

Question 2

Regression:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import random
6 import pandas as pd
7 import time
```

```
1 from google.colab import files
2 uploaded = files.upload()
```

Choose Files 2 files

- **A2Q2Data_test.csv**(text/csv) - 403470 bytes, last modified: 3/13/2024 - 100% done
 - **A2Q2Data_train.csv**(text/csv) - 8069086 bytes, last modified: 3/13/2024 - 100% done
- Saving A2Q2Data_test.csv to A2Q2Data_test.csv
Saving A2Q2Data_train.csv to A2Q2Data_train.csv

Managing Train and test datas:

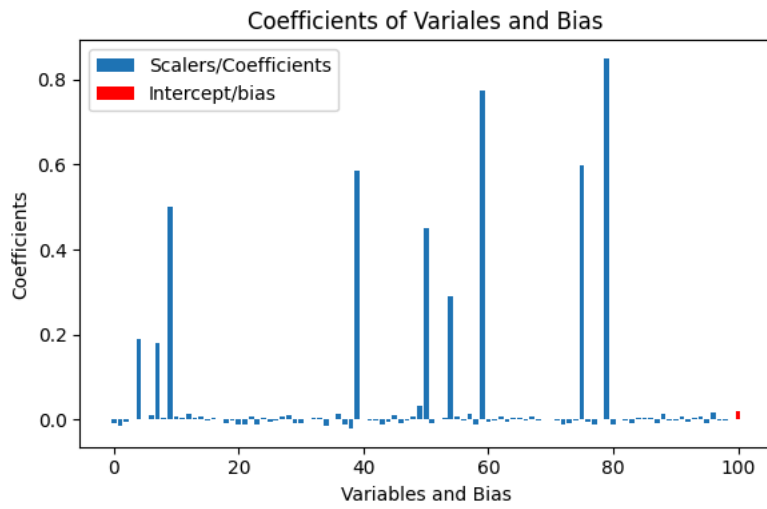
```
1 train = pd.read_csv('A2Q2Data_train.csv')
2 test = pd.read_csv('A2Q2Data_test.csv')
3
4 df2 = pd.DataFrame([dict(zip([col for col in train.columns], [col for col in train.columns]))])
5 train = pd.concat([df2,train], ignore_index = True)
6
7 X_train = train.iloc[:, 0:100].to_numpy()
8 y_train = train.iloc[:, 100].to_numpy()
9 X_train[0] = np.array([float(x) for x in X_train[0]])
10 y_train [0] = float(y_train[0])
11 X_train = X_train.astype('float32')
12 y_train = y_train.astype('float32')
13 ones = np.ones((X_train.shape[0],1))
14 X_train = np.hstack((X_train,ones))
15 y_train = y_train.reshape(10000,1)
```

```
1 df3 = pd.DataFrame([dict(zip([col for col in test.columns], [col for col in test.columns]))])
2
3 test = pd.concat([df3,test], ignore_index = True)
4
5 X_test = test.iloc[:, 0:100].to_numpy()
6 y_test = test.iloc[:, 100].to_numpy()
7 X_test[0] = np.array([float(x) for x in X_test[0]])
8 y_test[0] = float(y_test[0])
9 X_test = X_test.astype('float32')
10 y_test = y_test.astype('float32')
11 ones_test = np.ones((X_test.shape[0],1))
12 X_test = np.hstack((X_test,ones_test))
13 y_test = y_test.reshape(X_test.shape[0],1)
```

Using basic matrix Equation [*Analytical Solution using Derivation of W_{ml}*]

```
1 w_ml = ((np.linalg.inv(X_train.T @ X_train)) @ X_train.T ) @ y_train
```

```
1 fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True)
2 ax.bar(range(0,99),[w[0] for w in w_ml[0:99].tolist()])
3 ax.bar(100,w_ml[100].tolist(), color = 'red')
4 plt.title("Coefficients of Variables and Bias")
5 ax.set_xlabel('Variables and Bias')
6 ax.set_ylabel('Coefficients')
7 plt.legend({'Scalars/Coefficients':'0', 'Intercept/bias':'1'})
```



```

1 def cost(X,w,y):                                # Defining cost
2     return (np.linalg.norm(X @ w - y))**2

```

```

1 def R_sq(X,w,y):
2     cost_r = cost(X,w,y)
3     y_mean = np.mean(y)
4     mean_sq_error = 0
5
6     for i in range(y.shape[0]):
7         mean_sq_error += (y[i] - y_mean)**2
8
9     print(mean_sq_error)
10    print(cost_r)
11    return ((mean_sq_error - cost_r)/mean_sq_error) * 100

```

▽ Gradient Descent

```

1 def gradient(X,w,y):
2     n = np.shape(X)[0]
3     gradient_val = (2/n)*((X.T @ X) @ w - X.T @ y)
4     return(gradient_val)

```

