

A Note of Diffusion Map

Singyuan Yeh

2020 NCTS Mini-Course

July, 2020

Table of Contents

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map
- 7 Reference

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map
- 7 Reference

Definition

Given n data points $\{x_1, \dots, x_n\}$ in \mathbb{R}^d . Consider a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

Write $f_i = f(x_i)$, $1 \leq i \leq n$. We hope to minimize

$$E(f) = \sum_{i,j=1}^n A_{ij} (f_i - f_j)^2.$$

The matrix A is symmetric transition matrix, i.e. $A_{ij} \geq 0$ and

$$\sum_i A_{ij} = \sum_j A_{ij} = 1.$$

Definition

Let's look at the quadratic form $A_{ij} (f_i - f_j)^2$.

$$\begin{aligned}\sum_{i,j=1}^n A_{ij} (f_i - f_j)^2 &= \sum_{i=1}^n f_i^2 \left(\sum_{j=1}^n A_{ij} \right) - 2 \sum_{i,j=1}^n f_i A_{ij} f_j + \sum_{j=1}^n f_j^2 \left(\sum_{i=1}^n A_{ij} \right) \\ &= \sum_{i=1}^n f_i^2 + 2 \sum_{i,j=1}^n f_i A_{ij} f_j + \sum_{j=1}^n f_j^2 \\ &= 2f^T (I - A)f\end{aligned}$$

Define $L = I - A$. Hope to minimize $E(f) = 2f^T Lf$, subject to some constraints that is not mentioned here.

Now we turn our attention to Euclidean \mathbb{R}^1 Laplacian operator.

$$L = \begin{bmatrix} \ddots & & & \\ & 1 & -2 & 1 \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$$

Then, matrix $A = \begin{bmatrix} \ddots & & & \\ & 1 & 0 & 1 \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$ and $D = \begin{bmatrix} \ddots & & & \\ & -2 & & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$, where A is adjacent matrix in Euclidean space and $D_{ii} = -\sum_j A_{ij}$ shows degree of point i . Therefore $-L = D - A$ is roughly called **unnormalized** Laplacian operator.

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map
- 7 Reference

Affinity graph

Given a pair of graph $G = (V, E)$. If we allow the **affinity** function $\omega : E \rightarrow \mathbb{R}_+$, then we called **affinity** graph $G = (V, E, \omega)$.

Remark

A function ω can regarded as some kind of “distance function” between two vertices. We can convert a Euclidean discrete space into an affinity graph by setting a function $\omega(i, j) = 1$ for all $(i, j) \in E$.

Affinity matrix

- Given a graph $G = (V, E)$ and $|V| = n$, the adjacency matrix of G is matrix $W \in \mathbb{R}^{n \times n}$ defined by

$$W_{i,j} = \begin{cases} 1 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}.$$

- Given a **affinity** graph $G = (V, E, W)$ and $|V| = n$, the adjacency matrix of G is matrix $W \in \mathbb{R}^{n \times n}$ defined by

$$W_{i,j} = \begin{cases} \omega_{ij} & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}.$$

Degree matrix

Let G is affinity graph. The degree function $d : V \rightarrow \mathbb{R}_+$ is defined by

$$d(i) = \sum_{(i,j) \in E} W_{ij}.$$

The degree matrix $D \in \mathbb{R}^{n \times n}$ is defined by a diagonal matrix

$$D = \begin{bmatrix} d(1) & & \\ & \ddots & \\ & & d(n) \end{bmatrix}.$$

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
 - Definition
 - Spectral Propositions
 - Diagonalization of GL
 - Diffusion map
- 4 Example
- 5 Comparison to Kernel PCA

Definition

We focus on **undirected** graph $G = (V, E)$ with n vertices. Let $G = (V, E, \omega)$ be an undirected affinity graph.

- The unnormalized graph Laplacian (GL) is defined as $\tilde{L} = D - W$.
- If there is no isolated vertex, the normalized graph Laplacian (NGL) is defined as $L = I_n - D^{-1}W$. (NOT necessary symmetric)
- The **symmetrized** normalized graph Laplacian is defined as $\mathcal{L} = I_n - D^{-\frac{1}{2}}WD^{-\frac{1}{2}}$.

