

# PCA

May 24, 2021

*If you plan to run the assignment locally:* You can download the assignments and run them locally, but please be aware that as much as we would like our code to be universal, computer platform differences may lead to incorrectly reported errors even on correct solutions. Therefore, we encourage you to validate your solution in Coursera whenever this may be happening. If you decide to run the assignment locally, please: 1. Try to download the necessary data files from your home directory one at a time, 2. Don't update anything other than this Jupyter notebook back to Coursera's servers, and 3. Make sure this notebook maintains its original name after you upload it back to Coursera.

**Note:** You need to submit the assignment to be graded, and passing the validation button's test does not grade the assignment. The validation button's functionality is exactly the same as running all cells.

```
[3]: %matplotlib inline
      %load_ext autoreload
      %autoreload 2

      import matplotlib.pyplot as plt
      import numpy as np
      import seaborn as sns
      import pandas as pd
      import time
      import os
      from sklearn.decomposition import TruncatedSVD

      from aml_utils import test_case_checker, perform_computation
```

## 0.1 Attention:

This assignment is computationally heavy, and inefficient implementations may not pass the autograding even if they technically produce the correct results. To avoid this, make sure you read and understand all the instructions before starting to implement the tasks. Failure to follow the instructions closely will most likely cause timeouts.

It is **your responsibility** to make sure your implementation is not only **correct**, but also as **efficient** as possible. If you follow all the instructions provided, you should be able to have all the cells evaluated in under 10 minutes.

# 1 \*Assignment Summary

CIFAR-10 is a dataset of 32x32 images in 10 categories, collected by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. It is often used to evaluate machine learning algorithms. You can download this dataset from <https://www.cs.toronto.edu/~kriz/cifar.html>.

- For each category, compute the mean image and the first 20 principal components. Plot the error resulting from representing the images of each category using the first 20 principal components against the category.
- Compute the distances between mean images for each pair of classes. Use principal coordinate analysis to make a 2D map of the means of each categories. For this exercise, compute distances by thinking of the images as vectors.
- Here is another measure of the similarity of two classes. For class A and class B, define  $E(A | B)$  to be the average error obtained by representing all the images of class A using the mean of class A and the first 20 principal components of class B. Now define the similarity between classes to be  $(1/2)(E(A | B) + E(B | A))$ . If A and B are very similar, then this error should be small, because A's principal components should be good at representing B. But if they are very different, then A's principal components should represent B poorly. In turn, the similarity measure should be big. Use principal coordinate analysis to make a 2D map of the classes. Compare this map to the map in the previous exercise? are they different? why?

## References:

- Textbook section 5.1 [https://link.springer.com/chapter/10.1007/978-3-030-18114-7\\_5](https://link.springer.com/chapter/10.1007/978-3-030-18114-7_5)

## 2 0. Data

### 2.1 0.1 Description

CIFAR-10 is a dataset of 32x32 images in 10 categories, collected by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. It is often used to evaluate machine learning algorithms. You can download this dataset from <https://www.cs.toronto.edu/~kriz/cifar.html>.

### 2.2 0.2 Information Summary

- **Input/Output:** This data has a set of 32 pixel rows, 32 pixel columns, and 3 color channels. Therefore, each single image, is vectorized, will consist of  $32 \times 32 \times 3$  elements (i.e., each image has 3072 dimensions). There are a total of 60000 samples labelled from 10 class. The data set is balanced with each class having exactly 6000 samples.
- **Missing Data:** There is no missing data.
- **Final Goal:** We want to understand the data using multi-dimensional scaling methods.

### 2.3 0.3 Loading The Data

If you are curious how the original data was obtained, we used the torchvision API to download and pre-process it. The ready-to-use data is stored in numpy format for easier access.

```
[4]: if os.path.exists('../PCA-lib/cifar10.npz'):
      np_file = np.load('../PCA-lib/cifar10.npz')
```

```

train_images_raw = np_file['train_images_raw']
train_labels = np_file['train_labels']
eval_images_raw = np_file['eval_images_raw']
eval_labels = np_file['eval_labels']

else:
    import torchvision
    download_ = not os.path.exists('../PCA-lib/cifar10/')
    data_train = torchvision.datasets.CIFAR10('../PCA-lib/cifar10', train=True,
    ↪transform=None, target_transform=None, download=download_)
    data_eval = torchvision.datasets.CIFAR10('../PCA-lib/cifar10', train=False,
    ↪transform=None, target_transform=None, download=download_)
    train_images_raw = data_train.data
    train_labels = np.array(data_train.targets)
    eval_images_raw = data_eval.data
    eval_labels = np.array(data_eval.targets)
    np.savez('../PCA-lib/cifar10.npz', train_images_raw=train_images_raw,
    ↪train_labels=train_labels,
        eval_images_raw=eval_images_raw, eval_labels=eval_labels)

```

```

[5]: class_to_idx = {'airplane': 0,
                    'automobile': 1,
                    'bird': 2,
                    'cat': 3,
                    'deer': 4,
                    'dog': 5,
                    'frog': 6,
                    'horse': 7,
                    'ship': 8,
                    'truck': 9}

```

```

[6]: images_raw = np.concatenate([train_images_raw, eval_images_raw], axis=0)
labels = np.concatenate([train_labels, eval_labels], axis=0)
images_raw.shape, labels.shape

```

```

[6]: ((60000, 32, 32, 3), (60000,))

```

### 3 1. Principal Component Analysis

0. Let's say we have Data Matrix  $X$  with  $N$  rows (i.e., data points) and  $d$  columns (i.e., features).

$$X = [\cdots]_{N \times d}$$

1. Let's perform SVD on the  $X$ .

$$X = U_x S_x V_x^T$$

Let's assume  $N > d$  (We have 6000 data points per class, which is more than the 3072 dimensions).

