/Cmpsc 448/Homework/hw2/data.npy") # Load dataset (100 x 2)
18
19     X = data[:, 0] # Extract feature column
20     y = data[:, 1] # Extract target column
21
22     X = np.column_stack((np.ones(X.shape[0]), X)) # Add intercept column
23
24     return X, y
25
26 if __name__ == '__main__':
27     np.random.seed(448)
28     X,y= load_data() #load the features and the target values
29     w = np.zeros(X.shape[1]) # initialize the weights
30
31
32     w1, hist1 = bgd_l2(X,y,w,0.05,0.1,0.001,50)
33     plot_loss_history(hist1, [0.05,0.1,0.001,50])
34
35     w2, hist2 = bgd_l2(X,y,w,0.1,0.01,0.001,50)
36     plot_loss_history(hist2,[0.1,0.01,0.001,50])
37
38     w3, hist3 = bgd_l2(X,y,w,0.1,0,0.001,100)
39     plot_loss_history(hist3,[0.1,0.01,0.001,50])
40
41     w4, hist4 = bgd_l2(X,y,w,0.1,0,0,100)
42     plot_loss_history(hist4,[0.1,0.01,0.001,50])
43
44     #Stochastic Gradient Descent
45     w5, hist5 = sgd_l2(X,y,w,1,0.1,0.5,800)
46     plot_loss_history(hist5, [1,0.1,0.5,800], "Stochastic")
47
48     w6, hist6 = sgd_l2(X,y,w,1,0.01,0.1,800)
49     plot_loss_history(hist6,[1,0.01,0.1,800],"Stochastic")
50
51     w7, hist7 = sgd_l2(X,y,w,1,0,0,40)
52     plot_loss_history(hist7,[1,0,0,40],"Stochastic")
53
54     w8, hist8 = sgd_l2(X,y,w,1,0,0,800)
55     plot_loss_history(hist8,[1,0,0,800],"Stochastic")
56

```

```

import math
import random
import numpy as np

def bgd_l2(data, y, w, eta, delta, lam, num_iter):

    n, d = data.shape # extract size of dataset and number of features
    new_w = np.copy(w) # create new weight vector
    history_fw = [] # Store function values

    for _ in range(num_iter):
        gradient = np.zeros(d) # Initialize gradient
        loss = 0 # Initialize loss function value

        for i in range(n):
            x_i = data[i] # Feature vector
            y_i = y[i] # Target value
            pred = np.dot(new_w, x_i) # Compute prediction w * x

            # Follow function definition
            if y_i >= pred + delta:
                gradient += -2 * (y_i - pred - delta) * x_i #derivative of function with respect to x_i
                loss += (y_i - pred - delta) ** 2
            elif y_i <= pred - delta:
                gradient += -2 * (y_i - pred + delta) * x_i #derivative of function with respect to x_i
                loss += (y_i - pred + delta) ** 2

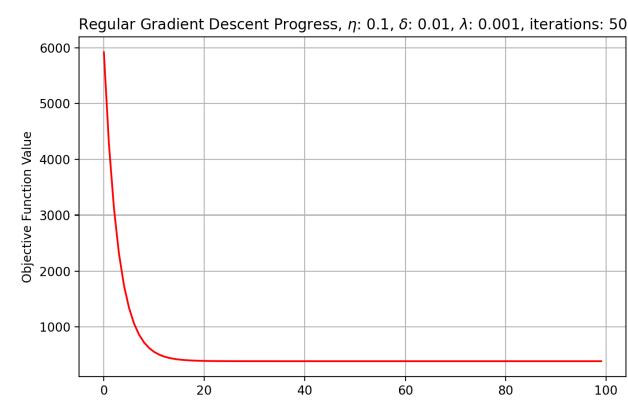
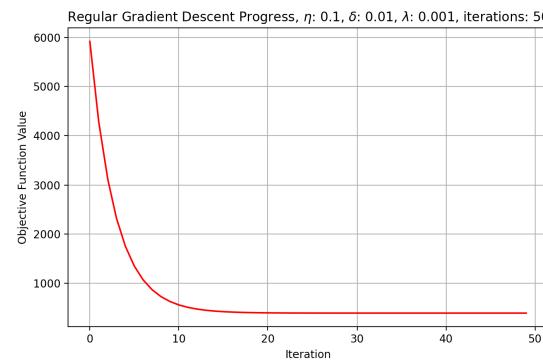
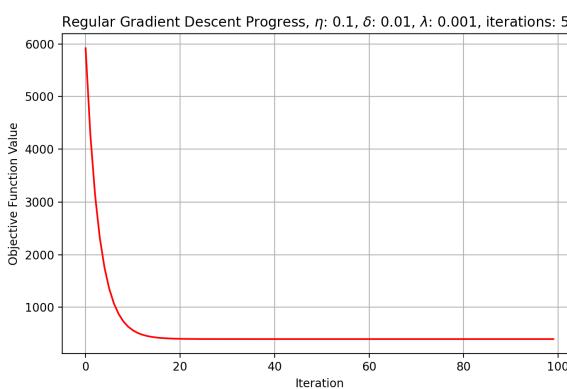
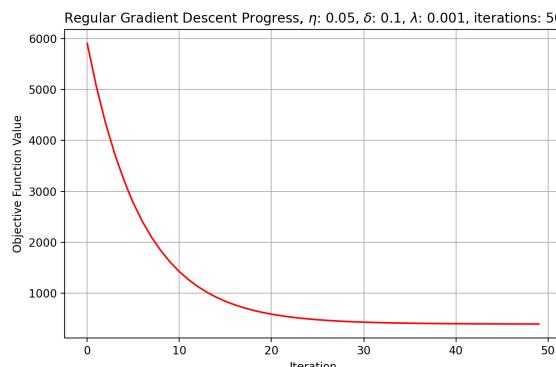
        # Add regularization term to gradient
        gradient = (gradient / n) + (2 * lam * new_w)

        # Update weights using gradient descent
        new_w = new_w - eta * gradient

        # Compute full objective function value and store
        reg_term = lam * np.sum(new_w ** 2) # regularization term
        fw = (loss / n) + reg_term # Compute full loss function
        history_fw.append(fw)

    return new_w, history_fw # Return updated weights and loss history

```



```

def sgd_l2(data, y, w, eta, delta, lam, num_iter, i=-1):
    n, d = data.shape
    new_w = np.copy(w)
    history_fw = []

    for t in range(1, num_iter + 1): # Start iterations from 1
        # Select data point
        if i == -1:
            ind = np.random.randint(0, n) # select random point
        else:
            ind = i # Use the specific data point

        x_i = data[ind] #feature vector
        y_i = y[ind] #target vector
        pred = np.dot(new_w, x_i) # Compute prediction w * x

        # Compute learning rate
        learn_rate = eta / np.sqrt(t)

        # Follow function definition
        gradient = np.zeros(d)
        loss = 0

        if y_i >= pred + delta:
            gradient = -2 * (y_i - pred - delta) * x_i #derivative of function with respect to x_i
            loss = (y_i - pred - delta) ** 2
        elif y_i <= pred - delta:
            gradient = -2 * (y_i - pred + delta) * x_i #derivative of function with respect to x_i
            loss = (y_i - pred + delta) ** 2

        # Add regularization term
        gradient += 2 * lam * new_w

        # Update weights
        new_w = new_w - learn_rate * gradient

        # Compute objective function value
        reg_term = lam * np.sum(new_w ** 2)
        fw = loss + reg_term
        history_fw.append(fw)

        # dont iterate again if i != -1
        if i != -1:
            break

    return new_w, history_fw

```

