



VISVESVARAYA NATIONAL INSTITUTE
OF TECHNOLOGY (VNIT), NAGPUR

Machine Learning with Python (ECL443)

Lab Report

Submitted by :

Prajyot Jadhav (BT20ECE046)

Semester 7

Submitted to :

Dr. Saugata Sinha

(Course Instructor)

Department of Electronics and Communication Engineering,
VNIT Nagpur

Contents

1 Experiment-6 2

Experiment-6

Aim: To compress the ovarian cancer dataset using PCA(Principal Component Analysis) and Autoencoder and to evaluate the effectiveness by comparing the reconstructed data with the original data.

Abstract: This assignment focuses on the reconstruction and evaluation of data from Principal Component (PC) space and comparing it to the performance of a trained Autoencoder. In the first part of the assignment, we reconstruct the original data from PC space, and compute the Mean Squared Error (MSE) between the original data and the reconstructed data. This process aims to assess the efficiency of PC space in capturing and representing the original data's variance. In the second part, we compare the MSE obtained from the Autoencoder's reconstructed data, with that of PCA, aiming to discern the differences in their data representation capabilities.

Introduction: Principal Component Analysis (PCA) is a widely used dimensionality reduction technique that helps in capturing the most important patterns or features within a dataset. It achieves this by transforming the data into a new coordinate system represented by Principal Components (PCs), which are orthogonal and ordered by the amount of variance they explain.

In the first part, we aim to assess the quality of data reconstruction from PC space. First the principal components for which the ratio of the sum of the eigen values corresponding to the principle components to be used to the sum of all eigen values is 0.95. Next, by transforming the original data into PC space while considering all Principal Components, we intend to capture as much variance as possible. We then reconstruct the data and calculate the Mean Squared Error (MSE) between the original data and the reconstructed data. A lower MSE indicates a more faithful reconstruction of the original data.

In the second part, we extend our evaluation to include the Autoencoder. We compute the MSE between the data reconstructed by the Autoencoder and the original data and compare this MSE to the one obtained using PCA.

Method:

- The ovarian cancer dataset is loaded from a .mat file using the `scipy.io.loadmat` function. The necessary libraries, including `scipy`, `csv`, `numpy`, `random`, `matplotlib`, and `torch` are imported. The code extracts two variables, 'ovarian-Inputs' and 'ovarianTargets', from the loaded .mat file. It then saves these variables as CSV files, 'data_1.csv' and 'data_2.csv', respectively.

- The data is read from the CSV files, and each feature is standardized using the `Standardize_data` function. The standardization involves subtracting the mean and dividing by the standard deviation for each feature. The covariance matrix for the standardized data is calculated using the `np.cov` function.
- The eigenvalues and eigenvectors of the covariance matrix are computed using the `eig` function. The eigenvectors are adjusted to have positive values. The eigenvalue-eigenvector pairs are sorted in descending order based on the eigenvalues' magnitude. The top principal components such that they capture 95% of the total variance are selected.
- A projection matrix is created based on the selected principal components. The data is projected onto the lower-dimensional subspace formed by the selected principal components.
- The data is reconstructed from this compressed data and then MSE loss is calculated between the original data and this reconstructed data.
- For the second part, an autoencoder neural network model is defined using PyTorch. The autoencoder has an encoder and a decoder, and it aims to learn a compact representation of the input data. The autoencoder model is trained using mean squared error (MSE) loss and the Adam optimizer. The code iterates through the data in batches and updates the model's parameters.
- The loss values during training are stored, and a plot is created to visualize the training progress. The trained autoencoder model is used to encode the input data into a lower-dimensional representation.
- Similar to the first part, data is reconstructed from a trained autoencoder and then the MSE loss is calculated between original and reconstructed data.

Results:

- The autoencoder was trained for 25 epochs with a batch size of 36. The following loss plot was obtained during the training of the autoencoder.