By the way SVD works, we should have

$$U_x = [\cdots]_{N \times d}$$

$$S_x = [\cdots]_{d \times d}$$

$$V_x = [\cdots]_{d \times d}$$

and

$$U_x^T U_x = I_{d \times d}$$

$$V_x^T V_x = I_{d \times d}$$

2. The textbook says we need the following decomposition for the covariance matrix  $\Sigma$ :

$$\Sigma \mathcal{U} = \mathcal{U} \Lambda$$

3. We assume that  $X$  has mean zero (i.e., we already subtracted the feature averages). If  $X$  has  $N$  rows (i.e., data items), we have

$$\Sigma = \frac{1}{N} X^T X$$

4. Let's find  $\Sigma$  in terms of  $U_x$ ,  $S_x$ , and  $V_x$

$$\Sigma = \frac{1}{N} X^T X = \frac{1}{N} V_x S_x U_x^T U_x S_x V_x^T = V_x \frac{S_x^2}{N} V_x^T$$

$$\Rightarrow \Sigma V_x = V_x \frac{S_x^2}{N}$$

5. By comparison, we have

$$\mathcal{U} = V_x$$

$$\Lambda = \frac{S_x^2}{N}$$

### 3.0.1 Considering the above:

1. **There is no need to compute the covariance matrix  $\Sigma$**  and then find its diagonalization; You can easily perform SVD on the data matrix  $X$ , and get what you need!
2. In fact, you do not even need the matrices  $V_x$  and  $U_x$  for computing the mean squared error; You can infer the mean squared error using only the  $S_x$  matrix.
  - Numpy's SVD function `np.linalg.svd` has an argument `compute_uv` that turns off returning the  $U$  and  $V$  matrices for better efficiency. Therefore, you may be able to save some runtime in large data sets if you only care about the mean squared error!

## 4 Task 1

Write a function `pca_mse` that takes two arguments as input

1. **data\_raw**: a numpy array with the shape  $(N, \cdots)$ , where  $N$  is the number of samples, and there may be many excess dimensions denoted by  $\cdots$ . You will have to reshape this input **data\_raw** matrix to obtain a shape of  $(N, d)$ , where  $d$  is the vectorized data's dimension. For example, **data\_raw** could have an input shape of  $(6000, 50, 50, 3)$ . In this case you will have to reshape the input data to have a shape of  $(6000, 7500)$ .

2. `num_components`: This is the number of PCA components that we want to retain. This variable is denoted by  $r$  in the PCA definition in the textbook.

and returns the variable `mse` which is the mean squared error of the PCA projection into the designated number of principal components.

**Important Note:** Make sure you use `np.linalg.svd` for the SVD operation. Do not use any other matrix factorization function for this question (such as `np.linalg.eig`).

**Important Note:** Make sure you read and understand the notes from the previous cells before you start implementing. Failing to properly set the arguments for `np.linalg.svd` or trying to find the mean squared error by calculating the reconstruction error may cause extreme delays and timeouts for your implementation.

**Hint:** If you don't know how to extract the mean squared error of the PCA projection, or don't have a fresh probability and statistics memory, take a look at the Principal Component Analysis chapter of the most recent version of the textbook; the subsection titled "The error in a low-dimensional representation" explains how to find the mean squared error of the PCA projection as a function of the eigenvalues that were dropped.

```
[7]: def pca_mse(data_raw, num_components=20):  
  
    # your code here  
    shapetuple = data_raw.shape  
    nn = len(shapetuple)  
    r=1  
    for i in range(1,nn):  
        r = r * shapetuple[i]  
    data_raw.reshape(shapetuple[0],r)  
    new_data_raw = data_raw - np.mean(data_raw,axis=0)  
    s = np.linalg.svd(new_data_raw, full_matrices=True,compute_uv=False)  
  
    lam = (s**2)/shapetuple[0]  
    mse = sum(lam[num_components:r])  
  
    return np.float64(mse)
```

```
[8]: # Performing sanity checks on your implementation  
some_data = (np.arange(35).reshape(5,7) ** 13) % 20  
some_mse = pca_mse(some_data, num_components=2)  
assert some_mse.round(3) == 37.903  
  
# Checking against the pre-computed test database  
test_results = test_case_checker(pca_mse, task_id=1)  
assert test_results['passed'], test_results['message']
```

```
[31]: #Task 1 Test Cell
```

```
[9]: if perform_computation:
    class_names = []
    class_mses = []
    for cls_name, cls_label in class_to_idx.items():
        data_raw = images_raw[labels == cls_label,:,:,]
        start_time = time.time()
        print(f'Processing Class {cls_name}', end='')
        cls_mse = pca_mse(data_raw, num_components=20)
        print(f' (The SVD operation took %.3f seconds)' % (time.
→time()-start_time))
        class_names.append(cls_name)
        class_mses.append(cls_mse)
```

Processing Class airplane (The SVD operation took 0.609 seconds)  
Processing Class automobile (The SVD operation took 0.605 seconds)  
Processing Class bird (The SVD operation took 0.619 seconds)  
Processing Class cat (The SVD operation took 0.597 seconds)  
Processing Class deer (The SVD operation took 0.608 seconds)  
Processing Class dog (The SVD operation took 0.601 seconds)  
Processing Class frog (The SVD operation took 0.609 seconds)  
Processing Class horse (The SVD operation took 0.607 seconds)  
Processing Class ship (The SVD operation took 0.597 seconds)  
Processing Class truck (The SVD operation took 0.600 seconds)

```
[10]: if perform_computation:
    fig, ax = plt.subplots(figsize=(9,4.), dpi=120)
    sns.barplot(class_names, class_mses, ax=ax)
    ax.set_title('The Mean Squared Error of Representing Each Class by the
→Principal Components')
    ax.set_xlabel('Class')
    _ = ax.set_ylabel('Mean Squared Error')
```

```

↳ -----
TypeError                                Traceback (most recent call
↳last)

<ipython-input-10-9fa5465e8c7e> in <module>
      1 if perform_computation:
      2     fig, ax = plt.subplots(figsize=(9,4.), dpi=120)
----> 3     sns.barplot(class_names, class_mses, ax=ax)
      4     ax.set_title('The Mean Squared Error of Representing Each Class
↳by the Principal Components')
      5     ax.set_xlabel('Class')
```

```

/opt/conda/lib/python3.8/site-packages/seaborn/_decorators.py in
↳ inner_f(*args, **kwargs)
    44         )
    45         kwargs.update({k: arg for k, arg in zip(sig.parameters,
↳ args)})
---> 46         return f(**kwargs)
    47     return inner_f
    48

