

Assignment 1

Methods of PCA, MDS and Isomap

DD2434 ADVANCED MACHINE LEARNING
FILIP BERGENTOFT, BERGENTO@KTH.SE

Problem 1

It is assumed throughout this problem that whenever a matrix A is referenced, it is a real and symmetric matrix of size $n \times n$.

Part (i): *Prove that a real symmetric matrix has real eigenvalues*

Let λ be an eigenvalue to A with a corresponding vector v , which is by the definition of an eigenvector non-zero.

The norm of a complex vector z is given by $\|z\| = \sqrt{z^T \bar{z}}$ and is non-zero for all $z \in \mathbb{C}$ except $z = 0$.

Given that A is real we can use that

$$\overline{Av} = \overline{\lambda v} = A\bar{v} = \overline{\lambda v} \quad (1)$$

Then we can expand $v^T A\bar{v}$ in the two following ways

$$v^T A\bar{v} = v^T (A\bar{v}) = \left\{ \text{Using equation (1)} \right\} = v^T \bar{\lambda} v = \bar{\lambda} v^T v = \bar{\lambda} \|v\|^2 \quad (2)$$

$$v^T A\bar{v} = (A^T v)^T \bar{v} = \left\{ A = A^T \right\} = (Av)^T \bar{v} = \lambda v^T \bar{v} = \lambda \|v\|^2 \quad (3)$$

Thus since equation (2) and (3) are equal (from the left hand side) we get that

$$\bar{\lambda} \|v\|^2 = \lambda \|v\|^2 \Rightarrow \lambda = \bar{\lambda} \quad (4)$$

since v is non-zero by the definition of an eigenvector. Thus are the eigenvalues of a real and symmetric matrix real.

Part (ii): qqq: Detta måste fixas!!

Part (iii): *Prove that a positive semi definite matrix has non negative eigenvalues*

Let B be a complex $n \times n$ positive semi-definite matrix accompanied by an eigenvalue λ and an associated eigenvector v . By the definition of a positive semi-definite matrix we know that for any vector $z \in \mathbb{C}$ it holds that $\bar{z}^T B z \geq 0$. The following thus holds

$$\bar{v}^T B v = \bar{v}^T \lambda v = \lambda \bar{v}^T v = \lambda \|v\|^2 \geq 0 \quad (5)$$

where $\|v\|^2$ is necessarily positive since eigenvectors are by definition non-zero. This implies that $\lambda \geq 0$, which was to be proven.

Part (iv): Let $A \in \mathbb{R}^{n \times n}$ symmetric and positive semi-definite matrix. Define a matrix $D = \{D_{ij} | D_{ij} = A_{ii} + A_{jj} - 2A_{ij}\}$. Show that there exists n vectors v_1, \dots, v_n , $v_i \in \mathbb{R}^n \forall i$ such that $D_{ij} = \|v_i - v_j\|_2^2$.

$$D_{ij} = \|v_i - v_j\|_2^2 = v_i^T v_i + v_j^T v_j - 2v_i^T v_j \quad (6)$$

Thus, if we can show that any matrix $A \in \mathbb{R}^{n \times n}$ that fulfils the given conditions, can have each of its elements expressed as a dot product between n given vectors v_1, \dots, v_n , $v_i \in \mathbb{R}^n \forall i$ we have shown what is asked.

Given that A is symmetric and positive semi-definite it can be decomposed into the following eigen-decomposition $A = Q\Lambda Q^T$ where Λ is a diagonal matrix of real eigenvalues (since A is PSD) and Q is an orthogonal matrix of eigenvectors of A . The decomposition can then be rewritten in the following manner

$$A = Q\Lambda Q^T = (\Lambda^{\frac{1}{2}} Q^T)^T (\Lambda^{\frac{1}{2}} Q^T) = V^T V \quad (7)$$

If we thus choose every v_i , $i = 1, \dots, n$ to be the i :th column in the matrix $(\Lambda^{\frac{1}{2}} Q^T)$ we get that $A_{ij} = (V^T V)_{ij} = v_i^T v_j$ which yields the final result

$$D_{ij} = A_{ii} + A_{jj} - 2A_{ij} = v_i^T v_i + v_j^T v_j - 2v_i^T v_j = \|v_i - v_j\|_2^2 \quad (8)$$

which was to be proven.