Some relation between definition:

- $L = D^{-1}\tilde{L}$. i.e. Normalized \tilde{L} .
- $\mathcal{L} = D^{\frac{1}{2}}(I - D^{-1}W)D^{-\frac{1}{2}} = D^{\frac{1}{2}}LD^{-\frac{1}{2}}$. i.e. L is similar to \mathcal{L} .

Transition matrix

Definition

Define the transition matrix of the random walk on the graph as $A = D^{-1} W$. It is **NOT** necessary symmetric.

Proposition

- $\sum_j A_{ij} = 1$
- $\sum_{j=1}^n (A^k)_{ij} = 1$

Remark: . A entry A_{ij} can be thought of as the probability of moving from i to j in one step of a random walk on G .

Proof of proposition

Proposition

- $\sum_j A_{ij} = 1$
- $\sum_{j=1}^n (A^k)_{ij} = 1$

Proof:

- $\sum_j A_{ij} = \sum_j \frac{1}{d_i} \delta_{ik} W_{kj} = \sum_j \frac{1}{d_i} W_{ij} = 1$

-

$$\begin{aligned} \sum_{j=1}^n (A^k)_{ij} &= \sum_{j_1, \dots, j_k, j=1}^n A_{ij_1} A_{j_1 j_2} \cdots A_{j_{k-1} j_k} A_{j_k j} \\ &= \sum_{j_1=1}^n A_{ij_1} \sum_{j_2=1}^n A_{j_1 j_2} \cdots \sum_{j_k=1}^n A_{j_{k-1} j_k} \sum_{j=1}^n A_{j_k j} \\ &= 1 \end{aligned}$$

Notation

Before some basic spectral properties of the GL are provided, we introduced some notation.

- Denote $\sigma(M)$ to be the spectrum of a given matrix M .
- Denote $\rho(M)$ to be the associated spectral radius

$$\rho(M) = \max_{f \neq 0} \frac{\|f^T M f\|}{f^T f} = \max\{|\lambda| : \lambda \in \sigma(M)\}$$

Definition

The Rayleigh quotient of a matrix $M \in \mathbb{R}^{n \times n}$ is defined as

$$RM(v) = \frac{\langle v, Mv \rangle}{\langle v, v \rangle}$$

Nonnegative definite GL (1)

Proposition of GL

The unnormalized graph Laplacian $\tilde{L} = D - W$ is nonnegative definite and $\sigma(\tilde{L}) \subset [0, 2\rho(D)]$.

Proof: Let $f \in \mathbb{R}^n$ and $d_i = \sum_j W_{ij}$. Since $W_{ij} \geq 0$ and $W_{ij} = W_{ji}$,

$$\begin{aligned} f^T \tilde{L} f &= f^T (D - W) f = \sum_{i,j=1}^n f_i (d_i \delta_{ij} - W_{ij}) f_j = \sum_{i=1}^n d_i f_i^2 - \sum_{i,j=1}^n f_i W_{ij} f_j \\ &= \sum_{i=1}^n \sum_{j=1}^n W_{ij} f_i^2 - \sum_{i,j=1}^n f_i W_{ij} f_j \\ &= \frac{1}{2} \left(\sum_{i=1}^n \sum_{j=1}^n W_{ij} f_i^2 + \sum_{j=1}^n \sum_{i=1}^n W_{ji} f_j^2 - \sum_{i,j=1}^n 2 f_i W_{ij} f_j \right) \\ &= \frac{1}{2} \sum_{i,j=1}^n W_{ij} (f_i - f_j)^2 \geq 0 \end{aligned}$$

Nonnegative definite GL (2)

Proof: Since above equality holds if $f_1 = \dots = f_n$, which implies the smallest eigenvalue of \tilde{L} is 0 w.r.t. eigenvector $\mathbf{1}$.