```

```

/opt/conda/lib/python3.8/site-packages/seaborn/categorical.py in
↳ barplot(x, y, hue, data, order, hue_order, estimator, ci, n_boot, units, seed,
↳ orient, color, palette, saturation, errcolor, errwidth, capsize, dodge, ax,
↳ **kwargs)
    3177 ):
    3178
-> 3179     plotter = _BarPlotter(x, y, hue, data, order, hue_order,
    3180                             estimator, ci, n_boot, units, seed,
    3181                             orient, color, palette, saturation,

```

```

/opt/conda/lib/python3.8/site-packages/seaborn/categorical.py in
↳ __init__(self, x, y, hue, data, order, hue_order, estimator, ci, n_boot,
↳ units, seed, orient, color, palette, saturation, errcolor, errwidth, capsize,
↳ dodge)
    1582         errwidth, capsize, dodge):
    1583         """Initialize the plotter."""
-> 1584         self.establish_variables(x, y, hue, data, orient,
    1585                                 order, hue_order, units)
    1586         self.establish_colors(color, palette, saturation)

```

```

/opt/conda/lib/python3.8/site-packages/seaborn/categorical.py in
↳ establish_variables(self, x, y, hue, data, orient, order, hue_order, units)
    154
    155         # Figure out the plotting orientation
--> 156         orient = infer_orient(
    157             x, y, orient, require_numeric=self.require_numeric
    158         )

```

```

/opt/conda/lib/python3.8/site-packages/seaborn/_core.py in
↳ infer_orient(x, y, orient, require_numeric)
    1343     elif require_numeric and "numeric" not in (x_type, y_type):
    1344         err = "Neither the `x` nor `y` variable appears to be
↳ numeric."

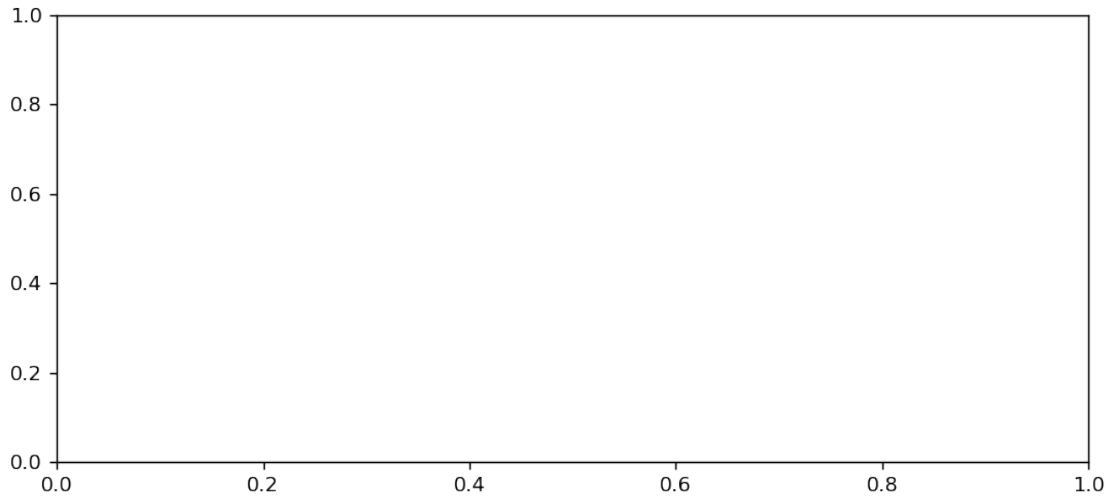
```

```

-> 1345         raise TypeError(err)
    1346
    1347     else:

```

TypeError: Neither the `x` nor `y` variable appears to be numeric.



## 5 2. Principal Coordinate Analysis

```

[11]: class_mean_list = []
      for cls_label in sorted(class_to_idx.values()):
          data_raw = images_raw[labels == cls_label, :, :, :]
          class_mean = np.mean(data_raw, axis=0).reshape(1,-1)
          class_mean_list.append(class_mean)
      class_means = np.concatenate(class_mean_list, axis=0)

```

## 6 Task 2

Write a function `mean_image_squared_distances` that takes the matrix `class_means` as an input and return the `SquaredDistances` matrix as output.

`class_means` is a numpy array like a traditional data matrix; it has a shape of  $(N, d)$  where there are  $N$  individual data-points where each is stored in a single  $d$  dimensional row.  $(N, d)$  could be anything, so do not make assumptions about it.

Your job is to produce the numpy array `SquaredDistances` whose  $i^{th}$  row and  $j^{th}$  column is the **squared** Euclidean distance between the  $i^{th}$  row of `class_means` and  $j^{th}$  row of `class_means`. Obviously \* The diagonal elements should be zero. \* The `SquaredDistances` should be symmetric.



```
[12]: def mean_image_squared_distances(class_means):
```

```
    # your code here
```

```
    M = class_means.shape[0]
```

```
    A_dots = (class_means*class_means).sum(axis=1).reshape((M,1))*np.  
    ↪ones(shape=(1,M))
```

```
    B_dots = (class_means*class_means).sum(axis=1)*np.ones(shape=(M,1))
```

```
    D_squared = A_dots + B_dots -2*class_means.dot(class_means.T)
```

```
    return D_squared
```

```
[13]: #Performing sanity checks
```

```
some_data = ((np.arange(35).reshape(5,7) ** 13) % 20) / 7.
```

```
some_dist = mean_image_squared_distances(some_data)
```

```
assert np.array_equal(some_dist.round(3), np.array([[ 0.    ,  4.551, 18.204,  8.  
    ↪306, 14.041],
```

```
                                                    [ 4.551,  0.    , 12.714,  3.  
    ↪918, 12.551],
```

```
                                                    [18.204, 12.714,  0.    ,  8.  
    ↪633,  8.735],
```

```
                                                    [ 8.306,  3.918,  8.633,  0.  
    ↪    ,  7.49 ],
```

```
                                                    [14.041, 12.551,  8.735,  7.  
    ↪49 ,  0.    ]]))
```

```
# Checking against the pre-computed test database
```

```
test_results = test_case_checker(mean_image_squared_distances, task_id=2)
```

```
assert test_results['passed'], test_results['message']
```

```
[14]: # Task 2 Test Cell
```

## 7 Task 3

Read and implement the Principal Coordinate Analysis procedure from your textbook by writing the function PCoA which takes the following arguments: 1. **SquaredDistances**: A numpy array which is square in shape, symmetric, and is the square of a distance matrix of some unknown set of points. The output of the `mean_image_squared_distances` function you wrote previously will be fed as this argument.