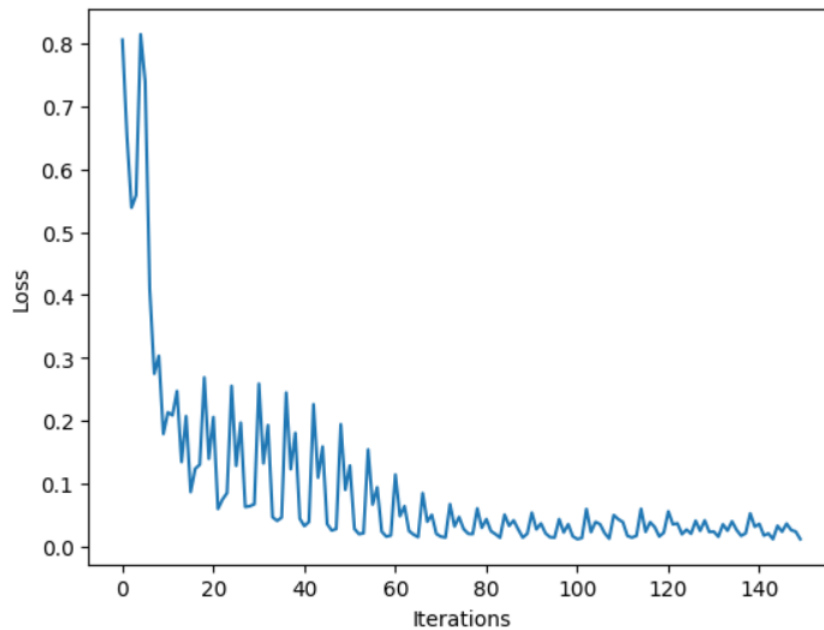


Figure 1: Loss vs Iterations

- After data compression using a PCA, using the principal components for which the ratio of the sum of the eigen values corresponding to the principle components to be used to the sum of all eigen values is 0.95. The data was reconstructed from PC space and the following results were obtained:

MSE between original and reconstructed data: 0.0489

- The above steps were repeated considering all PC components and the following result was obtained:

MSE between original and reconstructed data: 5.691665334111982e-27

- After data compression using an autoencoder, following MSE loss was obtained between reconstructed data and original data:

MSE between original and reconstructed data: 0.04032338038086891

Discussion:

- The results obtained in the first part of the assignment, where we reconstructed the original data from PC space, provided valuable insights into the efficiency of PCA in data representation. The MSE between the original data and the reconstructed data when 8 principal components were considered is significantly more than the MSE loss when all principal components were considered. When all principal components were considered the MSE loss was nearly equal to zero.
- In the second part of the assignment, the Autoencoder's performance was evaluated and compared to PCA. The MSE between the data reconstructed by the Autoencoder with latent space of 8, and the original data is less when compared to the MSE with the one obtained from PCA. A lower MSE for the Autoencoder might indicate its superior ability to model intricate relationships in the data. Autoencoders are capable of capturing intricate patterns, but their performance can be affected by the quality of training data and network architecture.

Conclusion: The ovarian cancer dataset was successfully compressed using PCA(Principal Component Analysis) and Autoencoder and its effectiveness was evaluated by comparing the reconstructed and original data.

Appendix:

```
1 import scipy.io
2 import csv
3 import numpy as np
4 import random
5 import matplotlib.pyplot as plt
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9
10 # Load .mat file
11 mat = scipy.io.loadmat('/content/ovarian_dataset.mat')
12
13 # Specify the variable name to convert to CSV
14 variable_name1 = 'ovarianInputs'
15 variable_name2 = 'ovarianTargets'
16 # Get the data from the loaded .mat file
17 #print(mat)
18 data1 = mat[variable_name1]
19 data2 = mat[variable_name2]
20
21 # Specify the CSV file name
```

```
22 csv_file_1 = '/content/data_1.csv'
23 csv_file_2 = '/content/data_2.csv'
24
25 # Write the data to CSV
26 with open(csv_file_1, 'w', newline='') as csvfile:
27     csvwriter = csv.writer(csvfile)
28     #for row in data1:
29     #    csvwriter.writerow(row)
30     for idx, row in enumerate(data1):
31         csvwriter.writerow(row)
32
33 with open(csv_file_2, 'w', newline='') as csvfile:
34     csvwriter = csv.writer(csvfile)
35     #for row in data1:
36     #    csvwriter.writerow(row)
37     for idx, row in enumerate(data2):
38         csvwriter.writerow(row)
39
40 #Reading Data from .csv file
41 with open('/content/data_1.csv', 'r') as f:
42     reader = csv.reader(f)
43     data_features = list(reader)
44
45 data_features = np.array(data_features, dtype=np.float32)
46
47 with open('/content/data_2.csv', 'r') as f:
48     reader = csv.reader(f)
49     data_labels = list(reader)
50
51 data_labels = np.array(data_labels, dtype=np.float32)
52 #data_labels = data_labels[0,:]
53 #data_labels = data_labels.reshape((1,data_labels.shape[0]))
54 #print(data_array.shape)
55 #print(data_array)
56
57
58
59 print(data_features)
60
61 def mean(x): # np.mean(X, axis = 0)
62     return sum(x)/len(x)
63
64 def std(x): # np.std(X, axis = 0)
65     return (sum((i - mean(x))**2 for i in x)/len(x))**0.5
66
67 def Standardize_data(X):
68     xx = X.transpose()
69     #print(xx)
70     summ = np.sum(xx, axis=0)
```

```

71     print(xx.shape[0])
72     mean = summ/xx.shape[0]
73     print(summ.shape)
74     stdd = np.std(xx,axis=0)
75     print(stdd.shape)
76     #X = X.transpose()
77     xx = (xx-mean)/stdd
78     #for i in range(X.shape[1]):
79     #    X[:,i] = (X[:,i]-summ[i])/stdd[i]
80
81     return xx
82 print(data_features.shape)
83 data_features = Standardize_data(data_features)
84 print(data_features.shape)
85 #data_features_2 = data_features.transpose()
86 #print(data_features)
87
88 def covariance(x):
89     #print(x.shape[0])
90     return (x.T @ x)/(x.shape[0]-1)
91
92 cov_mat = np.cov(data_features, rowvar=False)
93
94 # cov_mat = covariance(data_features) # np.cov(X_std.T)
95 print(cov_mat.shape)
96
97 from numpy.linalg import eig
98
99 # Eigendecomposition of covariance matrix
100 eig_vals, eig_vecs = eig(cov_mat)
101
102 # Adjusting the eigenvectors (loadings) that are largest in ...
103 # absolute value to be positive
104 max_abs_idx = np.argmax(np.abs(eig_vecs), axis=0)
105 signs = np.sign(eig_vecs[max_abs_idx, range(eig_vecs.shape[0])])
106 eig_vecs = eig_vecs*signs[np.newaxis,: ]
107 eig_vecs = eig_vecs.T
108
109 print('Eigenvalues \n', eig_vals)
110 print('Eigenvectors \n', eig_vecs)
111
112 # We first make a list of (eigenvalue, eigenvector) tuples
113 eig_pairs = [(np.abs(eig_vals[i]), eig_vecs[i,:]) for i in ...
114               range(len(eig_vals))]
115
116 # Then, we sort the tuples from the highest to the lowest based ...
117 # on eigenvalues magnitude
118 eig_pairs.sort(key=lambda x: x[0], reverse=True)