Problem 2

Given that we want to do PCA using k components on a matrix $Y \in \mathbb{R}^{p \times n}$, $p \leq n$ where the columns correspond to data points, we can decompose Y in the following manner using SVD.

$$Y = U\Sigma V^T \quad (9)$$

Where U contains the left singular vectors of Y and V contains the right singular vectors of Y .

The transformation matrix W is then given by the first k columns of the matrix U .

If we instead want to perform PCA using k components on Y^T we can use that

$$Y^T = (U\Sigma V^T)^T = V\Sigma U^T$$

The transformation matrix \widetilde{W} can thus be chosen as the first k columns in V .

Thus is a single SVD computation sufficient for computing PCA on both columns and rows.

Problem 3

Want to maximise the expression $tr(Y^T W W^T Y)$ where $Y \in \mathbb{R}^{d \times n}$, $W \in \mathbb{R}^{d \times k}$, $k < d$ under the condition that W has orthonormal columns which implies that $W^T W = I_k$. In order to achieve this we want to use Lagrange multipliers, and for mathematical convenience we will use the cyclic property of trace to shift around the matrices in the expression in the following manner.

$$tr(Y^T W W^T Y) = tr(Y Y^T W W^T) = tr(W^T Y Y^T W) \quad (10)$$

We thus want to maximise $tr(W^T Y Y^T W)$ subject to $W^T W - I_k = 0$ which yields the following Lagrangian function

$$L = tr(W^T Y Y^T W) - \sum_{i=1}^k \sum_{j=1}^k \lambda_{ij} (w_i^T w_j - \mathbf{1}\{i = j\}) \quad (11)$$

$$= \sum_{i=1}^k w_i^T Y Y^T w_i - \sum_{i=1}^k \sum_{j=1}^k \lambda_{ij} (w_i^T w_j - \mathbf{1}\{i = j\}) \quad (12)$$

where $\mathbf{1}$ is the indicator function.

Taking the derivative with respect to the λ_{lj} only results in the initially stated condition that

$$W^T W = I_k \quad (13)$$

Taking the derivative with respect to w_l , $l = 1, \dots, k$ yields

$$\frac{\partial L}{\partial w_l} = 2Y Y^T w_l - \sum_{j=1}^k (\lambda_{lj} w_j + \lambda_{jl} w_j) = 0 \quad (14)$$

By letting Λ be a matrix with elements $(\Lambda)_{lj} = \lambda_{lj}$ equation (14) can be expressed using matrices for all $l = 1, \dots, n$.

$$2YY^TW - W(\Lambda^T + \Lambda) = 0 \Rightarrow YY^TW = W \frac{\Lambda^T + \Lambda}{2} \quad (15)$$

Noting that $\frac{\Lambda^T + \Lambda}{2}$ is trivially symmetric it can be diagonalised in the following manner: $\frac{\Lambda^T + \Lambda}{2} = PDP^{-1}$. Substituting this into equation (15) results in

$$YY^TW = W \frac{\Lambda^T + \Lambda}{2} = WPD P^{-1} \quad (16)$$

$$\Rightarrow W^T YY^TW = W^T WPD P^{-1} = \left\{ W^T W = I_k \text{ by (13)} \right\} = PDP^{-1} \quad (17)$$

$$\Rightarrow P^{-1}W^T YY^TW P = P^{-1}PDP^{-1}P \quad (18)$$

$$\Rightarrow P^{-1}W^T YY^TW P = D \quad (19)$$

Thus is YY^T diagonalised by WP and D thus consist of k eigenvalues of YY^T and by substituting equation (18) into the original trace expression we get

$$tr(W^T YY^TW) = tr(PDP^{-1}) = tr(P^{-1}PD) = tr(D) = \sum_{i=1}^k d_i \quad (20)$$

where d_i are eigenvalues of YY^T . Thus in order to maximise this we just choose $d_i, i = 1, 2, \dots, k$ to be the k largest eigenvalues of YY^T .