2. **r**: This is the dimension of the visualization space, and corresponds to the same  $r$  variable in the textbook procedure.

Things to keep in mind: 1. There is an erratum in the textbook's description of the PCoA procedure. There is a missing negative sign when computing the matrix  $\mathcal{W}$ ; the correct definition of  $\mathcal{W}$  is  $\mathcal{W} := -\frac{1}{2}\mathcal{A}\mathcal{D}^{(2)}\mathcal{A}^T$ . 2. It is **vital** to make sure that eigenvalues are sorted as the textbook mentioned, and the eigenvectors are also ordered accordingly. Some decomposition functions such as numpy's `np.linalg.eig` do not guarantee to return the eigenvalues and eigenvectors in any sorted way, and `np.linalg.eigh` guarantees to return them in ascending order; you will have to make sure they are sorted as the textbook says.

**Note:** You should only use `np.linalg.eigh` for matrix factorization in this question since we're dealing with a symmetric matrix; do not use `np.linalg.eig`, `np.linalg.svd`, or any other matrix decomposition function in this question.

```
[29]: def PCoA(SquaredDistances, r=2):
    assert SquaredDistances.shape[0] == SquaredDistances.shape[1]
    num_points = SquaredDistances.shape[0]
    A = np.identity(num_points) - (np.ones([num_points, num_points])/num_points)
    W = -(0.5 * np.dot(np.dot(A, SquaredDistances), A.T))
    eigval, eigvec = np.linalg.eigh(W)
    eigval_mat = np.diag(np.flip(eigval))
    eigvec_f = np.flip(eigvec, axis=1)
    eigvec_r = eigvec_f[:, :r]
    eigval_matr = np.sqrt(eigval_mat[:r, :r])
    VT = np.dot(eigvec_r, eigval_matr)

    assert VT.shape[0] == num_points
    assert VT.shape[1] == r
    return VT
```

```
[30]: some_data = ((np.arange(35).reshape(5,7) ** 13) % 20) / 7.
some_dist = mean_image_squared_distances(some_data)
some_pcoa = PCoA(some_dist, r=2)
assert np.array_equal(some_pcoa.round(3), np.array([[-1.974,  0.421],
                                                    [-1.285, -0.646],
                                                    [ 1.98 , -1.137],
                                                    [-0.091, -0.266],
                                                    [ 1.369,  1.628]]))

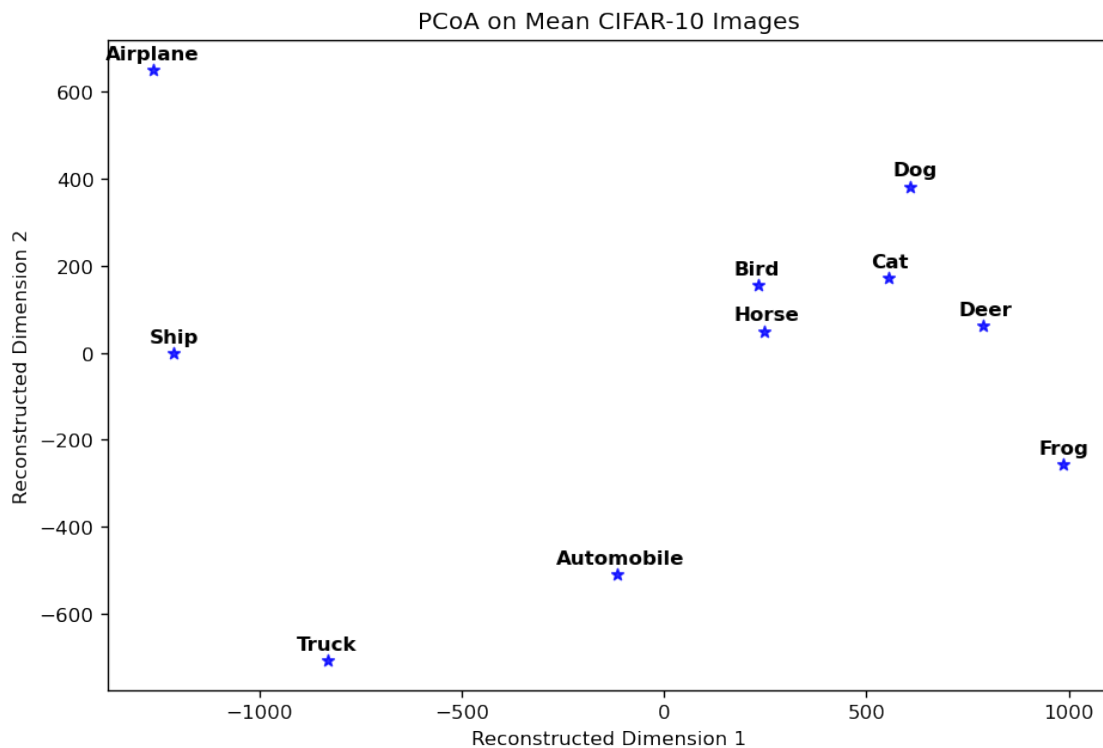
# Checking against the pre-computed test database
test_results = test_case_checker(lambda *args, **kwargs: PCoA(*args, **kwargs).
    ↳ astype(np.complex128), task_id=3)
assert test_results['passed'], test_results['message']
```

```
[31]: #Task 3 Test Cell
```

```
[32]: if perform_computation:
    SquaredDistances = mean_image_squared_distances(class_means)
```

```
VT = PCoA(SquaredDistances, r=2)
```

```
[33]: if perform_computation:
    class_names_list = sorted(list(class_to_idx.keys()))
    fig, ax = plt.subplots(figsize=(9,6.), dpi=120)
    x_components = VT[:,0]
    y_components = VT[:,1]
    sns.regplot(x=x_components, y=y_components, fit_reg=False, marker="*",
    ↪color="Blue", ax=ax)
    for class_idx in range(VT.shape[0]):
        num_letters = len(class_names_list[class_idx])
        ax.text(x_components[class_idx]-num_letters*15,
    ↪y_components[class_idx]+25,
                class_names_list[class_idx].capitalize(),
                horizontalalignment='left', size='medium', color='black',
    ↪weight='semibold')
    ax.set_xlabel('Reconstructed Dimension 1')
    ax.set_ylabel('Reconstructed Dimension 2')
    _ = ax.set_title('PCoA on Mean CIFAR-10 Images')
```