```



```
117 # For further usage
118 eig_vals_sorted = np.array([x[0] for x in eig_pairs])
119 eig_vecs_sorted = np.array([x[1] for x in eig_pairs])
120
121 #print(eig_pairs)
122
123 eig_vals_total = sum(eig_vals)
124 i=0
125 cum_sum = 0
126 threshold = 95
127 while(cum_sum<threshold):
128     cum_sum = cum_sum + eig_vals_sorted[i]/eig_vals_total*100
129     print(cum_sum)
130     i+=1
131 print(i)
132 explained_variance = [(i / eig_vals_total)*100 for i in ...
    eig_vals_sorted]
133 explained_variance = np.round(explained_variance, 2)
134 cum_explained_variance = np.cumsum(explained_variance)
135
136 print('Explained variance: {}'.format(explained_variance))
137 print('Cumulative explained variance: ...
    {}'.format(cum_explained_variance))
138
139 #plt.plot(np.arange(1,n_features+1), cum_explained_variance, '-o')
140 #plt.xticks(np.arange(1,n_features+1))
141 #plt.xlabel('Number of components')
142 #plt.ylabel('Cumulative explained variance');
143 #plt.show()
144
145
146 # Select top k eigenvectors
147 k = i
148 W1 = eig_vecs_sorted[:k, :] # Projection matrix
149
150 print(W1.shape)
151
152 X_proj1 = data_features.dot(W1.T)
153
154 print(X_proj1.shape)
155
156 recons_data_features1 = X_proj1.dot(W1)
157
158 MSE_Loss = ((data_features - ...
    recons_data_features1)**2).mean(axis=None)
159 print(f"MSE between original and reconstructed data: ...
    {MSE_Loss:.4f}")
160
161 # for all the principal components
```

```
162
163 W2 = eig_vecs_sorted[:, :] # Projection matrix
164
165 print(W2.shape)
166
167 X_proj2 = data_features.dot(W2.T)
168
169 print(X_proj2.shape)
170
171 recons_data_features2 = X_proj2.dot(W2)
172
173 MSE_Loss2 = ((data_features - ...
               recons_data_features2)**2).mean(axis=None)
174 print(f"MSE between original and reconstructed data: {MSE_Loss2}")
```

```
1 import scipy.io
2 import csv
3 import numpy as np
4 import random
5 import matplotlib.pyplot as plt
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9
10 # Load .mat file
11 mat = scipy.io.loadmat('/content/ovarian_dataset.mat')
12
13 # Specify the variable name to convert to CSV
14 variable_name1 = 'ovarianInputs'
15 variable_name2 = 'ovarianTargets'
16 # Get the data from the loaded .mat file
17 #print(mat)
18 data1 = mat[variable_name1]
19 data2 = mat[variable_name2]
20
21 # Specify the CSV file name
22 csv_file_1 = '/content/data_1.csv'
23 csv_file_2 = '/content/data_2.csv'
24
25 # Write the data to CSV
26 with open(csv_file_1, 'w', newline='') as csvfile:
27     csvwriter = csv.writer(csvfile)
28     #for row in data1:
29     #    csvwriter.writerow(row)
30     for idx, row in enumerate(data1):
31         csvwriter.writerow(row)
32
```

```
33 with open(csv_file_2, 'w', newline='') as csvfile:
34     csvwriter = csv.writer(csvfile)
35     #for row in data1:
36     #    csvwriter.writerow(row)
37     for idx, row in enumerate(data2):
38         csvwriter.writerow(row)
39
40 #Reading Data from .csv file
41 with open('/content/data_1.csv', 'r') as f:
42     reader = csv.reader(f)
43     data_features = list(reader)
44
45 data_features = np.array(data_features, dtype=np.float32)
46
47 with open('/content/data_1.csv', 'r') as f:
48     reader = csv.reader(f)
49     data_labels = list(reader)
50
51 data_labels = np.array(data_labels, dtype=np.float32)
52 #data_labels = data_labels[0,:]
53 #data_labels = data_labels.reshape((1, data_labels.shape[0]))
54 #print(data_array.shape)
