

and `class_B_data`.

- `class_A_data` is a numpy arrays with the shape (N, \dots) , where N is the number of samples, and there may be many excess dimensions denoted by \dots . You will have to reshape this input matrix to obtain a shape of (N, d) , where d is the vectorized data's dimension.
- `class_B_data` has the same data structure as `class_A_data`.

To compute $E[A|B]$:

- First, do whatever reshaping you have to do.
- Subtract Class A's mean from its data
- Use the `principal_components` function you wrote before to extract the 20 principal components of `class_B_data`.
- Project Class A's data onto the mentioned principal components and get back to the original space.
- Compute Class A's residuals (i.e., the difference between the original and the projection).
- Find the squared residual sizes for **each sample**, and then return their mean as the `E_A_cond_B` scalar. In other words, square class A's residuals, sum them over each sample (which should reduce the squared residual matrix to only N elements), and then report the mean of them as `E_A_cond_B`.

```
In [30]: def E_A_given_B(class_A_data, class_B_data):  
    # No code 0  
    # To compute E[A|B]:  
  
    # First, do whatever reshaping you have to do.  
    # reshape n dim to 2 dim by multiply dim 2 * dim 3 * ... * dim n  
  
    n = np.prod(class_A_data.shape[1:])  
    d = class_A_data.shape[0]  
    class_A_data = class_A_data.reshape(N, d)  
  
    d = np.prod(class_B_data.shape[1:])  
    N = class_B_data.shape[0]  
    class_B_data = class_B_data.reshape(N, d)  
  
    # Subtract Class A's mean from its data  
    #calculate the mean of each column with axis = 0 to get a column vector and reshape  
    mean_row = np.mean(class_A_data, axis = 0).reshape(1, d)  
    class_A_data = class_A_data - mean_row  
  
    # Use the principal_components function you wrote before to extract the 20 principal components  
    V = principal_components(class_B_data, num_components=20)  
  
    # Project Class A's data onto the mentioned principal components and get back to  
    projAonV_B = class_A_data @ V_B @ V_B.T  
  
    # Compute Class A's residuals (i.e., the difference between the original and the  
    e = class_A_data - projAonV_B  
  
    # Find the squared residual sizes for each sample, and then return their mean as  
    ## That is:  
    ## 1. square class A's residuals,  
    e_squared = e**2  
  
    ## 2. sum them over each sample (which should reduce the squared residual matrix  
    ## No note: since e.shape[0] = N, ie each row is a sample, must sum of rows, ie  
    sum_e_squared = np.sum(e_squared, axis = 1)  
  
    ## 3. E_A_cond_B = their mean  
    E_A_cond_B = np.mean(sum_e_squared, axis = 0)  
  
    # No code 0  
    return E_A_cond_B
```

```
In [31]: # Performing sanity checks on your implementation  
some_data = ((np.arange(35).reshape(5,7) ** 13) % 20) / 7.  
some_data = np.repeat(some_data, 8, axis=1)  
some_E = E_A_given_B(some_data, (some_data**1.02))  
assert some_E.round(3)==0.001  
  
# Checking against the pre-computed test database  
test_results = test_case_checker(E_A_given_B, task_id=6)  
assert test_results['passed'], test_results['message']
```

```
In [32]: # This cell is left empty as a separator. You can leave this cell as it is, and you should
```

```
In [33]: # This cell is left empty as a separator. You can leave this cell as it is, and you should
```

```
In [34]: if perform_computation:  
    num_classes = class_means.shape[0]  
    SimilarityMatrix = np.zeros((num_classes, num_classes))  
    for row in range(num_classes):  
        print('Row %i', end='')  
        row_at_time = time.time()  
        for col in range(row+1):  
            class_B_data = images_raw[labels == row, :, :, :]  
            class_B_data = images_raw[labels == col, :, :, :]  
            E_A_cond_B = E_A_given_B(class_A_data, class_B_data)  
            E_B_cond_A = E_A_given_B(class_B_data, class_A_data)  
            SimilarityMatrix[col, row] = (E_A_cond_B + E_B_cond_A)/2.  
            SimilarityMatrix[row, col] = (E_A_cond_B + E_B_cond_A)/2.  
        print(' (This row took %.3f seconds to finish)'%(time.time() - row_at_time))  
  
Row 0 (This row took 2.356 seconds to finish)  
Row 1 (This row took 4.652 seconds to finish)  
Row 2 (This row took 6.995 seconds to finish)  
Row 3 (This row took 9.379 seconds to finish)  
Row 4 (This row took 11.717 seconds to finish)  
Row 5 (This row took 14.006 seconds to finish)  
Row 6 (This row took 16.134 seconds to finish)  
Row 7 (This row took 18.470 seconds to finish)  
Row 8 (This row took 20.873 seconds to finish)  
Row 9 (This row took 23.274 seconds to finish)
```

```
In [35]: # This cell is left empty as a separator. You can leave this cell as it is, and you should
```

If you apply any general `SimilarityMatrix` variable to the previously defined `PCoA` function, you may get `NaN` entries due to the fact that they may not generally be a metric distance matrix (i.e, having non-zero diagonal elements and the triangle inequality not always holding).