```

1 def gradient_descent(X,y,no_of_epochs, alpha):
2     w = np.zeros((X.shape[1],1))
3     cost_gradient_descent = []
4     diff_from_w_ml = []
5
6     for epoch in range(no_of_epochs):
7         grad = gradient(X_train, w, y_train)
8         w = w - alpha*grad
9         cost_gradient_descent.append(cost(X_train,w,y_train))
10        diff_from_w_ml.append(np.linalg.norm(w - w_ml))
11
12    #Plot for w-w_ml vs epoch
13    fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True)
14    ax.plot(diff_from_w_ml)
15    ax.set_xlabel('iteration')
16    ax.set_ylabel('||MLE||2')
17
18    #Plot for error in each iteration
19    fig2, ax2 = plt.subplots(figsize=(6, 4), tight_layout=True)
20    ax2.plot(cost_gradient_descent)
21    ax2.set_xlabel('iteration')
22    ax2.set_ylabel('Δ Error')
23    print(diff_from_w_ml[-1])
24    return w

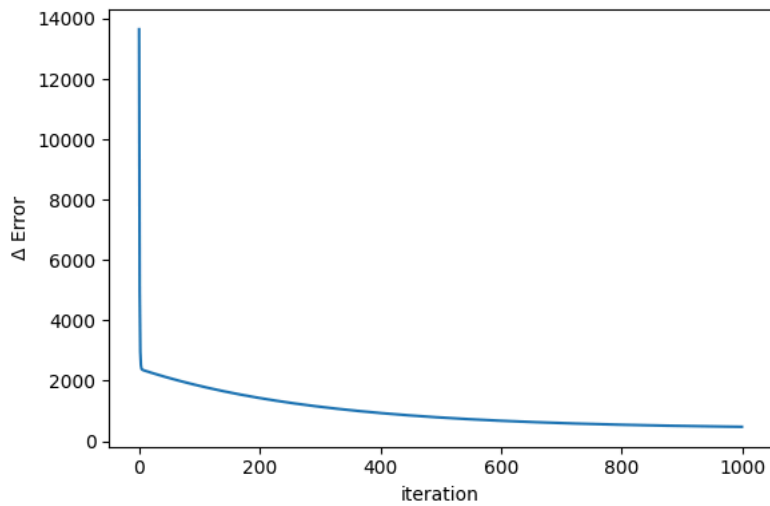
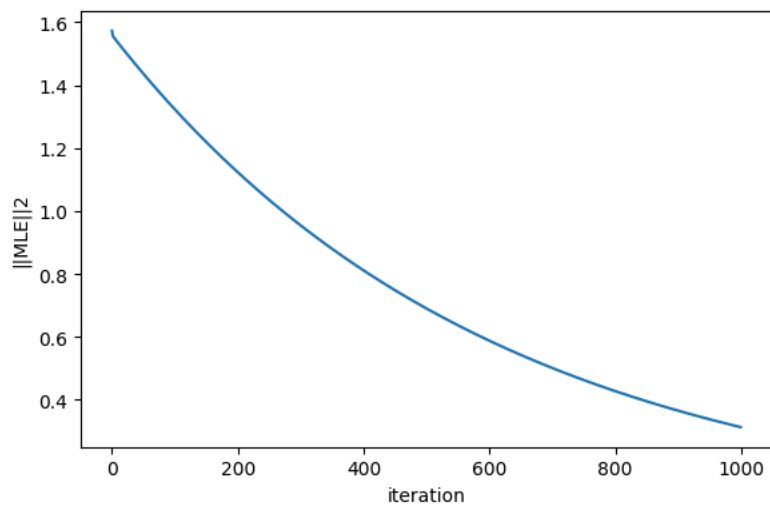
```

```

1 w_gd = gradient_descent(X_train, y_train, 1000, 1e-2)

```

0.3134321943474898



```
1 print("cost mlt train:",cost(X_train, w_ml, y_train))
2 print("cost ml test:",cost(X_test, w_ml, y_test))
3 print("R sq err ml test:",R_sq(X_test, w_ml, y_test))
4 print(" ")
5 print("cost gd:",cost(X_train, w_gd, y_train))
6 print("cost gd:",cost(X_test, w_gd, y_test))
7 print("R sq err gd test:",R_sq(X_test, w_gd, y_test))
```

```
cost mlt train: 396.85210962198335
cost ml test: 185.375750537584
[113.77624]
185.375750537584
R sq err ml test: [-62.930107]
```

```
cost gd: 472.51052247767205
cost gd: 155.55303127788213
[113.77624]
155.55303127788213
R sq err gd test: [-36.718376]
```

✓ Stochastic Gradient Descent