## 8 3. Generalized PCoA with Non-Metric Similarities

### 9 Task 4

Write a function `principal_components_precise_svd` that returns the principal components of a data matrix and takes the following arguments as input

1. `data_raw`: a numpy array 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 `data_raw` matrix to obtain a shape of  $(N, d)$ , where  $d$  is the vectorized data's dimension. For example, `data_raw` could have an input shape of  $(6000, 50, 50, 3)$ . In this case you will have to reshape the input data to have a shape of  $(6000, 7500)$ .
2. `num_components`: This is the number of PCA components that we want to retain. This variable is denoted by  $r$  in the PCA definition in the textbook.

and returns the variable `V_x` which is a numpy array with the shape  $(d, \text{num\_components})$ . The columns are the unitary principal components sorted descendingly with respect to the eigenvalues.

**Important Note:** Do not try to recover the covariance matrix  $\Sigma$  and then find its eigenvalues. This can prove to be both inefficient and unnecessary. As the theoretical review before the first task concluded, **There is no need to compute the covariance matrix  $\Sigma$** . Instead, all you need to do is to find the SVD of the data matrix, and extract the principal components from it.

**Important Note:** Do not use any matrix factorization function other than `np.linalg.svd` for this task; incorporating any other matrix factorization function (such as `np.linalg.eig`) may not be compatible with the results we expect and may even be inefficient.

```
[38]: def principal_components_precise_svd(data_raw, num_components=20):
    # your code here
    shapetuple = data_raw.shape
    nn = len(shapetuple)
    r=1
    for i in range(1,nn):
        r = r * shapetuple[i]
    data_raw.reshape(shapetuple[0],r)
    new_data_raw = data_raw - np.mean(data_raw,axis=0)
    U,sigma,vh = np.linalg.svd(new_data_raw, full_matrices=True)
    vh = vh.T
    V_x = vh[:, :num_components]

    # Don't mind/change the following lines.
    # This is for mitigating the ambiguity up to -/+1 factor in PCs.
    # (i.e., if x is a unitary PC, then -x is also a unitary PC).
    # We multiply each column by the sign of the largest element (in absolute
    ↪value) of that column
    sign_unambiguity = np.sign(V_x[np.abs(V_x).argmax(axis=0), np.arange(V_x.
    ↪shape[1])]).reshape(1,-1)
    V_x *= sign_unambiguity
```

```
return V_x
```

```
some_data = (np.arange(35).reshape(5,7) ** 13) % 20
some_pcs = principal_components_precise_svd(some_data, num_components=2)
assert np.array_equal(some_pcs.round(3), np.array([[ -0.123, -0.114],
                                                    [ -0.43  ,  0.119],
                                                    [ -0.021,  0.41  ],
                                                    [ -0.603, -0.164],
                                                    [  0.084,  0.491],
                                                    [ -0.223,  0.724],
                                                    [  0.616,  0.109]]))

# Checking against the pre-computed test database
test_results = test_case_checker(principal_components_precise_svd, task_id=4)
assert test_results['passed'], test_results['message']
```

```
#Task 4 Test Cell
```

The following cell will run your `principal_components_precise_svd` function on a single class of data, and provide you with some running time estimate.

```
if perform_computation:
    first_class_features = images_raw[labels == 0, :, :, :]

    starting_time = time.time()
    first_class_pcs = principal_components_precise_svd(first_class_features,
↪ num_components=20)
    end_time = time.time()

    print('Finding the principal components on a single class took %.3f seconds.
↪'%(end_time-starting_time))
```

```

IndexError                                Traceback (most recent call
last)

<ipython-input-41-11e8b4af82da> in <module>
      3
      4     starting_time = time.time()
----> 5     first_class_pcs =
principal_components_precise_svd(first_class_features, num_components=20)
      6     end_time = time.time()
      7

```

```

<ipython-input-38-6b427f5b9728> in_
↳ principal_components_precise_svd(data_raw, num_components)
    17     # (i.e., if x is a unitary PC, then -x is also a unitary PC).
    18     # We multiply each column by the sign of the largest element (in_
↳ absolute value) of that column
    ---> 19     sign_unambiguity = np.sign(V_x[np.abs(V_x).argmax(axis=0), np.
↳ arange(V_x.shape[1])]).reshape(1,-1)
    20     V_x *= sign_unambiguity
    21     return V_x

```

```

IndexError: shape mismatch: indexing arrays could not be broadcast_
↳ together with shapes (3,32,6000) (3,)

```

Although, this performance is extremely hardware-dependent, it's certainly not negligible. Keep in mind that we will have to run this function about 100 times on data of the same size to construct a similarity matrix in later tasks; any speedup may very well be appreciated.

Most of the computation time in the previous task was spent on the SVD factorization. Essentially, we found all the singular values and directions, ignored most of them, and only kept the top 20. This can be a very good place to start saving on computation; if only there was an SVD variant which you could tell in advance that you're only interested in the top 20 components, so that it wouldn't waste your time computing non-important singular values and directions...

**Spoiler Alert:** Such an efficient SVD variant exists, and sometimes is referred to as the “Truncated SVD” in application. Next task will be a redo of the previous task using this fast factroization.

