

# Asmt 5: Dimensionality Reduction

---

Turn in through GradeScope by 1:00pm:

Wednesday, November 12

100 points

## Overview

In this assignment you will explore regression techniques on high-dimensional data.

You will use two data sets for this assignment:

- <http://www.cs.utah.edu/~jeffp/teaching/DM/A5/A.csv>
- <http://www.cs.utah.edu/~jeffp/teaching/DM/A5/D.csv>

For python, you can use the following approach to load the data:

```
df = pd.read_csv('A.csv')
```

*As usual, it is recommended that you use LaTeX for this assignment. If you do not, you may lose points if your assignment is difficult to read or hard to follow. Find a sample form in this directory: <http://www.cs.utah.edu/~jeffp/teaching/latex/>*

Click [here](#) for an example template specifically created for this assignment.

## 1 Singular Value Decomposition, PCA, and LDA (70 points)

First we will compute the SVD of the matrix  $A$  we have loaded

```
import numpy as np
from scipy import linalg as LA
U, s, Vt = LA.svd(A, full_matrices=False)
```

Then take the top  $k$  components of  $A$  for values of  $k = 1$  through  $k = 6$  using

```
Uk = U[:, :k]
Sk = S[:k, :k]
Vtk = Vt[:k, :]
Ak = Uk @ Sk @ Vtk
```

**A (10 points):** Compute and report the  $L_2$  norm of the difference between  $A$  and  $A_k$  for each value of  $k$  (6 values) using

```
LA.norm(A-Ak, 2)
```

**B (10 points):** Find the smallest value  $k$  so that the  $L_2$  norm of  $A-A_k$  is less than 10% that of the  $L_2$  norm of  $A$ ;  $k$  might or might not be larger than 6. Report (i) the  $L_2$  norm for  $A$ , (ii) the  $L_2$  norm for your choice of  $A-A_k$ , and (iii) your choice of  $k$ .

**C (10 points):** Plot the points in 2 dimensions by projecting  $A$  onto the top 2 right singular values.

**D (10 points):** Now repeat (**B**) for PCA. First center the data to get  $\tilde{A}$ , and then find the value  $k$  where the  $L_2$  norm of  $\tilde{A}-A_k$  is less than 10% that of the  $L_2$  norm of  $\tilde{A}$ .

**E (10 points):** Plot the points in 2 dimensions by projecting  $A$  onto the top 2 principal components.

**F (20 points):** Next repeat the creation of 2-dimensional representation of creating the dataset but with LDA. Treat the data as in 3 clusters of 100 points each. The first 100 rows are the first cluster, the next 100 the second cluster, and the third 100 the third cluster. Run Latent Discriminant Analysis with these labels. The output should be a size=300 point set in 2 dimensions.

Plot the points in 2 dimensions by projecting A onto the top 2 LDA components.

## 2 Multidimensional Scaling (30 points)

You will apply multidimensional scaling on an all-pairs distance among US airports, stored as matrix D.

**A (10 points):** Transform the matrix into a  $D^{(2)}$  matrix where each element is squared. Report the Frobenius norm of  $D^{(2)}$ .

**B (10 points):** Double center the matrix so  $M = -\frac{1}{2}C_n D^{(2)} C_n$ , and report the Frobenius norm of  $M$ .

**C (10 points):** Plot the data in 2 dimensions on the top 2 eigenvectors of  $M$ .

## 3 BONUS (10 points):

Professor Phillips recently found the following method for distance metric learning, without much explanation. Let  $C = \{x_i, x'_i\}$  be a set of  $n_C$  points we would like to be close from each other, each  $x_i$  in  $\mathbb{R}^d$ . Let  $F = \{x_j, x'_j\}$  be a set of  $n_F$  points we would like to far from each other, each  $x_j$  in  $\mathbb{R}^d$ . Define an  $d \times d$  matrix

$$M = \left( \alpha I + \frac{\beta}{n_C} \sum_{\{x_i, x'_i\} \in C} (x_i - x'_i)(x_i - x'_i)^T - \frac{\gamma}{n_F} \sum_{\{x_j, x'_j\} \in F} (x_j - x'_j)(x_j - x'_j)^T \right)^{-1},$$

for some values of  $\alpha, \beta, \gamma > 0$ .

Describe why, or under which scenarios, or for which  $\alpha, \beta, \gamma$  that the resulting Mahalanobis distance  $d_M(p, q) = \sqrt{(p - q)^T M (p - q)}$  respects the close and far input pairs. Or show that it does not generally work (very well). Maybe trying it out on some data would be a good place to start.