This issue can be best seen when having a similarity measure that is extremely uneven (i.e., when the small entries are extremely small and the large entries are extremely large). This will make it difficult for the triangle inequality to hold. It is a good idea to amend the PCoA in a way that can deal with such non-metric similarity measures.

```
In [36]: VT = None  
if perform_computation:  
    VT = PCoA(SimilarityMatrix**40, r=10)  
VT
```

```
AssertionError                                Traceback (most recent call last)  
/tmp/ipykernel_83/3393768783.py in <module>  
      1 VT = None  
----> 2 if perform_computation:  
      3     VT = PCoA(SimilarityMatrix**40, r=10)  
      4 VT  
  
/tmp/ipykernel_83/3291470657.py in PCoA(SquaredDistances, r)  
     26  
     27     assert VT.shape[0] == num_points  
--> 28     assert VT.shape[1] == r  
     29     return VT  
AssertionError:
```

Task 7

Write the function `Lingoes_PreProcessing` that does some pre-processing to the `SimilarityMatrix` to make it have the Euclidean property and the triangles to close.

Here is a very brief and to the point description from the r documentation page (<https://www.rdocumentation.org/packages/spe/versions/5.2/topics/pcoa>).

"In the Lingoes (1971) procedure, a constant c1, equal to twice absolute value of the largest negative eigenvalue of the original principal coordinate analysis, is added to each original squared distance in the distance matrix, except the diagonal values. A new principal coordinate analysis, performed on the modified distances, has at most (n-2) positive eigenvalues, at least 2 null eigenvalues, and no negative eigenvalue."

If you're interested, you can read more about correction for negative eigenvalues in http://biol09.biol.umontreal.ca/PLCources/Ordination_sections_1.3+1.4_PCoA_Eng.pdf.

The function `Lingoes_PreProcessing` takes the numpy array `SimilarityMatrix` as input, and returns `ProcessedSimilarityMatrix` based on the following condition:

- If all eigenvalues computed during PCoA are non-negative, then `ProcessedSimilarityMatrix` should be the same as the `SimilarityMatrix`.
- Otherwise, follow the instructions to perform the Lingoes correction on the `SimilarityMatrix` and return `ProcessedSimilarityMatrix`.

In other words, this is what you're supposed to do:

- Perform the PCoA analysis on `SimilarityMatrix`, right up to the point where you find the eigenvalues. Do not go any further. More precisely, you should only find the eigenvalues of the matrix W corresponding to `SimilarityMatrix` in the PCoA analysis.
- Find the minimum eigenvalue and call it λ_{\min} .
- If $\lambda_{\min} \geq 0$, then stop and return `SimilarityMatrix` as it was without any change.
- If $\lambda_{\min} < 0$, then add $2|\lambda_{\min}|$ to all the non-diagonal elements of `SimilarityMatrix` and return the resulting matrix.

Important Note: Do not call the PCoA function on `SimilarityMatrix`. You should not call the whole `PCoA` function on `SimilarityMatrix`, as you do not care about the output reconstructions of `PCoA`. Instead, you need the eigenvalues for further processing, which are not returned by the `PCoA` function.

Note: You do not need a `for` loop for adding a scalar to the non-diagonal elements of a matrix; you can add the scalar to all the elements of the matrix, and then subtract it from the same scalar multiple of the identity matrix (i.e., using a function like `np.eye` for instance).

```
Procedure: 6.2 Principal Coordinate Analysis  
  
Assume we have a matrix  $D^{(n)}$  consisting of the squared differences between each pair of  $N$  points. (We do not need to know the points.) We wish to compute a set of points in  $s$  dimensions, such that the distances between these points are as similar as possible to the distances in  $D^{(n)}$ .  
  
1. Form  $A = [I - \frac{1}{n}11^T]$   
2. Form  $W = -\frac{1}{2}A^{(2)}A^{(2)T}$   
3. Form  $U, A$ , such that  $WU = UA$ ,  $U$ : eigenvectors of  $W$  and  $A$ : eigenvalues of  $W$ . Ensure that the entries of  $A$  are sorted in the decreasing order.  
  
The entries of  $U$  are scaled such that the trace of  $U^T U$  is equal to the trace of  $W$ .
```

