



Sheet#1 PCA

Data Matrix

1. Given the Data Matrix on the right answer the following questions
 - a. What is number of dimensions?
 - b. What are the types of the attributes?
 - c. What is the distance between x_1 and x_3 ?
 - d. What is the length of x_2 ?
 - e. What is the $\cos(\text{angle})$ between x_2 and x_4 ?
 - f. Do we need attribute scaling?
 - g. Compute the attribute scaled data matrix after scaling each attribute linearly between 0 and 1
 - h. Repeat parts c,d,e on the scaled data matrix in part (g)
2. Given the Data Matrix on the right submit your python code and its output that will do the following
 - a. Compute the norm of each instance. (5x1)
 - b. Compute the Cosine similarity matrix (5x5) matrix
 - c. Compute the Euclidean Distance matrix of the instances (5x5)

| ID | a1 | a2 | a3 | a4 |
|----------------------|----|----|----|----|
| 1 | 10 | 60 | 10 | 90 |
| 2 | 20 | 50 | 40 | 70 |
| 3 | 30 | 50 | 30 | 40 |
| 4 | 20 | 50 | 20 | 60 |
| 5 | 10 | 60 | 30 | 10 |
| DATA MATRIX D | | | | |

Principal Component Analysis

3. Given Data matrix above. Consider a_1 , a_2 and a_4 only
 - a. Write down the new data matrix **D3** (5x3)
 - b. Plot the data using 3d scatter plots
 - c. Compute the **mean** vector (3x1)
 - d. Compute centered data matrix **Z** by subtracting mean vector from the Data Matrix. (5x3)
 - e. Compute Covariance matrix **COV** (3x3)
 - f. Use python solvers to find eigenvalues (Diagonal 3x3 matrix) and eigen vectors (3x3) matrix. **Take care of the eigenvalues order.**

- g. Verify $U^T \Lambda U = \mathbf{COV}$.
- h. Compute the explained variance by the eigenvector corresponding to the largest eigenvalue. Do you think one eigenvector is good enough?
- i. Compute the projection matrix \mathbf{P} to go to 2-dimensions. Consider the top two eigenvectors of matrix \mathbf{U} according to eigenvalues. (3x2)
- j. Project the instances into a 2-Dimension space. $\mathbf{x}_j = \mathbf{P}^T \mathbf{x}$
- k. Plot the resulting Data matrix $\mathbf{D2}$ using scatter plots.

4. Given the data below , answer the following questions

- A. Compute 3x3 Covariance matrix of the 5 tuples dataset we have.
- B. The trace of the covariance matrix is the sum of the eigenvalues of the matrix.

1. **Compute** the three eigenvalues of the covariance matrix if

$$\frac{\lambda_a}{\lambda_b} = 0.505 \text{ and } \frac{\lambda_b}{\lambda_c} = 0.647, \text{ where } \lambda_a < \lambda_b < \lambda_c$$

2. **Determine** the explained variance using only λ_b, λ_c

| | X_1 | X_2 | X_3 |
|----------------|-------|-------|-------|
| \mathbf{x}_1 | 0.5 | 4.5 | 2.5 |
| \mathbf{x}_2 | 2.2 | 1.5 | 0.1 |
| \mathbf{x}_3 | 3.9 | 3.5 | 1.1 |
| \mathbf{x}_4 | 2.1 | 1.9 | 4.9 |
| \mathbf{x}_5 | 0.5 | 3.2 | 1.2 |

Notes and hints

1. NO hand written reports are allowed.
2. Make sure you did everything on your own.
3. Review slicing operators `[:]` in python/numpy. They are helpful with matrices and vectors.
4. Try to avoid loops unless you can't. Numpy optimized many operations on vectors and matrices.
5. Find many useful functions in numpy library. Always check the documentation and/ or stackoverflow.com
 - `numpy.linalg.eigh(A)` A is a matrix. Computes eigenvectors and eigenvalues of symmetric matrix A
 - `numpy.dot(A,B)` A, B are matrices/vectors. Computes dot product
 - `numpy.mean(A, axis=0)` A is matrix, axis =0 will average over the columns
 - `numpy.diag(A)` converts vector A into a diagonal matrix.
 - `numpy.vstack((A,B))` expand matrices A,B into one wider matrix → number of rows of A and B must match.
 - `numpy.transpose(A)` will compute the transpose of a matrix/vector A.