On the other hands, since $(f_i - f_j)^2 \leq 2(f_i^2 + f_j^2)$ and $\rho(D) = \max_i d_i$,

$$\begin{aligned} f^T \tilde{L} f &= \frac{1}{2} \sum_{i=1, j=1}^n W_{ij} (f_i - f_j)^2 \leq \sum_{i=1, j=1}^n W_{ij} (f_i^2 + f_j^2) \\ &= \sum_{i=1, j=1}^n W_{ij} f_i^2 + \sum_{i=1, j=1}^n W_{ij} f_j^2 \\ &= 2 \sum_i d_i f_i^2 \leq 2\rho(D) f^T f \end{aligned}$$

Some remarks of eigenvalue of NG

- From above, we know $\mathbf{1}^T \tilde{L} \mathbf{1} = 0$. Furthermore, $\tilde{L} \mathbf{1} = \mathbf{0} \mathbf{1} = \mathbf{0}$
- Since $\tilde{L} = DL$ and D is invertible, $\mathbf{0} = L \mathbf{1} = (I - D^{-1}W) \mathbf{1}$. Hence, $D^{-1}W \mathbf{1} = \mathbf{1}$. Thus, 1 is eigenvalue of A .
- Note that A is not necessary symmetric. Since A is similar to $D^{\frac{1}{2}} A D^{-\frac{1}{2}} = D^{\frac{1}{2}} (D^{-1}W) D^{-\frac{1}{2}} = D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$, which is symmetric.
- As mentioned above, \mathcal{L} is similar to $L = I - D^{-1}W$. If λ is eigenvalue of A , then $1 - \lambda$ is eigenvalue of L and \mathcal{L} .

Some spectrum properties

Lemma

$\rho(A) = 1$, $\sigma(A) \subset [-1, 1]$ and $\sigma(L) = \sigma(\tilde{L}) \subset [0, 2]$.

Diagonalization of GL (1)

- 1 As mentioned above, A is similar to symmetric matrix $D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$.
- 2 Since $D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$ is symmetric, exist orthonormal matrix O (i.e. $O^T O = 1$) such that

$$D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = O \Lambda O^T$$

where diagonal matrix $\Lambda = \begin{bmatrix} \mu_1 & & \\ & \ddots & \\ & & \mu_n \end{bmatrix}$ and

$$1 = |\mu_1| \geq |\mu_2| \geq \dots \geq |\mu_n|.$$

- 3 The eigenvalue of L and \mathcal{L} is $\{\lambda_i = 1 - \mu_i\}$

Diagonalization of GL (2)

- ④ Since $D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = O \Lambda O^T$, we can build a relation to A

$$\begin{aligned} A &= D^{-1} W = D^{-\frac{1}{2}} D^{-\frac{1}{2}} W D^{-\frac{1}{2}} D^{\frac{1}{2}} = D^{-\frac{1}{2}} O \Lambda O^T D^{\frac{1}{2}} \\ &= U \Lambda V^T, \end{aligned}$$

where $U = D^{-\frac{1}{2}} O$ and $V = D^{\frac{1}{2}} O$.

Properties of $U\Lambda V^T$ (1)

Proposition of $U\Lambda V^T$

- ① $UV^T = VU^T = U^T V = V^T U = I$
- ② $AU = U\Lambda$ and $V^T A = \Lambda V^T$
- ③ Denote two vectors $u = \frac{1}{n}\mathbf{1}$ and $v = \frac{1}{\sum d_i} [d_1, \dots, d_n]^T$, which are normalized by 1-norm, i.e. $\|\cdot\|_1$. Then, $Au = u$ and $v^T A = v^T$.

Proof:

- ① Plug $U = D^{-\frac{1}{2}} O$ and $V = D^{\frac{1}{2}} O$ into equation.
- ② By $A = U\Lambda V^T$, the following can be computed directly

$$AU = U\Lambda V^T U = U\Lambda.$$

Properties of $U \Lambda V^T$ (2)

- ③ First, we know $A \mathbf{1} = \mathbf{1}$, so $Au = u$ done!.

Second, since $A = D^{-1} W$, we can get

$$(v^T A)_j = \sum_i v_i A_{ij} = \frac{\sum_i d_i \frac{w_{ij}}{d_i}}{\sum_l d_l} = \frac{\sum_i W_{ij}}{\sum_l d_l} = \frac{d_j}{\sum_l d_l} = v_j$$

Now, it's sufficient to introduce **diffusion map**.