## 10 Task 5

Similar to `principal_components_precise_svd`, write a function `principal_components` that uses `scikit-learn`'s `TruncatedSVD` decomposition instead of the precise `np.linalg.svd` decomposition that was used in the previous task. As in the previous task, `principal_components` should return the principal components of a data matrix and take the following arguments as input

1. `data_raw`: a numpy array 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 `data_raw` matrix to obtain a shape of  $(N, d)$ , where  $d$  is the vectorized data's dimension. For example, `data_raw` could have an input shape of  $(6000, 50, 50, 3)$ . In this case you will have to reshape the input data to have a shape of  $(6000, 7500)$ .
2. `num_components`: This is the number of PCA components that we want to retain. This variable is denoted by  $r$  in the PCA definition in the textbook.

`principal_components` should return the variable `V_x` which is a numpy array with the shape  $(d, \text{num\_components})$ . The columns are the unitay principal components sorted descendingly with respect to the eigenvalues.

**Important Note:** You should only use `scikit-learn`'s `TruncatedSVD` decomposition for this task. You can read about this function at <https://scikit->

[learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html](http://learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html).

- You must use the `randomized` algorithm implementation as it is more efficient.
- Since this heuristic is stochastic, you must provide `random_state=12345` as an input argument to this object's constructor for reproducibility.
- Use exactly 5 iterations for this heuristic (i.e., specify `n_iter` to be exactly 5).

**Important Note:** Do not try to recover the covariance matrix  $\Sigma$  and then find its eigenvalues. This can prove to be both inefficient and unnecessary. As the theoretical review before the first task concluded, **There is no need to compute the covariance matrix  $\Sigma$ .** Instead, all you need to do is to find the SVD of the data matrix, and extract the principal components from it.

```
[48]: def principal_components(data_raw, num_components=20):

    shapetuple = data_raw.shape
    nn = len(shapetuple)
    r=1
    for i in range(1,nn):
        r = r * shapetuple[i]
    data_raw.reshape(shapetuple[0],r)
    new_data_raw = data_raw - np.mean(data_raw,axis=0)
    svd = TruncatedSVD(n_components = num_components, algorithm='randomized',
    ↪n_iter=5, random_state=12345 )
    transformed = svd.fit_transform(new_data_raw)
    V_x = svd.components_.T

    assert V_x.ndim==2
    # Don't mind/change the following lines.
    # This is for mitigating the ambiguity up to -/+1 factor in PCs.
    # (i.e., if x is a unitary PC, then -x is also a unitary PC).
    # We multiply each column by the sign of the largest element (in absolute
    ↪value) of that column
    sign_unambiguity = np.sign(V_x[np.abs(V_x).argmax(axis=0), np.arange(V_x.
    ↪shape[1])]).reshape(1,-1)
    V_x *= sign_unambiguity
    return V_x
```

```
[49]: some_data = (np.arange(35).reshape(5,7) ** 13) % 20
some_pcs = principal_components(some_data, num_components=2)
assert np.array_equal(some_pcs.round(3), np.array([[ -0.123, -0.114],
                                                    [ -0.43 ,  0.119],
                                                    [ -0.021,  0.41 ],
                                                    [ -0.603, -0.164],
                                                    [  0.084,  0.491],
                                                    [ -0.223,  0.724],
                                                    [  0.616,  0.109]]))
```

```
# Checking against the pre-computed test database
test_results = test_case_checker(principal_components, task_id=5)
assert test_results['passed'], test_results['message']
```

[50]: #Task 5 Test Cell

```
[51]: if perform_computation:
        first_class_features = images_raw[labels == 0, :, :]

        starting_time = time.time()
        first_class_pcs = principal_components(first_class_features,
        num_components=20)
        end_time = time.time()

        print('Finding the principal components on a single class took %.3f seconds.
        '%(end_time-starting_time))
```

```

-----

ValueError                                Traceback (most recent call
last)

<ipython-input-51-b137cd82a5bb> in <module>
      3
      4     starting_time = time.time()
----> 5     first_class_pcs = principal_components(first_class_features,
num_components=20)
      6     end_time = time.time()
      7

<ipython-input-48-74a798387a83> in principal_components(data_raw,
num_components)
      9     new_data_raw = data_raw - np.mean(data_raw,axis=0)
     10     svd = TruncatedSVD(n_components = num_components,
algorithm='randomized', n_iter=5, random_state=12345 )
--> 11     transformed = svd.fit_transform(new_data_raw)
     12     V_x = svd.components_.T
     13

/opt/conda/lib/python3.8/site-packages/sklearn/decomposition/
truncated_svd.py in fit_transform(self, X, y)
    162         Reduced version of X. This will always be a dense array.
```



```

163         """
--> 164         X = self._validate_data(X, accept_sparse=['csr', 'csc'],
165                                 ensure_min_features=2)
166         random_state = check_random_state(self.random_state)

/opt/conda/lib/python3.8/site-packages/sklearn/base.py in
↪_validate_data(self, X, y, reset, validate_separately, **check_params)
419         out = X
420         elif isinstance(y, str) and y == 'no_validation':
--> 421         X = check_array(X, **check_params)
422         out = X
423         else:

/opt/conda/lib/python3.8/site-packages/sklearn/utils/validation.py in
↪inner_f(*args, **kwargs)
61         extra_args = len(args) - len(all_args)
62         if extra_args <= 0:
---> 63         return f(*args, **kwargs)
64
65         # extra_args > 0

/opt/conda/lib/python3.8/site-packages/sklearn/utils/validation.py in
↪check_array(array, accept_sparse, accept_large_sparse, dtype, order, copy,
↪force_all_finite, ensure_2d, allow_nd, ensure_min_samples,
↪ensure_min_features, estimator)
657         "into decimal numbers with dtype='numeric'")
↪from e
658         if not allow_nd and array.ndim >= 3:
--> 659         raise ValueError("Found array with dim %d. %s expected
↪<= 2."
660                             % (array.ndim, estimator_name))
661

```

ValueError: Found array with dim 4. Estimator expected <= 2.