```

1 def SGD(X, y, batch_size, no_of_epochs, learning_rate):
2     data = np.hstack((X, y))
3     n = data.shape[0]
4     d = X.shape[1]
5
6     n_batches = int(n/batch_size)
7     w_iter = np.zeros((no_of_epochs,d)) #List containing w_i in each iteration
8     cost_list = []
9     diff_from_wml = []
10    w_sgd = np.zeros((X.shape[1],1))
11    for i in range(no_of_epochs):
12        idx = np.random.randint(data.shape[0], size = batch_size) #Choosing random list of Batch size
13        batch = data[idx,:]
14        X_batch = batch[:, 0:data.shape[1]-1]
15        y_batch = batch[:, data.shape[1]-1]
16        y_batch = y_batch.reshape((batch_size, 1))
17        grad = gradient(X_batch, w_sgd, y_batch)
18        w_sgd = w_sgd - (learning_rate)*grad
19        diff_from_wml.append((np.linalg.norm(w_sgd - w_ml)**2))
20        cost_list.append(cost(X_batch, w_sgd, y_batch))
21        w_iter[i] = w_sgd.T
22
23    #Plot for w-w_ml vs epoch
24    fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True)
25    ax.plot(diff_from_wml)
26    ax.set_xlabel('Iteration')
27    ax.set_ylabel('||MLE||2')
28
29    #Plot for error in each iteration
30    fig2, ax2 = plt.subplots(figsize=(6, 4), tight_layout=True)
31    ax2.plot(cost_list)
32    ax2.set_xlabel('iteration')
33    ax2.set_ylabel('Δ Error')
34    w_sgd_avg = np.mean(w_iter, axis = 0)
35    return w_sgd

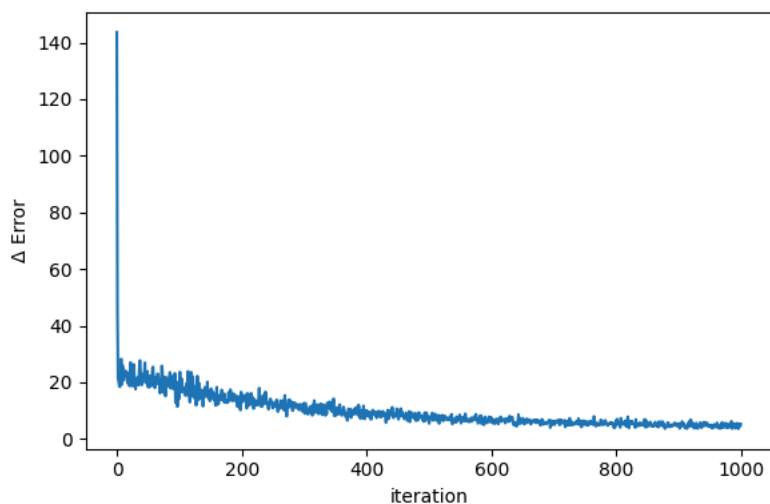
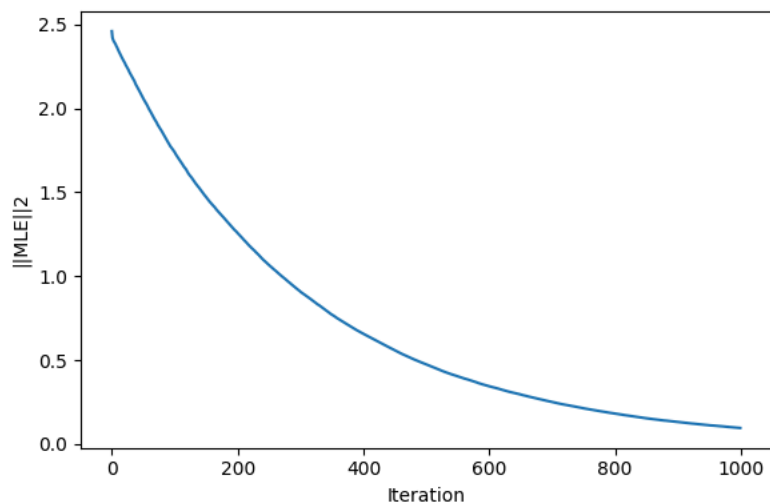
```