Eigenmap (1)

As setting above, let $A = D^{-1}W = U\Lambda V^T$ where $\Lambda = \text{diag}(\mu_1, \dots, \mu_n)$ with $1 = \mu_1 \geq \mu_2 \geq \dots \geq \mu_n$. Take m with $m+1 \leq n$. The m -dimension eigenmap for i th-vertex is defined as

$$\text{Eig}_m(i) = [u_2(i), \dots, u_{m+1}(i)]^T$$

Remark

It map vertex i in \mathbb{R}^n to \mathbb{R}^m , where $m \leq n$.

Eigenmap (2)

Note that the $n \times m$ (n data reduced in m dimension) matrix

$$\begin{bmatrix} \text{Eig}_m(1)^T \\ \text{Eig}_m(2)^T \\ \vdots \\ \text{Eig}_m(n)^T \end{bmatrix} = \begin{bmatrix} u_2(1) & u_3(1) & \cdots & u_{m+1}(1) \\ u_2(2) & u_3(2) & \cdots & u_{m+1}(2) \\ \vdots & \vdots & & \vdots \\ u_2(n) & u_3(n) & \cdots & u_{m+1}(n) \end{bmatrix} = \begin{bmatrix} | & | & & | \\ u_2 & u_3 & \cdots & u_{m+1} \\ | & | & & | \end{bmatrix}$$

Definition

As setting above, let $A = D^{-1}W = U\Lambda V^T$ where $\Lambda = \text{diag}(\mu_1, \dots, \mu_n)$ with $1 = \mu_1 \geq \mu_2 \geq \dots \geq \mu_n$. Take diffusion time $t > 0$. The diffusion map (DM) $\Phi_t: V \rightarrow \mathbb{C}^{n-1}$ is defined by

$$\Phi_t(i) = [\mu_2^t u_2(i), \mu_3^t u_3(i), \dots, \mu_n^t u_n(i)]^T.$$

Furthermore, if all eigenvalue are nonnegative, then the embedding is into \mathbb{R}^n .

Truncated diffusion map (tDM)

Definition

As setting above, let $A = D^{-1}W = U\Lambda V^T$ where $\Lambda = \text{diag}(\mu_1, \dots, \mu_n)$ with $1 = \mu_1 \geq \mu_2 \geq \dots \geq \mu_n$. Take diffusion time $t > 0$. The truncated diffusion map with time t and threshold δ is a map $\Phi_t^\delta: V \rightarrow \mathbb{C}^{m(t,\delta)-1}$ is defined by

$$\Phi_t^\delta(i) = \left[\mu_2^t u_2(i), \dots, \mu_{m(t,\delta)}^t u_{m(t,\delta)}(i) \right]^T,$$

where $m(t, \delta) := \max \{i : |\mu_i|^t > \delta |\mu_2|^t\}$.

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 **Example**
 - Data set
 - Result
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map

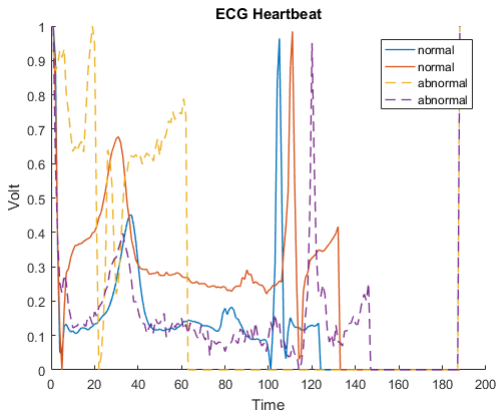
Introduction to dataset

This dataset is from website kaggle, called The PTB Diagnostic ECG Database. This dataset has been used in exploring heartbeat classification. The signals correspond to electrocardiogram (ECG) shapes of heartbeats for the normal case and the cases affected by PTB diagnostic (lung problem).

Because it is a large dataset, I just use 1000 normal and 500 abnormal to train the model.

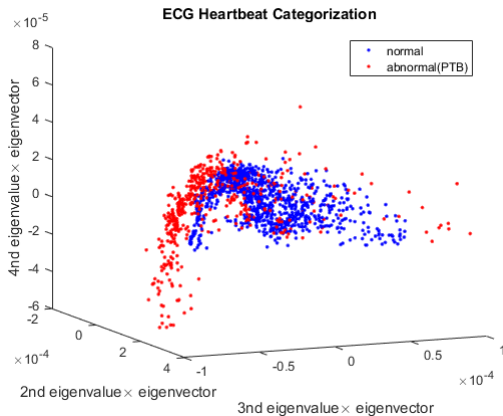
Plot the data

The following figure is about ECG heartbeats with PTB diagnostic or without PTB diagnostic.