55 #print(data_array)
56
57
58
59 print(data_features.shape)
60
61 class AE(torch.nn.Module):
62     def __init__(self):
63         super().__init__()
64
65         # Building an linear encoder with Linear
66         # layer followed by Relu activation function
67         # 784 ==> 9
68         self.encoder = torch.nn.Sequential(
69             torch.nn.Linear(100, 64),
70             torch.nn.ReLU(),
71             torch.nn.Linear(64, 32),
72             torch.nn.ReLU(),
73             torch.nn.Linear(32, 8),
74         )
75
76         # Building an linear decoder with Linear
77         # layer followed by Relu activation function
78         # The Sigmoid activation function
79         # outputs the value between 0 and 1
80         # 9 ==> 784
81         self.decoder = torch.nn.Sequential(
```

```
82         torch.nn.Linear(8, 32),
83         torch.nn.ReLU(),
84         torch.nn.Linear(32, 64),
85         torch.nn.ReLU(),
86         torch.nn.Linear(64, 100),
87         #torch.nn.Sigmoid()
88     )
89
90     def forward(self, x):
91         encoded = self.encoder(x)
92         decoded = self.decoder(encoded)
93         return encoded, decoded
94
95
96 # In[16]:
97
98
99 # Model Initialization
100 model = AE()
101
102 # Validation using MSE Loss function
103 loss_function = torch.nn.MSELoss()
104
105 # Using an Adam Optimizer with lr = 0.1
106 optimizer = torch.optim.Adam(model.parameters(),
107                               lr = 0.005,
108                               weight_decay = 1e-8)
109
110
111 # In[17]:
112
113
114 epochs = 25
115 batch_size = 36
116 outputs = []
117 losses = []
118 for epoch in range(epochs):
119     for i in range(data.features.shape[1]//batch_size):
120
121         batch = data.features[:,i*batch_size:(i+1)*batch_size]
122         # Reshaping the image to (-1, 784)
123         #image = image.reshape(-1, 28*28)
124
125         # Output of Autoencoder
126         batch = batch.transpose()
127         batch = torch.Tensor(batch)
128         encoded,reconstructed = model(batch)
129
130         # Calculating the loss function
```

```
131     #print(reconstructed)
132     loss = loss_function(reconstructed, batch)
133
134     # The gradients are set to zero,
135     # the gradient is computed and stored.
136     # .step() performs parameter update
137     optimizer.zero_grad()
138     loss.backward()
139     optimizer.step()
140
141     # Storing the losses in a list for plotting
142     losses.append(loss.detach())
143     outputs.append((epochs, batch, reconstructed))
144
145 # Defining the Plot Style
146 plt.style.use('fivethirtyeight')
147 plt.xlabel('Iterations')
148 plt.ylabel('Loss')
149
150 # Plotting the last 100 values
151 plt.plot(losses)
152
153
154 # In[18]:
155
156
157 features = torch.Tensor(data.features.transpose())
158 encoded, _ = model(features)
159 latent = encoded.detach().numpy()
160 print(latent.shape)
161
162 # Reconstruct the data using the trained Autoencoder
163 with torch.no_grad():
164     _, reconstructed_data = ...
165         model(torch.Tensor(data.features.transpose()))
166
167 # Convert the reconstructed_data to a NumPy array
168 reconstructed_data = reconstructed_data.numpy()
169
170 # Calculate the MSE between the original and reconstructed data
171 mse = ((data.features.transpose() - reconstructed_data)**2).mean()
172
173 print(f"MSE between original and reconstructed data: {mse}")
```