SquaredDistances: The output of the `mean_image_squared_distances` PCoA takes `SquaredDistances`

```
In [37]: def Lingoes_PreProcessing(SimilarityMatrix):  
    assert SimilarityMatrix.shape[0] == SimilarityMatrix.shape[1]  
    num_points = SimilarityMatrix.shape[0]  
  
    # No code 0  
    # Perform the PCoA analysis on SimilarityMatrix right up to the point where you find  
    # More precisely, you should only find the eigenvalues of the matrix W corresponding  
    #  
    N = num_points  
    I = np.identity(N)  
    A = I - np.ones((N,N)) * (1/N)  
    D2 = SimilarityMatrix**2  
    W = -0.5 * A @ D2 @ A.T * N * N  
    # by default the evalues & evectors are returned in ascending order  
    lam, V = np.linalg.eigh(W)  
    # Find the minimum eigenvalue and call it Amin  
    lam_min = lam[0]  
    # If Amin >= 0, then stop and return SimilarityMatrix as it was without any change  
    if lam_min >= 0:  
        ProcessedSimilarityMatrix = SimilarityMatrix  
    # If Amin < 0, then add 2|Amin| to all the non-diagonal elements of Similarity  
    else:  
        ProcessedSimilarityMatrix = SimilarityMatrix + 2*abs(lam_min) * (np.eye(N)*2)  
    # No code 1  
    return ProcessedSimilarityMatrix
```

```
In [38]: # Performing sanity checks on your implementation  
some_data = ((np.arange(35).reshape(5,7) ** 13) % 20) / 7.  
some_dist = mean_image_squared_distances(some_data)**5.  
assert np.array_equal(some_lingoes.round(1), np.array([[ 898987.0, 898987.1, 2896177.1, 1229280.0, 944977.7, 1208489.7, 13442744.7],  
               [ 898987.1, 898987.0, 1229280.0, 944977.7, 1208489.7, 13442744.7, 2896177.1],  
               [2896177.1, 1229280.0, 944977.7, 1208489.7, 13442744.7, 2896177.1, 898987.1],  
               [ 944977.7, 1208489.7, 13442744.7, 2896177.1, 898987.1, 1229280.0, 898987.1],  
               [1208489.7, 13442744.7, 2896177.1, 898987.1, 1229280.0, 898987.1, 898987.1],  
               [13442744.7, 2896177.1, 898987.1, 1229280.0, 898987.1, 898987.1, 898987.1]]))  
  
# Checking against the pre-computed test database  
test_results = test_case_checker(Lingoes_PreProcessing, task_id=7)  
assert test_results['passed'], test_results['message']
```

```
In [39]: # This cell is left empty as a separator. You can leave this cell as it is, and you should
```

```
In [40]: def PCoA_lingoes(SimilarityMatrix, r=2):  
    ProcessedSimilarityMatrix = Lingoes_PreProcessing(SimilarityMatrix)  
    return PCoA(ProcessedSimilarityMatrix, r=r)
```

```
In [41]: VT = None  
if perform_computation:  
    VT = PCoA_lingoes(SimilarityMatrix, r=2)  
VT
```

```
Out[41]: array([[ -119.98786571, 355.24236881],  
               [ -407.7649715 , -194.47481966],  
               [ 200.74413716, 69.91510988],  
               [ 228.74930337, -150.11666377],  
               [ 34.48666889, 156.6042217 ],  
               [ 343.06053853, -112.89238769],  
               [ 112.75711242, -44.2970879 ],  
               [ 147.41580395, 23.49865708],  
               [-290.00131424, 180.45366172],  
               [-249.45943287, -284.34308071]])
```

```
In [42]: if perform_computation:  
    class_names_list = sorted(list(class_to_idx.keys()))  
    fig, ax = plt.subplots(figsize=(9,6.), dpi=10)  
    x_components = VT[:,0]  
    y_components = VT[:,1]  
    sns.regplot(x=x_components, y=y_components, fit_reg=False, marker="x", color="Blue",  
               class_names_list(class_idx).capitalize(),  
               ax.text(x_components[class_idx]-num_letters*8, y_components[class_idx]+10,  
                       class_names_list(class_idx).capitalize(),  
                       horizontalalignment="left", size="medium", color="black", weight="semibold",  
                       dx=5, dy=5)  
    ax.set_xlabel("Reconstructed Dimension 1")  
    ax.set_ylabel("Reconstructed Dimension 2")  
    _ = ax.set_title("Generalized PCoA on CIFAR-10 Images")
```



```
In [ ]:
```

```
In [ ]:
```