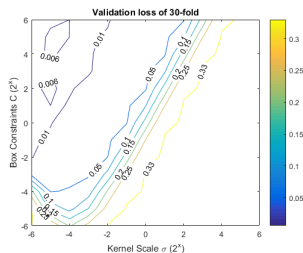
Diffusion map

Choose kernel scale $\sigma = 0.005$. Then,

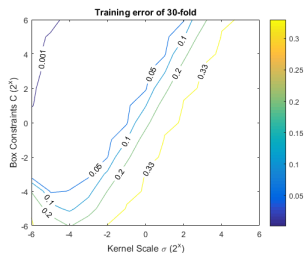


Classification by SVM

Tuning procedure between box constraints C and kernel scale σ by 30-fold.



(a) Validation loss of SVM (σ, C)



(b) Training error of SVM (σ, C)

The best cross validation loss is 0.0047 with $(\sigma, C) = (0.03125, 4)$.

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
 - Classical PCA
 - Kernel PCA
 - Dimension reduction by Kernel PCA
 - Result

Classical PCA (1)

Given sample points $\{x_1, \dots, x_n\} \subset \mathbb{R}^m$, project to axis vector $v \in \mathbb{R}^m$, and the projection points are $\{v^T x_1, \dots, v^T x_n\}$.

Assume that the mean of every components are zero. The variance on such axis v is

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (v^T x_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^n (v^T x_i - 0)^2 = \frac{1}{n} \sum_{i=1}^n (v^T x_i)^2$$

Hence, let $C = \frac{1}{n} X X^T$ be covariance matrix,

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (v^T x_i) (v^T x_i)^T = v^T \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) v = v^T C v$$

where $X = [x_1 \dots x_n]$ is $m \times n$ matrix. Note that $\sum_{i=1}^n x_i x_i^T$ is a $m \times m$ matrix.

Classical PCA (2)

Goal

The goal of PCA is to find vector v such that it maximum the variance of projection.

$$v = \arg \max_{v \in R^m, \|v\|=1} v^T C v$$

The maximum variance is the maximum eigenvalue λ and corresponding eigenvector is the solution projection axis v .

Let kernel matrix $K_{ij} = k(x_i, x_j)$. Suppose exist $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^l$ such that $K_{ij} = \phi(x_i)^T \phi(x_j)$.

Problem (1)

How to get zero mean kernel matrix \tilde{K} by using $\tilde{\phi}(x_i)$, which is removed mean from $\phi(x_i)$?

Problem (2)

How to find covariance matrix $C = \tilde{\phi}(X)\tilde{\phi}(X)^T$?

Find covariance matrix

Because ϕ is not specifically defined, $\tilde{\phi}(X)\tilde{\phi}(X)^T$ cannot be solved. Now, we rewrite the equation $\tilde{\phi}(X)\tilde{\phi}(X)^T v = \lambda v$,

$$\tilde{\phi}(X)^T \tilde{\phi}(X) \tilde{\phi}(X)^T \tilde{\phi}(X) \bar{v} = \lambda \tilde{\phi}(X)^T \tilde{\phi}(X) \bar{v}$$

$$\tilde{K} \tilde{K} \bar{v} = \lambda \tilde{K} \bar{v}$$

$$\tilde{K} \bar{v} = \lambda \bar{v},$$

where $\tilde{K} = \tilde{\phi}(x_i)^T \tilde{\phi}(x_i)$ is zero mean kernel matrix and $v = \tilde{\phi}(X) \bar{v}$ is projection axis such that it has maximum variance.

