Remember that only PDF submissions are accepted.

1. A dataset D has 64 points. Each point in D is a vector of length 6830, in other words there are6830 features. If we do Principal Component Analyses (PCA) of this dataset D, how many principal components with non-zero variance would we get? Explain why?

   To make my argument clear: N= 64, P=6830

   The rule: A sample of size N with P dimensions has at most N−1 principal components if N≤P.

   This means that the answer to our question is 63 principal components.

   Why N-1 principal components?

      Principal components are the axes of most variance in a dataset. Let's imagine that we have 2 points in that dataset, that would mean that there is only 1 axis of variance between them. This example scales up to data sets with a higher dimension as well. Therefore, the answer is N-1 or 63.

2. (Programming) Read this tech blog (https://sebastianraschka.com/Articles/2014_pca_step_by_step. html) (also uploaded to Canvas) and implement a PCA on your own. Feel free to call any command to compute eigenvector, eigenvalue, covariance matrix, SVD, etc.

   Just hand-in your code, and you don't have to play with data for now. (But you are encouraged to play with some toy data you find on your own to validate if your implementation is correct or not). You can reuse the code to this question in your mini-project, and play with data then. **See next page**

```python
# pca.py > ...
1    #author: Colin Vande Vijvere
2    #assignment: HW00
3    #---------------------------------------
4    #PRINCIPAL COMPONENT ANALYSIS (PCA)
5    #---------------------------------------
6    import numpy as np
7
8    matrix = np.array([[1,2,3],
9                       [4,5,6],
10                      [7,8,9],
11                      [9,10,11]])
12   print('input data:\n', matrix, '\n')
13
14   #getting mean of each column in the matrix
15   mean = np.mean(matrix.T, axis = 1)
16   print ('mean of each collumn:\n', mean,'\n')
17
18   #center columns by subtracting column means
19   matrix_C = matrix - mean
20   print ('matrix_C:\n', matrix_C,'\n')
21
22   #covariance matrix
23   matrix_cov = np.cov(matrix_C.T)
24   print ('covariance matrix:\n', matrix_cov,'\n')
25
26   #eigendecomposition of covariance matrix
27   values, vectors = np.linalg.eig(matrix_cov)
28   print(vectors)
29   print(values)
30   print('\n')
31
32
33   #projecting the data
34   projector = vectors.T.dot(matrix_C.T)
35   print(projector.T)
```