Using this `principal_components` function, and the `images_raw` array, you could reconstruct an arbitrary image using a small number of components, see the effect of the number of components on the reconstructed image's quality, and share your results on Piazza!

## 11 Task 6

Write the function `E_A_given_B` that computes the  $E[A|B]$  and takes the two matrices `class_A_data` and `class_B_data`.

1. `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.
2. `class_B_data` has the same data structure as `class_A_data`.

To compute  $E[A|B]$ : 1. First, do whatever reshaping you have to do. 2. Subtract Class A's mean from its data 3. Use the `principal_components` function you wrote before to extract the 20 principal components of `class_B_data`. 4. Project Class A's data onto the mentioned principal components and get back to the original space. 5. Compute Class A's residuals (i.e., the difference between the original and the projection). 5. 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`.

```
[55]: def E_A_given_B(class_A_data, class_B_data):
    shapetuple = class_A_data.shape
    nn = len(shapetuple)
    r=1
    for i in range(1,nn):
        r = r * shapetuple[i]
    class_A_data.reshape(shapetuple[0],r)
    shapetupleII = class_B_data.shape
    mm = len(shapetupleII)
    f=1
    for i in range(1,mm):
        f = f * shapetupleII[i]
    class_B_data.reshape(shapetupleII[0],f)
    new_cAd = class_A_data - np.mean(class_A_data,axis=0)
    VxB = principal_components(class_B_data,num_components=20)
    class_A_projection = np.dot( VxB, np.dot(VxB.T, new_cAd.T)).T
    E_A_cond_B = np.mean(np.sum(np.square(class_A_projection - new_cAd),
    ↪axis=1))

    return E_A_cond_B
```

```
[56]: 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']
```

```
[57]: #Task 6 Test Cell
```

```
[58]: if perform_computation:
    num_classes = class_means.shape[0]
    SimilarityMatrix = np.zeros((num_classes, num_classes))
    for row in range(num_classes):
        print(f'Row {row}', end='')
        row_st_time = time.time()
        for col in range(row+1):
            class_A_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(f' (This row took %.3f seconds to finish)'%(time.time() -
↪row_st_time))
```

Row 0

```
↪-----
```

```
ValueError                                Traceback (most recent call↪
↪last)
```

```
<ipython-input-58-c9e323725cb3> in <module>
      8         class_A_data = images_raw[labels == row, :, :, :]
      9         class_B_data = images_raw[labels == col, :, :, :]
----> 10         E_A_cond_B = E_A_given_B(class_A_data, class_B_data)
      11         E_B_cond_A = E_A_given_B(class_B_data, class_A_data)
      12         SimilarityMatrix[col, row] = (E_A_cond_B + E_B_cond_A)/2.
```

```
<ipython-input-55-ed65683ec264> in E_A_given_B(class_A_data,↪
↪class_B_data)
      13     class_B_data.reshape(shapetupleII[0],f)
      14     new_cAd = class_A_data - np.mean(class_A_data,axis=0)
----> 15     VxB = principal_components(class_B_data,num_components=20)
      16     class_A_projection = np.dot( VxB, np.dot(VxB.T, new_cAd.T)).T
      17     E_A_cond_B = np.mean(np.sum(np.square(class_A_projection -↪
↪new_cAd), axis=1))
```

```
<ipython-input-48-74a798387a83> in principal_components(data_raw,↪
↪num_components)
```

```

    9     new_data_raw = data_raw - np.mean(data_raw,axis=0)
    10     svd = TruncatedSVD(n_components = num_components,
↪algorithm='randomized', n_iter=5, random_state=12345 )
    --> 11     transformed = svd.fit_transform(new_data_raw)
    12     V_x = svd.components_.T
    13

/opt/conda/lib/python3.8/site-packages/sklearn/decomposition/
↪_truncated_svd.py in fit_transform(self, X, y)
    162         Reduced version of X. This will always be a dense array.
    163         """
    --> 164         X = self._validate_data(X, accept_sparse=['csr', 'csc'],
    165                                     ensure_min_features=2)
    166         random_state = check_random_state(self.random_state)

/opt/conda/lib/python3.8/site-packages/sklearn/base.py in
↪_validate_data(self, X, y, reset, validate_separately, **check_params)
    419         out = X
    420         elif isinstance(y, str) and y == 'no_validation':
    --> 421         X = check_array(X, **check_params)
    422         out = X
    423         else:

/opt/conda/lib/python3.8/site-packages/sklearn/utils/validation.py in
↪inner_f(*args, **kwargs)
    61         extra_args = len(args) - len(all_args)
    62         if extra_args <= 0:
    --> 63             return f(*args, **kwargs)
    64
    65         # extra_args > 0

/opt/conda/lib/python3.8/site-packages/sklearn/utils/validation.py in
↪check_array(array, accept_sparse, accept_large_sparse, dtype, order, copy,
↪force_all_finite, ensure_2d, allow_nd, ensure_min_samples,
↪ensure_min_features, estimator)
    657         "into decimal numbers with dtype='numeric'")
↪from e
    658         if not allow_nd and array.ndim >= 3:
    --> 659             raise ValueError("Found array with dim %d. %s expected
↪<= 2."
    660                                     % (array.ndim, estimator_name))
    661

```

ValueError: Found array with dim 4. Estimator expected <= 2.

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.