Remove mean

Remove mean of $\phi(x_i)$,

$$\tilde{\phi}(x_i) = \phi(x_i) - \frac{1}{n}\phi(X)\mathbf{1},$$

where $\mathbf{1}$ is column vector. Now,

$$\begin{aligned}\tilde{\phi}(X) &= \left[\phi(x_1) - \frac{1}{n}\phi(X)\mathbf{1} \cdots \phi(x_n) - \frac{1}{n}\phi(X)\mathbf{1} \right] \\ &= \phi(X) - \frac{1}{n}\phi(X)\mathbf{1}\mathbf{1}^T\end{aligned}$$

Write $\mathbf{1}\mathbf{1}^T$ as $\mathbf{1}_n$ It's sufficient to compute zero mean kernel matrix \tilde{K} ,

$$\begin{aligned}\tilde{K} &= \tilde{\phi}(X)^T \tilde{\phi}(X) = \left[\phi(X) - \frac{1}{n}\phi(X)\mathbf{1}_n \right]^T \left[\phi(X) - \frac{1}{n}\phi(X)\mathbf{1}_n \right] \\ &= K - K\mathbf{1}_n - \mathbf{1}_n K + \mathbf{1}_n K \mathbf{1}_n\end{aligned}$$

- **Normalized**

$$v^T v = \bar{v}^T \tilde{\phi}(X)^T \tilde{\phi}(X) \bar{v} = \bar{v}^T K \bar{v} = \lambda \bar{v}^T \bar{v}$$

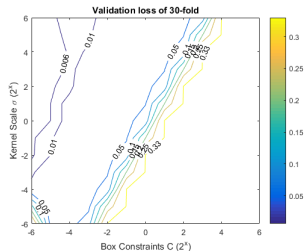
Hence, we have to assign $\frac{\bar{v}}{\sqrt{\lambda}}$ to \bar{v}

- **Projection**

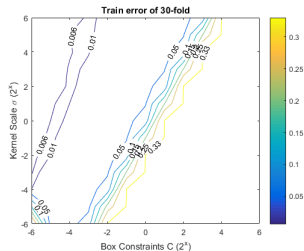
In classical PCA $v^T x_i$ is projection of x_i to v . Now, let each components of column vector P is projection of $\phi(x_i)$ to v ,

$$P = \tilde{\phi}(X)^T v = \tilde{\phi}(X)^T \tilde{\phi}(X) \bar{v} = K \bar{v} = \lambda \bar{v}.$$

Tuning procedure between box constraints C and kernel scale σ by 30-fold.



(c) Validation loss of SVM (σ , C)



(d) Training error of SVM (σ , C)

Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map
 - Nystrom Method
 - Diffusion map by Nystrom method

Nystrom method approximate Gram matrix (1)

Suppose a sample set $\mathcal{X} = \{x_i\}$ with corresponding $n \times n$ kernel matrix K , i.e. $K_{ij} = k(x_i, x_j)$. Then the subset $\mathcal{Z} = \{z_i\} \subset \mathcal{X}$, which contains landmark points, with corresponding $k \times k$ kernel matrix H , i.e. $H_{ij} = k(z_i, z_j)$.

Theorem (Williams and Seeger [3])

With above notation,

$$K \approx \hat{K} = EH^{-1}E^T,$$

where $E_{ij} = k(x_i, z_j)$.

Remark: The matrix E could be seemed as extrapolation matrix, which is submatrix of K .

Nystrom method approximate Gram matrix (2)

The eigen-system of the kernel matrix is $K\Phi_K = \Phi_K\Lambda_K$

$$\Phi_K \approx \sqrt{\frac{k}{n}} E \Phi_Z \Lambda_Z^{-1}, \quad \Lambda_K \approx \frac{n}{k} \Lambda_Z$$

Then,

$$\begin{aligned} K &\simeq \left(\sqrt{\frac{k}{n}} E \Phi_Z \Lambda_Z^{-1} \right) \left(\frac{n}{k} \Lambda_Z \right) \left(\sqrt{\frac{k}{n}} E \Phi_Z \Lambda_Z^{-1} \right) \\ &= E H^{-1} E^T \end{aligned}$$

Nystrom decomposition (Zhang [4])

Given the low rank approximation $K \approx \hat{K} = GG^T$ where $G \in \mathbb{R}^{n \times k}$ and $k < n$, the top k eigenvector U of K can be obtained as $U \approx GV\hat{\Lambda}^{-\frac{1}{2}}$, where $V, \hat{\Lambda} \in \mathbb{R}^{k \times k}$ are from eigenvalue decomposition of the $k \times k$ matrix $S = G^T G = V\hat{\Lambda}V^T$