Now to show that if we select $W = U_k$ we get the same maximum value. Using the skinny SVD of $Y = U\Sigma V^T$ and substituting into the original trace expression we get

$$tr(Y^T WW^T Y) = tr(V\Sigma^T U^T WW^T U\Sigma V^T) \quad (21)$$

$$= tr(V^T V\Sigma U^T WW^T U\Sigma) = tr(\Sigma U^T WW^T U\Sigma) \quad (22)$$

$$= tr(\Sigma U^T U_k U_k^T U\Sigma) = tr((\Sigma U^T U_k)(\Sigma U^T U_k)^T) \quad (23)$$

$$= tr((\Sigma I_{d \times k})(\Sigma I_{d \times k})^T) = tr\Sigma_{k \times k}^2 \quad (24)$$

$$= \sum_{i=1}^k \sigma_i^2 = \left\{ k \text{ largest singular values} \right\} = \sum_{i=1}^k d_i \quad (25)$$

Thus is $\text{tr}(W^T Y Y^T W)$ subject to $W^T W - I_k = 0$ maximised by choosing $W = U_k$ and the maximum value is given by $\sum_{i=1}^k \sigma_i^2$.

Problem 4

In order to show that it provides a correct estimation of the Gram matrix S we need to prove that we can derive a similarity matrix S from a matrix D where $(D)_{ij} = d_{ij} = (z_i - z_j)^T (z_i - z_j)$ such that

1. $S = Z^T Z$
2. Z gives 0 reconstruction error for the matrix D

Normally $s_{ij} = y_i^T y_j$, we claim that it is valid to express

$$s_{ij} = -\frac{1}{2}(d_{ij}^2 - d_{1j}^2 - d_{1i}^2) \quad (26)$$

where

$$d_{ij}^2 = (y_i - y_j)^2 = y_i^2 + y_j^2 - 2y_i y_j \quad (27)$$

.

Substituting (27) into (26) yields the following

$$s_{ij} = -\frac{1}{2}(d_{ij}^2 - d_{1j}^2 - d_{1i}^2) \quad (28)$$

$$= -\frac{1}{2}(y_i^2 + y_j^2 - 2y_i^T y_j - y_1^2 - y_i^2 + 2y_1^T y_i - y_1^2 - y_j^2 + 2y_1^T y_j) \quad (29)$$

$$= -\frac{1}{2}(-2y_i^T y_j - 2y_1^2 + 2y_1^T y_i + 2y_1^T y_j) \quad (30)$$

$$= (y_i - y_1)^T (y_j - y_1) = z_i^T z_j \quad (31)$$

Let T be a matrix of the same size as Y , with all columns as y_1 . Then the Gram matrix S can be written as

$$S = (Y - T)^T (Y - T) = Z^T Z \quad (32)$$

Thus is 1. shown. Now we want to show 2.

We know that $d_{i,j} = (z_i - z_j)^T (z_i - z_j)$ and by substituting the result from equation (31) we get the following

$$d_{ij}^2 = (z_i - z_j)^T (z_i - z_j) \quad (33)$$

$$= (y_i - y_1 - (y_j - y_1))^T (y_i - y_1 - (y_j - y_1)) \quad (34)$$

$$= (y_i - y_j)^T (y_i - y_j) \quad (35)$$

Which shows that Z gives 0 reconstruction error for the distance matrix D .

Problem 5

We are considering the classical MDS when $Y \in \mathbb{R}^{d \times n}$ is known which implies that we can construct a similarity matrix $S = Y^T Y$. An MDS embedding can then be obtained by performing eigen-decomposition on S which yields the following

$$S = Y^T Y = Q \Lambda Q^T = (\Lambda^{-\frac{1}{2}} Q^T)^T (\Lambda^{-\frac{1}{2}} Q^T) \quad (36)$$

where Λ is a diagonal matrix with eigenvalues of $S = Y^T Y$ sorted in descending order.

The latent matrix X can then be chosen as

$$X = I_{k \times n} \Lambda^{-\frac{1}{2}} Q^T, \quad k < d \quad (37)$$

However, this can also be written in terms of the SVD of $Y = U \Sigma V^T$

$$S = Y^T Y = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T = (\Sigma V^T)^T (\Sigma V^T) \quad (38)$$

$$= (\Lambda^{-\frac{1}{2}} Q^T)^T (\Lambda^{-\frac{1}{2}} Q^T) \quad (39)$$

Since Λ and Σ has ordered diagonals by magnitude and V contains the right singular vectors we know that $V = Q$ which gives that

$$\Lambda^{-\frac{1}{2}} Q^T = \Sigma V^T \quad (40)$$

Thus can the latent matrix also be expressed as

$$X = I_{k \times n} \Sigma V^T, \quad k < d \quad (41)$$

In PCA we decompose Y into its SVD, yielding

$$Y = U \Sigma V^T \quad (42)$$

The latent matrix is then chosen as

$$X = (U I_{n \times k})^T Y = I_{k \times n} U^T Y \quad (43)$$

$$= \left\{ Y = U \Sigma V^T \Rightarrow U^T Y = \Sigma V^T \right\} \quad (44)$$

$$= I_{k \times n} \Sigma V^T, \quad k < d \quad (45)$$

Comparing equations (41) and (45) we see that the results are equivalent.