```

1 w_sgd_avg = SGD(X_train, y_train, 100, 1000, 1e-2)
2
3 print("cost sgd train :",cost(X_train, w_sgd_avg, y_train))
4 print("cost sgd test:" ,cost(X_test,w_sgd_avg,y_test))
5 print("R sq err sgd test:",R_sq(X_test, w_sgd_avg, y_test))

```

cost sgd train : 471.0785349470669
cost sgd test: 155.54034805441407
[113.77624]
155.54034805441407
R sq err sgd test: [-36.70723]



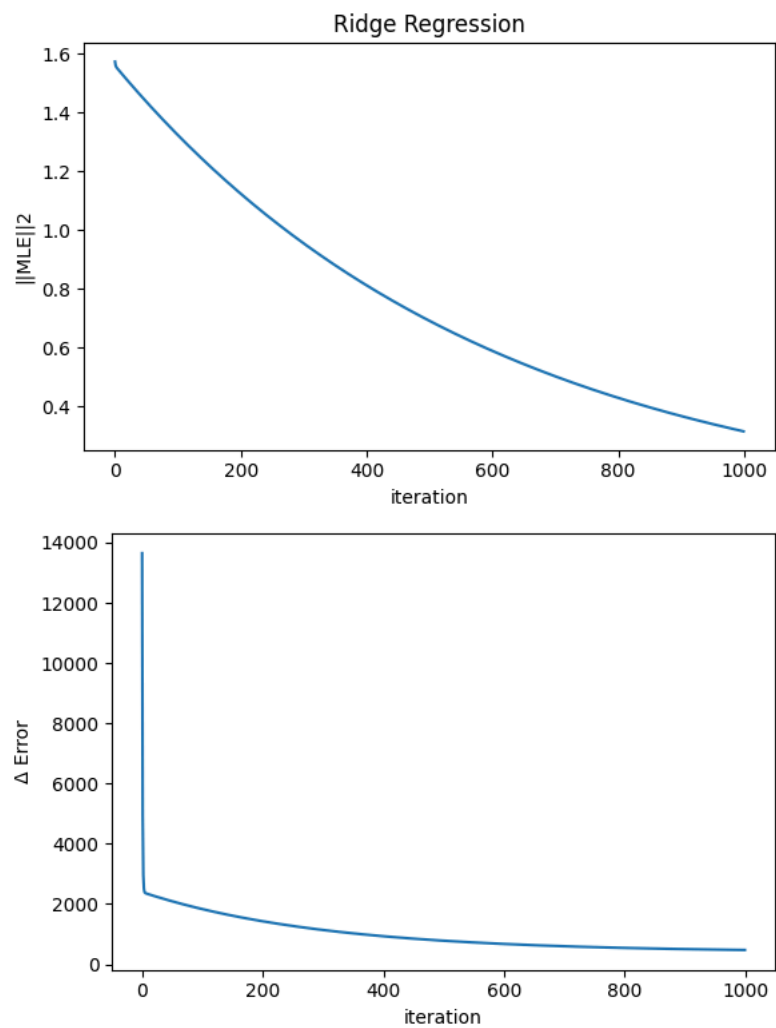
✓ Ridge Regression

```
1 #New gradient term due to regularisation
2 def l2_gradient(X, y, w, l2_penalty):
3     n = X.shape[0]
4     d = X.shape[1]
5     return (2/n)*((X.T @ X + np.identity(d)*l2_penalty) @ w) - (X.T @ y)
```

```
1 def ridge_regression(X,y,l2_penalty, learning_rate, number_of_epochs):
2     w = np.zeros((X.shape[1],1))
3     diff_from_wml = []
4     cost_ridge = []
5
6     for epoch in range(number_of_epochs):
7         grad = l2_gradient(X = X, y = y, w = w, l2_penalty = l2_penalty)
8         w = w - learning_rate*grad
9         diff_from_wml.append(np.linalg.norm(w - w_ml))
10        cost_ridge.append(cost(X = X, y = y, w = w))
11
12    #mod MLE vs episodes
13    fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True)
14    ax.plot(diff_from_wml)
15    ax.set_title('Ridge Regression')
16    ax.set_xlabel('iteration')
17    ax.set_ylabel('||MLE||2')
18
19
20    fig2, ax2 = plt.subplots(figsize=(6, 4), tight_layout=True)
21    ax2.plot(cost_ridge)
22
23    ax2.set_xlabel('iteration')
24    ax2.set_ylabel('Δ Error')
25    return w
```

```
1 wrr = ridge_regression(X = X_train, y = y_train, learning_rate = 1e-2, l2_penalty = 2.028, number_of_epochs = 1000)
2 print(cost(X = X_train, y = y_train, w = wrr))
3 print(cost(X = X_test, y = y_test, w = wrr))
```

473.4368228883739
155.38913236299499



K Fold Method - Cross validation

```
1 def kfoldMethod(X,y,number_of_epochs, k):  
2     n = X.shape[0]  
  
1 optimum_penalty = kfoldMethod(X_train, y_train, 1000, 10)  
2 print(optimum_penalty)  
3 wrr = ridge_regression(X = X_train, y = y_train, learning_rate = 1e-2, l2_penalty = optimum_penalty, number_of_epochs = 1000)
```

1.7254545454545454