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

```
[59]: array([[0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
           [0., 0., 0., 0., 0., 0., 0., 0., 0., 0.]])
```

## 12 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/ape/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/PLcourses/Ordination\\_sections\\_1.3+1.4\\_PCoA\\_Eng.pdf](http://biol09.biol.umontreal.ca/PLcourses/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: 1. If all eigenvalues computed during PCoA are non-negative, then `ProcessedSimilarityMatrix` should be the same as

the `SimilarityMatrix`. 2. 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: 1. 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. 2. Find the minimum eigenvalue and call it  $\lambda_{\min}$ . 3. If  $\lambda_{\min} \geq 0$ , then stop and return `SimilarityMatrix` as it was without any change. 4. 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).

```
[62]: def Lingoes_PreProcessing(SimilarityMatrix):

    assert SimilarityMatrix.shape[0] == SimilarityMatrix.shape[1]
    num_points = SimilarityMatrix.shape[0]

    A = np.identity(num_points) - (np.ones([num_points, num_points])/num_points)
    W = -(0.5 * np.dot(np.dot(A, SimilarityMatrix), A.T))
    eigval, eigvec = np.linalg.eigh(W)
    eigval_min = np.amin(eigval)
    if (eigval_min >= 0):
        ProcessedSimilarityMatrix = SimilarityMatrix
    elif (eigval_min < 0):
        ProcessedSimilarityMatrix = SimilarityMatrix.copy()
        icr = 2 * (abs(eigval_min))
        for i in range(0, num_points):
            for j in range(0, num_points):
                if (i != j):
                    ProcessedSimilarityMatrix[i][j] += icr

    return ProcessedSimilarityMatrix
```

```
[63]: some_data = ((np.arange(35).reshape(5,7) ** 13) % 20) / 7.
some_dist = mean_image_squared_distances(some_data)**5.
some_lingoes = Lingoes_PreProcessing(some_dist)
assert np.array_equal(some_lingoes.round(1), np.array([[ 0. , 898987.1,
↪ 2896177.9, 936570.7, 1442744.7],
                                                         [ 898987.1, 0. ,
↪ 1229280.9, 897958.5, 1208489.7],
```

```

                                [2896177.9, 1229280.9,
↪    0. , 944977.4, 947878.7],
                                [ 936570.7, 897958.5,
↪944977.4, 0. , 920604.3],
                                [1442744.7, 1208489.7,
↪947878.7, 920604.3, 0. ]]))

# Checking against the pre-computed test database
test_results = test_case_checker(Lingoes_PreProcessing, task_id=7)
assert test_results['passed'], test_results['message']

```

```
[64]: # Task 7 Test Cell
```

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

```

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

```

```
[66]: array([[0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.],
             [0., 0.]])

```

```
[67]: if perform_computation:
        class_names_list = sorted(list(class_to_idx.keys()))
        fig, ax = plt.subplots(figsize=(9,6.), dpi=120)
        x_components = VT[:,0]
        y_components = VT[:,1]
        sns.regplot(x=x_components, y=y_components, fit_reg=False, marker="*",
↪color="Blue", ax=ax)
        for class_idx in range(VT.shape[0]):
            num_letters = len(class_names_list[class_idx])
            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')

```

```
ax.set_xlabel('Reconstructed Dimension 1')
ax.set_ylabel('Reconstructed Dimension 2')
_ = ax.set_title('Generalized PCoA on CIFAR-10 Images')
```

```

↳
↳ -----
ValueError                                Traceback (most recent call↳
↳ last)

/opt/conda/lib/python3.8/site-packages/IPython/core/formatters.py in↳
↳ __call__(self, obj)
    339         pass
    340     else:
--> 341         return printer(obj)
    342         # Finally look for special method names
    343         method = get_real_method(obj, self.print_method)

/opt/conda/lib/python3.8/site-packages/IPython/core/pylabtools.py in↳
↳ <lambda>(fig)
    246
    247     if 'png' in formats:
--> 248         png_formatter.for_type(Figure, lambda fig: print_figure(fig,↳
↳ 'png', **kwargs))
    249     if 'retina' in formats or 'png2x' in formats:
    250         png_formatter.for_type(Figure, lambda fig:↳
↳ retina_figure(fig, **kwargs))

/opt/conda/lib/python3.8/site-packages/IPython/core/pylabtools.py in↳
↳ print_figure(fig, fmt, bbox_inches, **kwargs)
    130         FigureCanvasBase(fig)
    131
--> 132     fig.canvas.print_figure(bytes_io, **kw)
    133     data = bytes_io.getvalue()
    134     if fmt == 'svg':

/opt/conda/lib/python3.8/site-packages/matplotlib/backend_bases.py in↳
↳ print_figure(self, filename, dpi, facecolor, edgecolor, orientation, format,↳
↳ bbox_inches, pad_inches, bbox_extra_artists, backend, **kwargs)
    2208
    2209         try:
-> 2210             result = print_method(
    2211                 filename,
```



2212 dpi=dpi,

```
    /opt/conda/lib/python3.8/site-packages/matplotlib/backend_bases.py in
↳ wrapper(*args, **kwargs)
    1637         kwargs.pop(arg)
    1638
-> 1639         return func(*args, **kwargs)
    1640
    1641     return wrapper
```

```
    /opt/conda/lib/python3.8/site-packages/matplotlib/backends/backend_agg.
↳ py in print_png(self, filename_or_obj, metadata, pil_kwargs, *args)
    507         *metadata*, including the default 'Software' key.
    508         """
--> 509         FigureCanvasAgg.draw(self)
    510         mpl.image.imsave(
    511             filename_or_obj, self.buffer_rgba(), format="png",
↳ origin="upper",
```

```
    /opt/conda/lib/python3.8/site-packages/matplotlib/backends/backend_agg.
↳ py in draw(self)
    400     def draw(self):
    401         # docstring inherited
--> 402         self.renderer = self.get_renderer(cleared=True)
    403         # Acquire a lock on the shared font cache.
    404         with RendererAgg.lock, \
```

```
    /opt/conda/lib/python3.8/site-packages/matplotlib/backends/backend_agg.
↳ py in get_renderer(self, cleared)
    416             and getattr(self, "_lastKey", None) == key)
    417         if not reuse_renderer:
--> 418             self.renderer = RendererAgg(w, h, self.figure.dpi)
    419             self._lastKey = key
    420         elif cleared:
```

```
    /opt/conda/lib/python3.8/site-packages/matplotlib/backends/backend_agg.
↳ py in __init__(self, width, height, dpi)
    94         self.width = width
    95         self.height = height
---> 96         self._renderer = _RendererAgg(int(width), int(height), dpi)
    97         self._filter_renderers = []
    98
```

```
ValueError: Image size of 609169x49774 pixels is too large. It must be  
↳ less than 216 in each direction.
```

```
<Figure size 1080x720 with 1 Axes>
```

```
[ ]:
```