Diffusion map by Nystrom method (1)

- 1 Given a Gaussian kernel matrix K , let $G = EH^{-\frac{1}{2}}$, where H is kernel matrix of landmark points \mathcal{Z} . That is, $K \approx GG^T = EH^{-1}E^T$.
- 2 Now, approximate degree matrix $D_{ii} = (K\mathbf{1})_i$ by $D_{ii} \approx (GG^T\mathbf{1})_i$.
- 3 Since $D^{-\frac{1}{2}}WD^{-\frac{1}{2}} = O\Lambda O^T$, let $\tilde{G} = D^{-\frac{1}{2}}G$. That is, $D^{-\frac{1}{2}}WD^{-\frac{1}{2}} \approx \tilde{G}\tilde{G}^T$.
- 4 The eigen-decomposition is $\tilde{G}^T\tilde{G} = V\hat{\Lambda}V^T$, where V is eigenvector and $\hat{\Lambda}$ is eigenvalue matrix, approximating to k eigenvalue of $D^{-\frac{1}{2}}WD^{-\frac{1}{2}}$.

Diffusion map by Nystrom method (2)

- 5 It's sufficient to approximate eigenvector matrix O by $O \approx \hat{O} = \bar{G} V \hat{\Lambda}^{-\frac{1}{2}}$.

Nystrom method approximate eigen-system

The eigen-system of $D^{-\frac{1}{2}} W D^{-\frac{1}{2}} = O \Lambda O^T$.

- $\hat{O} = \bar{G} V \hat{\Lambda}^{-\frac{1}{2}}$ approximate k eigenvector O .
- $\hat{\Lambda}$ approximate k eigenvalue Λ .

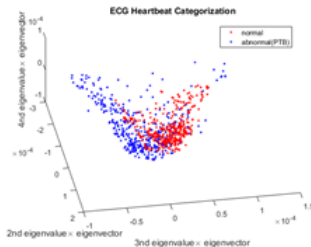
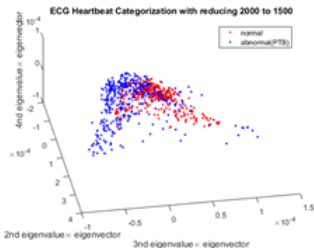
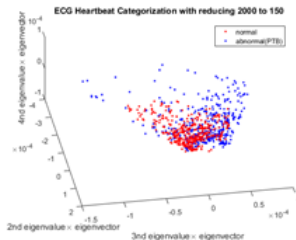
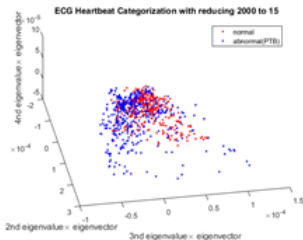
Diffusion map by Nystrom method (3)

Nystrom method approximate diffusion map

As mentioned above, let $A = D^{-1}W = U\Lambda V^T$, where $U = D^{-\frac{1}{2}}O$. Now, it's sufficient to approximate top k eigenvalue by $\hat{\Lambda}$ and to approximate top k eigenvector by $D^{-\frac{1}{2}}\hat{O} = D^{-\frac{1}{2}}\bar{G}V\hat{\Lambda}^{-\frac{1}{2}}$.

Diffusion map by Nystrom method Result

Take 15/150/1500 from 2000 samples.



Outline

- 1 Motivation
- 2 Affinity Graph
- 3 Graph Laplacian
- 4 Example
- 5 Comparison to Kernel PCA
- 6 Reduced Diffusion Map
- 7 Reference

- [1] MAOPEI TSUI, *2020 NCTS Mini-Course on Manifold Learning Lecture Notes*.
- [2] HAUTIENG WU, *Mathematics of Massive Data Analysis*.
- [3] CHRISTOPHER WILLIAMS AND MATTHIAS SEEGER, *Using the Nystrom Method to Speed Up Kernel Matrix*.
- [4] KAI ZHANG, IVOR TSANG AND JAMES KWOK, *Improved Nyström Low-Rank Approximation and Error Analysis*.