Regarding the computational efficiency I believe that there are two areas to consider: computational complexity and numerical accuracy. The complexities of computing the SVD and eigen-decomposition are similar and somewhat dependant on the properties of the input matrix, eigen-decomposition often performing a bit better. However, for the MDS we are required to perform a matrix multiplication beforehand which both introduces additional cost in time but can also introduce issues with numerical accuracy causing the resulting eigen-decomposition to have imaginary parts. For this reason the SVD should be preferred.

Problem 6

We want to argue that the process to obtain a neighbourhood graph G in the Isomap method may yield a disconnected graph. Given that we have a set of points that corresponds to those visualised in 1 the process to obtain a neighbourhood graph G will yield a disconnected graph if we only connect a point to its 2 nearest neighbours. This as for any given point in cluster A , its two nearest neighbours both belong to cluster A and for any given point in cluster B , its two nearest neighbours both belong to cluster B . Thus are cluster A and B not connected and the graph is disconnected, causing the Isomap algorithm to fail.

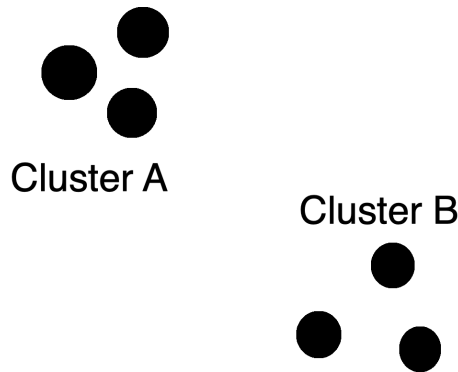


Figure 1: Example of 2 clusters

An example of heuristic to path this problem is to use create artificial data points that bridge the gap between the different clusters, resulting in the

graph to become connected again. This could be achieved by using the fact that a disconnected neighbourhood graph implies that there exists clusters within the data and use the following scheme:

1. Identify all C_i , $i = 1, 2, \dots, k$ clusters by using the neighbourhood function G
2. Pick a random point p_i , $i = 1, 2, \dots, k$ of each cluster (the mean of the cluster would also work but harder to justify)
3. Connect all points p_i with each other using $\text{linspace}(p_i, p_j, n)$ where n is chosen such that the distance between the artificial data points is equal to the smallest distance of the original data set.

This would ensure that all clusters would be tied together as the distances between the artificial points equals the smallest distance of the original data set, ensuring that the artificial points will always belong to the neighbourhood sets of each other. The expected value of the distance between two clusters using this method will roughly correspond to the distance between the two cluster means.

Problem 7

The goal of this problem is to visualise how similar different animals at a zoo are by projecting the given data from \mathbb{R}^{16} to

$$\mathbb{R}^2$$

using three different embeddings; PCA, MDS and isomap.

Preprocessing the data

In order to enable the applications of the embeddings the data has to be preprocessed by first removing the columns 'type' and 'animal name' from the data set. These can later be used in the visualisation.

The remaining attributes are all boolean with values in $\{0, 1\}$ except the attribute 'legs' which takes values in $\{0, 2, 4, 6, 8\}$. This attribute thus takes on values up to 8 times the magnitude in comparison to the remaining attributes which corresponds to it having an importance weight of 8 and the other a weight of 1 when taking the distance between points. For this reason the magnitude of 'legs' will be scaled down by 8 times, causing it to instead take values in $\{0, \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1\}$.

Visualisation of data

In order to project the data in 2D such that similar animals are projected close to one another the first two resulting latent variables from a method will be plotted against each other as they contain the most information about the data set as a result of the sorting of the singular values and eigenvalues in descending order.

PCA

The implementation of the PCA method is short using Numpy's Singular Value Decomposition function and following the exact method presented in the lecture. The only concern I had was the centring of the data matrix as it removes the boolean nature of the matrix. However I argue that this only alter the values and does not remove the information of the attributes rendering the action as viable.

MDS

Isomap

Comparison